

Reference Data For Engineers

*Radio, Electronics,
Computer, and
Communications*

NINTH EDITION



**Reference Data for
Engineers:
Radio, Electronics,
Computer, and
Communications**

Reference Data for Engineers: Radio, Electronics, Computer, and Communications

Ninth Edition

Mac E. Van Valkenburg, Editor-in-Chief, Eighth Edition
Wendy M. Middleton, Editor-in-Chief, Ninth Edition


Originally titled, *Reference Data for Radio Engineers*



Newnes

Boston Oxford Johannesburg Melbourne New Delhi


Newnes is an imprint of Butterworth-Heinemann.

 A member of the Reed Elsevier group

Copyright © 1998, 2002, by Butterworth-Heinemann
Copyright © 1993 by SAM, Prentice Hall Computer Publishing
Copyright © 1968, 1975, and 1985 by SAMS
Copyright © 1956 by International Telephone and Telegraph Corp.
Copyright © 1943, 1946, and 1949 by Federal Telephone and Radio Corp.

All rights reserved.

No part of this publication may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, electronic, mechanical, photocopying, recording, or otherwise, without the prior written permission of the publisher.

-  Recognizing the importance of preserving what has been written, Butterworth-Heinemann prints its books on acid-free paper whenever possible.



Butterworth-Heinemann supports the efforts of American Forests and the Global ReLeaf program in its campaign for the betterment of trees, forests, and our environment.

ISBN: 0-7506-7291-9

The publisher offers special discounts on bulk orders of this book.
For information, please contact:

Manager of Special Sales
Butterworth-Heinemann
225 Wildwood Avenue
Woburn, MA 01801-2041
Tel: 781-904-2500
Fax: 781-904-2620

For information on all Newnes publications available, contact our World Wide Web home page at: <http://www.newnespress.com>

10 9 8 7 6 5 4 3 2 1

Printed in the United States of America

DEDICATION

IN MEMORIAM
Mac E. Van Valkenburg
1921–1997

IN MEMORIAM
Edward C. Jordan
1910–1991

Preface

The ninth edition provides new knowledge in the fields of satellite technology, space communication, microwave science, telecommunication, global positioning systems, frequency data, and radar. Also revised and updated were the chapters on intellectual property and patents, probability and design, power electronics—rectifiers, filter, and power supplies, properties of materials, units, constants and conversion factors, active filter design, antennas, integrated circuits, digital signal processing, and surface acoustic wave devices.

I am sure that as you read this text you will also appreciate the hard work and expertise of these authors. This edition of *Reference Data for Electrical Engineers* will continue its tradition of being an important reference text for all electrical engineers.

Highlights of Updated Chapters

Chapter 1: Radio Frequency Spectrum Management and Time Frequency Standards - by F. Matos

- ITU information is available on the web for download.
- While the Omnibus Budget Reconciliation Act did open the door for auctions, it only authorized them for a short time. A better reference is: *Radiofrequency Spectrum Management: Background, Status, and Current Issues*, January 30, 2001, Richard M. Nunno, Congressional Research Service, Washington, DC Order Code RL30829.
- The GPS Time Scale began on January 6, 1980 and is not adjusted for leap seconds. Therefore, it differs from UTC by the number of leap seconds that have occurred January 6, 1980. As of June 2001, GPS Time is 13 seconds ahead of UTC.
- Much of the data in this chapter cross-references with Chapter 37 on GPS.
- The updates on GPS before and after S/A was turned off are quite good, as is the discussion of the merits of different types of atomic clocks.
- For up-to-date GPS information, the web address for the United States Coast Guard Navigation Center is: <http://www.navcen.uscg.gov/>
- For the readers' reference: A nanometer is 10⁻⁹ meter. The *meter* is defined as: the length of the path traveled by light in vacuum during a time interval of 1/299 792 458 of a second.
- The constant for the speed of light can be found at: http://www.physics.nist.gov/cgi-bin/cuu/Value?csearch_for=light

Chapter 3: Units, Constants, and Conversion Factors - by E. Mechtly

- Temperatures are discussed in the context of Table 2, SI Base and Derived Units. See entries for *Kelvin* and for *degree Celsius* in Table 2.
- Many tables were updated. The style information on how to correctly express different units is a resource that many will be likely to use in the future.
- The values for constants have been updated and citations including worldwide web addresses have been added.

Chapter 4: Properties of Materials - by E. Mechtly

- The most current versions of the table of Superconducting Materials can be found at internet sources such as the NIST Superconducting Materials Database at <http://www.ceramics.nist.gov/srd> or by the "High-T (critical) Update," sponsored by the Argonne National Laboratory at <http://www.iitap.iastate.edu>.
- The tables on wires, gauges, and drills sizes lie on the interface between engineering and shop technology. Here, again, the internet offers a multitude of equivalent sources. For example, a web search for "wire gauge" or "AWG" yields many sources. Likewise, does a search for "drills."

Chapter 10: Active Filter Design - by R. Schaumann

- GHz applications and three amplifier filters sections have been updated.
- The chapter now has additional detail on switched filters.

Chapter 14: Power Electronics—Rectifiers, Filters, and Power Supplies - by O. Pauk

- Manufacturer's data sheets are now available for free download from the web as PDF files.
- Information on switching power supplies is of high quality and is up-to-date.

Chapter 16: Electron Tubes - by D. Abe, R. Abrams, B. Danly, H. Freund, K. Jensen, B. Levush, R. Myers, R. Parker, J. Shaw, A. Shih, and J. Yater

- For more than fifty years, the venerable vacuum tube has proved to be the technology of choice for many applications. The advent of solid-state power amplifiers has also made a profound impact. This healthy technological competition has no doubt

been to the benefit of both technologies, and the reader should not mistake the proud history of vacuum electronics as an indication of a technology at its sunset.

- This chapter describes the numerous advances in the area of vacuum electronics, which continue to make the technology the appropriate choice for many applications. Present-day applications for microwave tubes include satellite and terrestrial communications, industrial and scientific applications from heating materials to heating fusion reactors, and continuing military applications, such as radar and electronic warfare.
- Driven by these continuing needs, this technology area has continued to see impressive advances in recent times, as described in this chapter.

Chapter 20: Integrated Circuits - by H. Hosack

- Updated information and figures add to this essential chapter.
- New information has been added on MEM's and submicron topics.

Chapter 27: Satellite and Space Communications - by G. Hyde

- A discussion of Iridium is included in the chapter. However, Iridium's fate is far from clear.
- BSS definition: Broadcast Satellite Services.
- DARS definition: Digital Audio Radio System.
- ALOHA definition: A made-up name, not an acronym - made up by Prof. Abramson (of the University of Hawaii, of course) Somebody may have fitted words to the letters, but the author did not hear it defined by Abramson when he lectured on it at Comsat Labs many years ago.
- Multi-beam antennas of 15 meters diameter for MSS are mentioned at the end of the Spacecraft Antenna section.
- Note that 3G wireless is a future system and it is not clear when it will supersede 2.5G, which seems to be growing rapidly.

Chapter 28: Discrete-Time Signal Processing - by K. Jenkins and D. Jones

- The SAW filter data has been updated.
- Advances in Digital signal processing have been presented.
- Technological improvements in Surface Acoustic design techniques are covered.

Chapter 32: Antennas - by R. Hansen

- Updated data on antennas.

Chapter 36: Radar - by M. Skolnik

- Information of the newest advances in Radar.

Chapter 37: Radio Navigation Systems - by D. Jablonski

- The most recent data on GPS.
- Information on the remaining uses for LORAN.
- Insights into upcoming changes in the field.

Chapter 38: Common Carrier Transmission - by L. Sibley

- Provides good information on DSL.

Chapter 44: Probability and Statistics - by S. Gupta

- Gives concise clear data for design use.

Chapter 46: Cellular Telecommunications Systems - by W. Lee

- Updates on 3G and IMT-2000.

Chapter 49: Miscellaneous Data - by D. Bender

- Copyright law changes highlighted.
- Intellectual property protection strategies outlined.

Contents

Chapter 1	Radio Frequency Spectrum Management and Time Frequency Standards
Chapter 2	International Telephony Recommendations
Chapter 3	Units, Constants, and Conversion Factors
Chapter 4	Properties of Materials
Chapter 5	Components or Parts
Chapter 6	Fundamentals of Networks
Chapter 7	Fourier Waveform Analysis
Chapter 8	Filters, Simple Bandpass Design
Chapter 9	Filters, Modern-Network-Theory Design
Chapter 10	Active Filter Design
Chapter 11	Attenuators
Chapter 12	Measurements and Analysis
Chapter 13	Magnetic-Core Transformers and Reactors
Chapter 14	Power Electronics—Rectifiers, Filters, and Power Supplies
Chapter 15	Feedback Control Systems
Chapter 16	Electron Tubes
Chapter 17	Power Grid-Tube Circuits
Chapter 18	Semiconductors and Transistors
Chapter 19	Transistor Circuits
Chapter 20	Integrated Circuits
Chapter 21	Optoelectronics
Chapter 22	Optical Communications
Chapter 23	Analog Communications
Chapter 24	Digital Communications
Chapter 25	Information Theory and Coding
Chapter 26	Computer Communications Networks
Chapter 27	Satellite and Space Communications
Chapter 28	Discrete-Time Signal Processing
Chapter 29	Transmission Lines
Chapter 30	Waveguides and Resonators
Chapter 31	Scattering Matrices
Chapter 32	Antennas
Chapter 33	Electromagnetic-Wave Propagation
Chapter 34	Radio Noise and Interference
Chapter 35	Broadcasting, Cable Television, and Recording System Standards
Chapter 36	Radar
Chapter 37	Radio Navigation Systems
Chapter 38	Common Carrier Transmission
Chapter 39	Switching Networks and Traffic Concepts
Chapter 40	Electroacoustics
Chapter 41	Lasers
Chapter 42	Computer Organization and Programming
Chapter 43	Logic Design
Chapter 44	Probability and Statistics
Chapter 45	Reliability and Life Testing
Chapter 46	Cellular Telecommunications Systems
Chapter 47	Mathematical Equations
Chapter 48	Mathematical Tables
Chapter 49	Miscellaneous Data

Contributors

David K. Abe (Ch. 16)

Electronics Engineer, Vacuum Electronics Branch, Electronic Systems and Technology Division, Naval Research Laboratory.
Member IEEE, APS.

Richard H. Abrams, Jr. (Ch. 16)

Head, Emission Sources Section, Vacuum Electronics Branch, Naval Research Laboratory.
Member IEEE, APS, AAAS, AOC.

Richard W. Avery (Ch. 38)

Member of Technical Staff. AT&T - Bell Laboratories.

Fred J. Banzi (Ch. 13)

Member of Technical Staff. AT&T - Bell Laboratories.
Member IEEE.

Pier L. Bargellini (Ch. 27)

Senior Scientist (retired), COMSAT Laboratories.
Fellow IEEE.
Assoc. Fellow AIAA.

Stanley Baron (Ch. 35)

Managing Director Technical Development, National Broadcasting Co. (NBC).
Fellow SMPTE.
Member Royal Television Society.
Past Engineering Vice-President SMPTE.
Chair, CCIR TG 11/3 - Digital Terrestrial Television Broadcasting.
David Sarnoff Gold Medal.

David Bender (Ch. 49)

Counsel, White & Case (New York).
American Bar Assn.
Computer Law Assn. Board of Directors.
American Intellectual Property Law Assn.
Member ACM.

Richard E. Blahut (Ch. 25)

IBM Fellow, IBM Corporation.
Fellow IEEE.
Past-President IEEE/GIT.
National Academy of Engineering.

Robert Blaszczczyk (Ch. 5)

Applications Engineer, VRN Division of Vemitron Corp.

John Bohrer (Ch. 5)

Consultant, Bohrcon.
Past-Chairman ECC.
Past-President AI Chem.

John F. X. Browne (Ch. 35)

President, John F. X. Browne and Associates.
Fellow SMPTE.
Governor SMPTE.
Senior Member IEEE.
Past-President Assn. of Federal Communications Consulting Engineers.

Charles A. Cain (Ch. 49)

Professor of Electrical Engineering and Biocengineering, University of Michigan.
Senior Member IEEE.
Member of Directors, Bioelectromagnetic Society.

Bernard B. Carnigila (Ch. 13)

Member of Technical Staff. AT&T - Bell Laboratories.

Dorothy H. Cerni (Ch. 2)

Technical Information Specialist, Institute for Telecommunication Sciences, National Telecommunications and Information Administration (NITA).

Pallab K. Chatterjee (Ch. 20)

Vice-President and Chief Technical Officer, Components Section, Director Research and Development, Texas Instruments, Inc.
Fellow IEEE.
Member APS.

Marvin Chodorow (Ch. 16 & 17)

Professor Emeritus of Applied Physics and Electrical Engineering, Ginzton Laboratory, Stanford University.
Fellow IEEE.
Lamme Medal (1982).
W. R. G. Baker Award (1962).

Erich Christian (Ch. 9)

Master of Advanced Engineering, ITT Telecom (retired).
Professor North Carolina State University.
Senior Member IEEE.

Robert Coakley (Ch. 12)

President, Elenex Corp.
Fellow (British) IEEE.

Douglass D. Crombie (Ch. 2 & 33)

Senior Editor, The Acrospace Corporation.
Former Director of the Institute of Telecommunication Science, National Telecommunications And Information Administration (NITA).
Chief Scientist (retired) NITA.
National Academy of Engineering.
Fellow IEEE.

Bruce G. Danly (Ch. 16)

Supervisory Research Physicist, Head, High Power Devices Section, Vacuum Electronics Branch, Naval Research Laboratory.
Member IEEE, APS.
Member of the Team awarded DoD Robert L. Woods Award for Development of the High-power 94 GHz Gyro-Klystron technology.

Thijs de Hass (Ch. 2)

Networks and Systems Division, ITS, National Telecommunications And Information Administration (NITA).

Georges A. Deschamps (Ch. 31 & 44)

Professor Emeritus of Electrical and Computer Engineering, UIUC.
National Academy of Engineering.
Fellow IEEE, AAAS.
Member APS.

Thomas DeTemple (Ch. 41)

Professor of Electrical and Computer Engineering, UIUC.
Senior Member IEEE.
Member APS.

Michiel deWitt (Ch. 20)

Senior Member, Technical Staff, Texas Instruments, Inc.
Member IEEE, APS.

John J. Dupre (Ch. 12)

R&D Section Manager, Lightwave Instruments, Hewlett-Packard Co.
Member IEEE.

John D. Dyson (Ch. 31)

Professor of Electrical and Computer Engineering, UIUC.
Fellow IEEE.
Member APS, MTT, I and M, EMC, Antenna Standards Committee.

Hermann Fickenscher (Ch. 13)

Supervisor, Magnetic Components Group, AT&T - Bell Laboratories.
Senior Member IEEE.
Member Electronics Transformers Technical Committee.

Arthur Fong (Ch. 12)

Manager (retired) Corporate Engineering Design, Hewlett-Packard Corp.
Fellow IEEE.

Henry P. Freund (Ch. 16)

Senior Research Physicist, Science Applications International Corp.
Member and Fellow APS.
New York Academy of Sciences.

Chester S. Gardner (Ch. 22)

Professor of Electrical and Computer Engineering, UIUC.
Fellow OSA, IEEE.
Vice-Chairman for Remote Sensing OS Technical Council.

Mario Gerla (Ch. 26)

Professor and Vice-Chair of Computer Science, University of California, Los Angeles.
Member IEEE.
Fulbright Fellow.

Stephen J. Greer (Ch. 12)

Interface Engineer, Hewlett-Packard Co.
Member IEEE.
Secretary IEC/SC/G5C Working Group.

Surendra M. Gupta (Ch. 44)

Director of Laboratory for Responsible Manufacturing, Northeastern University, Boston, MA.
Professor of Mechanical, Industrial and Manufacturing Engineering.
Member IIE, INFORMS, DSI, POMS, APICS.
Elected Member of Alpha Pi Mu Honor Society, Delta Mu Delta Honor Society, Sigma Xi Honor Society, Tau Beta Pi Honor Society.
2000 Outstanding Research Award, Northwestern University.
Registered Professional Engineer, MA.

Robert C. Hansen (Ch. 32)

President, R. C. Hansen, Inc.
Phi Kappa Phi, Tau Beta Pi, Eta Kappa Nu.
Fellow IEE and IEEE.
National Academy of Engineering.

Mark Harward (Ch. 20)

Member of Technical Staff, Texas Instruments, Inc.
Member IEEE.
Member Eta Kappa Nu.

Charles F. Hempstead (Ch. 13)

Supervisor (retired), Magnetics and Characterization, AT&T - Bell Laboratories.

Harold H. Hosack (Ch. 20)

Associate Director, Materials and Process Science, Semiconductor Research Corp.

Bill J. Hunsinger (Ch. 28)

Co-Chair, Comlinear Corp.
Fellow IEEE.

Geoffrey Hyde (Ch. 27)

Editor-in-Chief, Space Communications Journal (IOS Press).
Fellow and Life Fellow IEEE.
Member IEEE-USA APC, IEEE (AP-S, Comm. Soc.), USNC/URSI Comm. B and F, AIAA.
IEEE AP-S Best Paper Award 1967.
RCA Sarnoff Fellowships 1964, 1965.

Tatsuo Itoh (Ch. 29)

Hayden Head Professor in Engineering, University of Texas.
Fellow IEEE.
Editor IEEE Transactions, MMT.
Chairman Microwave Field Theory Comm.

Daniel G. Jablonski (Ch. 37)

Member Principal Professional Staff, Johns Hopkins University Applied Physics.
Instructor, Whiting School of Engineering, Johns Hopkins University.
Adjunct Professor, Capitol College, Laurel, MD.
Senior Member IEEE.
Member American Physical Society.
Licensed Professional Engineer, MD.

Horace G. Jackson (Ch. 19)

Lecturer, Dept. of Electrical Engineering and Computer Science, University of California, Berkeley.

W. Kenneth Jenkins (Ch. 28)

Head of Electrical Engineering, Pennsylvania State College.

Fellow IEEE.

IEEE Millennium Medal, 2000.

IEEE Circuits and Systems Society Golden Jubilee Medal, 2000.

Co-winner of the 2000 International Award of the George Montefiore Foundation (U. of Liege, Belgium) for outstanding career contributions to the field of Electrical Engineering.

Recipient of the IEEE Circuits and Systems Society Meritorious Service Award (1990).

Past President of IEEE Circuits and Systems Society (1985).

Kevin Lynn Jensen (Ch. 16)

Research Physicist, Electron Science and Technology Division, Naval Research Laboratory.

Member APS, IEEE, NYAS, MRS, ECS.

Amos E. Joel (Ch. 39)

Switching Consultant, Bell Telephone Laboratories (retired).

National Academy of Engineering.

Fellow IEEE.

IEEE Medal of Honor (1992).

ITU Centenary Award (1983).

Alexander Graham Bell Medal (1976).

Douglas L. Jones (Ch. 28)

Professor of Electrical and Computer Engineering, UIUC.

Fellow IEEE.

Fulbright Fellow (1987–88).

Durwood R. Kressler (Ch. 13)

Member of Technical Staff, AT&T - Bell Laboratories.

Kenneth R. Laker (Ch. 10)

Alfred Fittler Moore Professor of Electrical Engineering, University of Pennsylvania.

Formerly, Supervisor, Signal Processing Subsystems Group, AT&T - Bell Laboratories.

Fellow IEEE.

Past-President Circuits and Systems Society.

President, Division 1 of IEEE.

William C. Y. Lee (Ch. 46)

Chairman of LinkAir Communications, Inc. and developer of LAS-CDMA—a new patented technology for wireless telecom systems.

Baruch Levush (Ch. 16)

Head of Theory and Design Section, Vacuum Electronics Branch, Electronics Science and Technology division, Naval Research Laboratories.

Member IEEE, APS.

Fellow IEEE.

Member of the Team awarded DoD Robert L. Woods Award for Development of the High-power 94 GHz Gyro-Klystron technology.

Sigfried Linkwitz (Ch. 12)

EMC Specialist, Hewlett Packard Company.

Satwinder D. S. Malhi (Ch. 20)

Member Technical Staff, Texas Instruments, Inc.

Member IEEE.

Frederick Matos (Ch. 1)

Telecommunication policy analyst, National Telecommunications and Information Administration (NITA).
Former IEEE Congressional Fellow.

Edward McClusky (Ch. 43)

Professor and Director, Center for Reliable Computing, Stanford University.
Fellow IEEE.
Member ACM.
SIGARCH Director 1980–83.
Past-President IEEE Computer Society.

Eugene A. Mechtly (Ch. 3 & 4)

Professor, Former, 2-term, Assistant Dean College of Engineering, UIUC.
Member American Society for Testing and Materials.
Member National Conference on Weights and Measures.

Douglas L. Merriott (Ch. 45)

Associate Professor of Mechanical and Industrial Engineering, UTUC.

Edward Mette (Ch. 5)

Resistive Products Division, TRW, Inc.

Mehrdad Moslehi (Ch. 20)

Member of Technical Staff, Texas Instruments, Inc.
Member IEEE.

David C. Munson Jr. (Ch. 7)

Professor of Electrical and Computer Engineering, UIUC.
Member IEEE.

Robert E. Myers (Ch. 16)

Program Manager, Vacuum Electronics Branch, Vacuum Electronics Engineering Facility, Naval Research Laboratory.
Member IEEE, ASM.

Arthur Olsen Jr. (Ch. 13)

Supervisor, Magnetics and Data Group, AT&T - Bell Laboratories.

John R. Page Jr. (Ch. 12)

R&D Manager, RF and Microwave Instruments, Hewlett-Packard Company.
Member IEEE.

George Papen (Ch. 22)

Professor of Electrical and Computer Engineering, UIUC.
Fellow OSA.
Vice-Chairman for Remote Sensing OS Technical Council.
Fellow IEEE.

Robert K. Parker (Ch. 16)

Head, Vacuum Electronics Branch, Naval Research Laboratory.
Fellow IEEE, APS.

Ondrej Pauk (Ch. 14)

Applications Engineer, Motorola Transportation Systems Group, SMARTMOS Product Technology Center.

William R. Perkins (Ch. 15)

Professor of Electrical and Computer Engineering, UIUC.
Fellow IEEE.
Editor-in-Chief IEEE Press.
IEEE Centennial Medal.
Distinguished Member, IEEE Controls Systems Society (President 1985).
Halliburton Education Leadership Award.

Ronald E. Pratt (Ch. 12)

R&D Project Manager, Hewlett-Packard Company.
Manager IEEE.

Rick Price (Ch. 5)

Capacitor Division, TRW, Inc.

Donald H. Priest (Ch. 17)

Consultant.
Former Chief Advisory Engineer (retired), EIMAC Division of Varian Assoc.

Michael B. Pursley (Ch. 23 & 24)

Professor of Electrical and Computer Engineering, UIUC.
Fellow IEEE.
Past-President IEEE/GIT.

Frederick M. Renley Jr. (Ch. 35)

Staff Vidco Specialist, Center for Information Technology Integration, University of Michigan.
Fellow, Honorary Member SMPTE.
Chairman Working Group on Advanced Television Production.
Member IEEE.
Fellow SBE.

Eugene A. Robinson (Ch. 37)

Senior Fellow, Texas Instruments, Inc.
Engineering Director, Missile Systems.

Lawrence C. J. Roscoe (Ch. 38)

Department Head, AT&T - Bell Laboratories.
Member IEEE.

Douglas K. Rytting (Ch. 12)

R&D Section Manager, High Performance Network Analyzers, Hewlett-Packard Company.

Gary D. Sasaki (Ch. 12)

R&D Section Manager, Frequency and Time Instruments, Hewlett-Packard Company.
Member IEEE.

Rolf Schaumann (Ch. 10)

Professor Emeritus, Former Chair of Electrical and Computer Engineering at Portland State University.
Fellow IEEE.
Member HKN, TBP.
Millennium Medal, IEEE, 2000.
Golden Jubilee Medal, IEEE Circuits and Systems Society, 1999.
Meritorious Service Award, IEEE Circuits and Systems Society, 1993.

Paul D. Schomer (Ch. 40)

Team Leader - Acoustics, U. S. Army Construction Engineering Research Laboratory.
Fellow ASA.
Member IEEE, Institute of Noise Control Engineering.
Chairman, NAE.
U. S. National Counterpart Committee to CIB-W-51 Acoustics.

Jonathan Shaw (Ch. 16)

Electronics Engineer, Naval Research Laboratories.
Member AVS, MRS, IEEE.

W. Ford Shepherd (Ch. 8 & 9)

Director, R&D Network Systems Division, ITT Telcomm.

Arnold Shih (Ch. 16)

Research Physicist, Naval Research Laboratories.
Member APS.

Ludwell A Sibley (Ch. 38)

Consultant.
Formerly Director, Transmission Engineering and Network Technical Support, Telcordia Technologies, Inc.
Senior Member IEEE.

Merrill I. Skolnik (Ch. 36)

Superintendent Emeritus, Radar Division, Naval Research Laboratory.
Fellow IEEE.

Kent Sternstrom (Ch. 13)

Member of Technical Staff, AT&T - Bell Laboratories.

Gregory E. Stillman (Ch. 21)

Professor of Electrical and Computer Engineering UIUC.
Fellow IEEE.
Past-President Electron Device Society.
Jack Morton Award 1999.
Gallium Arsenide Society, Heinrich Welker Medal 1999.
Chairman, National Research Council, Center for Electronic Engineering Assessment, National Institute of Standards and Technology.
National Academy of Engineering.

Ben G. Streetman (Ch. 18)

Professor and J. S. Cockrell Centennial Chair, Dept. of Electrical Engineering, University of Texas, Austin.
Fellow IEEE.
Member Electrochemical Society.
ASEE Terman Award.

George W. Swenson (Ch. 34 & 40)

Professor Emeritus of Electrical Engineering and Astronomy, UIUC.
Member of the Acoustics Team, Construction Engineering Research Laboratory, U. S. Army Corps. of Engineers.

Robert Symons (Ch. 16)

Technical Director, Litton Electron Device Division.
Fellow IEEE.
Former Associate Editor, IEEE Transactions on Electron Devices.

Richard Thomas (Ch. 34)

Senior Staff Scientist, National Radio Astronomy Observatory.
Fellow IEEE.

Fouad A. Tobagi (Ch. 26)

Associate of Electrical Engineering, Stanford University.
Fellow IEEE.
Member ACM.

Joseph A. Toro (Ch. 5)

Manager, Special Products, TWR Capacitor Division.
Member IEEE.

Fred G. Turnbull (Ch. 14)

Electronics Engineer, General Electric Co.
Fellow IEEE.

Joseph T. Verleyen (Ch. 41)

Professor of Electrical and Computer Engineering, UIUC.
Senior Member IEEE.
Member APS.
Chairman, Gaseous Electronics Conference.

Theodore F. VonKampen (Ch. 5)

Engineering Manager, TWR Capacitor Division.

John Wakerly (Ch. 42)

Vice-President of Engineering, Alantec, San Jose, CA.
Consulting Associate Professor, Stanford University.
Member IEEE.
IEEE Centennial Award 1984.

Richard Walker (Ch. 13)

Member of Technical Staff, AT&T - Bell Laboratories.

Alva L. Wallis (Ch. 48)

Meteorologist, National Climatic Data Center, Asheville, NC.

Bruno Weinschel (Ch. 11)

Chief Engineer, Weinschel Associates.
Fellow and Past-President IEEE.
Past-Chairman CPEM.
Past-Chairman U. S. Commission I of URSI.
Fellow (British) IEE.

Joan E. Yater (Ch. 16)

Research Physicist, Naval Research Laboratories.
Member APS.

Acknowledgments

The search for new knowledge springs from current knowledge. When we have outgrown the current we look for new frontiers. The previous edition of *Reference Data for Electrical Engineers* has provided a foundation for this new edition and gave it the nutrients to grow. I would like to thank the Editor-in-Chief of the eighth edition, Mac E. Van Valkenburg for establishing a solid foundation for new growth.

The backbone of the ninth edition is comprised of authors from the eighth edition who revised and updated their chapters. This gives the ninth edition endurance and stability. I would like to thank Geoffrey Hyde, Merrill Skolnik, Ludwell Sibley, Rolf Schumann, Fredrick Matos, Eugene Mechtly, Harold Hosack, David Bender, Kevin Jenkins, Douglas Jones, William Lee, and Robert Hansen for participating in this new edition and for their continuing generous sharing of knowledge. Their advice and assistance to me was of immeasurable value.

Adding to the expansion of our knowledge were several new authors. Ondrej Pauk wrote on power electronics—rectifiers, filters, and power supplies. Surendra Gupta provided a much-needed update of the probability and statistics chapter. Daniel Jablonski expanded upon the GPS chapter to give us an up-to-date direction for the ninth edition. Robert Parker, Kevin Jensen, and their associates collaborated to provide an excellent source of knowledge on microwave technology.

Finally, I would like to acknowledge the strength and support of my husband, James; my children, Rebecca and David; my granddaughter, Aislinn (who provided distraction when needed); and last but not least, my best friend, an electrical engineer and electronic journalist, Lee Goldberg.

Wendy M. Middleton
2/4/01

1 Radio Frequency Spectrum Management and Time and Frequency Standards

Revised by Fredrick Matos

Wavelength–Frequency Conversion 1-2

Conversion Equations

Nomenclature of Radio Frequency Bands

Radio Frequency Spectrum Management 1-2

International

US Radio Frequency Spectrum Management

Spectrum Engineering Software

Frequency Allocations

Frequency Allotments

Time Scales 1-142

International Definition of the Second

Time Zones

Time and Frequency Standards 1-144

Quartz Oscillators

Rubidium Oscillators

Cesium Oscillators

Hydrogen Masers

Time and Frequency Broadcasts 1-148

Short-Wave Services—WWV and WWVH

Low-Frequency Services—WWVB

WAVELENGTH-FREQUENCY CONVERSION

Fig. 1 permits conversion between frequency and wavelength. The f scale may be multiplied by a power of 10 if the λ scale is divided by the same power of 10.

Conversion Equations

Propagation velocity

$$c \approx 3 \times 10^8 \text{ meters/second}$$

Wavelength in meters

$$\lambda_m = \frac{300\,000}{f \text{ in kilohertz}} = \frac{300}{f \text{ in megahertz}}$$

Wavelength in centimeters

$$\lambda_{cm} = \frac{30}{f \text{ in gigahertz}}$$

Wavelength in feet

$$\lambda_{ft} = \frac{984\,000}{f \text{ in kilohertz}} = \frac{984}{f \text{ in megahertz}}$$

Wavelength in inches

$$\lambda_{in} = \frac{11.8}{f \text{ in gigahertz}}$$

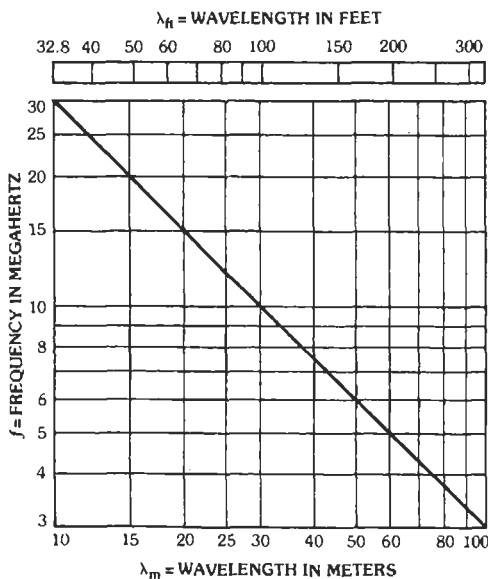


Fig. 1. Wavelength-frequency conversion.

- 1 angstrom unit, Å = 3.937×10^{-9} inch
- = 1×10^{-10} meter
- = 1×10^{-4} micrometer
- 1 micrometer, μm = 3.937×10^{-5} inch
- = 1×10^{-6} meter
- = 1×10^4 angstrom units

(Note that the term "micrometer" has superseded the term "micron.")

Nomenclature of Radio Frequency Bands

Table 1 is adapted from the Radio Regulations of the International Telecommunication Union, Geneva, 1990. Table 2 shows an alternative and more widely used set of designations.

RADIO FREQUENCY SPECTRUM MANAGEMENT

International

The distribution of radio frequencies for various telecommunications services is established by international treaty by the International Telecommunication Union (ITU), an organ of the United Nations. Frequency allocations, procedures for international registration and coordination of frequencies and satellites, frequency sharing, and other regulations, procedures, and technical standards are contained in the ITU *Radio Regulations*. Copies of the *Radio Regulations*, technical standards, and related information may be obtained from the International Telecommunication Union, Place des Nations, CH-1211, Geneva 20, Switzerland.

The ITU *Radio Regulations* have the status of international law and are amended at a treaty conference referred to as a World Radiocommunication Conference (WRC), normally convened every two or three years to address various topics. The tables of frequency allocations, regulations, and procedures of individual nations generally follow the ITU *Radio Regulations* with nations amending their own domestic rules and regulations following those of the WRC.

US Radio Frequency Spectrum Management

The Communications Act of 1934 provides the foundation for US spectrum rules and regulations, management, and usage. The basic authority is delineated in Sections 303, 304, and 305 of the Act. Section 303 presents the general powers of the Federal Communications Commission (FCC) regarding transmitting stations; 304 deals with waiving frequency claims; and 305 provides that Federal Government owned stations shall be assigned frequencies by the President (delegated to the Department of Commerce National Communications and Information Administration [NTIA]).

TABLE 1. NOMENCLATURE OF FREQUENCY BANDS

Band Number*	Frequency Range	Metric Subdivision	Adjectival Designation	
4	3 to 30 kilohertz	Myriametric waves	VLF	Very low frequency
5	30 to 300 kilohertz	Kilometric waves	LF	Low frequency
6	300 to 3000 kilohertz	Hectometric waves	MF	Medium frequency
7	3 to 30 megahertz	Decametric waves	HF	High frequency
8	30 to 300 megahertz	Metric waves	VHF	Very high frequency
9	300 to 3000 megahertz	Decimetric waves	UHF	Ultrahigh frequency
10	3 to 30 gigahertz	Centimetric waves	SHF	Superhigh frequency
11	30 to 300 gigahertz	Millimetric waves	EHF	Extremely high frequency
12	300 to 3000 gigahertz or 3 terahertz	Decimillimetric waves		

* "Band Number N" extends from 0.3×10^N to 3×10^N hertz. The upper limit is included in each band; the lower limit is excluded.

Based on ITU Radio Regulations, 1998.

TABLE 2. LETTER DESIGNATIONS FOR FREQUENCY BANDS

Letter	Frequency Band
L	1000–2000 MHz
S	2000–4000 MHz
C	4000–8000 MHz
X	8000–12 000 MHz
K _u	12–18 GHz
K	18–27 GHz
K _a	27–40 GHz
V	40–75 GHz
W	75–110 GHz

Source: IEEE Standard Letter Designations for Radar Frequency Bands, IEEE Std 521–1984, reaffirmed 1989.

Section 305 is particularly significant because it provides for the separation of authority between the Federal Government and the non-Federal Government, or private, sector. Section 305 has resulted in two US spectrum regulatory bodies: the FCC regulating the non-Federal Government sector, and the NTIA regulating the Federal Government sector. Section 305 has also resulted in agreements between the Federal Government and non-Government sectors that essentially divide the spectrum usage into three parts: exclusive Federal Government use, exclusive non-Federal Government use, and use shared between the two sectors.

The NTIA is aided by other federal agencies and departments through an advisory group, the Interdepartment Radio Advisory Committee (IRAC). IRAC carries out frequency coordination for the Federal Government agencies, recommends technical standards, and reviews major Federal Government systems to assure spectrum availability. The IRAC also provides advice to the NTIA on spectrum policy issues.

Although the NTIA and FCC generally operate independently of each other, they coordinate closely on spectrum matters. An FCC liaison representative participates in the IRAC, and the NTIA participates in the rule-making process of the FCC with the advice of the IRAC. FCC and NTIA spectrum sharing coordination is also carried out daily as required.

Spectrum Engineering Software

Spectrum engineering software is available free of charge from the National Telecommunications and Information Administration, Washington, DC. The Microcomputer Spectrum Analysis Models (MSAM) software contains 13 models, including models for intermodulation and propagation.

Frequency Allocations

Frequencies are allocated into tables via international treaty set by the ITU, and domestically in the United States by the Federal Communications Commission (FCC) and the National Telecommunications and Information Administration (NTIA). Many other nations have their own domestic tables. Table 3 presents the International Table of Frequency Allocations and the United States National Table of Frequency Allocations from the *Manual of Regulations and Procedures for Federal Radio Frequency Management* published by the NTIA or simply the *NTIA Manual*. (The latest edition of the NTIA Manual can be obtained from the Website www.ntia.doc.gov.) A very similar table is provided in Part 2 of the FCC Rules and Regulations. The international table contained in the following is obtained from the ITU's *Radio Regulations*, Edition of 1998.

The footnotes to the tables are important and need to be reviewed when considering frequency allocations. There are four different types of footnotes: 1) international;

2) United States comprehensive with the prefix "US;" 3) US Federal Government applying to only US Federal Government users, with the prefix "G;" and 4) Non-Federal-Government footnotes, applying to the US private sector, including state and local governments, abbreviated "NG."

ITU Table of Frequency Allocations—The ITU Table of Frequency Allocations (Table 3) is that table contained in Article S5 of the ITU Radio Regulations, Edition of 1998.

National Table of Frequency Allocations—The National Table of Frequency Allocations is comprised of the US Government Table of Frequency Allocations and the FCC Table of Frequency Allocations. The National Table indicates the normal national frequency allocation planning and the degree of conformity with the ITU Table. When required in the national interest and consistent with national rights, as well as obligations undertaken by the United States to other countries that may be affected, additional uses of frequencies in any band may be authorized to meet service needs other than those provided for in the National Table.

Specific exceptions to the National Table of Frequency Allocations are as follows:

A Government frequency assignment may be authorized in a non-Government band, as an exception, provided a) the assignment is coordinated with the FCC and b) no harmful interference will be caused to the service rendered by non-Government stations, present or future.

A non-Government frequency assignment may be authorized in a Government band, as an exception, provided a) the assignment is coordinated with the IRAC and b) no harmful interference will be caused to the service rendered by Government stations, present or future.

In the case of bands shared by Government and non-Government services, frequency assignments therein shall be subject to coordination between the IRAC and the FCC and no priority is recognized unless the terms of such priority are specifically defined in the National Table of Frequency Allocations or unless they are subject to mutually agreed arrangements in specific cases.

US Government Table of Frequency Allocations—The US Government Table of Frequency Allocations shall be used as a guide in the assignment of radio frequencies to Government radio stations in the United States and Possessions. Exceptions to the Table may be made by the IRAC after careful consideration to avoid harmful interference and to ensure compliance with the ITU Radio Regulations.

For the use of frequencies by Government radio stations outside the United States and Possessions, Government agencies shall be guided insofar as practicable by the ITU Table of Frequency Allocations and, where applicable, by the authority of the host government. Maximum practicable effort should be made to avoid

the possibility of harmful interference to other authorized US operations. If harmful interference is considered likely, it is incumbent upon the agency conducting the operation to coordinate with other US Flag users, as provided for in Section 8.3.11 of the NTIA Manual.

Application of the US Government Table is subject to the recognition that:

below 25 000 kHz the Table is only applicable in the assignment of frequencies after 5 September 1961;

under Article 48 of the International Telecommunication Constitution, administrations "retain their entire freedom with regard to military radio installations of their army, naval and air forces;" and under No. S4.4 of the ITU Radio Regulations, administrations may assign frequencies in derogation of the ITU Table of Frequency Allocations "on the express condition that harmful interference shall not be caused to services carried on by stations operating in accordance with the provisions of the Convention and of these Regulations."

Some frequency assignments below 25 000 kHz that were made before 5 September 1961, are not in conformity with the Government Table. Because of the exception mentioned in the first subparagraph above, the status of these assignments can be determined only on a case-by-case basis. With this exception, the rules pertaining to the relative status between radio services are as follows:

Stations of secondary services are on a noninterference basis to the primary service; they

- (a) shall not cause harmful interference to stations of primary services to which frequencies are already assigned or to which frequencies may be assigned at a later date;
- (b) cannot claim protection from harmful interference from stations of a primary service to which frequencies are already assigned or may be assigned at a later date;
- (c) can claim protection, however, from harmful interference from stations of the same or other secondary service(s) to which frequencies may be assigned at a later date.

Additional allocation—where a band is indicated in a footnote of the Table as "also allocated" to a service in an area smaller than a Region, or in a particular country. For example, an allocation which is added in this area or in this country to the service or services that are indicated in the Table.

Alternative allocation—where a band is indicated in a footnote of the Table as "allocated" to one or more services in an area smaller than a Region, or in a particular country. For example, an allocation that replaces, in this area or in this country, the allocation indicated in the Table.

Different category of service—where the allocation category (primary or secondary) of the service in

the Table is changed. For example, the Table reflects the allocation as Fixed, Mobile, and RADIOLOCATION, the category of these services are changed by the footnote to FIXED, MOBILE, and Radiolocation.

An allocation or a footnote to the Government Table denoting relative status between radio services automatically applies to each assignment in the band to which the footnote or allocation pertains, unless at the time of a particular frequency assignment action a different provision is decided upon for the assignment concerned.

A priority note reflecting the same provisions as an allocation or an applicable footnote to the Government Table is redundant and shall not be applied to frequency assignments.

An assignment that is in conformity with the service allocation (as amplified by pertinent footnotes) for the band in which it is contained takes precedence over assignments therein that are not in conformity unless, at the time of the frequency assignment action, a different provision is decided upon.

Where in this Table a band is indicated as allocated to more than one service, such services are listed in the following order:

- (a) services, the names of which are printed in all capital letters (example: FIXED); these services are called "primary" services;
- (b) services, the names of which are printed in "normal characters" (example: Mobile); these are "secondary" services.

The columns to the right of the double line show the national provisions; those to the left show the provisions of the ITU Table of Frequency Allocations.

Column 1 indicates the national band limits.

Column 2 indicates the Government allocation, including all "US" and "G" footnotes considered to be applicable to the Government nationally. Where the allocated service is followed by a function in parentheses, e.g., SPACE (space-to-Earth), the allocation is limited to the function shown.

Column 3 indicates the non-Government allocation including all "US" footnotes and certain "NG" footnotes as contained in Part 2 of the FCC Rules and Regulations. Where the allocated service is followed by a function in parentheses, e.g., SPACE (space-to-Earth), the allocation is limited to the function shown. These data have been included in the Government Table for information purposes only.

Column 4 contains such remarks as serve to amplify the Government and non-Government allocations or point to an understanding between the FCC and IRAC/NTIA in respect thereof.

The international footnotes shown in the columns to the left of the double line are applicable only in the relationships between the United States and other countries. An international footnote is applicable to the US Table of Allocations if the number also appears in

Columns 2 and 3 of the US Table. The international footnote is then applicable to both Government and non-Government use.

The texts of footnotes in this Table are listed in numerical order at the end of the Table, in sections headed Government Footnotes, US Footnotes, International Footnotes, and NG Footnotes.

Chart of Regions as Defined in Table of Frequency Allocations—For the allocation of frequencies the world has been divided into three Regions* as shown in Fig. 2 and described in Nos. S5.3 to S5.9.

Government (G) Footnotes

(These footnotes, each consisting of the letter "G" followed by one or more digits, denote stipulations applicable only to the Government.)

G2—In the bands 216–225, 420–450 (except as provided by US217), 890–902, 928–942, 1300–1400, 2310–2390, 2417–2450, 2700–2900, 5650–5925, and 9000–9200 MHz, the Government radiolocation is limited to the military services.

G5—In the bands 162.0125–173.2, 173.4–174, 406.1–410, and 410–420 MHz, the fixed and mobile services are all allocated on a primary basis to the Government nonmilitary agencies.

G6—Military tactical fixed and mobile operations may be conducted nationally on a secondary basis; (1) to the meteorological aids service in the band 403–406 MHz; and (2) to the radio astronomy service in the band 406.1–410 MHz. Such fixed and mobile operations are subject to local coordination to ensure that harmful interference will not be caused to the services to which the bands are allocated.

G8—Low-power Government radio control operations are permitted in the band 420–450 MHz.

G11—Government fixed and mobile radio services including low power radio control operations are permitted in the band 902–928 MHz on a secondary basis.

G15—Use of the band 2700–2900 MHz by the military fixed and shipborne air defense radiolocation installations will be fully coordinated with the meteorological aids and aeronautical radionavigation services. The military air defense installations will be moved from the band 2700–2900 MHz at the earliest practicable date. Until such time as military air defense installations can be accommodated satisfactorily elsewhere in the spectrum, such operations will, insofar as practicable, be adjusted to meet the requirements of the aeronautical radionavigation service.

G19—Use of the band 9000–9200 MHz by military fixed and shipborne air defense radiolocation installations will be fully coordinated with the aeronautical radionavigation service, recognizing fully the safety aspects of the latter. Military air defense installations will be accommodated ultimately outside this band. Until such time as military defense installations can be accommodated satisfactorily elsewhere in the spectrum, such operations will, in-so-far as practicable, be adjusted to meet the requirements of the aeronautical radionavigation service.

* It should be noted that where the words "regions" or "regional" are without a capital "R" in these Regulations, they do not relate to the three Regions here defined for purposes of frequency allocation.

TABLE 3. TABLE OF FREQUENCY ALLOCATIONS*

0-130 kHz (VLF/LF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
Below 9 (Not Allocated) S5.53.S5.54			Below 9 (Not Allocated) S5.53.S5.54		
9-14 RADIONAVIGATION			9-14 RADIONAVIGATION US18 US294		
14-19.95 FIXED MARITIME MOBILE S5.57 S5.55 S5.56			14-19.95 FIXED MARITIME MOBILE S5.57 US294	14-19.95 Fixed US294	International Fixed (23)
19.95-20.05 STANDARD FREQUENCY AND TIME SIGNAL (20 kHz)			19.95-20.05 STANDARD FREQUENCY AND TIME SIGNAL (20 kHz) US294		
20.05-70 FIXED MARITIME MOBILE S5.57 S5.56 S5.58			20.05-59 FIXED MARITIME MOBILE S5.57 US294	20.05-59 FIXED US294	International Fixed (23)
			59-61 STANDARD FREQUENCY AND TIME SIGNAL (60 kHz) US294		
			61-70 FIXED MARITIME MOBILE S5.57 US294	61-70 FIXED US294	International Fixed (23)
70-72 RADIONAVIGATION S5.60 S5.56	70-90 FIXED MARITIME MOBILE S5.57 MARITIME RADIO- NAVIGATION S5.60 Radiolocation	70-72 RADIONAVIGATION S5.60 Fixed Maritime mobile S5.57 S5.59	70-90 FIXED MARITIME MOBILE S5.57 Radiolocation	70-90 FIXED Radiolocation	International Fixed (23) Private Land Mobile (90)
72-84 FIXED MARITIME MOBILE S5.57 RADIONAVIGATION S5.60 S5.56		72-84 FIXED MARITIME MOBILE S5.57 RADIONAVIGATION S5.60			

84-86 RADIONAVIGATION S5.60		84-86 RADIONAVIGATION S5.60 Fixed Maritime mobile S5.57 S5.59		
86-90 FIXED MARITIME MOBILE S5.57 RADIONAVIGATION S5.56	S5.61	86-90 FIXED MARITIME MOBILE S5.57 RADIONAVIGATION S5.60	S5.60 US294	S5.60 US294
90-110 RADIONAVIGATION S5.62 Fixed S5.64			90-110 RADIONAVIGATION S5.62 US18 US104 US294	Private Land Mobile (90)
110-112 FIXED MARITIME MOBILE RADIONAVIGATION S5.64	110-130 FIXED MARITIME MOBILE MARITIME RADIO- NAVIGATION S5.60 Radiolocation	110-112 FIXED MARITIME MOBILE RADIONAVIGATION S5.60 S5.64	110-130 FIXED MARITIME MOBILE Radiolocation	International Fixed (23) Maritime (80) Private Land Mobile (90)
112-115 RADIONAVIGATION S5.60		112-117.6 RADIONAVIGATION S5.60 Fixed Maritime mobile S5.64 S5.65		
115-117.6 RADIONAVIGATION S5.60 Fixed Maritime mobile S5.64 S5.66				
117.6-126 FIXED MARITIME MOBILE RADIONAVIGATION S5.60 S5.64		117.6-126 FIXED MARITIME MOBILE RADIONAVIGATION S5.60 S5.64		
126-129 RADIONAVIGATION S5.60	S5.61 S5.64	126-129 RADIONAVIGATION S5.60 Fixed Maritime mobile S5.64 S5.65	S5.60 S5.64 US294	
See next page for 129-130		See next page for 129-130		

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

130–505 kHz (LF/MF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
129–130 FIXED MARITIME MOBILE RADIONAVIGATION S5.60 S5.64	See previous page for 110–130 kHz	129–130 FIXED MARITIME MOBILE RADIONAVIGATION S5.60 S5.64	See previous page for 110–130 kHz		See previous page for 110–130 kHz
130–148.5 FIXED MARITIME MOBILE S5.64 S5.67	130–160 FIXED MARITIME MOBILE	130–160 FIXED MARITIME MOBILE RADIONAVIGATION	130–160 FIXED MARITIME MOBILE		International Fixed (23) Maritime (80)
148.5–255 BROADCASTING	S5.64	S5.64	S5.64 US294		International Fixed (23)
	160–190 FIXED	160–190 FIXED Aeronautical radionavigation	160–190 FIXED MARITIME MOBILE 459 US294	160–190 FIXED 459 US294	
S5.68 S5.69 S5.70	190–200 AERONAUTICAL RADIONAVIGATION		190–200 AERONAUTICAL RADIONAVIGATION US18 US226 US294		Aviation (87)
	200–275 AERONAUTICAL RADIONAVIGATION Aeronautical mobile	200–285 AERONAUTICAL RADIONAVIGATION Aeronautical mobile	200–275 AERONAUTICAL RADIONAVIGATION Aeronautical mobile US18 US294		
255–283.5 BROADCASTING AERONAUTICAL RADIONAVIGATION S5.70 S5.71	275–285 AERONAUTICAL RADIONAVIGATION Aeronautical mobile Maritime radionavigation (radiobeacons)		275–285 AERONAUTICAL RADIONAVIGATION Aeronautical mobile Maritime radionavigation (radiobeacons) US18 US294		
283.5–315 AERONAUTICAL RADIONAVIGATION MARITIME RADIONAVIGATION (radiobeacons) S5.73 S5.72 S5.74	285–315 AERONAUTICAL RADIONAVIGATION MARITIME RADIONAVIGATION (radiobeacons) S5.73		285–325 MARITIME RADIONAVIGATION (radiobeacons) S5.73 Aeronautical radionavigation (radiobeacons)		

315-325 AERONAUTICAL RADIONAVIGATION Maritime radionavigation (radiobeacons) S5.73 S5.72 S5.75	315-325 MARITIME RADIONAVIGATION (radiobeacons) S5.73 Aeronautical radionavigation	315-325 AERONAUTICAL RADIONAVIGATION MARITIME RADIONAVIGATION (radiobeacons) S5.73	US18 US294		
325-405 AERONAUTICAL RADIONAVIGATION S5.72	325-335 AERONAUTICAL RADIONAVIGATION Aeronautical mobile Maritime radionavigation (radiobeacons)	325-405 AERONAUTICAL RADIONAVIGATION Aeronautical mobile	325-335 AERONAUTICAL RADIONAVIGATION (radiobeacons) Aeronautical mobile Maritime radionavigation (radiobeacons) US18 US294		
	335-405 AERONAUTICAL RADIONAVIGATION Aeronautical mobile		335-405 AERONAUTICAL RADIONAVIGATION (radiobeacons) Aeronautical mobile US18 US294		
405-415 RADIONAVIGATION S5.76 S5.72	405-415 RADIONAVIGATION S5.76 Aeronautical mobile		405-415 RADIONAVIGATION S5.76 Aeronautical mobile US18 US294		Maritime (80) Aviation (87)
415-435 MARITIME MOBILE S5.79 AERONAUTICAL RADIONAVIGATION S5.72	415-495 MARITIME MOBILE S5.79 S5.79A Aeronautical radionavigation S5.80		415-435 MARITIME MOBILE S5.79 AERONAUTICAL RADIONAVIGATION S5.80 US294		
435-495 MARITIME MOBILE S5.79 S5.79A Aeronautical radionavigation S5.72 S5.81 S5.82	S5.77 S5.78 S5.81 S5.82		435-495 MARITIME MOBILE S5.79 Aeronautical radionavigation 471 472A US231 US294	435-495 MARITIME MOBILE S5.79 471 472A US231 US294	Maritime (80)
495-505 MOBILE (distress and calling) S5.83			495-505 MOBILE (distress and calling) 472		

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

505–2107 kHz (MF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
505–526.5 MARITIME MOBILE S5.79 S.579A S5.84 AERONAUTICAL RADIONAVIGATION	505–510 MARITIME MOBILE S5.79	505–526.5 MARITIME MOBILE S5.79 S.579A S5.84 AERONAUTICAL RADIONAVIGATION Aeronautical mobile Land mobile	505–510 MARITIME MOBILE S5.79		Maritime (80)
	S5.81		471		
	510–525 MOBILE S5.79A S5.84 AERONAUTICAL RADIONAVIGATION		510–525 MARITIME MOBILE (ships only) 474 AERONAUTICAL RADIONAVIGATION (radiobeacons)	US14 US18 US225	
S5.72 S5.81	525–535 BROADCASTING S5.86 AERONAUTICAL RADIONAVIGATION	S5.81	525–535 MOBILE US221 AERONAUTICAL RADIONAVIGATION (radiobeacons)		Aviation (87) Private Land Mobile (90)
526.5–1606.5 BROADCASTING	526.5–535 BROADCASTING S5.86 AERONAUTICAL RADIONAVIGATION	526.5–535 BROADCASTING Mobile	US18 US239		
	S5.88				
S5.87 S5.87A	535–1605 BROADCASTING	535–1605 BROADCASTING	535–1605	535–1605 BROADCASTING	Radio Broadcasting (AM) (73) Auxiliary Broadcasting (74) Alaska Fixed (80) Private Land Mobile (90)
	1605–1625 BROADCASTING S5.89		1605–1625 MOBILE US221	1605–1705 BROADCASTING 480	
1606.5–1625 FIXED MARITIME MOBILE S5.90 LAND MOBILE	S5.90	1606.5–1800 FIXED MOBILE RADIOLOCATION RADIONAVIGATION	US238		
			1615–1625		
S5.92			US238 US299		
1625–1635 RADIOLOCATION	1625–1705 FIXED MOBILE BROADCASTING S5.89 Radiolocation		1625–1705 Radiolocation		
S5.93					
1635–1800 FIXED MARITIME MOBILE S5.90 LAND MOBILE	S5.90		US238 US299	US238 US299 US321 NG128	

S5.92 S5.96	1705-1800 FIXED MOBILE RADIOLOCATION AERONAUTICAL RADIO NAVIGATION	S5.91	1705-1800 FIXED MOBILE RADIOLOCATION		International Fixed (23) Maritime (80) Private Land Mobile (90)
1800-1810 RADIOLOCATION S5.93	1800-1850 AMATEUR	1800-2000 AMATEUR FIXED MOBILE except aeronautical mobile RADIO NAVIGATION Radiolocation	1800-1900	1800-1900 AMATEUR	Amateur (97)
1810-1850 AMATEUR S5.98 S5.99 S5.100 S5.101					
1850-2000 FIXED MOBILE except aeronautical mobile S5.92 S5.96 S5.103	1850-2000 AMATEUR FIXED MOBILE except aeronautical mobile RADIOLOCATION RADIO NAVIGATION S5.102	S5.97	1900-2000 RADIOLOCATION		Private Land Mobile (90) Amateur (97)
2000-2025 FIXED MOBILE except aeronautical mobile (R) S5.92 S5.103	2000-2065 FIXED MOBILE		2000-2065 FIXED MOBILE	2000-2065 MARITIME MOBILE NG19	Maritime (80)
2025-2045 FIXED MOBILE except aeronautical mobile (R) Meteorological aids S5.104 S5.92 S5.103					
2045-2160 FIXED MARITIME MOBILE LAND MOBILE S5.92	2065-2107 MARITIME MOBILE S5.105 S5.106		2065-2107 MARITIME MOBILE S5.105 US296 US340		
	See next page for 2107-2170 kHz		See next page for 2107-2170 kHz		See next page for 2107-2170 kHz

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

2107--3230 kHz (MF/HF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
See previous page for 2045-2160 kHz	2107-2170 FIXED MOBILE		2107-2170 FIXED MOBILE	2107-2170 FIXED LAND MOBILE MARITIME MOBILE NG19	International Fixed (23) Maritime (80) Aviation (87) Private Land Mobile (90)
2160-2170 RADIOLOCATION S5.93 S5.107			US340	US340	
2170-2173.5 MARITIME MOBILE			2170-2173.5 MARITIME MOBILE (telephony)	2170-2173.5 MARITIME MOBILE	Maritime (80)
			US340	US340	
2173.5-2190.5 MOBILE (distress and calling)			2173.5-2190.5 MOBILE (distress and calling)		Maritime (80) Aviation (87)
S5.108 S5.109 S5.110 S5.111			S5.108 S5.109 S5.110 S5.111 US279 US340		
2190.5-2194 MARITIME MOBILE			2190.5-2194 MARITIME MOBILE (telephony)	2190.5-2194 MARITIME MOBILE	Maritime (80)
			US340	US340	
2194-2300 FIXED MOBILE except aeronautical mobile (R)	2194-2300 FIXED MOBILE		2194-2495 FIXED MOBILE	2194-2495 FIXED LAND MOBILE MARITIME MOBILE NG19	International Fixed (23) Maritime (80) Aviation (87) Private Land Mobile (90)
S5.92 S5.103 S5.112	S5.112				
2300-2498 FIXED MOBILE except aeronautical mobile (R) BROADCASTING S5.113	2300-2495 FIXED MOBILE BROADCASTING S5.113		US340	US340	
S5.103	2495-2501 STANDARD FREQUENCY AND TIME SIGNAL (2500 KHZ)		2495-2501 STANDARD FREQUENCY AND TIME SIGNAL (2500 KHZ)		
2498-2501 STANDARD FREQUENCY AND TIME SIGNAL (2500 KHZ)			US340		

2501–2502 STANDARD FREQUENCY AND TIME SIGNAL Space research		2501–2502 STANDARD FREQUENCY AND TIME SIGNAL US340 G106	2501–2502 STANDARD FREQUENCY AND TIME SIGNAL US340	
2502–2625 FIXED MOBILE except aeronautical mobile (R) S5.92 S5.103 S5.114	2502–2505 STANDARD FREQUENCY AND TIME SIGNAL	2502–2505 STANDARD FREQUENCY AND TIME SIGNAL US340		
	2505–2850 FIXED MOBILE	2505–2850 FIXED MOBILE	2505–2850 FIXED LAND MOBILE MARITIME MOBILE	International Fixed (23) Maritime (80) Aviation (87) Private Land Mobile (90)
		2625–2650 MARITIME MOBILE MARITIME RADIONAVIGATION S5.92		
2650–2850 FIXED MOBILE except aeronautical mobile (R) S5.92 S5.103		US285 US340	US285 US340	
2850–3025 AERONAUTICAL MOBILE (R) S5.111 S5.115		2850–3025 AERONAUTICAL MOBILE (R) S5.111 S5.115 US283 US340		Aviation (87)
3025–3155 AERONAUTICAL MOBILE (OR)		3025–3155 AERONAUTICAL MOBILE (OR) US340		
3155–3200 FIXED MOBILE except aeronautical mobile (R) S5.116 S5.117	3200–3230 FIXED MOBILE except aeronautical mobile (R) BROADCASTING S5.113 S5.116	3155–3230 FIXED MOBILE except aeronautical mobile (R)		International Fixed (23) Maritime (80) Aviation (87) Private Land Mobile (90)
		US340		

TABLE 3 (CONT.) TABLE OF FREQUENCY ALLOCATIONS*

1230-5000 KHz (HF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
1230-1400 FIXED MOBILE except aeronautical mobile BROADCASTING SS 117			1230-1400 FIXED MOBILE except aeronautical mobile Radioeducation		International Fixed (123) Maritime (300) Aviation (87) Private Land Mobile (90)
SS 116 SS 118			US 340		
1400-1500 AERONAUTICAL MOBILE (R)			1400-1500 AERONAUTICAL MOBILE (R)		Aviation (87)
			US 264 US 340		
1500-1750 AMATEUR SS 120 FIXED MOBILE except aeronautical mobile	1500-1750 AMATEUR SS 120 SS 119	1500-1900 AMATEUR SS 120 FIXED MOBILE	1500-4000	1500-4000 AMATEUR SS 120	Aviation (97)
N 002	1750-4000 AMATEUR SS 120				
1600-1600 FIXED AERONAUTICAL MOBILE (R) LAND MOBILE	FIXED MOBILE except aeronautical mobile (R)				
1600-1650 AERONAUTICAL MOBILE (R)		1600-1650 AERONAUTICAL MOBILE BROADCASTING			
SS 123					
1950-4000 FIXED BROADCASTING	SS 122 SS 124 SS 125	1950-4000 FIXED BROADCASTING SS 126	US 340	US 340	
4000-4363 FIXED MARTIME MOBILE SS 127			4000-4363 MARTIME MOBILE SS 127		International Fixed (124) Maritime (300)
SS 126			US 236 US 340		
4063-4438 MARTIME MOBILE SS 109A SS 109 SS 110 SS 130 SS 131 SS 132			4063-4438 MARTIME MOBILE SS 109 SS 110 SS 130 SS 132		
SS 128 SS 129			US 82 US 296 US 340		

4438-4650 FIXED MOBILE except aeronautical mobile (R)	4438-4650 FIXED MOBILE except aeronautical mobile (R)	4438-4650 FIXED MOBILE except aeronautical mobile (R) US340	International Fixed (23) Maritime (80) Aviation (87) Private Land Mobile (90)
4650-4700 AERONAUTICAL MOBILE (R)		4650-4700 AERONAUTICAL MOBILE (R) US282 US283 US340	Aviation (87)
4700-4750 AERONAUTICAL MOBILE (OR)		4700-4750 AERONAUTICAL MOBILE (OR) US340	
4750-4850 FIXED AERONAUTICAL MOBILE (OR) LAND MOBILE BROADCASTING S5.113	4750-4850 FIXED MOBILE except aeronautical mobile (R) BROADCASTING S5.113	4750-4850 FIXED BROADCASTING S5.113 Land mobile US340	International Fixed (23) Maritime (80) Aviation (87)
4850-4995 FIXED LAND MOBILE BROADCASTING S5.113		4850-4995 FIXED MOBILE US340	4850-4995 FIXED US340
4995-5003 STANDARD FREQUENCY AND TIME SIGNAL (5000 kHz)		4995-5003 STANDARD FREQUENCY AND TIME SIGNAL (5000 kHz) US340	
5003-5005 STANDARD FREQUENCY AND TIME SIGNAL Space research		5003-5005 STANDARD FREQUENCY AND TIME SIGNAL US340 G106	5003-5005 STANDARD FREQUENCY AND TIME SIGNAL US340
5005-5060 FIXED BROADCASTING S5.113		5005-5060 FIXED US340	
		International Fixed (23) Maritime (80) Aviation (87) Private Land Mobile (90)	

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

5060-9040 kHz (HF)					Remarks
International Table			United States Table		
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
5060-5250 FIXED Mobile except aeronautical mobile S5.133			5060-5450 FIXED Mobile except aeronautical mobile US212 US340		International Fixed (23) Maritime (80) Aviation (87) Private Land Mobile (90)
5250-5450 FIXED MOBILE except aeronautical mobile					
5450-5480 FIXED AERONAUTICAL MOBILE (OR) LAND MOBILE	5450-5480 FIXED AERONAUTICAL MOBILE (R)	5450-5480 FIXED AERONAUTICAL MOBILE (OR) LAND MOBILE	5450-5480 AERONAUTICAL MOBILE (R) US283 US340		Aviation (87)
5480-5680 AERONAUTICAL MOBILE (R) S5.111 S5.115			5480-5680 AERONAUTICAL MOBILE (R) S5.111 S5.115 US283 US340		
5680-5730 AERONAUTICAL MOBILE (OR) S5.111 S5.115			5680-5730 AERONAUTICAL MOBILE (OR) S5.111 S5.115 US340		
5730-5900 FIXED LAND MOBILE	5730-5900 FIXED MOBILE except aeronautical mobile (R)	5730-5900 FIXED Mobile except aeronautical mobile (R)	5730-5950 FIXED MOBILE except aeronautical mobile (R) US340		International Fixed (23) Maritime (80) Aviation (87)
5900-5950 BROADCASTING S5.134 S5.136					
5950-6200 BROADCASTING			5950-6200 BROADCASTING US340		Radio Broadcast (HF) (73)
6200-6525 MARITIME MOBILE S5.109 S5.110 S5.130 S5.132 S5.137			6200-6525 MARITIME MOBILE S5.109 S5.110 S5.130 S5.132 US82 US296 US340		Maritime (80)
6525-6685 AERONAUTICAL MOBILE (R)			6525-6685 AERONAUTICAL MOBILE (R) US283 US340		Aviation (87)

6685-6765 AERONAUTICAL MOBILE (OR)			6685-6765 AERONAUTICAL MOBILE (OR) US340		
6765-7000 FIXED Land mobile S5.139 S5.138			6765-7000 FIXED Mobile S5.139 S5.138 US340		ISM Equipment (18) International Fixed (23) Aviation (87)
7000-7100 AMATEUR S5.120 AMATEUR-SATELLITE S5.140 S5.141			7000-7100 US340	7000-7100 AMATEUR S5.120 AMATEUR-SATELLITE US340	Amateur (97)
7100-7300 BROADCASTING	7100-7300 AMATEUR S5.120 S5.142	7100-7300 BROADCASTING	7100-7300 US340	7100-7300 AMATEUR S5.120 S5.142 US340	
7300-7350 BROADCASTING S5.134 S5.143			7300-8100 FIXED Mobile US340		International Fixed (23) Maritime (80) Aviation (87) Private Land Mobile (90)
7350-8100 FIXED Land mobile S5.144					
8100-8195 FIXED MARITIME MOBILE			8100-8195 MARITIME MOBILE US236 US340		Maritime (80)
8195-8815 MARITIME MOBILE S5.109 S5.110 S5.132 S5.145 S5.111			8195-8815 MARITIME MOBILE S5.109 S5.110 S5.132 S5.145 S5.111 US82 US296 US340		
8815-8965 AERONAUTICAL MOBILE (R)			8815-8965 AERONAUTICAL MOBILE (R) US340		Aviation (87)
8965-9040 AERONAUTICAL MOBILE (OR)			8965-9040 AERONAUTICAL MOBILE (OR) US340		

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

9040-13410 kHz (HF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
9040-9400 FIXED			9040-9500 FIXED		International Fixed (23) Maritime (80) Aviation (87)
9400-9500 BROADCASTING S5.134 S5.146			US340		
9500-9900 BROADCASTING S5.147			9500-9900 BROADCASTING S5.147 S5.148 US235 US340		International Fixed (23) Radio Broadcast (HF) (73)
9900-9995 FIXED			9900-9995 FIXED US340		International Fixed (23) Aviation (87)
9995-10003 STANDARD FREQUENCY AND TIME SIGNAL (10000kHz) S5.111			9995-10003 STANDARD FREQUENCY AND TIME SIGNAL (10000kHz) S5.111 US340		
10003-10005 STANDARD FREQUENCY AND TIME SIGNAL Space research S5.111			10003-10005 STANDARD FREQUENCY AND TIME SIGNAL S5.111 US340 G106	10003-10005 STANDARD FREQUENCY AND TIME SIGNAL S5.111 US340	
10005-10100 AERONAUTICAL MOBILE (R) S5.111			10005-10100 AERONAUTICAL MOBILE (R) S5.111 US283 US340		Aviation (87)
10100-10150 FIXED Amateur S5.120			10100-10150 US247 US340	10100-10150 AMATEUR S5.120 US247 US340	Amateur (97)
10150-11175 FIXED Mobile except aeronautical mobile (R)			10150-11175 FIXED Mobile except aeronautical mobile (R) US340		International Fixed (23) Aviation (87)
11175-11275 AERONAUTICAL MOBILE (OR)			11175-11275 AERONAUTICAL MOBILE (OR) US340		

11275–11400 AERONAUTICAL MOBILE (R)	11275–11400 AERONAUTICAL MOBILE (R) US283 US340	Aviation (87)
11400–11600 FIXED	11400–11650 FIXED	International Fixed (23) Aviation (87)
11600–11650 BROADCASTING S5.134 S5.146	US340	
11650–12050 BROADCASTING S5.147	11650–12050 BROADCASTING US235 US340	International Fixed (23) Radio Broadcast (HF) (73)
12050–12100 BROADCASTING S5.134 S5.146	12050–12230 FIXED	International Fixed (23) Aviation (87)
12100–12230 FIXED	US340	
12230–13200 MARITIME MOBILE S5.109 S5.110 S5.132 S5.145	12230–13200 MARITIME MOBILE S5.109 S5.110 S5.132 S5.145 US82 US296 US340	International Fixed (23) Maritime (80)
13200–13260 AERONAUTICAL MOBILE (OR)	13200–13260 AERONAUTICAL MOBILE (OR) US340	
13260–13360 AERONAUTICAL MOBILE (R)	13260–13360 AERONAUTICAL MOBILE (R) US283 US340	Aviation (87)
13360–13410 FIXED RADIO ASTRONOMY S5.149	13360–13410 RADIO ASTRONOMY S5.149 G115	13360–13410 RADIO ASTRONOMY S5.149

TABLE 3 (CONT.) TABLE OF FREQUENCY ALLOCATIONS²

13410-17900 kHz (HF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
13410-13570 FIXED Mobile except aeronautical mobile (R)			13410-13570 FIXED Mobile except aeronautical mobile (R)	13410-13570 FIXED	ISM Equipment (18) International Fixed (23) Aviation (87)
SS.150			SS.150 US340	SS.150 US340	
13570-13600 BROADCASTING SS.134			13570-13600 FIXED Mobile except aeronautical mobile (R)	13570-13600 FIXED	International Fixed (23) Aviation (87)
SS.151			US340	US340	
13600-13800 BROADCASTING			13600-13800 BROADCASTING		International Fixed (23) Radio Broadcast (HF) (73)
			SS.148 US340		
13800-13870 BROADCASTING SS.134			13800-14000 FIXED Mobile except aeronautical mobile (R)	13800-14000 FIXED	International Fixed (23) Aviation (87)
SS.151					
13870-14000 FIXED Mobile except aeronautical mobile (R)			US340	US340	
14000-14250 AMATEUR SS.120 AMATEUR SATELLITE			14000-14350	14000-14250 AMATEUR SS.120 AMATEUR SATELLITE	Amateur (97)
				US340	
14250-14350 AMATEUR SS.120			US340	14250-14350 AMATEUR SS.120	
SS.152				US340	
14350-14990 FIXED Mobile except aeronautical mobile (R)			14350-14990 FIXED Mobile except aeronautical mobile (R)	14350-14990 FIXED	International Fixed (23) Aviation (87)
			US340	US340	
14990-15005 STANDARD FREQUENCY AND TIME SIGNAL (15000 kHz)			14990-15005 STANDARD FREQUENCY AND TIME SIGNAL (15000 kHz)		
SS.111			SS.111 US340		

15005-15010 STANDARD FREQUENCY AND TIME SIGNAL Space research	15005-15010 STANDARD FREQUENCY AND TIME SIGNAL US340 G106	15005-15010 STANDARD FREQUENCY AND TIME SIGNAL US340	
15010-15100 AERONAUTICAL MOBILE (OR)	15010-15100 AERONAUTICAL MOBILE (OR) US340		
15100-15600 BROADCASTING	15100-15600 BROADCASTING S5.148 US340	International Fixed (23) Radio Broadcast (HF) (73)	
15600-15800 BROADCASTING S5.134 S5.146	15600-16360 FIXED US340	International Fixed (23) Aviation (87)	
15800-16360 FIXED S5.153			
16360-17410 MARITIME MOBILE S5.109 S5.110 S.5132 S5.145	16360-17410 MARITIME MOBILE S5.109 S5.110 S.5132 S5.145 US82 US296 US340	Maritime (80)	
17410-17480 FIXED	17410-17550 FIXED US340	International Fixed (23) Aviation (87)	
17480-17550 BROADCASTING S5.134 S5.146			
17550-17900 BROADCASTING	17550-17900 BROADCASTING S5.148 US340	International Fixed (23) Radio Broadcast (HF) (73)	

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

17900-2285 kHz (HF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
17900-17970 AERONAUTICAL MOBILE (R)			17900-17970 AERONAUTICAL MOBILE (R) US283 US340		Aviation (87)
17970-18030 AERONAUTICAL MOBILE (OR)			17970-18030 AERONAUTICAL MOBILE (OR) US340		
18030-18052 FIXED			18030-18068 FIXED		International Fixed (23) Maritime (80)
18052-18068 FIXED Space research			US340		
18068-18168 AMATEUR S5.120 AMATEUR-SATELLITE S5.154			18068-18168 US340	18068-18168 AMATEUR S5.120 AMATEUR-SATELLITE US340	International Fixed (23) Amateur (97)
18168-18780 FIXED Mobile except aeronautical mobile			18168-18780 FIXED Mobile US340		International Fixed (23) Maritime (80) Aviation (87)
18780-18900 MARITIME MOBILE			18780-18900 MARITIME MOBILE US82 US296 US340		International Fixed (23) Maritime (80)
18900-19020 BROADCASTING S5.134 S5.146			18900-19680 FIXED Mobile		International Fixed (23) Aviation (87)
19020-19680 FIXED			US340		
19680-19800 MARITIME MOBILE S5.132			19680-19800 MARITIME MOBILE S5.132 US340		Maritime (80)
19800-19990 FIXED			19800-19990 FIXED US340		International Fixed (23) Aviation (87)

19990-19995 STANDARD FREQUENCY AND TIME SIGNAL Space research S5.111	19990-19995 STANDARD FREQUENCY AND TIME SIGNAL Space research G106 S5.111 US340	19990-19995 STANDARD FREQUENCY AND TIME SIGNAL Space research S5.111 US340	
19995-20010 STANDARD FREQUENCY AND TIME SIGNAL (20000 kHz) S5.111	19995-20010 STANDARD FREQUENCY AND TIME SIGNAL (20000 kHz) S5.111 US340 G106	19995-20010 STANDARD FREQUENCY AND TIME SIGNAL (20000 kHz) S5.111 US340	
20010-21000 FIXED Mobile	20010-21000 FIXED Mobile US340	20010-21000 FIXED US340	
21000-21450 AMATEUR S5.120 AMATEUR-SATELLITE	21000-21450 US340	21000-21450 AMATEUR S5.120 AMATEUR-SATELLITE US340	Amateur (97)
21450-21850 BROADCASTING	21450-21850 BROADCASTING S5.148 US340		International Fixed (23) Radio Broadcast (HF) (73)
21850-21870 FIXED S5.155A S5.155	21850-21924 FIXED US340		International Fixed (23) Aviation (87)
21870-21924 FIXED S5.155B			
21924-22000 AERONAUTICAL MOBILE (R)	21924-22000 AERONAUTICAL MOBILE (R) US340		Aviation (87)
22000-22855 MARITIME MOBILE S5.132 S5.156	22000-22855 MARITIME MOBILE S5.132 US82 US296 US340		International Fixed (23) Maritime (80)

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

22855-26175 kHz (HF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
22855-23000 FIXED S5.156			22855-23000 FIXED US340		International Fixed (23) Aviation (87)
23000-23200 FIXED Mobile except aeronautical mobile (R) S5.156			23000-23200 FIXED Mobile except aeronautical mobile (R) US340	23000-23200 FIXED US340	
23200-23350 FIXED S5.156A AERONAUTICAL MOBILE (OR)			23200-23350 AERONAUTICAL MOBILE (OR) US340		
23350-24000 FIXED MOBILE except aeronautical mobile S5.157			23350-24890 FIXED MOBILE except aeronautical mobile US340	23350-24890 FIXED US340	International Fixed (23) Aviation (87)
24000-24890 FIXED LAND MOBILE					
24890-24990 AMATEUR S5.120 AMATEUR-SATELLITE			24890-24990 US340	24890-24990 AMATEUR S5.120 AMATEUR-SATELLITE US340	Amateur (97)
24990-25005 STANDARD FREQUENCY AND TIME SIGNAL (25000 kHz)			24990-25005 STANDARD FREQUENCY AND TIME SIGNAL (25000 kHz) US340		
25005-25010 STANDARD FREQUENCY AND TIME SIGNAL Space research			25005-25010 STANDARD FREQUENCY AND TIME SIGNAL US340 G106	25005-25010 STANDARD FREQUENCY AND TIME SIGNAL US340	
25010-25070 FIXED MOBILE except aeronautical mobile			25010-25070 US340	25010-25070 LAND MOBILE US340 NG112	Private Land Mobile (90)

25070–25210 MARITIME MOBILE	25070–25210 MARITIME MOBILE US82 US281 US296 US340	25070–25210 MARITIME MOBILE US82 US281 US296 US340 NG112	Maritime (80) Private Land Mobile (90)
25210–25550 FIXED MOBILE except aeronautical mobile	25210–25330 US340	25210–25330 LAND MOBILE US340	Private Land Mobile (90)
	25330–25550 FIXED MOBILE except aeronautical mobile US340	25330–25550 US340	
25550–25670 RADIO ASTRONOMY S5.149	25550–25670 RADIO ASTRONOMY US74 S5.149		
25670–26100 BROADCASTING	25670–26100 BROADCASTING US25 US340		Radio Broadcast (HF) (73) Remote Pickup (74D)
26100–26175 MARITIME MOBILE S5.132	26100–26175 MARITIME MOBILE S5.132 US340		Auxiliary Broadcasting (74) Maritime (80)

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

26175 kHz–33MHz (HF/VHF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
26175–27500 FIXED MOBILE except aeronautical mobile			26175–26480 US340	26175–26480 LAND MOBILE US340	Auxiliary Broadcasting (74)
			26480–26950 FIXED MOBILE except aeronautical mobile US10 US340	26480–26950 US10 US340	
			26950–27410 S5.150 US340	26950–26960 FIXED S5.150 US340	ISM Equipment (18) International Fixed (23)
				26960–27230 MOBILE except aeronautical mobile S5.150 US340	ISM Equipment (18) Personal Radio (95)
			27230–27410 FIXED MOBILE except aeronautical mobile S5.150 US340	ISM Equipment (18) Private Land Mobile (90) Personal Radio (95)	
27500–28000 METEOROLOGICAL AIDS FIXED MOBILE			27410–27540 US340	27410–27540 FIXED LAND MOBILE US340	Private Land Mobile (90)
			27540–28000 FIXED MOBILE US298 US340	27540–28000 US298 US340	

28-29.7 AMATEUR AMATEUR-SATELLITE	28-29.89	28-29.7 AMATEUR AMATEUR-SATELLITE US340	Amateur (97)
29.7-30.005 FIXED MOBILE	US340	29.7-29.8 LAND MOBILE US340	Private Land Mobile (90)
	29.8-29.89 FIXED US340	29.8-29.89 FIXED US340	International Fixed (23) Aviation (87)
	29.89-29.91 FIXED MOBILE US340	29.89-29.91 US340	
	29.91-30 US340	29.91-30 FIXED US340	International Fixed (23) Aviation (87)
30.005-30.01 SPACE OPERATION (satellite identification) FIXED MOBILE SPACE RESEARCH	30-30.56 FIXED MOBILE	30-30.56	
30.01-37.5 FIXED MOBILE	30.56-32	30.56-32 FIXED LAND MOBILE NG124	Private Land Mobile (90)
	32-33 FIXED MOBILE	32-33	
	See next page for 33-37.5 MHz		See next page for 33-37.5 MHz

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

33-50 MHz (VHF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
See previous page for 30.01-37.5 MHz			33-34	33-34 FIXED LAND MOBILE NG124	Private Land Mobile (90)
			34-35 FIXED MOBILE	34-35	
			35-36	35-36 FIXED LAND MOBILE	Public Mobile (22) Private Land Mobile (90)
			36-37 FIXED MOBILE US220	36-37 US220	
			37-37.5	37-37.5 LAND MOBILE NG124	Private Land Mobile (90)
			37.5-38.25 FIXED MOBILE Radio astronomy	37.5-38 Radio astronomy S5.149	37.5-38 LAND MOBILE Radio astronomy S5.149 NG59 NG124
S5.149	38-38.25 FIXED MOBILE RADIO ASTRONOMY S5.149 US81	38-38.25 RADIO ASTRONOMY S5.149 US81			
38.25-39.986 FIXED MOBILE	38.25-39 FIXED MOBILE 39-40	38.25-39 39-40 LAND MOBILE NG124	Private Land Mobile (90)		
39.986-40.02 FIXED MOBILE Space research	40-42 FIXED MOBILE	40-40.98	ISM Equipment (18)		

40.02-40.98 FIXED MOBILE S5.150				S5.150 US210	
40.98-41.015 FIXED MOBILE Space research S5.160 S5.161				40.98-42	
41.015-44 FIXED MOBILE S5.160 S5.161			S5.150 US210 US220	US220	
44-47 FIXED MOBILE S5.162 S5.162A			42-46.6	42-43.69 FIXED LAND MOBILE NG124 NG141	Public Mobile (22) Private Land Mobile (90)
				43.69-46.6 LAND MOBILE NG124 NG141	Private Land Mobile (90)
47-68 BROADCASTING S5.162A S5.163 S5.164 S5.165 S5.169 S5.171			47-50 FIXED MOBILE	47-50 FIXED MOBILE BROADCASTING	
			47-49.6	47-49.6 LAND MOBILE NG124	Private Land Mobile (90)
			49.6-50 FIXED MOBILE	49.6-50	
See next page for 50-68 MHz			See next page for 50-73 MHz	See next page for 50-72 MHz	See next page for 50-72 MHz

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

50-128.8125 MHz (VHF)						
International Table			United States Table		Remarks	
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government		
See previous page for 47-68 MHz	50-54 AMATEUR S5.166 S5.167 S5.168 S5.170		50-73	50-54 AMATEUR	Amateur (97)	
See previous page for 47-68 MHz	54-68 BROADCASTING Fixed Mobile S5.172	54-68 FIXED MOBILE BROADCASTING		54-72 BROADCASTING NG128 NG149	Broadcast Radio (TV) (73) Auxiliary Broadcasting (74)	
68-74.8 FIXED MOBILE except aeronautical mobile S5.149 S5.174 S5.175 S5.177 S5.179	68-72 BROADCASTING Fixed Mobile S5.173	68-74.8 FIXED MOBILE S5.149 S5.176 S5.179		72-73 FIXED MOBILE NG3 NG49 NG56	Public Mobile (22) Private Land Mobile (90) Personal Radio (95)	
	72-73 FIXED MOBILE					
	73-74.6 RADIO ASTRONOMY S5.178		73-74.6 RADIO ASTRONOMY US74			
	74.6-74.8 FIXED MOBILE		74.6-74.8 FIXED MOBILE US273			Private Land Mobile (90)
74.8-75.2 AERONAUTICAL RADIONAVIGATION S5.180 S5.181			74.8-75.2 AERONAUTICAL RADIONAVIGATION S5.180		Aviation (87)	
75.2-87.5 FIXED MOBILE except aeronautical mobile	75.2-75.4 FIXED MOBILE S5.179		75.2-75.4 FIXED MOBILE US273		Private Land Mobile (90)	

S5.175 S5.179 S5.184 S5.187	75.4-76 FIXED MOBILE	75.4-87 FIXED MOBILE	75.4-88	75.4-76 FIXED MOBILE	Public Mobile (22) Private Land Mobile (90) Personal Radio (95)
	76-88 BROADCASTING Fixed Mobile	S5.149 S5.182 S5.183 S5.188		76-88 BROADCASTING	Broadcast Radio (TV) (73) Auxiliary Broadcasting (74)
87.5-100 BROADCASTING S5.190	S5.185	87-100 FIXED MOBILE BROADCASTING	88-108	88-108 BROADCASTING	Broadcast Radio (FM) (73) Auxiliary Broadcasting (74)
100-108 BROADCASTING S5.192 S5.194	88-100 BROADCASTING			US93	US93 NG2 NG128 NG129
108-117.975 AERONAUTICAL RADIONAVIGATION S5.197			108-117.975 AERONAUTICAL RADIONAVIGATION US93 G126	Note: Footnote G126 states DGPS stations may be authorized on a primary basis in the 108-117.975 MHz band, but the FCC has not yet adopted this footnote.	
117.975-137 AERONAUTICAL MOBILE (R) S5.111 S5.198 S5.199 S5.200 S5.201 S5.202 S5.203 S5.203A S5.203B	117.975-121.9375 AERONAUTICAL MOBILE (R) S5.111 S5.199 S5.200 591 US26 US28		121.9375-123.0875		121.9375-123.0875 AERONAUTICAL MOBILE
	591 US30 US31 US33 US80 US102 US213		591 US30 US31 US33 US80 US102 US213		
	123.0875-123.5875 AERONAUTICAL MOBILE S5.200 591 US32 US33 US112		123.5875-128.8125 AERONAUTICAL MOBILE (R) 591 US26		
	See next page for 128.8125-137 MHz				See next page for 128.8125-137 MHz

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

128.8125–148 MHz (VHF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
See previous page for 117.975–137 MHz			128.8125–132.0125	128.8125–132.0125 AERONAUTICAL MOBILE (R)	Aviation (87)
			591	591	
			132.0125–136.00 AERONAUTICAL MOBILE (R)		
			591 US26		
			136–137	136–137 AERONAUTICAL MOBILE (R)	Satellite Communications (25) Aviation (87)
			591 US244	591 US244	
137–137.025 SPACE OPERATION (space-to-Earth) METEOROLOGICAL-SATELLITE (space-to-Earth) MOBILE-SATELLITE (space-to-Earth) S5.208A S5.209 SPACE RESEARCH (space-to-Earth) Fixed Mobile except aeronautical mobile (R) S5.204 S5.205 S5.206 S5.207 S5.208			137–137.025 SPACE OPERATION (space-to-Earth) METEOROLOGICAL-SATELLITE (space-to-Earth) MOBILE-SATELLITE (space-to-Earth) 599B US318 US319 US320 SPACE RESEARCH (space-to-Earth) 599A		Satellite Communications (25)
137.025–137.175 SPACE OPERATION (space-to-Earth) METEOROLOGICAL-SATELLITE (space-to-Earth) SPACE RESEARCH (space-to-Earth) Fixed Mobile-satellite (space-to-Earth) S5.208A S5.209 Mobile except aeronautical mobile (R) S5.204 S5.205 S5.206 S5.207 S5.208			137.025–137.175 SPACE OPERATION (space-to-Earth) METEOROLOGICAL-SATELLITE (space-to-Earth) SPACE RESEARCH (space-to-Earth) Mobile-satellite (space-to-Earth) 599B US318 US319 US320 599A		
137.175–137.825 SPACE OPERATION (space-to-Earth) METEOROLOGICAL-SATELLITE (space-to-Earth) MOBILE-SATELLITE (space-to-Earth) S5.208A S5.209 SPACE RESEARCH (space-to-Earth) Fixed Mobile except aeronautical mobile (R) S5.204 S5.205 S5.206 S5.207 S5.208			137.175–137.825 SPACE OPERATION (space-to-Earth) METEOROLOGICAL-SATELLITE (space-to-Earth) MOBILE-SATELLITE (space-to-Earth) 599B US318 US319 US320 SPACE RESEARCH (space-to-Earth) 599A		

137.825–138 SPACE OPERATION (space-to-Earth) METEOROLOGICAL-SATELLITE (space-to-Earth) SPACE RESEARCH (space-to-Earth) Fixed Mobile-satellite (space-to-Earth) S5.208A S5.209 Mobile except aeronautical mobile (R)			137.825–138 SPACE OPERATION (space-to-Earth) METEOROLOGICAL-SATELLITE (space-to-Earth) SPACE RESEARCH (space-to-Earth) Mobile-satellite (space-to-Earth) 599B US318 US319 US320		
S5.204 S5.205 S5.206 S5.207 S5.208			599A		
138–143.6 AERONAUTICAL MOBILE (OR)	138–143.6 FIXED MOBILE RADIOLOCATION Space research (space-to-Earth)	138–143.6 FIXED MOBILE Space research (space-to-Earth)	138–144 FIXED MOBILE	138–144	
S5.210 S5.211 S5.212 S5.214	S5.207 S5.213	S5.207 S5.213			
143.6–143.65 AERONAUTICAL MOBILE (OR) SPACE RESEARCH (space-to-Earth)	143.6–143.65 FIXED MOBILE RADIOLOCATION SPACE RESEARCH (space-to-Earth)	143.6–143.65 FIXED MOBILE SPACE RESEARCH (space-to-Earth)			
S5.211 S5.212 S5.214	S5.207 S5.213	S5.207 S5.213	US10 G30	US10	
143.65–144 AERONAUTICAL MOBILE (OR)	143.65–144 FIXED MOBILE RADIOLOCATION Space research (space-to-Earth)	143.65–144 FIXED MOBILE Space research (space-to-Earth)			
S5.210 S5.211 S5.212 S5.214	S5.207 S5.213	S5.207 S5.213	144–148	144–146 AMATEUR S510 AMATEUR-SATELLITE	Amateur (97)
144–146 AMATEUR S5.120 AMATEUR-SATELLITE	S5.216	S5.216			
146–148 FIXED MOBILE except aeronautical mobile (R)	146–148 AMATEUR	146–148 AMATEUR FIXED MOBILE		146–148 AMATEUR	
	S5.217	S5.217			

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

148-162.0125 MHz (VHF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
148-149.9 FIXED MOBILE except aeronautical mobile (R) MOBILE-SATELLITE (Earth-to-space) S5.209 S5.218 S5.219 S5.221	148-149.9 FIXED MOBILE MOBILE-SATELLITE (Earth-to-space) S5.209 S5.218 S5.219 S5.221		148-149.9 FIXED MOBILE MOBILE-SATELLITE (Earth-to-space) 599B US319 US320 US323 US325 S5.218 608A US10 G30	148-149.9 MOBILE-SATELLITE (Earth-to-space) 599B US319 US320 US323 US325 S5.218 608A US10	Satellite Communications (25)
149.9-150.05 MOBILE-SATELLITE (Earth-to-space) S5.209 S5.224A RADIONAVIGATION-SATELLITE S5.224B S5.220 S5.222 S5.223			149.9-150.05 MOBILE-SATELLITE (Earth-to-space) 599B US319 US322 RADIONAVIGATION-SATELLITE S5.223 608B		
150.05-153 FIXED MOBILE except aeronautical mobile RADIO ASTRONOMY S5.149	150.05-156.7625 FIXED MOBILE S5.225 S5.226 S5.227		150.05-150.8 FIXED MOBILE US216 G30	150.05-150.8 US216	
153-154 FIXED MOBILE except aeronautical mobile (R) Meteorological aids			150.8-152.855 US216	150.8-152.855 FIXED LAND MOBILE US216 NG4 NG51 NG112 NG124	Public Mobile (22) Private Land Mobile (90)
154-156.7625 FIXED MOBILE except aeronautical mobile (R)			152.855-154 S5.226	152.855-154 LAND MOBILE NG4 NG124	Auxiliary Broadcasting (74) Private Land Mobile (90)
S5.226 S5.227			154-156.2475 S5.226	154-156.2475 FIXED LAND MOBILE S5.226 NG112 NG117 NG124 NG148	Maritime (80) Private Land Mobile (90)
			156.2475-157.0375	156.2475-157.0375 MARITIME MOBILE	

156.7625–156.8375 MARITIME MOBILE (distress and calling) S5.111 S5.226				
156.8375–174 FIXED MOBILE except aeronautical mobile	156.8375–174 FIXED MOBILE	S5.226 S5.227 US77 US106 US107 US266	S5.226 S5.227 US77 US106 US107 US266 NG117	
		157.0375–157.1875 MARITIME MOBILE	157.0375–157.1875	Private Land Mobile (90)
		S5.226 US214 US266 G109	S5.226 US214 US266	
		157.1875–157.45	157.1875–157.45 LAND MOBILE MARITIME MOBILE	Maritime (80) Private Land Mobile (90)
		S5.226 US223 US266	S5.226 US223 US266 NG111	
		157.45–161.575	157.45–161.575 FIXED LAND MOBILE	Public Mobile (22) Maritime (80) Private Land Mobile (90)
		S5.226 US266	S5.226 US266 NG6 NG28 NG70 NG111 NG112 NG124 NG148 NG155	
		161.575–161.625	161.575–161.625 MARITIME MOBILE	Public Mobile (22) Maritime (80)
S5.226 US77	S5.226 US77 NG6 NG17			
161.625–161.775	161.625–161.775 LAND MOBILE	Public Mobile (22) Auxiliary Broadcasting (74)		
S5.226	S5.226 NG6			
161.775–162.0125	161.775–162.0125 LAND MOBILE MARITIME MOBILE	Public Mobile (22) Maritime (80) Private Land Mobile (90)		
S5.226 US266	S5.226 US266 NG6			
S5.226 S5.229	S5.226 S5.230 S5.231 S5.232	See next page for 162.0125–174 MHz		See next page for 162.0125–174 MHz

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

162.0125–322 MHz (VHF/UHF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
See previous page for 156.8375–174 MHz			162.0125–173.2 FIXED MOBILE	162.0125–173.2	Auxiliary Broadcasting (74) Private Land Mobile (90)
			S5.226 US8 US11 US13 US216 US223 US300 US312 G5	S5.226 US8 US11 US13 US216 US223 US300 US312	
			173.2–173.4	173.2–173.4 FIXED Land mobile	Private Land Mobile (90)
			173.4–174 FIXED MOBILE G5	173.4–174	
174–223 BROADCASTING	174–216 BROADCASTING Fixed Mobile S5.234	174–223 FIXED MOBILE BROADCASTING	174–216	174–216 BROADCASTING NG115 NG128 NG149	Broadcast Radio (TV) (73) Auxiliary Broadcasting (74)
	216–220 FIXED MARITIME MOBILE Radiolocation S5.241		216–220 MARITIME MOBILE Fixed Radiolocation S5.241 G2 Aeronautical mobile Land mobile US210 US229 US274 US317	216–220 MARITIME MOBILE Fixed Aeronautical mobile Land mobile US210 US229 US274 US317 NG152	Maritime (80) Private Land Mobile (90) Personal Radio (95) Amateur (97) Note: 216–220 MHz will become a mixed-use band in January 2002
	220–225 AMATEUR FIXED MOBILE Radiolocation S5.241		220–222 FIXED LAND MOBILE Radiolocation S5.241 G2 US335	220–222 FIXED LAND MOBILE US335	Private Land Mobile (90)
	S5.235 S5.237 S5.243		S5.233 S5.238 S5.240 S5.245	222–225 Radiolocation S5.241 G2	222–225 AMATEUR

223-230 BROADCASTING Fixed Mobile		223-230 FIXED MOBILE BROADCASTING AERONAUTICAL RADIONAVIGATION Radiolocation			
S5.243 S5.246 S5.247	225-235 FIXED MOBILE	S5.250	225-235 FIXED MOBILE	225-235	
230-235 FIXED MOBILE		230-235 FIXED MOBILE AERONAUTICAL RADIONAVIGATION			
S5.247 S5.251 S5.252		S5.250			
235-267 FIXED MOBILE			235-267 FIXED MOBILE	235-267	
S5.111 S5.199 S5.252 S5.254 S5.256			S5.111 S5.199 S5.256 G27 G100	S5.111 S5.199 S5.256	
267-272 FIXED MOBILE Space operation (space-to-Earth)			267-322 FIXED MOBILE	267-322	
S5.254 S5.257					
272-273 SPACE OPERATION (space-to-Earth) FIXED MOBILE					
S5.254					
273-312 FIXED MOBILE					
S5.254					
312-315 FIXED MOBILE Mobile-satellite (Earth-to-space) S5.254 S5.255			G27 G100		
See next page for 315-322 MHz					

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

322-410 MHz (UHF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
315-322 FIXED MOBILE SS.254			See previous page for 267-322 MHz		See previous page for 267-322 MHz.
322-328.6 FIXED MOBILE RADIO ASTRONOMY SS.149			322-328.6 FIXED MOBILE SS.149 G27	322-328.6 SS.149	
328.6-335.4 AERONAUTICAL RADIONAVIGATION SS.258 SS.259			328.6-335.4 AERONAUTICAL RADIONAVIGATION SS.258		
335.4-387 FIXED MOBILE SS.254			335.4-399.9 FIXED MOBILE G27 G100	335.4-399.9	
387-390 FIXED MOBILE Mobile-satellite (space-to-Earth) SS.208A SS.254 SS.255					
390-399.9 FIXED MOBILE SS.254					
399.9-400.05 MOBILE-SATELLITE (Earth-to-space) SS.209 A.5.224A RADIONAVIGATION-SATELLITE SS.222 SS.224B SS.260 SS.220			399.9-400.05 MOBILE-SATELLITE (Earth-to-space) US319 US322 RADIONAVIGATION-SATELLITE SS.260		
400.05-400.15 STANDARD FREQUENCY AND TIME SIGNAL-SATELLITE (400.1 MHz) SS.261 SS.262			400.05-400.15 STANDARD FREQUENCY AND TIME SIGNAL-SATELLITE (400.1 MHz) SS.261		

<p>400.15-401 METEOROLOGICAL AIDS METEOROLOGICAL-SATELLITE (space-to-Earth) MOBILE-SATELLITE (space-to-Earth) S5.208A S5.209 SPACE RESEARCH (space-to-Earth) S5.263 Space operation (space-to-Earth)</p> <p>S5.262 S5.264</p>	<p>400.15-401 METEOROLOGICAL AIDS (radiosonde) METEOROLOGICAL-SATELLITE (space-to-Earth) MOBILE-SATELLITE (space-to-Earth) 599B US319 US320 US324 SPACE RESEARCH (space-to-Earth) S5.263 Space operation (space-to-Earth)</p> <p>647B US70</p>	<p>400.15-401 METEOROLOGICAL AIDS (radiosonde) MOBILE-SATELLITE (space-to-Earth) 599B US319 US320 US324 SPACE RESEARCH (space-to-Earth) S5.263 Space operation (space-to-Earth)</p> <p>647B US70</p>	<p>Satellite Communications (25)</p>
<p>401-402 METEOROLOGICAL AIDS SPACE OPERATION (space-to-Earth) EARTH EXPLORATION-SATELLITE (Earth-to-space) METEOROLOGICAL-SATELLITE (Earth-to-space) Fixed Mobile except aeronautical mobile</p>	<p>401-402 METEOROLOGICAL AIDS (radiosonde) SPACE OPERATION (space-to-Earth) Earth exploration-satellite (Earth-to-space) Meteorological-satellite (Earth-to-space)</p> <p>US70</p>		
<p>402-403 METEOROLOGICAL AIDS EARTH EXPLORATION-SATELLITE (Earth-to-space) METEOROLOGICAL-SATELLITE (Earth-to-space) Fixed Mobile except aeronautical mobile</p>	<p>402-403 METEOROLOGICAL AIDS (radiosonde) US70 Earth exploration-satellite (Earth-to-space) Meteorological-satellite (Earth-to-space)</p> <p>US345</p>		<p>Personal Radio (95)</p>
<p>403-406 METEOROLOGICAL AIDS Fixed Mobile except aeronautical mobile</p>	<p>403-406 METEOROLOGICAL AIDS (radiosonde) US70</p> <p>US345 G6</p>	<p>403-406 METEOROLOGICAL AIDS (radiosonde) US70</p> <p>US345</p>	
<p>406-406.1 MOBILE-SATELLITE (Earth-to-space)</p> <p>S5.266 S5.267</p>	<p>406-406.1 MOBILE-SATELLITE (Earth-to-space)</p> <p>S5.266 S5.267</p>		
<p>406.1-410 FIXED MOBILE except aeronautical mobile RADIO ASTRONOMY</p> <p>S5.149</p>	<p>406.1-410 FIXED MOBILE RADIO ASTRONOMY US74</p> <p>US13 US117 G5 G6</p>	<p>406.1-410 RADIO ASTRONOMY US74</p> <p>US13 US117</p>	

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

410-470 MHz (UHF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
410-420 FIXED MOBILE except aeronautical mobile SPACE RESEARCH (space-to-space) S5.268			410-420 FIXED MOBILE US13 G5	410-420 US13	
420-430 FIXED MOBILE except aeronautical mobile Radiolocation S5.269 S5.270 S5.271			420-450 RADIOLOCATION G2 S5.286 US7 US87 US217 US228 US230 G8	420-450 Amateur S5.282 S5.286 US7 US87 US217 US228 US230 NG135	Private Land Mobile (90) Amateur (97)
430-440 AMATEUR RADIOLOCATION S5.138 S5.271 S5.272 S5.273 S5.274 S5.275 S5.276 S5.277 S5.280 S5.281 S5.282 S5.283	430-440 RADIOLOCATION Amateur S5.271 S5.276 S5.277 S5.278 S5.279 S5.281 S5.282				
440-450 FIXED MOBILE except aeronautical mobile Radiolocation S5.269 S5.270 S5.271 S5.284 S5.285 S5.286					
450-455 FIXED MOBILE S5.209 S5.271 S5.286 S5.286A S5.286B S5.286C S5.286D S5.286E			450-454 S5.286 US87	450-454 LAND MOBILE S5.286 US87 NG112 NG124	Auxiliary Broadcasting (74) Private Land Mobile (90)
455-456 FIXED MOBILE S5.209 S5.271 S5.286A S5.286B S5.286C S5.286E			454-456	454-455 FIXED LAND MOBILE NG12 NG112 NG148	Public Mobile (22) Maritime (80)
455-456 FIXED MOBILE S5.209 S5.271 S5.286A S5.286B S5.286C S5.286E	455-456 FIXED MOBILE MOBILE-SATELLITE (Earth-to-space) S5.286A S5.286B S5.286C S5.209 S5.271	455-456 FIXED MOBILE S5.209 S5.271 S5.286A S5.286B S5.286C S5.286E		455-456 LAND MOBILE	Auxiliary Broadcasting (74)

456-459 FIXED MOBILE S5.271 S5.287 S5.288			456-460 S5.288 669	456-460 FIXED LAND MOBILE S5.288 669 NG112 NG124 NG148	Public Mobile (22) Maritime (80) Private Land Mobile (90)
459-460 FIXED MOBILE S5.209 S5.271 S5.286A S5.286B S5.286C S5.286E	459-460 FIXED MOBILE MOBILE-SATELLITE (Earth-to-space) S5.286A S5.286B S5.286C S5.209 S5.271	459-460 FIXED MOBILE S5.209 S5.271 S5.286A S5.286B S5.286C S5.286E			
460-470 FIXED MOBILE Meteorological-satellite (space-to-Earth) S5.287 S5.288 S5.289 S5.290			460-470 Meteorological-satellite (space-to-Earth) S5.288 S5.289 669 US201 US209 US216	460-462.5375 FIXED LAND MOBILE S5.289 US201 US209 NG124	Private Land Mobile (90)
				462.5375-462.7375 LAND MOBILE S5.289 US201	Personal Radio (95)
				462.7375-467.5375 FIXED LAND MOBILE S5.289 669 US201 US209 US216 NG124	Private Land Mobile (90)
				467.5375-467.7375 LAND MOBILE S5.289 669 US201	Personal Radio (95)
				467.7375-470 FIXED LAND MOBILE S5.288 S5.289 US201 US216 NG124	Private Land Mobile (90)

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

470-849 MHz (UHF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
470-790 BROADCASTING	470-512 BROADCASTING Fixed Mobile S5.292 S5.293	470-585 FIXED MOBILE BROADCASTING S5.291 S5.298	470-608	470-512 FIXED BROADCASTING LAND MOBILE NG66 NG114 NG127 NG128 NG149	Public Mobile (22) Broadcast Radio (TV) (73) Auxiliary Broadcasting (74) Private Land Mobile (90)
	512-608 BROADCASTING S5.297			512-608 BROADCASTING NG128 NG149	Broadcast Radio (TV) (73) Auxiliary Broadcasting (74)
	608-614 RADIO ASTRONOMY Mobile-satellite except aeronautical mobile-satellite (Earth-to-space)	S5.149 S5.305 S5.306 S5.307	608-614 RADIO ASTRONOMY US74		
	614-890 BROADCASTING Fixed Mobile	610-890 FIXED MOBILE BROADCASTING	US246		
			614-890	614-698 BROADCASTING NG128 NG149	Broadcast Radio (TV) (73) Auxiliary Broadcast. (74)
				698-746 BROADCASTING NG128 NG149	Broadcast Radio (TV) (73) Auxiliary Broadcast. (74)
				Note: Band to be reallocated and auctioned by Sept. 30, 2002.	

S5.149 S5.291A S5.294 S5.296 S5.300 S5.302 S5.304 S5.306 S5.311 S5.312				746-764 FIXED MOBILE BROADCASTING	Wireless Communications (27) Broadcast Radio (TV) (73) Auxiliary Broadcasting (74) Private Land Mobile (90)
				NG128 NG159	
790-862 FIXED BROADCASTING	S5.293 S5.309 S5.311			764-776 FIXED MOBILE	Auxiliary Broadcasting (74) Private Land Mobile (90)
				NG128 NG158 NG159	
	806-890 FIXED MOBILE BROADCASTING			776-794 FIXED MOBILE BROADCASTING	Wireless Communications (27) Broadcast Radio (TV) (73) Auxiliary Broadcast. (74) Private Land Mobile (90)
				NG128 NG159	
				794-806 FIXED MOBILE	Auxiliary Broadcasting (74) Private Land Mobile (90)
				NG128 NG158 NG159	
				806-821 FIXED LAND MOBILE	Public Mobile (22) Private Land Mobile (90)
				NG31 NG63	
				821-824 LAND MOBILE	Private Land Mobile (90)
				NG63	
See next page for 862-890 MHz	S5.317 S5.318		S5.149 S5.305 S5.306 S5.307 S5.311 S5.320	824-849 FIXED LAND MOBILE	Public Mobile (22)
				NG63 NG151	
				See next page for 849-894 MHz	See next page for 866-896 MHz

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

849-940 MHz (UHF)

International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
See previous pages for 470-862 MHz	See previous pages for 614-890 MHz	See previous pages for 585-890 MHz	See previous pages for 614-890 MHz	See previous pages for 614-849 MHz	See previous pages for 614-849 MHz
				849-851 AERONAUTICAL MOBILE NG63	Public Mobile (22)
862-890 FIXED MOBILE except aeronautical mobile BROADCASTING S5.322 S5.319 S5.323				851-866 FIXED LAND MOBILE NG31 NG63	Public Mobile (22) Private Land Mobile (90)
				866-869 LAND MOBILE NG63	Private Land Mobile (90)
890-942 FIXED MOBILE except aeronautical mobile BROADCASTING S5.322 Radiolocation	890-902 FIXED MOBILE except aeronautical mobile Radiolocation S5.318 S5.325	890-942 FIXED MOBILE BROADCASTING Radiolocation	890-902	869-894 FIXED LAND MOBILE	Public Mobile (22)
				US116 US268 NG63 NG151	
				894-896 AERONAUTICAL MOBILE US116 US268	
				896-901 FIXED LAND MOBILE US116 US268	Private Land Mobile (90)
			US116 US268 G2	901-902 FIXED MOBILE US116 US268	Personal Communication (24)

S5.323	902-928 FIXED Amateur Mobile except aeronautical mobile Radiolocation S5.150 S5.325 S5.326	S5.327	902-928 RADIOLOCATION G59 S5.150 US215 US218 US267 US275 G11	902-928 S5.150 US215 US218 US267 US275	ISM Equipment (18) Private Land Mobile (90) Amateur (97)					
	928-942 FIXED MOBILE except aeronautical mobile Radiolocation		928-932 US116 US215 US268 G2	928-929 FIXED US116 US215 US268 NG120	929-930 FIXED LAND MOBILE US116 US215 US268	930-931 FIXED MOBILE US116 US215 US268	931-932 FIXED LAND MOBILE US116 US215 US268	Public Mobile (22) Private Land Mobile (90) Fixed Microwave (101)		
								Private Land Mobile (90)		
								Personal Communications (24)		
								Public Mobile (22)		
								932-935 FIXED US215 US268 G2	932-935 FIXED US215 US268 NG120	Public Mobile (22) Fixed Microwave (101)
								935-940 US116 US215 US268 G2	935-940 FIXED LAND MOBILE US116 US215 US268	Private Land Mobile (90)
								940-941 US116 US268 G2	940-941 FIXED MOBILE US116 US268	Personal Communications (24)

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

941-1429 MHz (UHF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
See previous page for 890-942 MHz	See previous page for 928-942 MHz	See previous page for 890-942 MHz	941-944 FIXED	941-944 FIXED	Public Mobile (22) Fixed Microwave (101)
942-960 FIXED MOBILE except aeronautical mobile BROADCASTING S5.322	942-960 FIXED MOBILE	942-960 FIXED MOBILE BROADCASTING	US268 US301 US302 G2	US268 US301 US302 NG120	
S5.323		S5.320	944-960	944-960 FIXED NG120	Public Mobile (22) International Fixed (23) Auxiliary Broadcast. (74) Fixed Microwave (101)
960-1215 AERONAUTICAL RADIONAVIGATION S5.328			960-1215 AERONAUTICAL RADIONAVIGATION S5.328 US224		Aviation (87)
1215-1240 EARTH EXPLORATION-SATELLITE (active) RADIOLOCATION RADIONAVIGATION-SATELLITE (space-to-Earth) S5.329 SPACE RESEARCH (active) S5.330 S5.331 S5.332			1215-1240 RADIOLOCATION S5.33 G56 RADIONAVIGATION-SATELLITE (space-to-Earth)	1215-1240 S5.333	
1240-1260 EARTH EXPLORATION-SATELLITE (active) RADIOLOCATION RADIONAVIGATION-SATELLITE (space-to-Earth) S5.329 SPACE RESEARCH (active) Amateur S5.330 S5.331 S5.332 S5.334 S5.335			1240-1300 RADIOLOCATION S5.333 G56 S5.334	1240-1300 Amateur S5.282 S5.333 S5.334	Amateur (97)
1260-1300 EARTH EXPLORATION-SATELLITE (active) RADIOLOCATION SPACE RESEARCH (active) Amateur S5.282 S5.330 S5.331 S5.332 S5.334 S5.335					
1300-1350 AERONAUTICAL RADIONAVIGATION S5.337 Radiolocation S5.149			1300-1350 AERONAUTICAL RADIO-NAVIGATION S5.337 Radiolocation G2 S5.149	1300-1350 AERONAUTICAL RADIO-NAVIGATION S5.337 S5.149	Aviation (87)

<p>1350-1400 FIXED MOBILE RADIOLOCATION</p>	<p>1350-1400 RADIOLOCATION</p>	<p>1350-1390 RADIOLOCATION G2 Fixed Mobile</p> <p>S5.149 S5.334 S5.339 US311 G27 G114</p>	<p>1350-1390</p> <p>S5.149 S5.334 S5.339</p>	
<p>S5.149 S5.338 S5.339</p>	<p>S5.149 S5.334 S5.339</p>	<p>1390-1400 RADIOLOCATION G2 Fixed Mobile</p> <p>S5.149 US311 S5.339 G27 G114</p>	<p>1390-1400</p> <p>S5.149 S5.339</p>	<p>Note: 1390-1400 MHz became non-Federal government exclusive spectrum in January 1999</p>
<p>1400-1427 EARTH EXPLORATION-SATELLITE (passive) RADIO ASTRONOMY SPACE RESEARCH (passive)</p> <p>S5.340 S5.341</p>	<p>1400-1427 EARTH EXPLORATION-SATELLITE (passive) RADIO ASTRONOMY US74 SPACE RESEARCH (passive)</p> <p>S5.341 US246</p>			
<p>1427-1429 SPACE OPERATION (Earth-to-space) FIXED MOBILE except aeronautical mobile</p> <p>S5.341</p>	<p>1427-1429 SPACE OPERATION (Earth-to-space) FIXED MOBILE except aeronautical mobile</p> <p>S5.341 G30</p>	<p>1427-1429 SPACE OPERATION (Earth-to-space) Fixed (telemetry) Land mobile (telemetry and telecommand)</p> <p>S5.341</p>	<p>Satellite Communications (25) Private Land Mobile (90)</p> <p>Note: 1427-1429 MHz became non-Federal government exclusive spectrum in January 1999</p>	

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

1429-1610 MHz (UHF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
1429-1452 FIXED MOBILE except aeronautical mobile S5.341 S5.342	1429-1452 FIXED MOBILE S5.343 S5.341		1429-1435 FIXED MOBILE S5.341 G30	1429-1435 Fixed (telemetry) Land mobile (telemetry and telecommand) S5.341	Private Land Mobile (90) Note: in January 1999, 1429-1432 MHz became non-Federal government exclusive spectrum and 1432-1435 MHz became mixed-use spectrum
1452-1492 FIXED MOBILE except aeronautical mobile BROADCASTING S5.345 S5.347 BROADCASTING- SATELLITE S5.345 S5.347 S5.341 S5.342	1452-1492 FIXED MOBILE S5.343 BROADCASTING S5.345 S5.347 BROADCASTING-SATELLITE S5.345 S5.347 S5.341 S5.344		1435-1525 MOBILE (aeronautical telemetry)		Aviation (87)
1492-1525 FIXED MOBILE except aeronautical mobile S5.341 S5.342	1492-1525 FIXED MOBILE S5.343 MOBILE-SATELLITE (space-to-Earth) S5.348A S5.341 S5.344 S5.348	1492-1525 FIXED MOBILE S5.341 S5.348A	S5.341 US78		
1525-1530 SPACE OPERATION (space-to-Earth) FIXED MOBILE-SATELLITE (space-to-Earth) Earth exploration-satellite Mobile except aeronautical mobile S5.349 S5.341 S5.342 S5.350 S5.351 S5.352A S5.354	1525-1530 SPACE OPERATION (space-to-Earth) MOBILE-SATELLITE (space-to-Earth) Earth exploration-satellite Fixed Mobile S5.343 S5.341 S5.351 S5.354	1525-1530 SPACE OPERATION (space-to-Earth) FIXED MOBILE-SATELLITE (space-to-Earth) Earth exploration-satellite Mobile S5.349 S5.341 S5.351 S5.352A S5.354	1525-1530 MOBILE-SATELLITE (space-to-Earth) Mobile (aeronautical telemetry) S5.341 S5.351 US78		Satellite Communications (25) Aviation (87)

<p>1530-1535 SPACE OPERATION (space-to-Earth) MOBILE-SATELLITE (space-to-Earth) S5.35A Earth exploration-satellite Fixed Mobile except aeronautical mobile</p> <p>S5.341 S5.342 S5.351 S5.354</p>	<p>1530-1535 SPACE OPERATION (space-to-Earth) MOBILE-SATELLITE (space-to-Earth) S5.35A Earth exploration-satellite Fixed Mobile S5.343</p> <p>S5.341 S5.351 S5.354</p>	<p>1530-1535 MOBILE-SATELLITE (space-to-Earth) MARITIME MOBILE-SATELLITE (space-to-Earth) Mobile (aeronautical telemetry)</p> <p>S5.341 S5.351 US78 US315</p>	
<p>1535-1559 MOBILE-SATELLITE (space-to-Earth)</p> <p>S5.341 S5.351 S5.353A S5.354 S5.355 S5.356 S5.357 S5.357A S5.359 S5.362A</p>	<p>1535-1544 MOBILE-SATELLITE (space-to-Earth) MARITIME MOBILE-SATELLITE (space-to-Earth)</p> <p>S5.341 S5.351 US315</p> <p>1544-1545 MOBILE-SATELLITE (space-to-Earth)</p> <p>S5.341 S5.356</p> <p>1545-1549.5 AERONAUTICAL MOBILE-SATELLITE (R) (space-to-Earth) Mobile-satellite (space-to-Earth)</p> <p>S5.341 S5.351 US308 US309</p> <p>1549.5-1558.5 AERONAUTICAL MOBILE-SATELLITE (R) (space-to-Earth) MOBILE-SATELLITE (space-to-Earth)</p> <p>S5.341 S5.351 US308 US309</p> <p>1558.5-1559 AERONAUTICAL MOBILE-SATELLITE (R) (space-to-Earth)</p> <p>S5.341 S5.351 US308 US309</p>	<p>Satellite Communications (25) Maritime (80)</p> <p>Aviation (87)</p> <p>Note: Footnote G126 states DGPS stations may be authorized on a primary basis in the 1559-1610 MHz band, but the FCC has not yet adopted this footnote.</p>	
<p>1559-1610 AERONAUTICAL RADIONAVIGATION RADIONAVIGATION-SATELLITE (space-to-Earth)</p> <p>S5.341 S5.355 S5.359 S5.363</p>	<p>1559-1610 AERONAUTICAL RADIONAVIGATION RADIONAVIGATION-SATELLITE (space-to-Earth)</p> <p>S5.341 US208 US260 G126</p>		

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

1610-1670 MHz (UHF)					Remarks
International Table			United States Table		
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
1610-1610.6 MOBILE-SATELLITE (Earth-to-space) AERONAUTICAL RADIONAVIGATION \$5.341 \$5.355 \$5.359 \$5.363 \$5.364 \$5.366 \$5.367 \$5.368 \$5.369 \$5.371 \$5.372	1610-1610.6 MOBILE-SATELLITE (Earth-to-space) AERONAUTICAL RADIONAVIGATION RADIODETERMINATION- SATELLITE (Earth-to- space) \$5.341 \$5.364 \$5.366 \$5.367 \$5.368 \$5.370 \$5.372	1610-1610.6 MOBILE-SATELLITE (Earth-to-space) AERONAUTICAL RADIONAVIGATION Radiodetermination-satellite (Earth-to-space) \$5.341 \$5.355 \$5.359 \$5.364 \$5.366 \$5.367 \$5.368 \$5.369 \$5.372	1610-1610.6 MOBILE-SATELLITE (Earth-to-space) US319 AERONAUTICAL RADIONAVIGATION US260 RADIODETERMINATION-SATELLITE (Earth-to-space) \$5.341 \$5.364 \$5.366 \$5.367 \$5.368 \$5.372 US208	Satellite Communications (25) Aviation (87)	
1610.6-1613.8 MOBILE-SATELLITE (Earth-to-space) RADIO ASTRONOMY AERONAUTICAL RADIONAVIGATION \$5.149 \$5.341 \$5.355 \$5.359 \$5.363 \$5.364 \$5.366 \$5.367 \$5.368 \$5.369 \$5.371 \$5.372	1610.6-1613.8 MOBILE-SATELLITE (Earth-to-space) RADIO ASTRONOMY AERONAUTICAL RADIONAVIGATION RADIODETERMINATION- SATELLITE (Earth-to- space) \$5.149 \$5.341 \$5.364 \$5.366 \$5.367 \$5.368 \$5.370 \$5.372	1610.6-1613.8 MOBILE-SATELLITE (Earth-to-space) RADIO ASTRONOMY AERONAUTICAL RADIONAVIGATION Radiodetermination-satellite (Earth-to-space) \$5.149 \$5.341 \$5.355 \$5.359 \$5.364 \$5.366 \$5.367 \$5.368 \$5.369 \$5.372	1610-1613.8 MOBILE-SATELLITE (Earth-to-space) US319 RADIO ASTRONOMY AERONAUTICAL RADIONAVIGATION US260 RADIODETERMINATION-SATELLITE (Earth-to-space) \$5.341 \$5.364 \$5.366 \$5.367 \$5.368 \$5.372 US208		
1613.8-1626.5 MOBILE-SATELLITE (Earth-to-space) AERONAUTICAL RADIONAVIGATION Mobile-satellite (space-to-Earth) \$5.341 \$5.355 \$5.359 \$5.363 \$5.364 \$5.365 \$5.366 \$5.367 \$5.368 \$5.369 \$5.371 \$5.372	1613.8-1626.5 MOBILE-SATELLITE (Earth-to-space) AERONAUTICAL RADIONAVIGATION RADIODETERMINATION- SATELLITE (Earth-to- space) Mobile-satellite (space-to- Earth) \$5.341 \$5.364 \$5.365 \$5.366 \$5.367 \$5.368 \$5.370 \$5.372	1613.8-1626.5 MOBILE-SATELLITE (Earth-to-space) AERONAUTICAL RADIONAVIGATION Mobile-satellite (space-to- Earth) Radiodetermination- satellite (Earth-to-space) \$5.341 \$5.355 \$5.359 \$5.364 \$5.365 \$5.366 \$5.367 \$5.368 \$5.369 \$5.372	1613.8-1626.5 MOBILE-SATELLITE (Earth-to-space) US319 AERONAUTICAL RADIONAVIGATION US260 RADIODETERMINATION-SATELLITE (Earth-to-space) Mobile-satellite (space-to-Earth) \$5.341 \$5.364 \$5.365 \$5.366 \$5.367 \$5.368 \$5.372 US208		

1626.5–1660 MOBILE-SATELLITE (Earth-to-space)	1626.5–1645.5 MOBILE-SATELLITE (Earth-to-space) MARITIME MOBILE-SATELLITE (Earth-to-space) S5.341 S5.351 US315	Satellite Communications (25) Maritime (80)
	1645.5–1645.5 MOBILE-SATELLITE (Earth-to-space) S5.341 S5.375	
	1646.5–1651 AERONAUTICAL MOBILE-SATELLITE (R) (Earth-to-space) Mobile-satellite (Earth-to-space) S5.341 S5.351 US308 US309	Aviation (87)
	1651–1660 MOBILE-SATELLITE (Earth-to-space) AERONAUTICAL MOBILE-SATELLITE (R) (Earth-to-space) S5.341 S5.351 US308 US309	
S5.341 S5.351 S5.353A S5.354 S5.355 S5.357A S5.359 S5.362A S5.374 S5.375 S5.376		
1660–1660.5 MOBILE-SATELLITE (Earth-to-space) RADIO ASTRONOMY S5.149 S5.341 S5.351 S5.354 S5.362A S5.376A	1660–1660.5 AERONAUTICAL MOBILE-SATELLITE (R) (Earth-to-space) RADIO ASTRONOMY S5.149 S5.341 S5.351 US308 US309	
1660.5–1668.4 RADIO ASTRONOMY SPACE RESEARCH (passive) Fixed Mobile except aeronautical mobile S5.149 S5.341 S5.379 S5.379A	1660.5–1668.4 RADIO ASTRONOMY US74 SPACE RESEARCH (passive) S5.341 US246	
1668.4–1670 METEOROLOGICAL AIDS FIXED MOBILE except aeronautical mobile RADIO ASTRONOMY S5.149 S5.341	1668.4–1670 METEOROLOGICAL AIDS (radiosonde) RADIO ASTRONOMY US74 S5.149 S5.341 US99	

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

1670-2110 MHz (UHF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
1670-1675 METEOROLOGICAL AIDS FIXED METEOROLOGICAL-SATELLITE (space-to-Earth) MOBILE S5.380 S5.341			1670-1675 METEOROLOGICAL AIDS (radiosonde) METEOROLOGICAL-SATELLITE (space-to-Earth) S5.341 US211		Note: 1670-1675 MHz became mixed-use spectrum in January 1999
1675-1690 METEOROLOGICAL AIDS FIXED METEOROLOGICAL-SATELLITE (space-to-Earth) MOBILE except aeronautical mobile S5.341	1675-1690 METEOROLOGICAL AIDS FIXED METEOROLOGICAL-SATELLITE (space-to-Earth) MOBILE except aeronautical mobile MOBILE-SATELLITE (Earth-to-space) S5.341 S5.377	1675-1690 METEOROLOGICAL AIDS FIXED METEOROLOGICAL-SATELLITE (space-to-Earth) MOBILE except aeronautical mobile S5.341	1675-1690 METEOROLOGICAL AIDS (radiosonde) METEOROLOGICAL-SATELLITE (space-to-Earth)		
1690-1700 METEOROLOGICAL AIDS METEOROLOGICAL-SATELLITE (space-to-Earth) Fixed Mobile except aeronautical mobile S5.289 S5.341 S5.382	1690-1700 METEOROLOGICAL AIDS METEOROLOGICAL-SATELLITE (space-to-Earth) MOBILE-SATELLITE (Earth-to-space) S5.289 S5.341 S5.377 S5.381	1690-1700 METEOROLOGICAL AIDS METEOROLOGICAL-SATELLITE (space-to-Earth) S5.289 S5.341 S5.381	S5.289 S5.341 US211		
1700-1710 FIXED METEOROLOGICAL-SATELLITE (space-to-Earth) MOBILE except aeronautical mobile S5.289 S5.341	1700-1710 FIXED METEOROLOGICAL-SATELLITE (space-to-Earth) MOBILE except aeronautical mobile MOBILE-SATELLITE (Earth-to-space) S5.289 S5.341 S5.377	1700-1710 FIXED METEOROLOGICAL-SATELLITE (space-to-Earth) MOBILE except aeronautical mobile S5.289 S5.341 S5.384	1700-1710 FIXED G118 METEOROLOGICAL-SATELLITE (space-to-Earth)	1700-1710 METEOROLOGICAL-SATELLITE (space-to-Earth) Fixed S5.289 S5.341	
1710-1930 FIXED MOBILE S5.380			1710-1755 FIXED MOBILE S5.341 US256	1710-1755 S5.341 US256	Note: Proceeds from the auction of the 1710-1755 MHz-mixed-use band are to be deposited not later than September 30, 2002.

			1755–1850 FIXED MOBILE G42	1755–1850	
S5.149 S5.341 S5.385 S5.386 S5.387 S5.388			1850–1990	1850–1990 FIXED MOBILE	RF Devices (15) Personal Communications (24) Fixed Microwave (101)
1930–1970 FIXED MOBILE S5.388	1930–1970 FIXED MOBILE Mobile-satellite (Earth-to-space) S5.388	1930–1970 FIXED MOBILE S5.388			
1970–1980 FIXED MOBILE S5.388					
1980–2010 FIXED MOBILE MOBILE-SATELLITE (Earth-to-space) S5.388 S5.389A S5.389B S5.389F			1990–2025	1990–2025 MOBILE-SATELLITE (Earth-to-space)	Satellite Communications (25) Auxiliary Broadcasting (74) Cable TV Relay (78)
2010–2025 FIXED MOBILE S5.388	2010–2025 FIXED MOBILE MOBILE-SATELLITE (Earth-to-space) S5.388 S5.389C S5.389D S5.389E S5.390	2010–2025 FIXED MOBILE S5.388	US111	US111	
2025–2110 SPACE OPERATION (Earth-to-space) (space-to-space) EARTH EXPLORATION-SATELLITE (Earth-to-space) (space-to-space) FIXED MOBILE S5.391 SPACE RESEARCH (Earth-to-space) (space-to-space) S5.392			2025–2110 US90 US111 US219 US222	2025–2110 FIXED MOBILE US90 US111 US219 US222 NG23 NG118	Auxiliary Broadcasting (74) Cable TV Relay (78)

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

2110-2345 MHz (UHF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
2110-2120 FIXED MOBILE SPACE RESEARCH (deep space) (Earth-to-space) S5.388			2110-2130	2110-2130 FIXED MOBILE	Public Mobile (22) Auxiliary Broadcasting (74) Cable TV Relay (78) Fixed Microwavc (101) Note: 2110-2150 MHz must be auctioned by September 30, 2002.
2120-2160 FIXED MOBILE	2120-2160 FIXED MOBILE Mobile-satellite (space-to-Earth)	2120-2160 FIXED MOBILE	US111 US252	US111 US252 NG23 NG118	
S5.388			2130-2200	2130-2150 FIXED MOBILE NG23 NG153	Public Mobile (22) Fixed Microwave (101)
S5.388				2150-2160 FIXED NG23	Domestic Public Fixed (21) Fixed Microwave (101)
2160-2170 FIXED MOBILE	2160-2170 FIXED MOBILE MOBILE-SATELLITE (space-to-Earth) S5.388 S5.389C S5.389D S5.389E S5.390	2160-2170 FIXED MOBILE	2160-2165 FIXED MOBILE NG23 NG153		Domestic Public Fixed (21) Public Mobile (22) Fixed Microwave (101)
S5.388 S5.392A			2165-2200 MOBILE-SATELLITE (space-to-Earth)		Public Mobile (22) Satellite Communications (25) Fixed Microwave (101)
2170-2200 FIXED MOBILE MOBILE-SATELLITE (space-to-Earth) S5.388 S5.389A S5.389F S5.392A			NG23		
2200-2290 SPACE OPERATION (space-to-Earth) (space-to-space) EARTH EXPLORATION-SATELLITE (space-to-Earth) (space-to-space) FIXED MOBILE S5.391 SPACE RESEARCH (space-to-Earth) (space-to-space)			2200-2290 SPACE OPERATION (space-to-Earth) (space-to-space) EARTH EXPLORATION-SATELLITE (space-to-Earth) (space-to-space) FIXED (line-of-sight only)	2200-2290	

S5.392		MOBILE (line-of-sight only including aeronautical telemetry, but excluding flight testing of manned aircraft) SPACE RESEARCH (space-to-Earth) (space-to-space)	US303	
2290–2300 FIXED MOBILE except aeronautical mobile SPACE RESEARCH (deep space) (space-to-Earth)		2290–2300 FIXED MOBILE except aeronautical mobile SPACE RESEARCH (deep space) (space-to-Earth)	2290–2300 SPACE RESEARCH (deep space) (space-to-Earth)	
2300–2450 FIXED MOBILE Amateur Radiolocation	2300–2450 FIXED MOBILE RADIOLOCATION Amateur	2300–2305	2300–2305 Amateur	Amateur (97) Note: 2300–2305 MHz became non-Federal government exclusive spectrum in August 1995
		G123		
		2305–2310	2305–2310 FIXED MOBILE except aeronautical mobile RADIOLOCATION Amateur US338	Wireless Comm. (27) Amateur (97)
		US338 G123		
S5.150 S5.282 S5.395	S5.150 S5.282 S5.393 S5.394 S5.396	2310–2360 Fixed Mobile Radiolocation G2	2310–2320 FIXED MOBILE US339 RADIOLOCATION BROADCASTING-SATELLITE US327 S5.396 US338	Wireless Communications (27)
		S5.396 US327 US328 G120	2320–2345 BROADCASTING-SATELLITE US327 Mobile US276 US328 S5.396	
		See next page	See next page	See next page

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

2345–2655 MHz (UHF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
See previous page for 2300–2450 MHz			See previous page for 2310–2360 MHz	2345–2360 FIXED MOBILE US339 RADIOLOCATION BROADCASTING- SATELLITE US327 S5.396	Wireless Communications (27)
			2360–2385 MOBILE US276 RADIOLOCATION G2 Fixed G120	2360–2385 MOBILE US276	
			2385–2390 MOBILE US276 RADIOLOCATION G2 Fixed G120	2385–2390 MOBILE US276	Note: 2385–2390 MHz will become non-Federal government exclusive spectrum in January 2005
			2390–2400 G122	2390–2400 AMATEUR	RF Devices (15) Amateur (97)
			2400–2402 S5.150 G123	2400–2402 Amateur S5.150 S5.282	ISM Equipment (18) Amateur (97)
			2402–2417 S5.150 G122	2402–2417 AMATEUR S5.150 S5.282	RF Devices (15) ISM Equipment (18) Amateur (97)
			2417–2450 Radiolocation G2 S5.150 G124	2417–2450 Amateur S5.150 S5.282	ISM Equipment (18) Amateur (97)
			2450–2483.5 FIXED MOBILE Radiolocation S5.150 S5.397	2450–2483.5 FIXED MOBILE RADIOLOCATION S5.150 S5.394	2450–2483.5 S5.150 US41

<p>2483.5–2500 FIXED MOBILE MOBILE-SATELLITE (space-to-Earth) Radiolocation</p> <p>S5.150 S5.371 S5.397 S5.398 S5.399 S5.400 S5.402</p>	<p>2483.5–2500 FIXED MOBILE MOBILE-SATELLITE (space-to-Earth) RADIOLOCATION RADIODETERMINATION- SATELLITE (space-to- Earth) S5.398</p> <p>S5.150 S5.402</p>	<p>2483.5–2500 FIXED MOBILE MOBILE-SATELLITE (space-to-Earth) RADIOLOCATION Radiodetermination-satellite (space-to-Earth) S5.398</p> <p>S5.150 S5.400 S5.402</p>	<p>2483.5–2500 MOBILE-SATELLITE (space-to-Earth) US319 RADIODETERMINATION- SATELLITE (space-to- Earth) S5.398</p> <p>S5.150 753F US41</p>	<p>2483.5–2500 MOBILE-SATELLITE (space-to-Earth) US319 RADIODETERMINATION- SATELLITE (space-to- Earth) S5.398</p> <p>S5.150 753F US41 NG147</p>	<p>ISM Equipment (18) Satellite Communications (25) Private Land Mobile (90) Fixed Microwave (101)</p>
<p>2500–2520 FIXED S5.409 S5.410 S5.411 MOBILE except aeronautical mobile MOBILE-SATELLITE (space-to-Earth) S5.403</p> <p>S5.405 S5.407 S5.408 S5.412 S5.414</p>	<p>2500–2520 FIXED S5.409 S5.411 FIXED-SATELLITE (space-to-Earth) S5.415 MOBILE except aeronautical mobile MOBILE-SATELLITE (space-to-Earth) S5.403</p> <p>S5.404 S5.407 S5.414 S5.415A</p>	<p>2500–2655</p>	<p>2500–2655 FIXED S5.409 S5.411 US205 FIXED-SATELLITE (space-to-Earth) NG102 BROADCASTING- SATELLITE NG101</p>	<p>Domestic Public Fixed (21) Auxiliary Broadcasting (74)</p>	
<p>2520–2655 FIXED S5.409 S5.410 S5.411 MOBILE except aeronautical mobile BROADCASTING- SATELLITE S5.413 S5.416</p> <p>S5.339 S5.403 S5.405 S5.408 S5.412 S5.417 S5.418</p>	<p>2520–2655 FIXED S5.409 S5.411 FIXED-SATELLITE (space-to-Earth) S5.415 MOBILE except aeronautical mobile BROADCASTING- SATELLITE S5.413 S5.416</p> <p>S5.339 S5.403</p>	<p>2520–2535 FIXED S5.409 S5.411 FIXED-SATELLITE (space-to-Earth) S5.415 MOBILE except aeronautical mobile BROADCASTING- SATELLITE S5.413 S5.416</p> <p>S5.403 S5.415A</p> <p>2535–2655 FIXED S5.409 S5.411 MOBILE except aeronautical mobile BROADCASTING- SATELLITE S5.413 S5.416</p> <p>S5.339 S5.418</p>	<p>S5.339 US205 US269</p>	<p>S5.339 US269</p>	

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

2655-3700 MHz (UHF/SHF)					Remarks
International Table			United States Table		
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
2655-2670 FIXED S5.409 S5.410 S5.411 MOBILE except aeronautical mobile BROADCASTING-SATELLITE S5.413 S5.416 Earth exploration-satellite (passive) Radio astronomy Space research (passive)	2655-2670 FIXED S5.409 S5.411 FIXED-SATELLITE (Earth-to-space) (space-to-Earth) S5.415 MOBILE except aeronautical mobile BROADCASTING-SATELLITE S5.413 S5.416 Earth exploration-satellite (passive) Radio astronomy Space research (passive)	2655-2670 FIXED S5.409 S5.411 FIXED-SATELLITE (Earth-to-space) S5.415 MOBILE except aeronautical mobile BROADCASTING-SATELLITE S5.413 S5.416 Earth exploration-satellite (passive) Radio astronomy Space research (passive)	2655-2690 Earth exploration-satellite (passive) Radio astronomy Space research (passive)	2655-2690 FIXED US205 NG47 FIXED-SATELLITE (Earth-to-space) NG102 BROADCASTING-SATELLITE NG101 Earth exploration-satellite (passive) Radio astronomy Space research (passive)	
S5.149 S5.412 S5.417 S5.420	S5.149 S5.420	S5.149 S5.420			
2670-2690 FIXED S5.409 S5.410 S5.411 MOBILE except aeronautical mobile MOBILE-SATELLITE (Earth-to-space) Earth exploration-satellite (passive) Radio astronomy Space research (passive)	2670-2690 FIXED S5.409 S5.411 FIXED-SATELLITE (Earth-to-space) (space-to-Earth) S5.415 MOBILE except aeronautical mobile MOBILE-SATELLITE (Earth-to-space) Earth exploration-satellite (passive) Radio astronomy Space research (passive)	2670-2690 FIXED S5.409 S5.411 FIXED-SATELLITE (Earth-to-space) S5.415 MOBILE except aeronautical mobile MOBILE-SATELLITE (Earth-to-space) Earth exploration-satellite (passive) Radio astronomy Space research (passive)			
S5.149 S5.419 S5.420	S5.149 S5.419 S5.420	S5.149 S5.419 S5.420 S5.420A	US205 US269	US269	
2690-2700 EARTH EXPLORATION-SATELLITE (passive) RADIO ASTRONOMY SPACE RESEARCH (passive) S5.340 S5.421 S5.422			2690-2700 EARTH EXPLORATION-SATELLITE (passive) RADIO ASTRONOMY US74 SPACE RESEARCH (passive) US246		
2700-2900 AERONAUTICAL RADIONAVIGATION S5.337 Radiolocation			2700-2900 AERONAUTICAL RADIO-NAVIGATION S5.337 METEOROLOGICAL AIDS Radiolocation G2 S5.423 US18 G15	2700-2900	
S5.423 S5.424				S5.423 US18	

2900–3100 RADIONAVIGATION S5.426 Radiolocation S5.425 S5.427		2900–3100 MARITIME RADIONAVIGATION Radiolocation G56 S5.427 US44 US316		2900–3100 MARITIME RADIONAVIGATION Radiolocation S5.427 US44 US316		Maritime 980)	
3100–3300 RADIOLOCATION Earth exploration-satellite (active) Space research (active) S5.149 S5.428		3100–3300 RADIOLOCATION S5.333 US110 G59 S5.149		3100–3300 Radiolocation S5.333 US110 S5.149			
3300–3400 RADIOLOCATION S5.149 S5.429 S5.430		3300–3400 RADIOLOCATION Amateur Fixed Mobile S5.149 S5.430		3300–3400 RADIOLOCATION Amateur S5.149 S5.429		3300–3500 Amateur Radiolocation US108 Amateur (97)	
3400–3600 FIXED FIXED-SATELLITE (space-to-Earth) Mobile Radiolocation S5.431		3400–3500 FIXED FIXED-SATELLITE (space-to-Earth) Amateur Mobile Radiolocation S5.433 S5.282 S5.432		3400–3500 RADIOLOCATION US108 G31 S5.149		3400–3500 Amateur Radiolocation US108 S5.149 S5.282	
		3500–3700 FIXED FIXED-SATELLITE (space-to-Earth) MOBILE except aeronautical mobile Radiolocation S5.433		3500–3650 RADIOLOCATION US110 G59 AERONAUTICAL RADIONAVIGATION (ground-based) G110 US245		3500–3600 Radiolocation US110 3600–3650 FIXED-SATELLITE (space-to-Earth) US245 Radiolocation US110	
3600–4200 FIXED FIXED-SATELLITE (space-to-Earth) Mobile		S5.435		3650–3700 RADIOLOCATION US110 G59 AERONAUTICAL RADIONAVIGATION (ground-based) G110 US245		3650–3700 FIXED-SATELLITE (space-to-Earth) US245 Radiolocation US110 Note: 3650–3700 MHz became mixed-use spectrum in January 1999	
		See next page for 3700–4200 MHz		See next page for 3700–4200 MHz			

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

3700–5650 MHz (SHF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
See previous page for 3600–4200 MHz.	3700–4200 FIXED FIXED-SATELLITE (space-to-Earth) Mobile except aeronautical mobile		3700–4200	3700–4200 FIXED NG41 FIXED-SATELLITE (space-to-Earth)	International Fixed (23) Satellite Communications (25) Fixed Microwave (101)
4200–4400 AERONAUTICAL RADIONAVIGATION S5.438 S5.437 S5.439 S5.440			4200–4400 AERONAUTICAL RADIONAVIGATION S5.440 US261		Aviation (87)
4400–4500 FIXED MOBILE			4400–4500 FIXED MOBILE	4400–4500	
4500–4800 FIXED FIXED-SATELLITE (space-to-Earth) S5.441 MOBILE			4500–4800 FIXED MOBILE US245	4500–4800 FIXED-SATELLITE (space-to-Earth) US245	
4800–4990 FIXED MOBILE S5.442 Radio astronomy			4800–4940 FIXED MOBILE S5.149 US203	4800–4940 S5.149 US203	
S5.149 S5.339 S5.443			4940–4990 FIXED MOBILE S5.149 S5.339 US257	4940–4990 S5.149 S5.339 US257	Note: 4940–4990 MHz became non-Federal government exclusive spectrum in March 1999
4990–5000 FIXED MOBILE except aeronautical mobile RADIO ASTRONOMY Space research (passive) S5.149			4990–5000 RADIO ASTRONOMY US74 Space research (passive) US246		
5000–5150 AERONAUTICAL RADIONAVIGATION S5.367 S5.444 S5.444A			5000–5150 AERONAUTICAL RADIONAVIGATION US260		Satellite Communications (25) Aviation (87)

<p>5150-5250 AERONAUTICAL RADIONAVIGATION FIXED-SATELLITE (Earth-to-space) S5.447A</p> <p>S5.446 S5.447 S5.447B S5.447C</p>			<p>Note: Footnote G126 states DGPS stations may be authorized on a primary basis in the 5000-5150 MHz band, but the FCC has not yet adopted this footnote.</p>
<p>5250-5255 EARTH EXPLORATION-SATELLITE (active) RADIOLOCATION SPACE RESEARCH S5.447D</p> <p>S5.448 S5.448A</p>	<p>5250-5350 RADIOLOCATION S5.333 US110 G59</p>	<p>5250-5350 Radiolocation S5.333 US110</p>	
<p>5255-5350 EARTH EXPLORATION-SATELLITE (active) RADIOLOCATION SPACE RESEARCH (active) S5.448 S5.448A</p>			
<p>5350-5460 EARTH EXPLORATION-SATELLITE (active) S5.448B AERONAUTICAL RADIONAVIGATION S5.449 Radiolocation</p>	<p>5350-5460 AERONAUTICAL RADIO- NAVIGATION S5.449 RADIOLOCATION G56</p> <p>US48</p>	<p>5350-5460 AERONAUTICAL RADIO- NAVIGATION S5.449 Radiolocation</p> <p>US48</p>	<p>Aviation (87)</p>
<p>5460-5470 RADIONAVIGATION S5.449 Radiolocation</p>	<p>5460-5470 RADIONAVIGATION S5.449 Radiolocation G56</p> <p>US49 US65</p>	<p>5460-5470 RADIONAVIGATION S5.449 Radiolocation</p> <p>US49 US65</p>	
<p>5470-5650 MARITIME RADIONAVIGATION Radiolocation</p> <p>S5.450 S5.451 S5.452</p>	<p>5470-5600 MARITIME RADIONAVIGATION Radiolocation G56</p> <p>US50 US65</p>	<p>5470-5600 MARITIME RADIONAVIGATION Radiolocation</p> <p>US50 US65</p>	<p>Maritime (80)</p>
	<p>5600-5650 MARITIME RADIONAVIGATION METEOROLOGICAL AIDS Radiolocation US51 G56</p> <p>S5.452 US65</p>	<p>5600-5650 MARITIME RADIONAVIGATION METEOROLOGICAL AIDS Radiolocation US51</p> <p>S5.452 US65</p>	

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

5650–7250 MHz (SHF)					Remarks	
International Table			United States Table			
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government		
5650–5725 RADIOLOCATION Amateur Space research (deep space) S5.282 S5.451 S5.453 S5.454 S5.455			5650–5925 RADIOLOCATION G2	5650–5830 Amateur	ISM Equipment (18) Amateur (97)	
5725–5830 FIXED-SATELLITE (Earth-to-space) RADIOLOCATION Amateur S5.150 S5.451 S5.453 S5.455 S5.456		5725–5830 RADIOLOCATION Amateur S5.150 S5.453 S5.455		S5.150 S5.282		
5830–5850 FIXED-SATELLITE (Earth-to-space) RADIOLOCATION Amateur Amateur-satellite (space-to-Earth) S5.150 S5.451 S5.453 S5.455 S5.456		5830–5850 RADIOLOCATION Amateur Amateur-satellite (space-to-Earth) S5.150 S5.453 S5.455		5830–5850 Amateur Amateur-satellite (space-to-Earth) S5.150		
5850–5925 FIXED FIXED-SATELLITE (Earth-to-space) MOBILE S5.150		5850–5925 FIXED FIXED-SATELLITE (Earth-to-space) MOBILE Amateur Radiolocation S5.150		5850–5925 FIXED FIXED-SATELLITE (Earth-to-space) MOBILE Radiolocation S5.150		5850–5925 FIXED-SATELLITE (Earth-to-space) US245 MOBILE NG160 Amateur S5.150 US245
5925–6700 FIXED FIXED-SATELLITE (Earth-to-space) MOBILE			5925–6425	5925–6425 FIXED NG41 FIXED-SATELLITE (Earth-to-space)	International Fixed (23) Satellite Communications (25) Fixed Microwave (101)	

S5.149 S5.440 S5.458	6425–6525	6425–6525 FIXED-SATELLITE (Earth-to-space) MOBILE S5.440 S5.458	Auxiliary Broadcasting (74) Cable TV Relay (78) Fixed Microwave (101)
	6525–6875	6525–6875 FIXED FIXED-SATELLITE (Earth-to-space) S5.441	Satellite Communications (25) Fixed Microwave (101)
6700–7075 FIXED FIXED-SATELLITE Earth-to-space) (space-to-Earth) S5.441 MOBILE	S5.458	S5.458	
S5.458 S5.458A S5.458B S5.458C	6875–7125	6875–7075 FIXED FIXED-SATELLITE (Earth-to-space) S5.441 MOBILE S5.458 NG118	Auxiliary Broadcasting (74) Cable TV Relay (78)
		7075–7125 FIXED MOBILE S5.458 NG118	
7075–7250 FIXED MOBILE	S5.458	S5.458 NG118	
S5.458 S5.459 S5.460	7125–7190 FIXED	7125–7190 S5.458 US252 G116	
	7190–7235 FIXED SPACE RESEARCH (Earth-to-space) S5.458	7190–7250	
	7235–7250 FIXED S5.458	S5.458	

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

7250–8215 MHz (SHF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
7250–7300 FIXED FIXED-SATELLITE (space-to-Earth) MOBILE S5.461			7250–7300 FIXED FIXED-SATELLITE (space-to-Earth) MOBILE-SATELLITE (space-to-Earth) Fixed G117	7250–8025	
7300–7450 FIXED FIXED-SATELLITE (space-to-Earth) MOBILE except aeronautical mobile S5.461			7300–7450 FIXED-SATELLITE (space-to-Earth) MOBILE-SATELLITE (space-to-Earth) Fixed G117		
7450–7550 FIXED FIXED-SATELLITE (space-to-Earth) METEOROLOGICAL-SATELLITE (space-to-Earth) MOBILE except aeronautical mobile S5.461A			7450–7550 FIXED FIXED-SATELLITE (space-to-Earth) METEOROLOGICAL-SAT- ELLITE (space-to-Earth) Mobile-satellite (space-to-Earth) G104 G117		
7550–7750 FIXED FIXED-SATELLITE (space-to-Earth) MOBILE except aeronautical mobile			7550–7750 FIXED FIXED-SATELLITE (space-to-Earth) Mobile-satellite (space-to-Earth) G117		
7750–7850 FIXED METEOROLOGICAL-SATELLITE (space-to-Earth) S5.416B MOBILE except aeronautical mobile			7750–7900 FIXED		
7850–7900 FIXED MOBILE except aeronautical mobile					

<p>7900-8025 FIXED FIXED-SATELLITE (Earth-to-space) MOBILE</p> <p>S5.461</p>	<p>7900-8025 FIXED-SATELLITE (Earth-to-space) MOBILE-SATELLITE (Earth-to-space) Fixed</p> <p>G117</p>		
<p>8025-8175 EARTH EXPLORATION-SATELLITE (space-to-Earth) FIXED FIXED-SATELLITE (Earth-to-space) MOBILE S5.463</p> <p>S5.462A</p>	<p>8025-8175 EARTH EXPLORATION-SATELLITE (space-to-Earth) FIXED FIXED-SATELLITE (Earth-to-space) Mobile-satellite (Earth-to-space) (no airborne transmissions)</p> <p>US258 G117</p>	<p>8025-8175</p> <p>US258</p>	
<p>8175-8215 EARTH EXPLORATION-SATELLITE (space-to-Earth) FIXED FIXED-SATELLITE (Earth-to-space) METEOROLOGICAL-SATELLITE (Earth-to-space) MOBILE S5.463</p> <p>S5.462A</p>	<p>8175-8215 EARTH EXPLORATION-SATELLITE (space-to-Earth) FIXED FIXED-SATELLITE (Earth-to-space) METEOROLOGICAL-SATELLITE (Earth-to-space) Mobile-satellite (Earth-to-space) (no airborne transmissions)</p> <p>US258 G104 G117</p>	<p>8175-8215</p>	

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

8215-10000 MHz (SHF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
8215-8400 EARTH EXPLORATION-SATELLITE (space-to-Earth) FIXED FIXED-SATELLITE (Earth-to-space) MOBILE S5.463 S5.462A			8215-8400 EARTH EXPLORATION-SATELLITE (space-to-Earth) FIXED FIXED-SATELLITE (Earth-to-space) Mobile-satellite (Earth-to-space) (no airborne transmissions) US258 G117	8215-8400 US258	
8400-8500 FIXED MOBILE except aeronautical mobile SPACE RESEARCH (space-to-Earth) S5.465 S5.466 S5.467			8400-8450 FIXED SPACE RESEARCH (space-to-Earth) (deep space only) 8450-8500 FIXED SPACE RESEARCH (space-to-Earth)	8400-8450 8450-8500 SPACE RESEARCH (space-to-Earth)	
8500-8550 RADIOLOCATION S5.468 S5.469			8500-9000 RADIOLOCATION S5.333 US110 G59 US53	8500-9000 Radiolocation S5.333 US110	
8550-8650 EARTH EXPLORATION-SATELLITE (active) RADIOLOCATION SPACE RESEARCH (active) S5.468 S5.469 S5.469A					
8650-8750 RADIOLOCATION S5.468 S5.469					
8750-8850 RADIOLOCATION AERONAUTICAL RADIONAVIGATION S5.470 S5.471					
8850-9000 RADIOLOCATION MARITIME RADIONAVIGATION S5.472 S5.473					

<p>9000–9200 AERONAUTICAL RADIONAVIGATION S5.337 Radiolocation</p> <p>S5.471</p>	<p>9000–9200 AERONAUTICAL RADIO- NAVIGATION S5.337 Radiolocation G2</p> <p>US48 US54 G19</p>	<p>9000–9200 AERONAUTICAL RADIO- NAVIGATION S5.337 Radiolocation</p> <p>US48 US54</p>	<p>Aviation (87)</p>
<p>9200–9300 RADIOLOCATION MARITIME RADIONAVIGATION S5.472</p> <p>S5.473 S5.474</p>	<p>9200–9300 MARITIME RADIO- NAVIGATION S5.472 Radiolocation US110 G59</p> <p>S5.474</p>	<p>9200–9300 MARITIME RADIO- NAVIGATION S5.472 Radiolocation US110</p> <p>S5.474</p>	
<p>9300–9500 RADIONAVIGATION S5.476 Radiolocation</p> <p>S5.427 S5.474 S5.475</p>	<p>9300–9500 RADIONAVIGATION S5.476 US66 Radiolocation US51 G56 Meteorological aids</p> <p>S5.427 S5.474 US67 US71</p>	<p>9300–9500 RADIONAVIGATION S5.476 US66 Radiolocation US51 Meteorological aids</p> <p>S5.427 S5.474 US67 US71</p>	
<p>9500–9800 EARTH EXPLORATION-SATELLITE (active) RADIOLOCATION RADIONAVIGATION SPACE RESEARCH (active)</p> <p>S5.476A</p>	<p>9500–10000 RADIOLOCATION S5.333 US110</p>	<p>9500–10000 RADIOLOCATION S5.333 US110</p>	
<p>9800–10000 RADIOLOCATION Fixed</p> <p>S5.477 S5.478 S5.479</p>	<p>S5.479</p>	<p>S5.479</p>	

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

10-12.7 GHz (SHF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
10-10.45 FIXED MOBILE RADIOLOCATION Amateur S5.479	10-10.45 RADIOLOCATION Amateur S5.479 S5.480	10-10.45 FIXED MOBILE RADIOLOCATION Amateur S5.479	10-10.45 RADIOLOCATION S5.479 US58 US108 G32	10-10.45 Radiolocation Amateur S5.479 US58 US108 NG42	Private Land Mobile (90) Amateur (97)
10.45-10.5 RADIOLOCATION Amateur Amateur-satellite S5.481			10.45-10.5 RADIOLOCATION US58 US108 G32	10.45-10.5 Radiolocation Amateur Amateur-satellite US58 US108 NG42 NG134	
10.5-10.55 FIXED MOBILE Radiolocation	10.5-10.55 FIXED MOBILE RADIOLOCATION		10.5-10.55 RADIOLOCATION US59		Private Land Mobile (90)
10.55-10.6 FIXED MOBILE except aeronautical mobile Radiolocation			10.55-10.6	10.55-10.6 FIXED	Fixed Microwave (101)
10.6-10.68 EARTH EXPLORATION-SATELLITE (passive) FIXED MOBILE except aeronautical mobile RADIO ASTRONOMY SPACE RESEARCH (passive) Radiolocation S5.149 S5.482			10.6-10.68 EARTH EXPLORATION-SATELLITE (passive) SPACE RESEARCH (passive) US265 US277	10.6-10.68 EARTH EXPLORATION-SATELLITE (passive) FIXED SPACE RESEARCH (passive) US265 US277	
10.68-10.7 EARTH EXPLORATION-SATELLITE (passive) RADIO ASTRONOMY SPACE RESEARCH (passive) S5.340 S5.483			10.68-10.7 EARTH EXPLORATION-SATELLITE (passive) RADIO ASTRONOMY US74 SPACE RESEARCH (passive) US246		
10.7-11.7 FIXED FIXED-SATELLITE (space-to-Earth) S5.441 S5.484A (Earth-to-space) S5.484 MOBILE except aeronautical mobile	10.7-11.7 FIXED FIXED-SATELLITE (space-to-Earth) S5.441 S5.484A MOBILE except aeronautical mobile		10.7-11.7 US211	10.7-11.7 FIXED NG41 FIXED-SATELLITE (space-to-Earth) S5.441 US211 NG104	International Fixed (23) Satellite Communications (25) Fixed Microwave (101)

11.7–12.5 FIXED MOBILE except aeronautical mobile BROADCASTING BROADCASTING- SATELLITE S5.487 S5.487A S5.492	11.7–12.1 FIXED S5.486 FIXED-SATELLITE (space-to-Earth) S5.484A Mobile except aeronautical mobile S5.485 S5.488	11.7–12.2 FIXED MOBILE except aeronautical mobile BROADCASTING BROADCASTING- SATELLITE S5.487 S5.487A S5.492	11.7–12.1 S5.486	11.7–12.1 FIXED-SATELLITE (space- to-Earth) NG143 NG145 Mobile except aeronautical mobile S5.486 S5.488	Satellite Communications (25) Fixed Microwave (101)
	12.1–12.2 FIXED-SATELLITE (space-to-Earth) S5.484A S5.485 S5.488 S5.489		12.1–12.2 S5.488	12.1–12.2 FIXED-SATELLITE (space- to-Earth) NG143 NG145 Mobile except aeronautical mobile S5.488	
	12.2–12.7 FIXED MOBILE except aeronautical mobile BROADCASTING BROADCASTING- SATELLITE S5.487A S5.488 S5.490 S5.492		12.2–12.7 FIXED MOBILE except aeronautical mobile BROADCASTING BROADCASTING- SATELLITE S5.484A S5.487 S5.491	12.2–12.7 FIXED MOBILE except aeronautical mobile BROADCASTING BROADCASTING- SATELLITE S5.488 S5.490	
12.5–12.75 FIXED-SATELLITE (space-to-Earth) S5.484A (Earth-to-space) S5.494 S5.495 S5.496	12.5–12.75 FIXED FIXED-SATELLITE (space-to-Earth) S5.484A MOBILE except aeronautical mobile BROADCASTING- SATELLITE S5.493 See next page for 12.7–12.75 GHz	12.5–12.75 FIXED FIXED-SATELLITE (space-to-Earth) S5.484A MOBILE except aeronautical mobile BROADCASTING- SATELLITE S5.493 See next page for 12.7–12.75 GHz	See next page for 12.7–12.75 GHz	See next page for 12.7–12.75 GHz	

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

12.7–14.5 GHz (SHF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
See previous page for 12.5–12.75 GHz	12.7–12.75 FIXED FIXED-SATELLITE (Earth-to-space) MOBILE except aeronautical mobile	See previous page for 12.5–12.75 GHz	12.7–12.75	12.7–12.75 FIXED FIXED-SATELLITE (Earth-to-space) MOBILE NG53 NG118	Auxiliary Broadcasting (74) Cable TV Relay (78) Fixed Microwave (101)
12.75–13.25 FIXED FIXED-SATELLITE (Earth-to-space) S5.441 MOBILE Space research (deep space) (space-to-Earth)			12.75–13.25 US251	12.75–13.25 FIXED FIXED-SATELLITE (Earth-to-space) S5.441 NG104 MOBILE US251 NG53 NG118	
13.25–13.4 EARTH EXPLORATION-SATELLITE (active) AERONAUTICAL RADIONAVIGATION S5.497 SPACE RESEARCH (active) S5.498A S5.499			13.25–13.4 AERONAUTICAL RADIONAVIGATION S5.497 Space research (Earth-to-space)		Aviation (87)
13.4–13.75 EARTH EXPLORATION-SATELLITE (active) RADIOLOCATION SPACE RESEARCH S5.501A Standard frequency and time signal-satellite (Earth-to-space) S5.499 S5.500 S5.501 S5.501B			13.4–13.75 RADIOLOCATION S5.333 US110 G59 Space research Standard frequency and time signal-satellite (Earth-to-space)	13.4–13.75 Radiolocation S5.333 US110 Space research Standard frequency and time signal-satellite (Earth-to-space)	Private Land Mobile (90)
13.75–14 FIXED-SATELLITE (Earth-to-space) S5.484A RADIOLOCATION Standard frequency and time signal-satellite (Earth-to-space) Space research S5.499 S5.500 S5.501 S5.502 S5.503 S5.503A			13.75–14 RADIOLOCATION US110 G59 Standard frequency and time signal-satellite (Earth-to-space) Space research US337	13.75–14 RADIOLOCATION US337 Radiolocation US110 Standard frequency and time signal-satellite (Earth-to-space) Space research	Satellite Communications (25) Private Land Mobile (90)
			S5.502 S5.503 S5.503A	S5.502 S5.503 S5.503A	

<p>14-14.25 FIXED-SATELLITE (Earth-to-space) S5.484A S5.506 RADIONAVIGATION S5.504 Mobile-satellite (Earth-to-space) except aeronautical mobile-satellite Space research</p>			<p>14-14.2 RADIONAVIGATION US292 Space research</p>	<p>14-14.2 FIXED-SATELLITE (Earth-to-space) RADIONAVIGATION US292 Land mobile-satellite (Earth-to-space) Space research</p>	<p>Satellite Communications (25) Maritime (80) Aviation (87)</p>
<p>S5.505</p>			<p>14.2-14.4</p>	<p>14.2-14.4 FIXED-SATELLITE (Earth-to-space) Land mobile-satellite (Earth-to-space) Mobile except aeronautical mobile</p>	<p>Satellite Communications (25) Fixed Microwave (101)</p>
<p>14.25-14.3 FIXED-SATELLITE (Earth-to-space) S5.484A S5.506 RADIONAVIGATION S5.504 Mobile-satellite (Earth-to-space) except aeronautical mobile-satellite Space research S5.505 S5.508 S5.509</p>					
<p>14.3-14.4 FIXED FIXED-SATELLITE (Earth-to-space) S5.484A S5.506 MOBILE except aeronautical mobile Mobile-satellite (Earth-to-space) except aeronautical mobile-satellite Radionavigation-satellite</p>	<p>14.3-14.4 FIXED-SATELLITE (Earth-to-space) S5.484A S5.506 Mobile-satellite (Earth-to-space) except aeronautical mobile-satellite Radionavigation-satellite</p>	<p>14.3-14.4 FIXED FIXED-SATELLITE (Earth-to-space) S5.484A S5.506 MOBILE except aeronautical mobile Mobile-satellite (Earth-to-space) except aeronautical mobile-satellite Radionavigation-satellite</p>			
<p>14.4-14.47 FIXED FIXED-SATELLITE (Earth-to-space) S5.484A S5.506 MOBILE except aeronautical mobile Mobile-satellite (Earth-to-space) except aeronautical mobile-satellite Space research (space-to-Earth)</p>			<p>14.4-14.47 Fixed Mobile</p>	<p>14.4-14.47 FIXED-SATELLITE (Earth-to-space) Land mobile-satellite (Earth-to-space)</p>	<p>Satellite Communications (25)</p>
<p>14.47-14.5 FIXED FIXED-SATELLITE (Earth-to-space) S5.484A S5.506 MOBILE except aeronautical mobile Mobile-satellite (Earth-to-space) except aeronautical mobile-satellite Radio astronomy S5.149</p>			<p>14.47-14.5 Fixed Mobile S5.149 US203</p>	<p>14.47-14.5 FIXED-SATELLITE (Earth-to-space) Land mobile-satellite (Earth-to-space) S5.149 US203</p>	

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

14.5–18.6 GHz (SHF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
14.5–14.8 FIXED FIXED-SATELLITE (Earth-to-space) S5.510 MOBILE Space research			14.5–14.7145 FIXED Mobile Space research	14.5–15.1365	
14.8–15.35 FIXED MOBILE Space research			14.7145–15.1365 MOBILE Fixed Space research	14.7145–15.1365	
			US310	US310	
			15.1365–15.35 FIXED Mobile Space research	15.1365–15.35	
S5.339			S5.339 US211	S5.339 US211	
15.35–15.4 EARTH EXPLORATION-SATELLITE (passive) RADIO ASTRONOMY SPACE RESEARCH (passive)			15.35–15.4 EARTH EXPLORATION-SATELLITE (passive) RADIO ASTRONOMY US74 SPACE RESEARCH (passive)		
S5.340 S5.511			US246		
15.4–15.43 AERONAUTICAL RADIONAVIGATION			15.4–15.7 AERONAUTICAL RADIONAVIGATION US260		Aviation (87)
S5.511D					
15.43–15.63 FIXED-SATELLITE (space-to-Earth) (Earth-to-space) S5.511A AERONAUTICAL RADIONAVIGATION					
S5.511C					
15.63–15.7 AERONAUTICAL RADIONAVIGATION					
S5.511D			733 797 US211		
15.7–16.6 RADIOLOCATION			15.7–16.6 RADIOLOCATION US110 G59	15.7–17.2 Radiolocation US110	Private Land Mobile (90)
S5.512 S5.513					

16.6–17.1 RADIOLOCATION Space research (deep space) (Earth-to-space) S5.512 S5.513		16.6–17.1 RADIOLOCATION US110 G59 Space research (deep space) (Earth-to-space)			
17.1–17.2 RADIOLOCATION S5.512 S5.513		17.1–17.2 RADIOLOCATION US110 G59			
17.2–17.3 EARTH EXPLORATION-SATELLITE (active) RADIOLOCATION SPACE RESEARCH (active) S5.512 S5.513 S5.513A		17.2–17.3 RADIOLOCATION US110 G59 Earth exploration-satellite (active) Space research (active)		17.2–17.3 Radiolocation US110 Earth exploration-satellite (active) Space research (active)	
17.3–17.7 FIXED-SATELLITE (Earth-to-space) S5.516 Radiolocation S5.514	17.3–17.7 FIXED-SATELLITE (Earth-to-space) S5.516 BROADCASTING-SATELLITE Radiolocation S5.514 S5.515 S5.517	17.3–17.7 FIXED-SATELLITE (Earth-to-space) S5.516 Radiolocation S5.514	17.3–17.7 Radiolocation US259 G59 US271	17.3–17.7 FIXED-SATELLITE (Earth-to-space) US271 US259	
17.7–18.1 FIXED FIXED-SATELLITE (space-to-Earth) S5.484A (Earth-to-space) S5.516 MOBILE	17.7–17.8 FIXED FIXED-SATELLITE (space-to-Earth) (Earth-to-space) S5.516 BROADCASTING-SATELLITE MOBILE S5.518 S5.515 S5.517	17.7–18.1 FIXED FIXED-SATELLITE (space-to-Earth) S5.484A (Earth-to-space) S5.516 MOBILE	17.7–17.8 US271	17.7–17.8 FIXED FIXED-SATELLITE (space-to-Earth) (Earth-to-space) US271 MOBILE NG144	Auxiliary Broadcasting (74) Cable TV Relay (78) Fixed Microwave (101)
	17.8–18.1 FIXED FIXED-SATELLITE (space-to-Earth) S5.484A (Earth-to-space) S5.516 MOBILE		17.8–18.6 S5.519 US334 G117	17.8–18.6 FIXED FIXED-SATELLITE (space-to-Earth) MOBILE S5.519 US334 NG144	Domestic Public Fixed (21) Satellite Communications (25) Auxiliary Broadcasting (74) Cable TV Relay (78) Fixed Microwave (101)
See next page for 18.1–18.6 GHz					

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

18.6–22.5 GHz (SHF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
18.1–18.4 FIXED FIXED-SATELLITE (space-to-Earth) S5.484A (Earth-to-space) S5.520 MOBILE S5.519 S5.521			See previous page for 17.8–18.6 GHz		See previous page for 17.8–19.7 GHz
18.4–18.6 FIXED FIXED-SATELLITE (space-to-Earth) S5.484A MOBILE					
18.6–18.8 FIXED FIXED-SATELLITE (space-to-Earth) S5.523 MOBILE except aeronautical mobile Earth exploration-satellite (passive) Space research (passive) S5.522	18.6–18.8 EARTH EXPLORATION- SATELLITE (passive) FIXED FIXED-SATELLITE (space-to-Earth) S5.523 MOBILE except aeronautical mobile MOBILE except aeronautical mobile SPACE RESEARCH (passive) S5.522	18.6–18.8 FIXED FIXED-SATELLITE (space-to-Earth) S5.523 MOBILE except aeronautical mobile Earth exploration-satellite (passive) Space research (passive) S5.522	18.6–18.8 EARTH EXPLORATION- SATELLITE (passive) SPACE RESEARCH (passive) US254 US255 US334 G117	18.6–18.8 EARTH EXPLORATION- SATELLITE (passive) FIXED FIXED-SATELLITE (space-to-Earth) MOBILE except aeronautical mobile SPACE RESEARCH (passive) US254 US255 US334 NG114	
18.8–19.3 FIXED FIXED-SATELLITE (space-to-Earth) S5.523A MOBILE			18.8–20.2		18.8–19.7 FIXED FIXED-SATELLITE (space-to-Earth) MOBILE US334 NG144
19.3–19.7 FIXED FIXED-SATELLITE (space-to-Earth) (Earth-to-space) S5.523B S5.523C S5.523D S5.523E MOBILE					
19.7–20.1 FIXED-SATELLITE (space-to-Earth) S5.484A Mobile-satellite (space-to-Earth) S5.524	19.7–20.1 FIXED-SATELLITE (space-to-Earth) S5.484A MOBILE-SATELLITE (space-to-Earth) S5.524 S5.525 S5.526 S5.527 S5.528 S5.529	19.7–20.1 FIXED-SATELLITE (space-to-Earth) S5.484A Mobile-satellite (space-to-Earth) S5.524	19.7–20.1 FIXED-SATELLITE (space-to-Earth) MOBILE-SATELLITE (space-to-Earth) S5.525 S5.526 S5.527 S5.528 S5.529 US334		

20.1–20.2 FIXED-SATELLITE (space-to-Earth) S5.484A MOBILE-SATELLITE (space-to-Earth)			US334 G117	20.1–20.2 FIXED-SATELLITE (space-to-Earth) MOBILE-SATELLITE (space-to-Earth)	
S5.524 S5.525 S5.526 S5.527 S5.528				S5.525 S5.526 S5.527 S5.528 US334	
20.2–21.2 FIXED-SATELLITE (space-to-Earth) MOBILE-SATELLITE (space-to-Earth) Standard frequency and time signal-satellite (space-to-Earth)			20.2–21.2 FIXED-SATELLITE (space-to-Earth) MOBILE-SATELLITE (space-to-Earth) Standard frequency and time signal-satellite (space-to-Earth)	20.2–21.2 Standard frequency and time signal-satellite (space-to-Earth)	
S5.524			G117		
21.2–21.4 EARTH EXPLORATION-SATELLITE (passive) FIXED MOBILE SPACE RESEARCH (passive)			21.2–21.4 EARTH EXPLORATION-SATELLITE (passive) FIXED MOBILE SPACE RESEARCH (passive)		Fixed Microwave (101)
US263					
21.4–22 FIXED MOBILE BROADCASTING-SATELLITE S5.530	21.4–22 FIXED MOBILE	21.4–22 FIXED MOBILE BROADCASTING-SATELLITE S5.530 S5.531	21.4–22 FIXED MOBILE		
22–22.21 FIXED MOBILE except aeronautical mobile			22–22.21 FIXED MOBILE except aeronautical mobile		
S5.149			S5.149		
22.21–22.5 EARTH EXPLORATION-SATELLITE (passive) FIXED MOBILE except aeronautical mobile RADIO ASTRONOMY SPACE RESEARCH (passive)			22.21–22.5 EARTH EXPLORATION-SATELLITE (passive) FIXED MOBILE except aeronautical mobile RADIO ASTRONOMY SPACE RESEARCH (passive)		
S5.149 S5.532			S5.149 US263		

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

22.5–27.5 GHz (SHF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
22.5–22.55 FIXED MOBILE			22.5–22.55 FIXED MOBILE US211		See previous page for 22.21–22.55 GHz
22.55–23.55 FIXED INTER-SATELLITE MOBILE S5.149			22.55–23.55 FIXED INTER-SATELLITE MOBILE S5.149 US278		Satellite Communications (25) Fixed Microwave (101)
23.55–23.6 FIXED MOBILE			23.55–23.6 FIXED MOBILE		Fixed Microwave (101)
23.6–24 EARTH EXPLORATION-SATELLITE (passive) RADIO ASTRONOMY SPACE RESEARCH (passive) S5.340			23.6–24 EARTH EXPLORATION-SATELLITE (passive) RADIO ASTRONOMY US74 SPACE RESEARCH (passive) US246		
24–24.05 AMATEUR AMATEUR-SATELLITE S5.150			24–24.05 S5.150 US211	24–24.05 AMATEUR AMATEUR-SATELLITE S5.150 US211	ISM Equipment (18) Amateur (97)
24.05–24.25 RADIOLOCATION Amateur Earth exploration-satellite (active) S5.150			24.05–24.25 RADIOLOCATION US110 G59 Earth exploration-satellite (active) S5.150	24.05–24.25 Radiolocation US110 Amateur Earth exploration-satellite (active) S5.150	ISM Equipment (18) Private Land Mobile (90) Amateur (97)
24.25–24.45 FIXED	24.25–24.45 RADIONAVIGATION	24.25–24.45 RADIONAVIGATION FIXED MOBILE	24.25–24.45	24.25–24.45 RADIONAVIGATION FIXED	Aviation (87) Fixed Microwave (101)
24.45–24.75 FIXED INTER-SATELLITE	24.45–24.65 INTER-SATELLITE RADIONAVIGATION S5.533	24.45–24.65 FIXED INTER-SATELLITE MOBILE RADIONAVIGATION S5.533	24.45–24.65 INTER-SATELLITE RADIONAVIGATION S5.533		Satellite Communications (25)

	24.65–24.75 INTER-SATELLITE RADIOLOCATION-SATELLITE (Earth-to-space)	24.65–24.75 FIXED INTER-SATELLITE MOBILE S5.533 S5.534	24.65–24.75 INTER-SATELLITE RADIOLOCATION-SATELLITE (Earth-to-space)	
24.75–25.25 FIXED	24.75–25.25 FIXED-SATELLITE (Earth-to-space) S5.535	24.75–25.25 FIXED FIXED-SATELLITE (Earth-to-space) S5.535 MOBILE S5.534	24.75–25.05 RADIONAVIGATION	Aviation (87)
			25.05–25.25	25.05–25.25 RADIONAVIGATION FIXED Aviation (87) Fixed Microwave (101)
25.25–25.5 FIXED INTER-SATELLITE S5.536 MOBILE Standard frequency and time signal-satellite (Earth-to-space)		25.25–25.5 FIXED INTER-SATELLITE S5.536 MOBILE Standard frequency and time signal-satellite (Earth-to-space)		25.25–27 Standard frequency and time signal-satellite (Earth-to-space) Earth exploration-satellite (space-to-space) Note: In this manual under the government column, a primary allocation has been added for the inter-satellite service in the bands 25.25–25.5, 25.5–27, and 27–27.5 GHz and the direction indicator for the Earth exploration-satellite service has been modified from space-to-space to space-to-Earth. At this time, the FCC has not adopted these changes.
25.5–27 EARTH EXPLORATION-SATELLITE (space-to-Earth) S5.536A S5.536B FIXED INTER-SATELLITE S5.536 MOBILE Standard frequency and time signal-satellite (Earth-to-space)		25.5–27 FIXED INTER-SATELLITE S5.536 MOBILE Earth exploration-satellite (space-to-Earth) Standard frequency and time signal-satellite (Earth-to-space)		
27–27.5 FIXED INTER-SATELLITE S5.536 MOBILE	27–27.5 FIXED FIXED-SATELLITE (Earth-to-space) INTER-SATELLITE S5.536 S5.537 MOBILE	27–27.5 FIXED INTER-SATELLITE S5.536 MOBILE		27–27.5 Earth exploration-satellite (space-to-space)

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

27.5–32 GHz (SHF/EHF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
27.5–28.5 FIXED FIXED-SATELLITE (Earth-to-space) S5.484A S5.539 MOBILE S5.538 S5.540			27.5–30 27.5–29.5 FIXED FIXED-SATELLITE (Earth-to-space) MOBILE		Satellite Communications (25) Fixed Microwave (101)
28.5–29.1 FIXED FIXED-SATELLITE (Earth-to-space) S5.484A S5.523A S5.539 MOBILE Earth exploration-satellite (Earth-to-space) S5.541 S5.540					
29.1–29.5 FIXED FIXED-SATELLITE (Earth-to-space) S5.523C S5.523E S5.535A S5.539 S5.541A MOBILE Earth exploration-satellite (Earth-to-space) S5.541 S5.540					
29.5–29.9 FIXED-SATELLITE (Earth-to-space) S5.484A S5.539 Earth exploration-satellite (Earth-to-space) S5.541 Mobile-satellite (Earth-to-space) S5.540 S5.542	29.5–29.9 FIXED-SATELLITE (Earth-to-space) S5.484A S5.539 MOBILE-SATELLITE (Earth-to-space) Earth exploration-satellite (Earth-to-space) S5.541 S5.525 S5.526 S5.527 S5.529 S5.540 S5.542	29.5–29.9 FIXED-SATELLITE (Earth-to-space) S5.484A S5.539 Earth exploration-satellite (Earth-to-space) S5.541 Mobile-satellite (Earth-to-space) S5.540 S5.542	29.5–29.9 FIXED-SATELLITE (Earth-to-space) MOBILE-SATELLITE (Earth-to-space) S5.525 S5.526 S5.527 S5.529		Satellite Communications (25)
29.9–30 FIXED-SATELLITE (Earth-to-space) S5.484A S5.539 MOBILE-SATELLITE (Earth-to-space) Earth exploration-satellite (Earth-to-space) S5.541 S5.543 S5.525 S5.526 S5.527 S5.538 S5.540 S5.542					

<p>30-31 FIXED-SATELLITE (Earth-to-space) MOBILE-SATELLITE (Earth-to-space) Standard frequency and time signal-satellite (space-to-Earth)</p>			<p>30-31 FIXED-SATELLITE (Earth-to-space) MOBILE-SATELLITE (Earth-to-space) Standard frequency and time signal-satellite (space-to- Earth)</p>	<p>30-31 Standard frequency and time signal-satellite (space-to- Earth)</p>	
<p>S5.542</p>			<p>G117</p>		
<p>31-31.3 FIXED MOBILE Standard frequency and time signal-satellite (space-to-Earth) Space research S5.544 S5.545</p>			<p>31-31.3 Standard frequency and time signal-satellite (space-to- Earth)</p>	<p>31-31.3 FIXED MOBILE Standard frequency and time signal-satellite (space- to-Earth)</p>	<p>Fixed Microwave (101)</p>
<p>S5.149</p>			<p>S5.149 US211</p>	<p>S5.149 US211</p>	
<p>31.3-31.5 EARTH EXPLORATION-SATELLITE (passive) RADIO ASTRONOMY SPACE RESEARCH (passive)</p>			<p>31.3-31.8 EARTH EXPLORATION-SATELLITE (passive) RADIO ASTRONOMY US74 SPACE RESEARCH (passive)</p>		
<p>S5.340</p>			<p>US246</p>		
<p>31.5-31.8 EARTH EXPLORATION- SATELLITE (passive) RADIO ASTRONOMY SPACE RESEARCH (passive) Fixed Mobile except aeronautical mobile</p>	<p>31.5-31.8 EARTH EXPLORATION- SATELLITE (passive) RADIO ASTRONOMY SPACE RESEARCH (passive)</p>	<p>31.5-31.8 EARTH EXPLORATION- SATELLITE (passive) RADIO ASTRONOMY SPACE RESEARCH (passive) Fixed Mobile except aeronautical mobile</p>			
<p>S5.149 S5.546</p>	<p>S5.340</p>	<p>S5.149</p>			
<p>31.8-32 FIXED RADIONAVIGATION SPACE RESEARCH (deep space) (space-to-Earth)</p>			<p>31.8-32 RADIONAVIGATION US69</p>		
<p>S5.547 S5.547B S5.548</p>			<p>US211 US262</p>		

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

32-40 GHz (EHF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
32-32.3 FIXED S5.547A INTER-SATELLITE RADIONAVIGATION SPACE RESEARCH (deep space) (space-to-Earth) S5.547 S5.547C S5.548			32-33 INTER-SATELLITE US278 RADIONAVIGATION US69		
32.3-33 FIXED S5.547A INTER-SATELLITE RADIONAVIGATION S5.547 S5.547D S5.548			S5.548 US262		
33-33.4 FIXED S5.547A RADIONAVIGATION S5.547 S5.547E			33-33.4 RADIONAVIGATION US69		
33.4-34.2 RADIOLOCATION S5.549			33.4-36 RADIOLOCATION US110 G34	33.4-36 Radiolocation US110	Private Land Mobile (90)
34.2-34.7 RADIOLOCATION SPACE RESEARCH (deep space) (Earth-to-space) S5.549					
34.7-35.2 RADIOLOCATION Space research S5.550 S5.549					
35.2-35.5 METEOROLOGICAL AIDS RADIOLOCATION S5.549					

<p>35.5–36 METEOROLOGICAL AIDS EARTH EXPLORATION-SATELLITE (active) RADIOLOCATION SPACE RESEARCH (active) S5.549 S5.551A</p>	<p>S5.551A US252</p>	<p>S5.551 US252</p>	
<p>36–37 EARTH EXPLORATION-SATELLITE (active) FIXED MOBILE SPACE RESEARCH (active) S5.149</p>	<p>36–37 EARTH EXPLORATION-SATELLITE (passive) FIXED MOBILE SPACE RESEARCH (passive) US263 US342</p>		
<p>37–37.5 FIXED MOBILE SPACE RESEARCH (space-to-Earth)</p>	<p>37–38 FIXED MOBILE SPACE RESEARCH (space-to-Earth)</p>	<p>37–37.6 FIXED MOBILE</p>	
<p>37.5–38 FIXED FIXED-SATELLITE (space-to-Earth) MOBILE SPACE RESEARCH (space-to-Earth) Earth exploration-satellite (space-to-Earth)</p>		<p>37.6–38.6 FIXED FIXED-SATELLITE (space-to-Earth) MOBILE</p>	<p>Satellite Communications (25)</p>
<p>38–39.5 FIXED FIXED-SATELLITE (space-to-Earth) MOBILE Earth exploration-satellite (space-to-Earth)</p>	<p>38–38.6 FIXED MOBILE</p>	<p>38.6–39.5 FIXED FIXED-SATELLITE (space-to-Earth) MOBILE</p>	<p>Auxiliary Broadcasting (74) Fixed Microwave (101)</p>
<p>39.5–40 FIXED FIXED-SATELLITE (space-to-Earth) MOBILE MOBILE-SATELLITE (space-to-Earth) Earth exploration-satellite (space-to-Earth)</p>	<p>39.5–40 FIXED-SATELLITE (space-to-Earth) MOBILE-SATELLITE (space-to-Earth) US291 G117</p>	<p>39.5–40 FIXED FIXED-SATELLITE (space-to-Earth) MOBILE MOBILE-SATELLITE (space-to-Earth) US291</p>	

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

40-50.2 GHz (EHF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
40-40.5 EARTH EXPLORATION-SATELLITE (Earth-to-space) FIXED FIXED-SATELLITE (space-to-Earth) MOBILE MOBILE-SATELLITE (space-to-Earth) SPACE RESEARCH (Earth-to-space) Earth exploration-satellite (space-to-Earth)			40-40.5 EARTH EXPLORATION-SATELLITE (Earth-to-space) FIXED-SATELLITE (space-to-Earth) MOBILE-SATELLITE (space-to-Earth) SPACE RESEARCH (Earth-to-space) Earth exploration-satellite (space-to-Earth) G117	40-40.5 FIXED-SATELLITE (space-to-Earth) MOBILE-SATELLITE (space-to-Earth)	Satellite Communications (25)
40.5-42.5 FIXED BROADCASTING BROADCASTING-SATELLITE Mobile S5.551B S5.551D	40.5-42.5 FIXED FIXED-SATELLITE (space-to-Earth) S5.551B S5.551E BROADCASTING BROADCASTING-SATELLITE Mobile S5.551C S5.551F		40.5-42.5 US211	40.5-41 FIXED-SATELLITE (space-to-Earth) BROADCASTING BROADCASTING-SATELLITE Mobile Fixed US211 41-42.5 FIXED BROADCASTING BROADCASTING-SATELLITE MOBILE US211	
42.5-43.5 FIXED FIXED-SATELLITE (Earth-to-space) S5.552 MOBILE except aeronautical mobile RADIO ASTRONOMY S5.149			42.5-43.5 FIXED FIXED-SATELLITE (Earth-to-space) MOBILE except aeronautical mobile RADIO ASTRONOMY US342	42.5-43.5 RADIO ASTRONOMY US342	

43.5-47 MOBILE S5.553 MOBILE-SATELLITE RADIONAVIGATION RADIONAVIGATION-SATELLITE S5.554	43.5-45.5 FIXED-SATELLITE (Earth-to-space) MOBILE-SATELLITE (Earth-to-space) G117	43.5-45.5	
	45.5-46.9 MOBILE MOBILE-SATELLITE (Earth-to-space) RADIONAVIGATION-SATELLITE S5.554		RF Devices (15)
	46.9-47 MOBILE MOBILE-SATELLITE (Earth- to-space) RADIONAVIGATION- SATELLITE S5.554	46.9-47 MOBILE MOBILE-SATELLITE (Earth-to-space) RADIONAVIGATION- SATELLITE FIXED S5.554	
47-47.2 AMATEUR AMATEUR-SATELLITE	47-48.2	47-47.2 AMATEUR AMATEUR-SATELLITE	Amateur (97)
47.2-50.2 FIXED FIXED-SATELLITE(Earth-to-space) S5.552 MOBILE S5.149 S5.340 S5.552A S5.555		47.2-48.2 FIXED FIXED-SATELLITE (Earth- to-space) US297 MOBILE US264 S5.555	
	48.2-50.2 FIXED FIXED-SATELLITE (Earth-to-space) US297 MOBILE US264 S5.555 US342		Satellite Communications (25)

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

50.2-65 GHz (EHF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
50.2-50.4 EARTH EXPLORATION-SATELLITE (passive) SPACE RESEARCH (passive) S5.340 S5.555A			50.2-50.4 EARTH EXPLORATION-SATELLITE (passive) FIXED MOBILE SPACE RESEARCH (passive) US263		
50.4-51.4 FIXED FIXED-SATELLITE (Earth-to-space) MOBILE Mobile-satellite (Earth-to-space)			50.4-51.4 FIXED FIXED-SATELLITE (Earth-to-space) MOBILE MOBILE-SATELLITE (Earth-to-space) G117	50.4-51.4 FIXED FIXED-SATELLITE (Earth-to-space) MOBILE MOBILE-SATELLITE (Earth-to-space)	
51.4-52.6 FIXED MOBILE S5.547 S5.556			51.4-54.25 EARTH EXPLORATION-SATELLITE (passive) SPACE RESEARCH (passive) RADIO ASTRONOMY		
52.6-54.25 EARTH EXPLORATION-SATELLITE (passive) SPACE RESEARCH (passive) S5.340 S5.556			US246		
54.25-55.78 EARTH EXPLORATION-SATELLITE (passive) INTER-SATELLITE S5.556A SPACE RESEARCH (passive) S5.556B			54.25-58.2 EARTH EXPLORATION-SATELLITE (passive) FIXED INTER-SATELLITE MOBILE 909 SPACE RESEARCH (passive)		
55.78-56.9 EARTH EXPLORATION-SATELLITE (passive) FIXED INTER-SATELLITE S5.556A MOBILE S5.558 SPACE RESEARCH (passive) S5.547 S5.557					

<p>56.9-57 EARTH EXPLORATION-SATELLITE (passive) FIXED INTER-SATELLITE S5.558A MOBILE S5.558 SPACE RESEARCH (passive) S5.547 S5.557</p>	<p>US263</p>	
<p>57-58.2 EARTH EXPLORATION-SATELLITE (passive) FIXED INTER-SATELLITE S5.556A MOBILE S5.558 SPACE RESEARCH (passive) S5.547 S5.557</p>		
<p>58.2-59 EARTH EXPLORATION-SATELLITE (passive) FIXED MOBILE SPACE RESEARCH (passive) S5.547 S5.556</p>	<p>58.2-59 EARTH EXPLORATION-SATELLITE (passive) SPACE RESEARCH (passive) RADIO ASTRONOMY US246</p>	
<p>59-59.3 EARTH EXPLORATION-SATELLITE (passive) FIXED INTER-SATELLITE S5.556A MOBILE S5.558 RADIOLOCATION S5.559 SPACE RESEARCH (passive)</p>	<p>59-64 FIXED INTER-SATELLITE MOBILE 909 RADIOLOCATION S5.138</p>	<p>RF Devices (15) ISM Equipment (18)</p>
<p>59.3-64 FIXED INTER-SATELLITE MOBILE S5.558 RADIOLOCATION S5.559 S5.138</p>		
<p>64-65 FIXED INTER-SATELLITE MOBILE except aeronautical mobile S5.547 S5.556</p>	<p>64-65 EARTH EXPLORATION-SATELLITE (passive) RADIO ASTRONOMY SPACE RESEARCH (passive) US246</p>	

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

65-95 GHz (EHF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
65-66 EARTH EXPLORATION-SATELLITE FIXED INTER-SATELLITE MOBILE except aeronautical mobile SPACE RESEARCH S5.547			65-66 EARTH EXPLORATION-SATELLITE SPACE RESEARCH Fixed Mobile		
66-71 INTER-SATELLITE MOBILE S5.553 S5.558 MOBILE-SATELLITE RADIONAVIGATION RADIONAVIGATION-SATELLITE S5.554			66-71 MOBILE S5.553 MOBILE-SATELLITE RADIONAVIGATION RADIONAVIGATION-SATELLITE S5.554		
71-74 FIXED FIXED-SATELLITE (Earth-to-space) MOBILE MOBILE-SATELLITE (Earth-to-space) S5.149 S5.556			71-74 FIXED FIXED-SATELLITE (Earth-to-space) MOBILE MOBILE-SATELLITE (Earth-to-space) US270		
74-75.5 FIXED FIXED-SATELLITE (Earth-to-space) MOBILE Space research (space-to-Earth)			74-75.5 FIXED FIXED-SATELLITE (Earth-to-space) US297 MOBILE		
75.5-76 AMATEUR AMATEUR-SATELLITE Space research (space-to-Earth)			75.5-76	75.5-76 AMATEUR AMATEUR-SATELLITE	Amateur (97)
76-81 RADIOLOCATION Amateur Amateur-satellite Space research (space-to-Earth)			76-81 RADIOLOCATION	76-77 RADIOLOCATION Amateur	RF Devices (15)
				77-77.5 RADIOLOCATION Amateur Amateur-satellite	Amateur (97)

<p>S5.560</p>	<p>S5.560</p>	<p>77.5-78 RADIOLOCATION AMATEUR AMATEUR-SATELLITE</p>	
<p>81-84 FIXED FIXED-SATELLITE (space-to-Earth) MOBILE MOBILE-SATELLITE (space-to-Earth) Space research (space-to-Earth)</p>	<p>81-84 FIXED FIXED-SATELLITE (space-to-Earth) MOBILE MOBILE-SATELLITE (space-to-Earth)</p>	<p>78-81 RADIOLOCATION Amateur Amateur-satellite S5.560</p>	
<p>84-86 FIXED MOBILE BROADCASTING BROADCASTING-SATELLITE</p> <p>S5.561</p>	<p>84-86 FIXED MOBILE</p> <p>S5.561 US211</p>	<p>84-86 FIXED MOBILE BROADCASTING BROADCASTING-SATELLITE</p> <p>S5.561 US211</p>	
<p>86-92 EARTH EXPLORATION-SATELLITE (passive) RADIO ASTRONOMY SPACE RESEARCH (passive)</p> <p>S5.340</p>	<p>86-92 EARTH EXPLORATION-SATELLITE (passive) RADIO ASTRONOMY US74 SPACE RESEARCH (passive)</p> <p>US246</p>		
<p>92-94 FIXED FIXED-SATELLITE (Earth-to-space) MOBILE RADIOLOCATION</p> <p>S5.149 S5.556</p>	<p>92-95 FIXED FIXED-SATELLITE (Earth-to-space) MOBILE RADIOLOCATION</p>		
<p>94-94.1 EARTH EXPLORATION-SATELLITE (active) RADIOLOCATION SPACE RESEARCH (active)</p> <p>S5.562</p>	<p>S5.149</p>		
<p>See next page for 94.1-95 GHz</p>			

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

95-150 GHz (EHF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
94.1-95 FIXED FIXED-SATELLITE (Earth-to-space) MOBILE RADIOLOCATION			See previous page for 92-95 GHz		See previous page for 92-95 GHz
95-100 MOBILE S5.553 MOBILE-SATELLITE RADIONAVIGATION RADIONAVIGATION-SATELLITE Radiolocation S5.149 S5.554 S5.555			95-100 MOBILE S5.553 MOBILE-SATELLITE RADIONAVIGATION RADIONAVIGATION-SATELLITE Radiolocation S5.149 S5.554		
100-102 EARTH EXPLORATION-SATELLITE (passive) FIXED MOBILE SPACE RESEARCH (passive) S5.341			100-102 EARTH EXPLORATION-SATELLITE (passive) SPACE RESEARCH (passive) S5.341 US246		
102-105 FIXED FIXED-SATELLITE (space-to-Earth) MOBILE S5.341			102-105 FIXED FIXED-SATELLITE (space-to-Earth) S5.341 US211		
105-116 EARTH EXPLORATION-SATELLITE (passive) RADIO ASTRONOMY SPACE RESEARCH (passive) S5.340 S5.341			105-116 EARTH EXPLORATION-SATELLITE (passive) RADIO ASTRONOMY US74 SPACE RESEARCH (passive) S5.341 US246		
116-119.98 EARTH EXPLORATION-SATELLITE (passive) FIXED INTER-SATELLITE MOBILE S5.558 SPACE RESEARCH (passive) S5.341			116-119.98 EARTH EXPLORATION-SATELLITE (passive) FIXED INTER-SATELLITE MOBILE S5.558 SPACE RESEARCH (passive) S5.341 US211 US263		

<p>119.98-120.02 EARTH EXPLORATION-SATELLITE (passive) FIXED INTER-SATELLITE MOBILE S5.558 SPACE RESEARCH (passive) Amateur S5.341</p>	<p>119.98-120.02 EARTH EXPLORATION-SATELLITE (passive) FIXED INTER-SATELLITE MOBILE S5.558 SPACE RESEARCH (passive) Amateur S5.341 US211 US263</p>		
<p>120.02-126 EARTH EXPLORATION-SATELLITE (passive) FIXED INTER-SATELLITE MOBILE S5.558 SPACE RESEARCH (passive) S5.138</p>	<p>120.02-126 EARTH EXPLORATION-SATELLITE (passive) FIXED INTER-SATELLITE MOBILE S5.558 SPACE RESEARCH (passive) S5.138 US211 US263</p>		ISM Equipment (18)
<p>126-134 FIXED INTER-SATELLITE MOBILE S5.558 RADIOLOCATION S5.559</p>	<p>126-134 FIXED INTER-SATELLITE MOBILE 909 RADIOLOCATION S5.559</p>		
<p>134-142 MOBILE S5.553 MOBILE-SATELLITE RADIONAVIGATION RADIONAVIGATION-SATELLITE Radiolocation S5.149 S5.340 S5.554 S5.555</p>	<p>134-142 MOBILE S5.553 MOBILE-SATELLITE RADIONAVIGATION RADIONAVIGATION-SATELLITE Radiolocation S5.149 S5.554 S5.555 917</p>		
<p>142-144 AMATEUR AMATEUR-SATELLITE</p>	<p>142-144</p>	<p>142-144 AMATEUR AMATEUR-SATELLITE</p>	Amateur (97)
<p>144-149 RADIOLOCATION Amateur Amateur-satellite S5.149 S5.555</p>	<p>144-149 RADIOLOCATION S5.149 S5.555</p>	<p>144-149 RADIOLOCATION Amateur Amateur-satellite S5.149 S5.555</p>	
<p>149-150 FIXED FIXED-SATELLITE (space-to-Earth) MOBILE</p>	<p>149-150 FIXED FIXED-SATELLITE (space-to-Earth) MOBILE</p>		

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

150-202 GHz (EHF)						
International Table			United States Table		Remarks	
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government		
150-151 EARTH EXPLORATION-SATELLITE (passive) FIXED FIXED-SATELLITE (space-to-Earth) MOBILE SPACE RESEARCH (passive) S5.149 S5.385			150-151 EARTH EXPLORATION-SATELLITE (passive) FIXED FIXED-SATELLITE (space-to-Earth) MOBILE SPACE RESEARCH (passive) S5.149 S5.385 US263			
151-156 FIXED FIXED-SATELLITE (space-to-Earth) MOBILE			151-164 FIXED FIXED-SATELLITE (space-to-Earth) US211			
156-158 EARTH EXPLORATION-SATELLITE (passive) FIXED FIXED-SATELLITE (space-to-Earth) MOBILE						
158-164 FIXED FIXED-SATELLITE (space-to-Earth) MOBILE						
164-168 EARTH EXPLORATION-SATELLITE (passive) RADIO ASTRONOMY SPACE RESEARCH (passive)			164-168 EARTH EXPLORATION-SATELLITE (passive) RADIO ASTRONOMY SPACE RESEARCH (passive) US246			
168-170 FIXED MOBILE			168-170 FIXED MOBILE			
170-174.5 FIXED INTER-SATELLITE MOBILE S5.558 S5.149 S5.385			170-174.5 FIXED INTER-SATELLITE MOBILE 909 S5.149 S5.385			

<p>174.5–176.5 EARTH EXPLORATION-SATELLITE (passive) FIXED INTER-SATELLITE MOBILE S5.558 SPACE RESEARCH (passive) S5.149 S5.385</p>	<p>174.5–176.5 EARTH EXPLORATION-SATELLITE (passive) FIXED INTER-SATELLITE MOBILE 909 SPACE RESEARCH (passive) S5.149 S5.385 US263</p>	
<p>176.5–182 FIXED INTER-SATELLITE MOBILE S5.558 S5.149 S5.385</p>	<p>176.5–182 FIXED INTER-SATELLITE MOBILE 909 S5.149 S5.385 US211</p>	
<p>182–185 EARTH EXPLORATION-SATELLITE (passive) RADIO ASTRONOMY SPACE RESEARCH (passive) S5.340 S5.563</p>	<p>182–185 EARTH EXPLORATION-SATELLITE (passive) RADIO ASTRONOMY SPACE RESEARCH (passive) US246</p>	
<p>185–190 FIXED INTER-SATELLITE MOBILE S5.558 S5.149 S5.385</p>	<p>185–190 FIXED INTER-SATELLITE MOBILE 909 S5.149 S5.385 US211</p>	
<p>190–200 MOBILE S5.553 MOBILE-SATELLITE RADIONAVIGATION RADIONAVIGATION-SATELLITE S5.341 S5.554</p>	<p>190–200 MOBILE S5.553 MOBILE-SATELLITE RADIONAVIGATION RADIONAVIGATION-SATELLITE S5.341 S5.554</p>	
<p>200–202 EARTH EXPLORATION-SATELLITE (passive) FIXED MOBILE SPACE RESEARCH (passive) S5.341</p>	<p>200–202 EARTH EXPLORATION-SATELLITE (passive) FIXED MOBILE SPACE RESEARCH (passive) S5.341 US263</p>	

TABLE 3 (CONT). TABLE OF FREQUENCY ALLOCATIONS*

202-400 GHz (EHF)					
International Table			United States Table		Remarks
Region 1	Region 2	Region 3	Federal Government	Non-Federal Government	
202-217 FIXED FIXED-SATELLITE (Earth-to-space) MOBILE S5.341			202-217 FIXED FIXED-SATELLITE (Earth-to-space) MOBILE S5.341		
217-231 EARTH EXPLORATION-SATELLITE (passive) RADIO ASTRONOMY SPACE RESEARCH (passive) S5.340 S5.341			217-231 EARTH EXPLORATION-SATELLITE (passive) RADIO ASTRONOMY US74 SPACE RESEARCH (passive) S5.341 US246		
231-235 FIXED FIXED-SATELLITE (Earth-to-space) MOBILE Radiolocation			231-235 FIXED FIXED-SATELLITE (space-to-Earth) MOBILE Radiolocation US211		
235-238 EARTH EXPLORATION-SATELLITE (passive) FIXED FIXED-SATELLITE (space-to-Earth) MOBILE SPACE RESEARCH (passive)			235-238 EARTH EXPLORATION-SATELLITE (passive) FIXED FIXED-SATELLITE (space-to-Earth) MOBILE SPACE RESEARCH (passive) US263		
238-241 FIXED FIXED-SATELLITE (space-to-Earth) MOBILE Radiolocation			238-241 FIXED FIXED-SATELLITE (space-to-Earth) MOBILE Radiolocation		
241-248 RADIOLOCATION Amateur Amateur-satellite S5.138			241-248 RADIOLOCATION S5.138	241-248 RADIOLOCATION Amateur Amateur-satellite S5.138	ISM Equipment (18) Amateur (97)

248-250 AMATEUR AMATEUR-SATELLITE	248-250	248-250 AMATEUR AMATEUR-SATELLITE	Amateur (97)
250-252 EARTH EXPLORATION-SATELLITE (passive) SPACE RESEARCH (passive) S5.149 S5.555	250-252 EARTH EXPLORATION-SATELLITE (passive) SPACE RESEARCH (passive) S5.149 S5.555		
252-265 MOBILE S5.553 MOBILE-SATELLITE RADIONAVIGATION RADIONAVIGATION-SATELLITE S5.149 S5.385 S5.554 S5.555 S5.564	252-265 MOBILE S5.553 MOBILE-SATELLITE RADIONAVIGATION RADIONAVIGATION-SATELLITE S5.149 S5.385 S5.554 S5.555 US211		
265-275 FIXED FIXED-SATELLITE (Earth-to-space) MOBILE RADIO ASTRONOMY S5.149	265-275 FIXED FIXED-SATELLITE (Earth-to-space) MOBILE RADIO ASTRONOMY S5.149		
275-400 (Not allocated) S5.565	275-300 FIXED MOBILE S5.565		
	300-400 (Not allocated) S5.565		Amateur (97)

* From *Manual of Regulations and Procedures for Federal Radio Frequency Management*, January 2000 Edition, as revised May and September 2000, US Department of Commerce, National Telecommunications and Information Administration, Washington, DC.

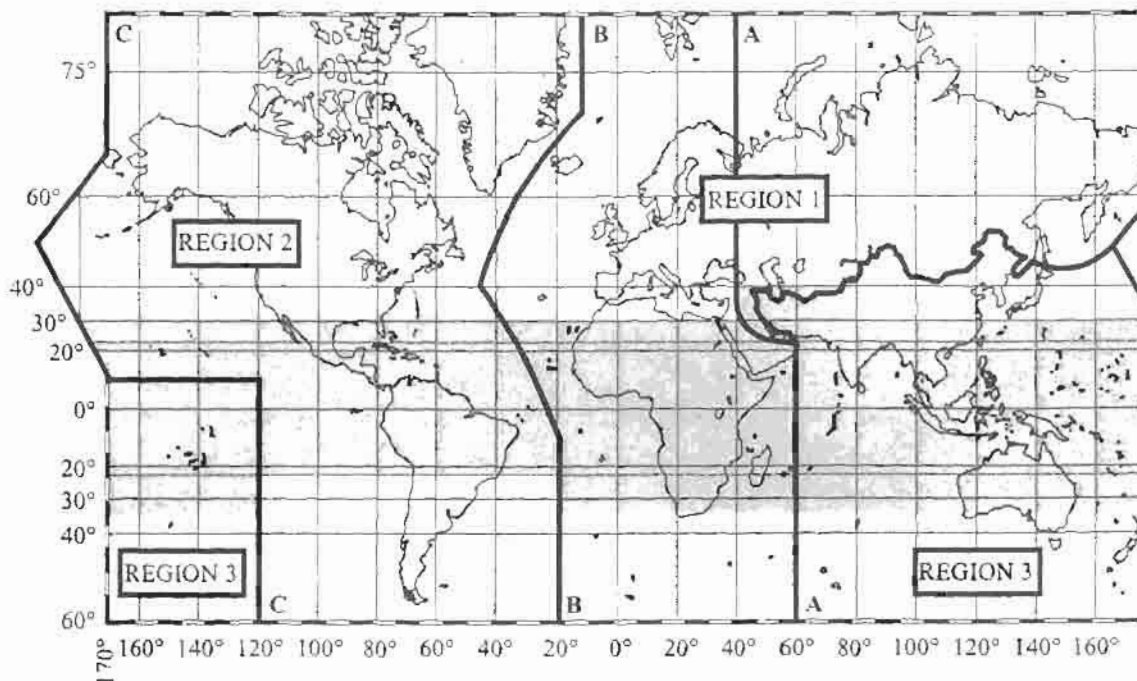


Fig. 2. The shaded part represents the Tropical Zones as defined in Nos. S5.16 to S5.20 and S5.21.

G27—In the bands 225–328.6, 335.4–399.9, and 1350–1400 MHz, the fixed and mobile services are limited to the military services.

G30—In the bands 138–144, 148–149.9, 150.05–150.8, 1427–1429, and 1429–1435 MHz, the fixed and mobile services are limited primarily to operations by the military services.

G31—In the bands 3300–3500 MHz, the Government radiolocation is limited to the military services, except as provided by footnote US108.

G32—Except for weather radars on meteorological-satellites in the band 9975–10025 MHz and for Government survey operations (see footnote US-108), Government radiolocation in the band 10000–10500 MHz is limited to the military services.

G34—In the band 34.4–34.5 GHz, weather radars on board meteorological satellites for cloud detection are authorized to operate on the basis of equality with military radiolocation devices. All other nonmilitary radiolocation in the band 33.4–36.0 GHz shall be secondary to the military services.

G42—Space command, control, range, and range rate systems for earth station transmission only (including installations on certain Navy ships) may be accommodated on a co-equal basis with the fixed and mobile services in the band 1761–1842 MHz. Specific frequencies required to be used at any location will be satisfied on a coordinated case-by-case basis.

G56—Government radiolocation in the bands 1215–1300, 2900–3100, 5350–5650, and 9300–9500 MHz is primarily for the military services; however, limited secondary use is permitted by other Government agencies in support of

experimentation and research programs. In addition, limited secondary use is permitted for survey operations in the band 2900–3100 MHz.

G59—In the bands 902–928 MHz, 3100–3300 MHz, 3500–3700 MHz, 5250–5350 MHz, 8500–9000 MHz, 9200–9300 MHz, 13.4–14.0 GHz, 15.7–17.7 GHz, and 24.05–24.25 GHz, all Government nonmilitary radiolocation shall be secondary to military radiolocation, except in the subband 15.7–16.2 GHz airport surface detection equipment (ASDE) is permitted on a co-equal basis subject to coordination with the military departments.

G100—The bands 235–322 MHz and 335.4–399.9 MHz are also allocated on a primary basis to the mobile-satellite service, limited to military operations.

G104—In the bands 7450–7550 and 8175–8215 MHz, it is agreed that although the military space radio communication systems, which include earth stations near the proposed meteorological-satellite, installations will precede the meteorological-satellite installations; engineering adjustments to either the military or the meteorological-satellite systems or both will be made as mutually required to assure compatible operations of the systems concerned.

G106—The bands 2501–2502 kHz, 5003–5005 kHz, 10003–10005 kHz, 15005–15010 kHz, 19990–19995 kHz, 20005–20010 kHz, and 25005–25010 kHz are also allocated, on a secondary basis, to the space research service. The space research transmissions are subject to immediate temporary or permanent shutdown in the event of interference to the reception of the standard frequency and time broadcasts.

G109—All assignments in the band 157.0375–157.1875 MHz are subject to adjustment to other frequencies in this band as long-term US maritime VHF planning develops, particularly that planning incident to support the National VHF-FM Radiotelephone Safety and Distress System (see Doc. 15624/1–1.9.111/1.9.125).

G110—Government ground-based stations in the aeronautical radionavigation service may be authorized between 3500 and 3700 MHz where accommodation in the 2700–2900 MHz band is not technically and/or economically feasible.

G114—In the band 1350–1400 MHz, the frequency 1381.05 MHz with emissions limited to ± 12 MHz is also allocated to Fixed and Mobile-Satellite Services (space-to-Earth) for the relay of nuclear burst data.

G115—In the band 13360–13410 kHz, the fixed service is allocated on a primary basis outside the conterminous United States. Within the conterminous United States, assignments in the fixed service are permitted, and will be protected for national defense purposes or if they are to be used only in an emergency jeopardizing life, public safety, or important property under conditions calling for immediate communication where other means of communication do not exist.

G116—The band 7125–7155 MHz is also allocated for Earth-to-space transmission in the Space Operations Service at a limited number of sites (not to exceed two), subject to established coordination procedures.

G117—In the bands 7250–7750 and 7900–8400 MHz, and 17.8–21.2, 30–31, 39.5–40.5, 43.5–45.5, and 50.4–51.4 GHz the Government fixed-satellite and mobile-satellite services are limited to military systems.

G118—Government fixed stations may be authorized in the band 1700–1710 MHz only if spectrum is not available in the band 1710–1850 MHz.

G120—Development of airborne primary radars in the band 2310–2390 MHz with peak transmitter power in excess of 250 watts for use in the United States is not permitted.

G121—In the band 285–325 kHz, the Maritime Radionavigation Service may also be used on a primary basis by Maritime Radionavigation Land Stations, and on a secondary basis by Radionavigation Land Stations, to transmit differential global positioning satellite (DGPS) information.

G122—The bands 2390–2400, 2402–2417, and 4660–4685 MHz were identified for immediate reallocation, effective 10 August 1994, for exclusive non-Government use under Title VI of the Omnibus Budget Reconciliation Act of 1993. Effective 10 August 1994, any Government operations in these bands are on a noninterference basis to authorized non-Government operations and shall not hinder the implementation of any non-Government operations.

G123—The bands 2300–2310 and 2400–2402 MHz were identified for reallocation, effective 10 August 1995, for exclusive non-Government use under Title VI of the Omnibus Budget Reconciliation Act of 1993. Effective 10 August 1995, any Government operations in these bands are on a noninterference basis to authorized non-Government operations and shall not hinder the implementation of any non-Government operations.

G124—The band 2417–2450 MHz was identified for reallocation, effective August 10, 1995, for mixed Government

and non-Government use under Title VI of the Omnibus Budget Reconciliation Act of 1993.

G125—The 4635–4660 MHz band was identified for reallocation, effective January 1, 1997, for exclusive non-Government use under Title VI of the Omnibus Budget Reconciliation Act of 1993. Effective January 1, 1997, any Government operations in this band are on a noninterference basis to authorized non-Government operations and shall not hinder the implementation of any non-Government operations. However, Government operation of mobile (including airborne) systems authorized as of March 22, 1995, within 80 km of Pico Del Este, PR (18°16' N, 65°46' W), Dam Neck, VA (36°46' N, 75°57' W), and St. Thomas, VI (18°21' N, 64°55' W) will be permitted on a fully protected basis until January 1, 2009.

G126—Differential-Global-Positioning-System (DGPS) Stations may be authorized on a primary basis in the bands 108–117.975 MHz, 1559–1610 MHz, and 5000–5150 MHz for the specific purpose of transmitting DGPS information intended for aircraft navigation.

G127—Federal Travelers Information Stations (TIS) on 1610 kHz have co-primary status with AM Broadcast assignments. Federal TIS authorized as of August 4, 1994, preclude subsequent assignment for conflicting allotments.

US Footnotes

(These footnotes, each consisting of the letters US followed by one or more digits, denote stipulations applicable to both Government and non-Government stations.)

US7—In the band 420–450 MHz and within the following areas, the peak envelope power output of a transmitter employed in the amateur service shall not exceed 50 watts, unless expressly authorized by the Commission after mutual agreement, on a case-by-case basis, between the Federal Communications Commission Engineer in Charge at the applicable district office and the military area frequency coordinator at the applicable military base. For areas (e) through (j), the appropriate military coordinator is located at Peterson AFB, CO.

- (a) Those portions of Texas and New Mexico bounded on the south by latitude 31°45' N, on the east by 104°00' W, on the north by latitude 34°30' N, and on the west by longitude 107°30' W.
- (b) The entire State of Florida including the Key West area and the areas enclosed within a 322-kilometer (200-mile) radius of Patrick Air Force Base, Florida (latitude 28°21' N, longitude 80°43' W), and within a 322-kilometer (200-mile) radius of Eglin Air Force Base, Florida (latitude 30°30' N, longitude 86°30' W).
- (c) The entire State of Arizona.
- (d) Those portions of California and Nevada south of latitude 37°10' N, and the areas enclosed within a 322-kilometer (200-mile) radius of the Pacific Missile Test Center, Point Mugu, California (latitude 34°09' N, longitude 119°11' W).
- (e) In the State of Massachusetts within a 160-kilometer (100-mile) radius around locations at Otis Air Force Base, Massachusetts (latitude 41°45' N, longitude 70°32' W).

- (f) In the State of California within a 240-kilometer (150-mile) radius around locations at Beale Air Force Base, California (latitude 39°08' N, longitude 121°26' W).
- (g) In the State of Alaska within a 160-kilometer (100-mile) radius of Clear, Alaska (latitude 64°17' N, longitude 149°10' W).
- (h) In the State of North Dakota within a 160-kilometer (100-mile) radius of Concrete, North Dakota (latitude 48°43' N, longitude 97°54' W).
- (i) In the States of Alabama, Florida, Georgia, and South Carolina within a 200-kilometer (124-mile) radius of Warner Robins Air Force Base, Georgia (latitude 32°38' N, longitude 83°35' W).
- (j) In the State of Texas within a 200-kilometer (124-mile) radius of Goodfellow Air Force Base, Texas (latitude 31°25' N, longitude 100°24' W).

US8—The use of frequencies 170.475, 171.425, 171.575, and 172.275 MHz east of the Mississippi River and 170.425, 170.575, 171.475, 172.225, and 172.375 MHz west of the Mississippi River may be authorized to fixed, land, and mobile stations operated by non-Federal forest firefighting agencies. In addition, land stations and mobile stations operated by non-Federal conservation agencies, for mobile relay operation only, may be authorized to use the frequency 172.275 MHz east of the Mississippi River and the frequency 171.475 MHz west of the Mississippi River. The use of any of the foregoing nine frequencies shall be on the condition that no harmful interference will be caused to Government stations.

US10—The use of the frequencies 26.62, 143.75, 143.90, and 148.15 MHz may be authorized to Civil Air Patrol land stations and Civil Air Patrol mobile stations.

US11—The use of the frequencies 166.250 and 170.150 MHz may be authorized to non-Government remote pickup broadcast base and land mobile stations and to non-Government base, fixed, and land mobile stations in the public safety radio services (the sum of the bandwidth of emission and tolerance is not to exceed 25 kHz, except that authorizations in existence as of December 20, 1974, using a larger bandwidth are permitted to continue in operation until December 20, 1979) in the continental United States (excluding Alaska) only, except within the area bounded on the west by the Mississippi River, on the north by the parallel of latitude 37°30' N, and on the east and south by that arc of the circle with centre at Springfield, Illinois, and radius equal to the airline distance between Springfield, Illinois, and Montgomery, Alabama, subtended between the foregoing west and north boundaries, on the conditions that harmful interference will not be caused to Government stations present or future in the Government band 162–174 MHz. The use of these frequencies by remote pickup broadcast stations will not be authorized for locations within 240 kilometers (150 miles) of New York City; and use of these frequencies by the public safety radio services will not be authorized except for locations within 240 kilometers (150 miles) of New York City.

US13—For the specific purpose of transmitting hydrological and meteorological data in cooperation with agencies of the Federal Government, the following frequencies may be

authorized to non-Government fixed stations on the condition that harmful interference will not be caused to Government stations:

MHz	MHz	MHz	MHz
169.425	170.275	171.125	406.175
169.450	170.300	171.825	409.675
169.475	170.325	171.850	409.725
169.500	171.025	171.875	412.625
169.525	171.050	171.900	412.675
170.225	171.075	171.925	412.725
170.250	171.100	406.125	412.775

Licensees holding a valid authorization on June 11, 1962, to operate on the frequencies 169.575, 170.375, or 171.975 MHz may continue to be authorized for such operations on the condition that harmful interference will not be caused to Government stations.

US14—When 500 kHz is being used for distress purposes, ship and coast stations using morse telegraphy may use 512 kHz for calling.

US18—Navigation aids in the United States and possessions in the bands 9–14 kHz, 90–110 kHz, 190–415 kHz, 510–535 kHz, and 2700–2900 MHz are normally operated by the US Government. However, authorizations may be made by the FCC for non-Government operation in these bands subject to the conclusion of appropriate arrangements between the FCC and the Government agencies concerned and upon special showing of need for service which the Government is not yet prepared to render.

US25—The use of frequencies in the band 25.85–26.1 MHz may be authorized in any area to non-Government remote pickup broadcast base and mobile stations on the condition that harmful interference is not caused to stations in the broadcasting service.

US26—The bands 117.975–121.4125 MHz, 123.5875–128.8125 MHz, and 132.0125–136 MHz are for air traffic control communications.

US28—The band 121.5875–121.9375 MHz is for use by aeronautical utility land and mobile stations, and for air traffic control communications.

US30—The band 121.9375–123.0875 MHz is available to FAA aircraft for communications pursuant to flight inspection functions in accordance with the Federal Aviation Act of 1958.

US31—Except as provided below the band 121.9375–123.0875 MHz is for use by private aircraft stations.

The frequencies 122.700, 122.725, 122.750, 122.800, 122.950, 122.975, 123.000, 123.050, and 123.075 MHz may be assigned to aeronautical advisory stations. In addition, at landing areas having a part-time or no airdrome control tower or FAA flight service station, these frequencies may be assigned on a secondary noninterference basis to aeronautical utility mobile stations and may be used by FAA ground vehicles for safety related communications during inspections conducted at such landing areas.

The frequencies 122.850, 122.900, and 122.925 MHz may be assigned to aeronautical multicom stations. In addition, 122.850 MHz may be assigned on a secondary noninterference basis to aeronautical utility mobile stations. In case of 122.925 MHz, US213 applies.

Air carrier aircraft stations may use 122,000 and 122,050 MHz for communication with aeronautical stations of the Federal Aviation Administration and 122,700, 122,800, 122,900, and 123,000 MHz for communications with aeronautical stations pertaining to safety of flight with and in the vicinity of landing areas not served by a control tower.

Frequencies in the band 121,9375–122,6875 MHz may be used by aeronautical stations of the Federal Aviation Administration for communication with private aircraft stations only, except that 122,000 and 122,050 MHz may also be used for communication with air carrier aircraft stations concerning weather information.

US32—Except for the frequencies 123.3 and 123.5 MHz, which are not authorized for Government use, the band 123,1125–123,5875 MHz is available for FAA Communications incident to flight test and inspection activities pertinent to aircraft and facility certification on a secondary noninterference basis.

US33—The band 123,1125–123,5875 MHz is for use by flight test and aviation instructional stations. The frequency 121,950 MHz is available for aviation instructional stations.

US41—The Government radiolocation service is permitted in the band 2450–2500 MHz on the condition that harmful interference is not caused to non-Government services.

US44—The non-Government radiolocation service may be authorized in the band 2900–3100 MHz on the condition that no harmful interference is caused to Government services.

US48—The non-Government radiolocation service may be authorized in the bands 5350–5460 MHz and 9000–9200 MHz on the condition that it does not cause harmful interference to the aeronautical radionavigation service or to the Government radiolocation service.

US49—The non-Government radiolocation service may be authorized in the band 5460–5470 MHz on the condition that it does not cause harmful interference to the aeronautical or maritime radionavigation services or to the Government radiolocation service.

US50—The non-Government radiolocation service may be authorized in the band 5470–5600 MHz on the condition that it does not cause harmful interference to the maritime radionavigation service or to the Government radiolocation service.

US51—In the bands 5600–5650 MHz and 9300–9500 MHz, the non-Government radiolocation service shall not cause harmful interference to the Government radiolocation service.

US53—In view of the fact that the band 13.25–13.4 GHz is allocated to doppler navigation aids, Government and non-Government airborne doppler radars in the aeronautical radionavigation service are permitted in the band 8750–8850 MHz only on the condition that they must accept any interference that may be experienced from stations in the radiolocation service in the band 8500–10000 MHz.

US54—Temporarily, and until certain operations of the radiolocation service in the band 9000–9200 MHz can be transferred to other appropriate frequency bands, the aeronautical radionavigation service may, in certain geographical areas, be subject to receiving some degree of interference from the radiolocation service.

US58—In the band 10000–10500 MHz, pulsed emissions are prohibited, except for weather radars on board meteorological-satellites in the band 10000–10025 MHz. The amateur service and the non-Government radiolocation service, which shall not cause harmful interference to the Government radiolocation service, are the only non-Government services permitted in this band. The non-Government radiolocation service is limited to survey operations as specified in footnote US108.

US59—The band 10.5–10.55 GHz is restricted to systems using type N0N emission with a power not to exceed 40 watts into the antenna.

US65—The use of the band 5460–5650 MHz by the maritime radionavigation service is limited to shipborne radars.

US66—The use of the band 9300–9500 MHz by the aeronautical radionavigation service is limited to airborne radars and associated airborne beacons. In addition, ground-based radar beacons in the aeronautical radionavigation service are permitted in the band 9300–9320 MHz on condition that harmful interference is not caused to the maritime radionavigation service.

US67—The use of the band 9300–9500 MHz by the meteorological aids service is limited to ground-based radars. Radiolocation installations will be coordinated with the meteorological aids service and, insofar as practicable, will be adjusted to meet the requirements of the meteorological aids service.

US69—In the band 31.8–33.4 GHz, ground-based radionavigation aids are not permitted except where they operate in cooperation with airborne or shipborne radionavigation devices.

US70—The meteorological aids service allocation in the band 400.15–406 MHz does not preclude the operation therein of associated ground transmitters.

US71—In the band 9300–9320 MHz, low-powered maritime radionavigation stations shall be protected from harmful interference caused by the operation of land-based equipment.

US74—In the bands 25.55–25.67, 73–74.6, 406.1–410, 608–614, 1400–1427, 1660.5–1670, 2690–2700, and 4990–5000 MHz and in the bands 10.68–10.7, 15.35–15.4, 23.6–24, 31.3–31.8, 86–92, 105–116, and 217–231 GHz, the radio astronomy service shall be protected from extraband radiation only to the extent such radiation exceeds the level that would be present if the offending station were operating in compliance with the technical standards or criteria applicable to the service in which it operates.

US77—Government stations may also be authorized for:

- (a) port operations use on a simplex basis by coast and ship stations on the frequencies 156.6 and 156.7 MHz;
- (b) duplex port operations use of the frequency 157.0 MHz for ship stations and 161.6 MHz for coast stations;
- (c) intership use of 156.3 MHz on a simplex basis; and
- (d) vessel traffic services under the control of the US Coast Guard on a simplex basis by coast and ship stations on the frequencies 156.25, 156.55, 156.6, and 156.7 MHz.

- (e) navigational bridge-to-bridge and navigational communications on a simplex basis by coast and ship stations on the frequencies 156.375 and 156.65 MHz.

US78—In the mobile service, the frequencies between 1435 and 1535 MHz will be assigned for aeronautical telemetry and associated telecommand operations for flight testing of manned or unmanned aircraft and missiles or for their major components. Permissible usage includes telemetry associated with launching and reentry into the earth's atmosphere as well as any incidental orbiting prior to reentry of manned objects undergoing flight tests. The following frequencies are shared with flight telemetering mobile stations: 1444.5, 1453.5, 1501.5, 1515.5, 1524.5, and 1525.5 MHz.

US80—Government stations may use the frequency 122.9 MHz subject to the following conditions:

- (a) All operations by Government stations shall be restricted to the purpose for which the frequency is authorized to non-Government stations and shall be in accordance with the appropriate provisions of the Commission's rules and regulations, Part 87, Aviation Services.
- (b) Use of the frequency is required for coordination of activities with Commission licensees operating on this frequency.
- (c) Government stations will not be authorized for operations at fixed locations.

US81—The band 38–38.25 MHz is used by both Government and non-Government radio astronomy observatories. No new fixed or mobile assignments are to be made and Government stations in the band 38–38.25 MHz will be moved to other bands on a case-by-case basis, as required, to protect radio astronomy observations from harmful interference. As an exception however, low-powered military transportable and mobile stations used for tactical and training purposes will continue to use the band. To the extent practicable, the latter operations will be adjusted to relieve such interference as may be caused to radio astronomy observations. In the event of harmful interference from such local operations, radio astronomy observatories may contact local military commands directly, with a view to effecting relief. A list of military commands, areas of coordination, and points of contact for purposes of relieving interference may be obtained upon request from the Office of the Chief Scientist, Federal Communications Commission, Washington, DC 20554.

US82—Until July 1, 1991, the assignable frequencies in the bands 4143.6–4146.6 kHz, 6218.6–6224.6 kHz, 8291.1–8297.3 kHz, 12429.2–12439.5 kHz, 16587.1–16596.4 kHz, and 22124–22139.5 kHz may be authorized on a shared non-priority basis to Government and non-Government ship and coast stations (SSB telephony, with peak envelope power not to exceed 1 kW). Effective July 1, 1991, the assignable frequencies in the bands 4146–4152 kHz, 6224–6233 kHz, 8294–8300 kHz, 12353–12368 kHz, 16528–16549 kHz, 18825–18846 kHz, 22159–22180 kHz, and 25100–25121 kHz may be authorized on a shared nonpriority basis to Government and non-Government ship and coast stations (SSB telephony, with peak envelope power not to exceed 1 kW).

US87—The frequency 450 MHz, with maximum emission bandwidth of 500 kHz, may be used by Government and non-Government stations for space telecommand at specific locations, subject to such conditions as may be applied on a case-by-case basis.

US90—In the band 2025–2110 MHz Earth-to-space and space-to-space transmissions may be authorized in the space research and earth exploration-satellite services subject to such conditions as may be applied on a case-by-case basis. Such transmissions shall not cause harmful interference to non-Government stations operating in accordance with the Table of Frequency Allocations. All space-to-space transmission reaching the earth's surface shall adhere to a power flux density of between –144 and –154 dBW/m²/4 kHz depending on the angle of arrival per ITU Radio Regulation 2557 and shall not cause harmful interference to the other space services.

US93—In the conterminous United States, the frequency 108.0 MHz may be authorized for use by VOR test facilities, the operation of which is not essential for the safety of life or property, subject to the condition that no interference is caused to the reception of FM broadcasting stations operating in the band 88–108 MHz. In the event that such interference does occur, the licensee or other agency authorized to operate the facility shall discontinue operation on 108 MHz and shall not resume operation until the interference has been eliminated or the complaint otherwise satisfied. VOR test facilities operating on 108 MHz will not be protected against interference caused by FM broadcasting stations operating in the band 88–108 MHz nor shall the authorization of a VOR test facility on 108 MHz preclude the Commission from authorizing additional FM broadcasting stations.

US99—In the band 1668.4–1670 MHz, the meteorological aids service (radiosonde) will avoid operations to the maximum extent practicable. Whenever it is necessary to operate radiosondes in the band 1668.4–1670 MHz within the United States, notification of the operations shall be sent as far in advance as possible to the Electromagnetic Spectrum Management Unit, National Science Foundation, Washington, DC 20550.

US102—In Alaska only, the frequency 122.1 MHz may also be used for air carrier air traffic control purposes at locations where other frequencies are not available to air carrier aircraft stations for air traffic control.

US104—The LORAN Radionavigation System has priority in the band 90–110 kHz in the United States and Possessions. Radiolocation land stations making use of LORAN type equipment may be authorized to both Government and non-Government on a Secondary Service basis for offshore radiolocation activities only at specific locations and subject to such technical and operational conditions (e.g., power, emission, pulse rate and phase code, hours of operation), including on-the-air testing, as may be required on a case-by-case basis to ensure protection of the LORAN Radionavigation System from harmful interference and to ensure mutual compatibility among radiolocation operators. Such authorizations to stations in the radiolocation service are further subject to showing of need for service which is not currently provided and which the Government is not yet prepared to render by way of the radionavigation service.

US106—The frequency 156.75 MHz is available for assignment to non-Government and Government stations for environmental communications in accordance with an agreed plan.

US107—The frequency 156.8 MHz is the international distress, safety and calling frequency for the maritime mobile VHF radiotelephone service for use by Government and non-Government ship and coast stations. Guard bands of 156.7625–156.7875 and 156.8125–156.8375 MHz are maintained.

US108—Within the bands 3300–3500 MHz and 10000–10500 MHz, survey operations, using transmitters with a peak power not to exceed 5 watts into the antenna, may be authorized for Government and non-Government use on a secondary basis to other Government radiolocation operations.

US110—In the frequency bands 3100–3300 MHz, 3500–3700 MHz, 5250–5350 MHz, 8500–9000 MHz, 9200–9300 MHz, 9500–10000 MHz, 13.4–14.0 GHz, 15.7–17.3 GHz, 24.05–24.25 GHz, and 33.4–36 GHz, the non-Government radiolocation service shall be secondary to the Government radiolocation service and to airborne doppler radars at 8800 MHz and shall provide protection to airport surface detection equipment (ASDE) operating between 15.7 and 16.2 GHz.

US111—In the band 1990–2120 MHz, Government space research earth stations may be authorized to use specific frequencies at specific locations for earth-to-space transmissions. Such authorizations shall be secondary to non-Government use of this band and subject to such other conditions as may be applied on a case-by-case basis.

- Corpus Christi, TX, 27°39' N, 097°23' W
- Fairbanks, AK, 64°59' N, 147°53' W
- Goldstone, CA, 35°18' N, 116°54' W
- Greenbelt, MD, 39°00' N, 076°50' W
- Guam, Mariana Is., 13°19' N, 144°44' E
- Kauai, HI, 22°08' N, 159°40' W
- Merritt Is., FL, 28°29' N, 080°35' W
- Roseman, NC, 35°12' N, 082°52' W
- Wallops Is., VA, 37°57' N, 075°28' W

US112—The frequency 123.1 MHz is for search and rescue communications. This frequency may be assigned for air traffic control communications at special aeronautical events on the condition that no harmful interference is caused to search and rescue communications during any period of search and rescue operations in the locale involved.

US116—In the bands 890–902 MHz, 928–932 MHz, and 935–941 MHz, no new assignments are to be made to Government radio stations after July 10, 1970, except, on a case-by-case basis, to experimental stations and to additional stations of existing networks in Alaska. Government assignments existing prior to July 10, 1970, to stations in Alaska may be continued. All other existing Government assignments shall be on a secondary basis to stations in the non-Government land mobile service and shall be subject to adjustment or removal from the bands 890–902 MHz, 928–932 MHz, and 935–941 MHz, at the request of the FCC.

US117—In the band 406.1–410 MHz, all new authorizations will be limited to a maximum 7 watts per kHz of necessary bandwidth; existing authorizations as of November 30, 1970, exceeding this power are permitted to continue in use.

New authorizations in this band for stations, other than mobile and transportable stations, within the following areas are subject to prior coordination by the applicant through the Electromagnetic Spectrum Management Unit, National Science Foundation, Washington, DC 20550 (202-357-9696):

Arecibo Observatory: Rectangle between latitudes 17°30' N and 19°00' N and between longitudes 65°10' W and 68°00' W.

Owens Valley Radio Observatory: Two contiguous rectangles, one between latitudes 36° N and 37° N and longitudes 117°40' W and 118°30' W and the second between latitudes 37° N and 38° N and longitudes 118° W and 118°50' W.

Sagamore Hill Radio Observatory: Rectangle between latitudes 42°10' N and 43°00' N and longitudes 70°31' W and 71°31' W.

Table Mountain Solar Observatory (NOAA) Boulder, Colorado (407–409 MHz only): Rectangle between latitudes 39°30' N and 40° 30' N and longitudes 104°30' W and 106°00' W or the Continental Divide whichever is farther east.

The non-Government use of this band is limited to the radio astronomy service and as provided by footnote US13.

US201—In the band 460–470 MHz, space stations in the earth exploration-satellite service may be authorized for space-to-Earth transmission on a secondary basis with respect to the fixed and mobile services. When operating in the meteorological-satellite service such stations shall be protected from harmful interference from other applications of the earth exploration-satellite service. The power flux density produced at the earth's surface by any space station in this band shall not exceed – 152 dBW/m²/4 kHz.

US203—Radio astronomy observations of the formaldehyde line frequencies 4825–4835 MHz and 14.470–14.500 GHz may be made at certain radio astronomy observatories as indicated below:

Bands to be Observed		Observatory
4 GHz	14 GHz	
X		National Astronomy and Ionosphere Center, Arecibo, Puerto Rico
X	X	National Radio Astronomy Observatory, Green Bank, West Virginia
X	X	National Radio Astronomy Observatory, Socorro, New Mexico
X	X	Hat Creek Observatory (U. of CA), Hat Creek, California
X	X	Haystack Radio Observatory (MIT-Lincoln Lab) Tyngsboro, Massachusetts
X	X	Owens Valley Radio Observatory (Cal. Tech.), Big Pine, California
	X	Five College Radio Astronomy Observatory, Quabbin Reservoir (near Amherst) Massachusetts

Every practicable effort will be made to avoid the assignment of frequencies to stations in the fixed or mobile services in these bands. Should such assignments result in harmful interference to these observations, the situation will be remedied to the extent practicable.

US205—Tropospheric scatter systems are prohibited in the band 2500–2690 MHz.

US208—Planning and use of the band 1559–1626.5 MHz necessitate the development of technical and/or operational sharing criteria to ensure the maximum degree of electromagnetic compatibility with existing and planned systems within the band.

US209—The use of frequencies 460.6625, 460.6875, 460.7125, 460.7375, 460.7625, 460.7875, 460.8125, 460.8375, 460.8625, 465.6625, 465.6875, 465.7125, 465.7375, 465.7625, 465.7875, 465.8125, 465.8375, and 465.8625 MHz may be authorized with 100 mW or less output power, to Government and non-Government radio stations for one-way, nonvoice biomedical telemetry operations in hospitals or in medical or convalescent centres.

US210—Use of frequencies in the bands 40.66–40.70 and 216–220 MHz may be authorized to Government and non-Government stations on a secondary basis for the tracking of, and telemetering of scientific data from, ocean buoys and wildlife. Airborne wildlife telemetry in the 216–220 MHz band will be limited to the 216.000–216.100 MHz portion of the band. Operation in these two bands is subject to the technical standards specified in (a) Section 8.2.42 of the NTIA Manual for Government use or in (b) Section 5.108 of the Commission's Rules for non-Government.

US211—In the bands 1670–1690 and 5000–5250 MHz and 10.7–11.7, 15.1365–15.35, 15.4–15.7, 22.5–22.55, 24–24.05, 31.0–31.3, 31.8–32, 40.5–42.5, 84–86, 102–105, 116–126, 151–164, 176.5–182, 185–190, 231–235, 252–265 GHz, applicants for airborne or space station assignments are urged to take all practicable steps to protect radio astronomy observations in the adjacent bands from harmful interference; however, US74 applies.

US212—In the State of Alaska, the carrier frequency 5167.5 kHz (assigned frequency 5168.9 kHz) is designated for emergency communications. This frequency may also be used in the Alaska-Private Fixed Service for calling and listening, but only for establishing communications before switching to another frequency. The maximum power is limited to 150 watts peak envelope power (PEP).

US213—The frequency 122.925 MHz is for use only for communications with or between aircraft when coordinating natural resources programs of Federal or State natural resources agencies, including forestry management and fire suppression, fish and game management and protection, and environmental monitoring and protection.

US214—The frequency 157.100 MHz is the primary frequency for liaison communications between ship stations and stations of the United States Coast Guard.

US215—Emissions from microwave ovens manufactured on and after January 1, 1980, for operation on the frequency 915 MHz, must be confined within the band 902–928 MHz. Emissions from microwave ovens manufactured prior to January 1, 1980, for operation on the frequency 915 MHz, must be confined within the band 902–940 MHz. Radiocommuni-

cations services operating within the band 928–940 MHz must accept any harmful interference that may be experienced from the operation of microwave ovens manufactured before January 1, 1980.

US216—The frequencies 150.775 and 150.790 and the bands 152–152.0150, 163.2375–163.2625, 462.9375–463.1875, and 467.9375–468.1875 MHz are authorized for Government and non-Government operations in medical radio communications systems.

US217—Pulse-ranging radiolocation systems may be authorized for Government and non-Government use in the 420–450 MHz band along the shorelines of Alaska and the contiguous 48 states. Spread-spectrum radiolocation systems may be authorized in the 420–435 MHz portion of the band for operation within the contiguous 48 states and Alaska. Authorizations will be granted on a case-by-case basis; however, operations proposed to be located within the zones set forth in US228 should not expect to be accommodated. All stations operating in accordance with this provision will be secondary to stations operating in accordance with the Table of Frequency Allocations.

US218—The band 902–928 MHz is available for Location and Monitoring Service (LMS) systems subject to their not causing harmful interference to the operation of all Government stations authorized in these bands. These systems must tolerate interference from the operation of industrial, scientific, and medical (ISM) devices and the operation of Government stations authorized in these bands.

US219—In the band 2025–2110 MHz Government Earth Resources Satellite Earth Stations in the Earth Exploration-Satellite Service may be authorized to use the frequency 2106.4 MHz for Earth-to-space transmissions for tracking, telemetry, and telecommand at the sites listed below. Such transmissions shall not cause harmful interference to non-Government operations:

Sioux Falls, SD, 43°32'03.1" N, 96°45'42.8" W
Fairbanks, AK, 64°58'36.6" N, 147°30'54.2" W

US220—The frequencies 36.25 and 41.71 MHz may be authorized to Government stations and non-Government stations in the Petroleum Radio Service, for oil spill containment and cleanup operations. The use of these frequencies for oil spill containment or cleanup operations is limited to the inland and coastal waterways regions.

US221—Use of the mobile service in the bands 525–535 kHz and 1605–1615 kHz is limited to distribution of public service information from Travelers Information stations operating on 530 kHz or 1610 kHz.

US222—In the band 2025–2035 MHz Geostationary Operational Environmental Satellite earth stations in the Space Research and Earth Exploration-Satellite Services may be authorized on a co-equal basis to use the frequency band 2025–2035 MHz for Earth-to-space transmissions for tracking, telemetry, and telecommand at the sites listed below:

Wallops Is., VA, 37°50'48" N, 75°27'33" W
Seattle, WA, 47°34'15" N, 122°33'10" W
Honolulu, HI, 21°21'12" N, 157°52'36" W

US223—Within 120 kilometers (75 miles) of the United States/Canada border on the Great Lakes, the Saint Lawrence Seaway, and the Puget Sound and the Strait of Juan de Fuca and its approaches, use of coast transmit frequency 162.025 MHz and ship station transmit frequency 157.425 MHz (VHF maritime mobile service channel 88) may be authorized by the maritime mobile service for public correspondence.

US224—Government systems utilizing spread-spectrum techniques for terrestrial communication, navigation, and identification may be authorized to operate in the band 960–1215 MHz on the condition that harmful interference will not be caused to the aeronautical radionavigation service. These systems will be handled on a case-by-case basis. Such systems shall be subject to a review at the national level for operational requirements and electromagnetic compatibility prior to development, procurement, or modification.

US225—In addition to its present Government use, the frequency band 510–525 kHz is available to Government and non-Government aeronautical radionavigation stations inland of the Territorial Base Line as coordinated with the military services. In addition, the frequency 510 kHz is available for non-Government ship-helicopter operations when beyond 185 kilometers (100 nautical miles) from shore and required for aeronautical radionavigation.

US226—In the State of Hawaii, stations in the aeronautical radionavigation service shall not cause harmful interference to US Navy reception from its station at Honolulu on 198 kHz.

US228—Applicants of operation in the band 420 to 450 MHz under the provisions of US217 should not expect to be accommodated if their area of service is within the following geographic areas:

- (a) Those portions of Texas and New Mexico bounded on the south by latitude 31°45' N, on the east by longitude 104°00' W, on the north by latitude 34°30' N, and on the west by longitude 107°30' W.
- (b) The entire State of Florida including the Key West area and the areas enclosed within a 322-kilometer (200-mile) radius of Patrick Air Force Base, Florida (latitude 28°21' N, longitude 80°43' W), and within a 322-kilometer (200-mile) radius of Eglin Air Force Base, Florida (latitude 30°30' N, longitude 86°30' W).
- (c) The entire State of Arizona.
- (d) Those portions of California and Nevada south of latitude 37°10' N, and the areas enclosed within a 322-kilometer (200-mile) radius of the Pacific Missile Test Center, Point Mugu, California (latitude 34°09' N, longitude 119°11' W).
- (e) In the State of Massachusetts within a 160-kilometer (100-mile) radius around locations at Otis Air Force Base, Massachusetts (latitude 41°45' N, longitude 70°32' W).
- (f) In the State of California within a 240-kilometer (150-mile) radius around locations at Beale Air Force Base, California (latitude 39°08' N, longitude 121°26' W).
- (g) In the State of Alaska within a 160-kilometer (100-mile) radius of Clear, Alaska (latitude 64°17' N, longitude 149°10' W).
- (h) In the State of North Dakota within a 160-kilometer (100-mile) radius of Concrete, North Dakota (latitude 48°43' N, longitude 97°54' W).
- (i) In the States of Alabama, Florida, Georgia, and South Carolina within a 200-kilometer (124-mile) radius of Warner Robins Air Force Base, Georgia (latitude 32°38' N, longitude 83°35' W).
- (j) In the State of Texas within a 200-kilometer (124-mile) radius of Goodfellow Air Force Base, Texas (latitude 31°25' N, longitude 100°24' W).

US229—Assignments to stations in the fixed and mobile services may be made on the condition that no harmful interference is caused to the Navy SPASUR system currently operating in the southern United States in the frequency band 216.88–217.08 MHz.

US230—Non-Government land mobile service is allocated on a primary basis in the bands 422.1875–425.4875 and 427.1875–429.9875 MHz within 80 kilometers (50 statute miles) of Detroit, MI, and Cleveland, OH, and in the bands 423.8125–425.4875 and 428.8125–429.9875 MHz within 80 kilometers (50 statute miles) of Buffalo, NY.

US231—When an assignment cannot be obtained in the bands between 200 and 525 kHz, which are allocated to aeronautical radionavigation, assignments may be made to aeronautical radiobeacons in the maritime mobile band 435–490 kHz, on a secondary basis, subject to the coordination and agreement of those agencies having assignments within the maritime mobile band that may be affected. Assignments to aeronautical radionavigation radiobeacons in the band 435–490 kHz shall not be a bar to any required changes to the Maritime Mobile Radio Service and shall be limited to Government not employing voice emissions.

US235—Until implementation procedures and schedules are determined by future conferences of the International Telecommunication Union, the bands 9775–9900, 11650–11700, 11975–12050, 13600–13800, 15450–15600, 17550–17700, and 21750–21850 kHz to be implemented by the broadcasting service are allocated as an alternative allocation to the fixed service. The bands 12230–12330, 16360–16460, 17360–17410, 18780–18900, 19680–19800, 22720–22855, 25110–25210, and 26100–26175 kHz to be implemented by the maritime mobile service are also allocated as an alternative allocation to the fixed service until July 1, 1991, when these bands are to be allocated exclusively to the maritime mobile service.

US236—Until implementation procedures and schedules are determined by future conferences of the International Telecommunication Union (see Resolution 319), the bands 4000–4063 and 8100–8195 kHz are also allocated on a primary basis to the fixed service.

US237—Until implementation procedures and schedules are determined by a future Regional Conference of the International Telecommunication Union, the band 1615–1625 kHz is also allocated on a primary basis to the radiolocation service.

US238—Until implementation procedures and schedules are determined by a future Regional Conference of the International Telecommunication Union, the band 1625–1705 kHz

is allocated to the radiolocation service on a primary basis as a different category of service.

US239—Aeronautical radionavigation stations (radio-beacons) may be authorized, primarily for offshore use, in the band 525–535 kHz on a noninterference basis to Travelers Information Stations.

US240—The bands 1715–1725 kHz and 1740–1750 kHz are allocated on a primary basis and the bands 1705–1715 kHz and 1725–1740 kHz on a secondary basis to the aeronautical radionavigation service (radiobeacons).

US244—The band 136–137 MHz is allocated to the non-Government aeronautical mobile (R) service on a primary basis and is subject to pertinent international treaties and agreements. The frequencies 136.000 MHz, 136.025 MHz, 136.050 MHz, 136.075 MHz, 136.125 MHz, 136.150 MHz, 136.175 MHz, 136.225 MHz, 136.250 MHz, 136.300 MHz, 136.325 MHz, 136.350 MHz, 136.400 MHz, 136.425 MHz, and 136.450 MHz are available on a shared basis to the Federal Aviation Administration for air traffic control purposes, such as automatic weather observation services, automatic terminal information services, and airport control tower communications. Stations licensed prior to January 2, 1990, using the 136–137 MHz band for space operations (space-to-Earth), meteorological-satellite service (space-to-Earth), and the space research service (space-to-Earth) may continue to use this band on a secondary basis to aeronautical mobile (R) service stations. No new assignments will be made to stations in the above space services.

US245—The Fixed-Satellite Service is limited to international intercontinental systems and subject to case-by-case electromagnetic compatibility analysis.

US246—No stations will be authorized to transmit in the bands 608–614 MHz, 1400–1427 MHz, 1660.5–1668.4 MHz, 2690–2700 MHz, 4990–5000 MHz, 10.68–10.70 GHz, 15.35–15.40 GHz, 23.6–24.0 GHz, 31.3–31.8 GHz, 51.4–54.25 GHz, 58.2–59.0 GHz, 64–65 GHz, 86–92 GHz, 100–102 GHz, 105–116 GHz, 164–168 GHz, 182–185 GHz, and 217–231 GHz.

US247—The band 10100–10150 kHz is allocated to the fixed service on a primary basis outside the United States and Possessions. Transmissions of stations in the amateur service shall not cause harmful interference to this fixed service use and stations in the amateur service shall make all necessary adjustments (including termination of transmission) if harmful interference is caused.

US251—The band 12.75–13.25 GHz is also allocated to the Space Research Service (Deep Space) (space-to-Earth) for reception only at Goldstone, California, 35°18' N, 116°54' W.

US252—The bands 2110–2120 MHz, 7145–7190 MHz, and 34.2–34.7 GHz are also allocated for Earth-to-space transmissions in the Space Research Service, limited to deep space communications at Goldstone, California.

US253—In the band 2300–2310 MHz, the fixed and mobile services shall not cause harmful interference to the amateur service.

US254—In the band 18.6–18.8 GHz, the fixed and mobile services shall be limited to a maximum equivalent isotropically radiated power of +35 dBW and the power delivered to the antenna shall not exceed –3 dBW.

US255—In the band 18.6–18.8 GHz, the fixed-satellite service shall be limited to a power flux density at the Earth's surface of –101 dBW/m² in a 200 MHz band for all angles of arrival.

US256—Radio astronomy observations may be made in the band 1718.8–1722.2 MHz on an unprotected basis. Agencies providing other services in this band in the geographic areas listed below should bear in mind that their operations may affect those observations, and those agencies are encouraged to minimize potential interference to the observations insofar as it is practicable.

Hat Creek Observatory, Hat Creek, California	Rectangle between latitudes 40°00'N and 42°00'N and between longitudes 120°15'W and 122°15'W.
Owens Valley Radio Observatory, Big Pine, California	Two contiguous rectangles, one between 36°00'N and 37°00'N and between longitudes 117°40'W and 118°30'W and the second between latitudes 37°00'N and 30°00'N and between longitudes 118°00'W and 118°50'W.
Haystack Radio Observatory, Tyngsboro, Massachusetts	Rectangle between latitudes 41°00'N and 43°00'N and between longitudes 71°00'W and 73°00'W.
National Astronomy and Ionosphere Center, Arecibo, Puerto Rico	Rectangle between latitudes 17°30'N and 19°00'N and between longitudes 65°10'W and 68°00'W.
National Radio Astronomy Observatory, Green Bank, West Virginia	Rectangle between latitudes 37°30'N and 39°15'N and between longitudes 78°30'W and 80°30'W.

US257—Radio astronomy observations may be made in the 4950–4990 MHz band at certain Radio Astronomy Observatories indicated below:

National Astronomy and Ionosphere Center, Arecibo, Puerto Rico	Rectangle between latitudes 17°30'N and 19°00'N and between longitudes 65°10'W and 68°00'W.
Haystack Radio Observatory, Tyngsboro, Massachusetts	Rectangle between latitudes 41°00'N and 43°00'N and between longitudes 71°00'W and 73°00'W.
National Radio Astronomy Observatory, Green Bank, West Virginia	Rectangle between latitudes 37°00'N and 39°15'N and between longitudes 78°30'W and 80°30'W.
National Radio Astronomy Observatory, Socorro, New Mexico	Rectangle between latitudes 32°30'N and 35°30'N and between longitudes 106°00'W and 109°00'W.

Owens Valley Radio Observatory, Big Pinc, California	Two contiguous rectangles, one between latitudes 36°00'N and 37°00'N and between longitudes 117°40'W and 118°30'W and the second between latitudes 37°00'N and 38°00'N and between longitudes 118°00'W and 118°50'W.
Hat Creek Observatory, Hat Creek, California	Rectangle between latitudes 40°00'N and 42°00'N and between longitudes 120°15'W and 122°15'W.

Every practicable effort will be made to avoid the assignment of frequencies in the band 4950–4990 MHz to stations in the fixed and mobile services within the geographic areas given above. In addition, every practicable effort will be made to avoid the assignment of frequencies in this band to stations in the aeronautical mobile service that operate outside of those geographic areas, but which may cause harmful interference to the listed observatories. Should such assignments result in harmful interference to these observatories, the situation will be remedied to the extent practicable.

US258—In the band 8025–8400 MHz, the non-Government earth exploration-satellite service (space-to-Earth) is allocated on a primary basis. Authorizations are subject to a case-by-case electromagnetic compatibility analysis.

US259—Stations in the radiolocation service in the band 17.3–17.7 GHz shall be restricted to operating powers of less than 51 dBW effective isotropic radiated power (eirp) after feeder link stations for the broadcasting-satellite service are authorized and brought into use.

US260—Aeronautical mobile communications that are an integral part of aeronautical radionavigation systems may be satisfied in the bands 1559–1626.5 MHz, 5000–5250 MHz, and 15.4–15.7 GHz.

US261—The use of the band 4200–4400 MHz by the Aeronautical Radionavigation service is reserved exclusively for airborne radio altimeters. Experimental stations will not be authorized to develop equipment for operational use in this band other than equipment related to altimeter stations. However, passive sensing in the Earth Exploration-Satellite and Space Research services may be authorized in this band on a secondary basis (no protection is provided from the radio altimeters).

US262—The band 31.8–32.3 GHz is also allocated for space-to-Earth transmissions in the Space Research Service, limited to deep space communications at Goldstone, California.

US263—In the frequency bands 21.2–21.4, 22.21–22.5, 36–37, 50.2–50.4, 54.25–58.2, 116–126, 150–151, 174.5–176.5, 200–202, and 235–238 GHz, the Space Research and the Earth Exploration-Satellite Services shall not receive protection from the Fixed and Mobile Services operating in accordance with the Table of Frequency Allocations.

US264—In the band 48.94–49.04 GHz, airborne stations shall not be authorized.

US265—In the band 10.6–10.68 GHz, the fixed service shall be limited to a maximum equivalent isotropically radiated power of 40 dBW and the power delivered to the antenna shall not exceed –3 dBW, per 250 kHz.

US266—Licensees in the Public Safety Radio Services holding a valid authorization on June 30, 1958, to operate in the frequency band 156.27–157.47 MHz or on the frequencies of 161.85, 161.91, or 161.97 MHz may, upon proper application, continue to be authorized for such operation, including expansion of existing systems, until such time as harmful interference is caused to the operation of any authorized station other than those licensed in the Public Safety Radio Service.

US267—In the band 902–928 MHz, amateur radio stations shall not operate within the States of Colorado and Wyoming, bounded by the area of: latitude 39° N to 42° N and longitude 103° W to 108° W.

US268—The bands 890–902 MHz and 928–942 MHz are also allocated to the radiolocation service for Government ship stations (offshore ocean areas) on the condition that harmful interference is not caused to non-Government land mobile stations. The provisions of footnote US116 apply.

US269—In the band 2500–2690 MHz, applicants for space station assignments are urged to take all practicable steps to protect radio astronomy observations in the adjacent band, 2690–2700 MHz, from harmful interference. Further, all applicants are urged to coordinate their proposed systems through the Electromagnetic Spectrum Management Unit, National Science Foundation, Washington, DC 20550, prior to systems development.

US270—The band 72.77–72.91 GHz is also allocated to the radio astronomy service. Applicants for frequency assignments in this band are urged to take all practicable steps to protect radio astronomy observations from harmful interference.

US271—The use of the band 17.3–17.8 GHz by the Fixed-Satellite Service (Earth-to-space) is limited to feeder links for Broadcasting-Satellite Service.

US273—In the 74.6–74.8 MHz and 75.2–75.4 MHz bands, stations in the fixed and mobile services are limited to a maximum power of 1 watt from the transmitter into the antenna transmission line.

US274—In the 216–220 MHz band, fixed, aeronautical mobile, and land mobile stations are limited to telemetering and associated telecommand operations.

US275—The band 902–928 MHz is allocated on a secondary basis to the amateur service subject to not causing harmful interference to the operations of Government stations authorized in this band or to Location and Monitoring Service (LMS) systems. Stations in the amateur service must tolerate any interference from the operations of industrial, scientific, and medical (ISM) devices, LMS systems, and the operations of Government stations authorized in this band. Further, the amateur service is prohibited in those portions of Texas and New Mexico bounded on the south by latitude 31°41' N, on the east by longitude 104°11' W, on the north by latitude 34°30' N, and on the west by longitude 107°30' W; in addition, outside this area but within 240 kilometers (150 miles) of these boundaries of White Sands Missile Range the service is

restricted to a maximum transmitter peak envelope power output of 50 watts.

US276—Except as otherwise provided for herein, use of the band 2310–2390 MHz by the mobile service is limited to aeronautical telemetering and associated telecommand operations for flight testing of manned or unmanned aircraft, missiles, or major components thereof. The following six frequencies are shared on a co-equal basis by Government and non-Government stations for telemetering and associated telecommand operations of expendable and reusable launch vehicles whether or not such operations involve flight testing: 2312.5, 2332.5, 2352.5, 2364.5, 2370.5, and 2382.5 MHz. All other mobile telemetering uses shall be secondary to the above uses.

US277—The band 10.6–10.68 GHz is also allocated on a primary basis to the radio astronomy service. However, the radio astronomy service shall not receive protection from stations in the Fixed Service, which is licensed to operate in the 100 most populous urbanized areas as defined by the US Census Bureau. The following radio astronomy sites have been coordinated for observations in this band: National Radio Astronomy Observatory, Green Bank, West Virginia (38°26'08" N; 79°49'42" W); National Radio Astronomy Observatory, Socorro, New Mexico (34°04'43" N; 107°37'04" W); Harvard Radio Astronomy Station, Fort Davis, Texas (30°38'08" N; 103°56'42" W); Hat Creek Observatory, Hat Creek, California (40°49'03" N; 121°28'24" W); Owens Valley Radio Observatory, Big Pine, California (37°13'54" N; 118°17'36" W); Naval Research Laboratory, Maryland Point, Maryland (38°22'26" N; 77°14'00" W).

US278—In the 22.55–23.55 and 32–33 GHz bands, non-geostationary inter-satellite links may operate on a secondary basis to geostationary inter-satellite links.

US279—The frequency 2182 kHz may be authorized to fixed stations associated with the maritime mobile service for the sole purpose of transmitting distress calls and distress traffic and urgency and safety signals and messages.

US281—In the band 25.07–25.11 MHz, non-Government stations in the Industrial Radio Services shall not cause harmful interference to, and must accept interference from, stations in the Maritime Mobile Service operating in accordance with the International Table of Frequency Allocations.

US282—In the band 4650–4700 kHz, frequencies may be authorized for non-Government communication with helicopters in support of offshore drilling operations on the condition that harmful interference will not be caused to services operating in accordance with the Table of Frequency Allocations.

US283—In the bands 2850–3025 kHz, 3400–3500 kHz, 4650–4700 kHz, 5450–5680 kHz, 6525–6685 kHz, 10005–10100 kHz, 11275–11400 kHz, 13260–13360 kHz, and 17900–17970 kHz frequencies in these bands may be authorized for non-Government flight test purposes on the condition that harmful interference will not be caused to services operating in accordance with the Table of Frequency Allocations.

US285—Under exceptional circumstances, the carrier frequencies 2635, 2638, and 2738 kHz may be authorized to coast stations.

US287—The band 14–14.5 GHz is also allocated to the non-Government land mobile-satellite service (Earth-to-space) on a secondary basis.

US290—In the band 1900–2000 kHz, amateur stations may continue to operate on a secondary basis to the Radiolocation Service, pending a decision as to their disposition through a future rule making proceeding in conjunction with implementation of the Standard Broadcasting Service in the 1625–1705 kHz band.

US291—Television pickup stations in the mobile service may be authorized to use frequencies in the band 38.6–40 GHz on a secondary basis to stations operating in accordance with the Table of Frequency Allocations.

US292—In the band 14.0–14.2 GHz, stations in the radiolocation service shall operate on a secondary basis to the fixed-satellite service.

US294—In the spectrum below 490 kHz electric utilities operate Power Line Carrier (PLC) systems on power transmission lines for communications important to the reliability and security of electric service to the public. These PLC systems operate under the provisions of Part 15 of the Federal Communication Commission's Rules and Regulations or Chapter 7 of the National Telecommunications and Information Administration's Manual of Regulations and Procedures for Federal Radio Frequency Management, on an unprotected and noninterference basis with respect to authorized radio users. Notification of intent to place new or revised radio frequency assignments or PLC frequency uses in the bands below 490 kHz is to be made in accordance with the Rules and Regulations of the FCC and NTIA, and users are urged to minimize potential interference to the degree practicable. This footnote does not provide any allocation status to PLC radio frequency uses.

US296—In the bands designated for ship wide-band telegraphy, facsimile, and special transmission systems, the following assignable frequencies are available to non-Federal Government stations on a shared basis with Government stations: 2070.5, 2072.5, 2074.5, 2076.5, 4154, 4170, 6235, 6259, 8302, 8338, 12370, 12418, 16551, 16615, 18848, 18868, 22182, 22238, 25123, and 25159 kHz.

US297—The bands 47.2–49.2 GHz and 74.0–75.5 GHz are also available for feeder links for the broadcasting-satellite service.

US298—Channels 27555, 27615, 27635, 27655, 27765, and 27860 kHz are available to eligibles in the Forest Products Radio Service on a secondary basis to Government operations including experimental stations. Operations in the Forest Products Radio Service on these channels will not exceed 150 watts and are limited to the States of Washington, Oregon, Maine, North Carolina, South Carolina, Tennessee, Georgia, Florida, Alabama, Mississippi, Louisiana, and Texas (eastern portion).

US299—Until implementation procedures and schedules are determined by a future Regional Conference of the International Telecommunication Union the frequency bands 1615–1625 and 1625–1705 kHz in Alaska are also allocated to the maritime mobile services and the Alaska fixed service.

US300—The frequencies 169.445, 169.505, 170.245, 170.305, 171.045, 171.105, 171.845, and 171.905 MHz are

available for wireless microphone operations on a secondary basis to Government and non-Government operations.

US301—Except as provided in US302, broadcast auxiliary stations licensed as of November 21, 1984, to operate in the band 942–944 MHz may continue to operate on a co-equal primary basis to other stations and services operating in the band in accordance with the Table of Frequency Allocations.

US302—The band 942–944 MHz in Puerto Rico is allocated as an alternative allocation to the fixed service for broadcast auxiliary stations only.

US303—In the band 2285–2290 MHz, non-Federal government space stations in the space research, space operations, and earth exploration-satellite services may be authorized to transmit to the Tracking and Data Relay Satellite System subject to such conditions as may be applied on a case-by-case basis. Such transmissions shall not cause harmful interference to authorized Federal Government stations. The power flux density at the Earth's surface from such non-Government stations shall not exceed -144 to -154 dBW/m²/4 kHz, depending on angle of arrival, in accordance with ITU Radio Regulation S21.16.

US307—The sub-band 5150–5216 MHz is also allocated for space-to-Earth transmissions in the fixed-satellite service for feeder links in conjunction with the radiodetermination-satellite service operating in the bands 1610–1626.5 MHz and 2483.5–2500 MHz. The total power flux density at the earth's surface shall in no case exceed -159 dBW/m²/4 kHz for all angles of arrival.

US308—In the frequency bands 1549.5–1558.5 MHz and 1651–1660 MHz, the Aeronautical Mobile-Satellite (R) requirements that cannot be accommodated in the 1545–1549.5 MHz, 1558.5–1559 MHz, 1646.5–1651 MHz, and 1660–1660.5 MHz bands shall have priority access with real-time preemptive capability for communications in the mobile-satellite service. Systems not interoperable with the aeronautical mobile-satellite (R) service shall operate on a secondary basis. Account shall be taken of the priority of safety-related communications in the mobile-satellite service.

US309—Transmissions in the bands 1545–1559 MHz from terrestrial aeronautical stations directly to aircraft stations, or between aircraft stations, in the aeronautical mobile (R) service are also authorized when such transmissions are used to extend or supplement the satellite-to-aircraft links. Transmissions in the band 1646.5–1660.5 MHz from aircraft stations in the aeronautical mobile (R) service directly to terrestrial aeronautical stations, or between aircraft stations, are also authorized when such transmissions are used to extend or supplement the aircraft-to-satellite links.

US310—In the band 14.896–15.121 GHz, non-Federal Government space stations in the space research service may be authorized on a secondary basis to transmit to Tracking and Data Relay Satellites subject to such conditions as may be applied on a case-by-case basis. Such transmissions shall not cause harmful interference to authorized Federal Government stations. The power flux density at the earth's surface from such non-Government stations shall not exceed -138 to -148 dBW/m²/4 kHz, depending on the angle of arrival, in accordance with Recommendation ITU-R SA.510-1(1994).

US311—Radio astronomy observations may be made in the 1350–1400 MHz band on an unprotected basis at certain radio astronomy observatories indicated below:

National Astronomy and Ionosphere Center, Arecibo, Puerto Rico	Rectangle between latitudes 17°30'N and 19°00'N and between longitudes 65°10'W and 68°00'W.	
National Radio Astronomy Observatory, Socorro, New Mexico	Rectangle between latitudes 32°30'N and 35°30'N and between longitudes 106°00'W and 109°00'W.	
National Radio Astronomy Observatory, Green Bank, West Virginia	Rectangle between latitudes 37°30'N and 39°15'N and between longitudes 78°30'W and 80°30'W.	
National Radio Astronomy Observatory, very Long Baseline Array Stations	80-kilometer (50-mile) radius centered on:	
	Latitude (North)	Longitude (West)
Pie Town, NM	34°18'	108°07'
Kitt Peak, AZ	31°57'	111°37'
Los Alamos, NM	35°47'	106°15'
Fort Davis, TX	30°38'	103°57'
North Liberty, IA	41°46'	91°34'
Brewster, WA	48°08'	119°41'
Owens Valley, CA	37°14'	118°17'
Saint Croix, VI	17°46'	64°35'
Mauna Kea, HI	19°48'	155°27'
Hancock, NH	42°56'	71°59'

Every practicable effort will be made to avoid the assignment of frequencies in the band 1350–1400 MHz to stations in the fixed and mobile services that could interfere with radio astronomy observations within the geographic areas given above. In addition, every practicable effort will be made to avoid assignment of frequencies in this band to stations in the aeronautical mobile service that operate outside of those geographic areas, but which may cause harmful interference to the listed observatories. Should such assignments result in harmful interference to these observatories, the situation will be remedied to the extent practicable.

US312—The frequency 173.075 MHz may also be authorized on a primary basis to non-Government stations in the Police Radio Service (with a maximum authorized bandwidth of 20 kHz) for stolen vehicle recovery systems.

US315—In the frequency bands 1530–1544 MHz and 1626.5–1645.5 MHz maritime mobile-satellite distress and safety communications, e.g., GMDSS, shall have priority access with real-time preemptive capability in the mobile-satellite service. Communications of mobile-satellite system stations not participating in the GMDSS shall operate on a secondary basis to distress and safety communications of sta-

tions operating in the GMDSS. Account shall be taken of the priority of safety-related communications in the mobile-satellite service.

US316—The band 2900–3000 MHz is also allocated on a primary basis to the Meteorological Aids Service. Operations in this service are limited to Government Next Generation Weather Radar (NEXRAD) systems where accommodation in the 2700–2900 MHz band is not technically practical and are subject to coordination with existing authorized stations.

US317—The band 218.0–219.0 MHz is allocated on a primary basis to the Interactive Video and Data Operations.

US318—Until January 1, 2000, the use of the 137–138 MHz band by the mobile-satellite service will be secondary to Government operations within the sub-bands: 137.333–137.367, 137.485–137.515, 137.605–137.635, and 137.753–137.787 MHz.

US319—In the 137–138, 148–149.9, 149.9–150.05, 399.9–400.05, 400.15–401, 1610–1626.5, and 2483.5–2500 MHz bands, Federal Government stations in the mobile-satellite service shall be limited to earth stations operating with non-Federal Government space stations.

US320—Use of the 137–138, 148–149.9, and 400.15–401 MHz bands by the mobile-satellite service is limited to non-voice, nongeostationary satellite systems and may include satellite links between land earth stations at fixed locations.

US322—Use of bands 149.9–150.05 MHz and 399.9–400.05 MHz by the mobile-satellite service (Earth-to-space) is limited to nonvoice, nongeostationary satellite systems, including satellite links between land earth stations.

US323—In the 148–149.9 MHz band, no individual mobile earth station shall transmit on the same frequency being actively used by fixed and mobile stations and shall transmit no more than 1% of the time during any 15 minute period; except, individual mobile earth stations in this band that do not avoid frequencies actively being used by the fixed and mobile services shall not exceed a power density of -16 dBW/4 kHz and shall transmit no more than 0.25% of the time during any 15 minute period. Any single transmission from any individual mobile earth station operating in this band shall not exceed 450 ms in duration and consecutive transmissions from a single mobile earth station on the same frequency shall be separated by at least 15 seconds. Land earth stations in this band shall be subject to electromagnetic compatibility analysis and coordination with fixed and mobile stations.

US324—Government and non-Government satellite systems in the 400.15–401 MHz band shall be subject to electromagnetic compatibility analysis and coordination.

US325—In the band 148–149.9 MHz fixed and mobile stations shall not claim protection from land earth stations in the mobile-satellite service that have been previously coordinated; Government fixed and mobile stations exceeding 27 dBW eirp, or an emission bandwidth greater than 38 kHz, will be coordinated with existing mobile-satellite service space stations.

US326—The 399.9–400.05 MHz band is allocated to the mobile-satellite service (Earth-to-space) on a primary basis after January 1, 1997, and shall be limited to nonvoice, nongeostationary satellite systems, including satellite links between land earth stations.

US327—The band 2310–2360 MHz is allocated to the broadcasting-satellite service (sound) and complementary terrestrial broadcasting service on a primary basis. Such use is limited to digital audio broadcasting and is subject to the provisions of Resolution 528.

US328—In the band 2310–2360 MHz, the mobile and radiolocation services are allocated on a primary basis until 1 January 1997 or until broadcasting-satellite (sound) service has been brought into use in such a manner as to affect or be affected by the mobile and radiolocation services in those service areas, whichever is later. The broadcasting-satellite (sound) service during implementation should also take cognizance of the expendable and reusable launch vehicle frequencies 2312.5, 2332.5, and 2352.5 MHz, to minimize the impact on this mobile service use to the extent possible.

US334—In the band 17.8–20.2 GHz, Government space stations and associated earth stations in the fixed-satellite (space-to-Earth) service may be authorized on a primary basis. For a Government geostationary satellite network to operate on a primary basis, the space station shall be located outside the arc measured from East to West, 70° W to 120° W. Coordination between Government fixed-satellite systems and non-Government systems operating in accordance with the US Table of Frequency Allocations is required.

US335—The primary Government and non-Government allocations for the various segments of the 220–222 MHz band are divided as follows:

- (1) The 220.0–220.55/221.0–221.55, 220.6–220.8/221.6–221.8, 220.85–220.90/221.85–221.90, and 220.925–221.0/221.925–222.0 MHz bands (Channels 1–110, 121–160, 171–180, and 186–200, respectively) are available for exclusive non-Government use.
- (2) The 220.55–220.60/221.55–221.60 MHz bands (Channels 111–120) are available for exclusive Government use.
- (3) The 220.80–220.85/221.80–221.85 and 220.900–220.925/221.900–221.925 MHz bands (Channels 161–170 and 181–185, respectively) are available for shared Government and non-Government use. The exclusive non-Government band segments are also available for temporary fixed geophysical telemetry operations on a secondary basis to the fixed and mobile services.

US337—In the band 13.75–13.80 GHz, earth stations in the fixed-satellite service shall be coordinated on a case-by-case basis through the frequency assignment subcommittee to minimize harmful interference to the Tracking and Data Relay Satellite System's forward space-to-space link (TDRSS forward link-to-LEO).

US338—In the 2305–2310 MHz band, space-to-Earth operations are prohibited. Additionally, in the 2305–2320 MHz band, all Wireless Communications Service (WCS) operations within 50 kilometers of $35^\circ 20' N$ and $116^\circ 53' W$ shall be coordinated through the Frequency Assignment Subcommittee of the Interdepartment Radio Advisory Committee to minimize harmful interference to NASA's Goldstone Deep Space facility.

US339—The bands 2310–2320 and 2345–2360 MHz are also available for aeronautical telemetering and associated telecommand operations for flight testing of manned or unmanned aircraft, missiles, or major components thereof on a secondary basis to the Wireless Communications Service. The following two frequencies are shared on a co-equal basis by Government and non-Government stations for telemetering and associated telecommand operations of expendable and reusable launch vehicles whether or not such operations involve flight testing: 2312.5 and 2352.5 MHz. Other mobile telemetering uses may be provided on a noninterference basis to the above uses. The broadcasting-satellite (sound) service during implementation should also take cognizance of the expendable and reusable launch vehicle frequencies 2312.5 and 2352.5 MHz, to minimize the impact on this mobile service use to the extent possible.

US340—The 2–30 MHz band is available on a secondary noninterference basis to Government and non-Government maritime and aeronautical stations for the purposes of measuring the quality of reception on radio channels. See 47 C.F.R. Sec. 87.149 for the list of protected frequencies and bands within this frequency range. Actual communications shall be limited to those frequencies specifically allocated to the maritime mobile and aeronautical mobile services.

US342—In making assignments to stations of other services to which are allocated the following bands:

13360–13410 kHz
 37.5–38.25 MHz
 322–328.6 MHz*
 1330–1400 MHz*
 1610.6–1613.8 MHz*
 1660–1670 MHz
 3260–3267 MHz*
 3332–3339 MHz*
 3345.8–3352.5 MHz*
 4825–4835 MHz*
 14.47–14.5 GHz*
 22.01–22.21 GHz*
 22.21–22.5 GHz
 22.81–22.86 GHz*
 23.07–23.12 GHz*
 31.2–31.3 GHz
 36.43–36.5 GHz*
 42.5–43.5 GHz
 48.94–49.04 GHz*
 97.88–98.08 GHz*
 140.69–140.98 GHz*
 144.68–144.98 GHz*
 145.45–145.75 GHz*
 146.82–147.12 GHz*
 262.24–262.76 GHz*
 265–275 GHz

(* indicates radio astronomy use for spectral line observations), all practicable steps shall be taken to protect the radio astronomy service from harmful interference. Emissions from spaceborne or airborne stations can be particularly serious sources of interference to the radio astronomy service (see Nos. 343/S4.5 and 344/S4.6 and Article 36/S29 of the ITU Radio Regulations).

International Footnotes

(These footnotes come from the ITU Radio Regulations, Geneva, 1998.)

Until such time as the FCC has considered the substantially revised international footnotes that have previously been adopted domestically, the old international footnotes shall apply in the United States. These old footnotes appear immediately after S5.565.

S5.3—Region 1: Region 1 includes the area limited on the east by line A (lines A, B, and C are defined below) and on the west by line B, excluding any of the territory of the Islamic Republic of Iran which lies between these limits. It also includes the whole of the territory of Armenia, Azerbaijan, Russian Federation, Georgia, Kazakstan, Mongolia, Uzbekistan, Kyrgyzstan, Tajikistan, Turkmenistan, Turkey, and Ukraine and the area to the north of Russian Federation that lies between lines A and C.

S5.4—Region 2: Region 2 includes the area limited on the east by line B and on the west by line C.

S5.5—Region 3: Region 3 includes the area limited on the east by line C and on the west by line A, except any of the territory of Armenia, Azerbaijan, Russian Federation, Georgia, Kazakstan, Mongolia, Uzbekistan, Kyrgyzstan, Tajikistan, Turkmenistan, Turkey, and Ukraine and the area to the north of Russian Federation. It also includes that part of the territory of the Islamic Republic of Iran lying outside of those limits.

S5.6—The lines A, B, and C are defined as follows:

S5.7—Line A: Line A extends from the North Pole along meridian 40° East of Greenwich to parallel 40° North; thence by great circle arc to the intersection of meridian 60° East and the Tropic of Cancer; thence along the meridian 60° East to the South Pole.

S5.8—Line B: Line B extends from the North Pole along meridian 10° West of Greenwich to its intersection with parallel 72° North; thence by great circle arc to the intersection of meridian 50° West and parallel 40° North; thence by great circle arc to the intersection of meridian 20° West and parallel 10° South; thence along meridian 20° West to the South Pole.

S5.9—Line C: Line C extends from the North Pole by great circle arc to the intersection of parallel 65°30′ North with the international boundary in Bering Strait; thence by great circle arc to the intersection of meridian 165° East of Greenwich and parallel 50° North; thence by great circle arc to the intersection of meridian 170° West and parallel 10° North; thence along parallel 10° North to its intersection with meridian 120° West; thence along meridian 120° West to the South Pole.

S5.53—Administrations authorizing the use of frequencies below 9 kHz shall ensure that no harmful interference is caused thereby to the services to which the bands above 9 kHz are allocated.

S5.54—Administrations conducting scientific research using frequencies below 9 kHz are urged to advise other administrations that may be concerned so that such research may be afforded all practicable protection from harmful interference.

S5.55—Additional allocation: in Armenia, Azerbaijan, Bulgaria, Russian Federation, Georgia, Kazakstan, Kyrgyzstan, Tajikistan, Turkmenistan, and Ukraine, the band 14–

17 kHz is also allocated to the radionavigation service on a primary basis (WRC-97).

S5.56—The stations of services to which the bands 14–19.95 kHz and 20.05–70 kHz and in Region 1 also the bands 72–84 kHz and 86–90 kHz are allocated may transmit standard frequency and time signals. Such stations shall be afforded protection from harmful interference. In Armenia, Azerbaijan, Belarus, Bulgaria, Georgia, Kazakstan, Mongolia, Uzbekistan, Kyrgyzstan, Slovakia, the Czech Republic, Russian Federation, Tajikistan, Turkmenistan, and Ukraine, the frequencies 25 kHz and 50 kHz will be used for this purpose under the same conditions (WRC-97).

S5.57—The use of the bands 14–19.95 kHz, 20.05–70 kHz, and 70–90 kHz (72–84 kHz and 86–90 kHz in Region 1) by the maritime mobile service is limited to coast radiotelegraph stations (A1A and F1B only). Exceptionally, the use of class J2B or J7B emissions is authorized subject to the necessary bandwidth not exceeding that normally used for class A1A or F1B emissions in the band concerned.

S5.58—*Additional allocation:* in Armenia, Azerbaijan, Bulgaria, Georgia, Kazakstan, Kyrgyzstan, Russian Federation, Tajikistan, Turkmenistan, and Ukraine, the band 67–70 kHz is also allocated to the radionavigation service on a primary basis (WRC-97).

S5.59—*Different category of service:* in Bangladesh, the Islamic Republic of Iran, and Pakistan, the allocation of the bands 70–72 kHz and 84–86 kHz to the fixed and maritime mobile services is on a primary basis (see No. S5.33).

S5.60—In the bands 70–90 kHz (70–86 kHz in Region 1) and 110–130 kHz (112–130 kHz in Region 1), pulsed radionavigation systems may be used on condition that they do not cause harmful interference to other services to which these bands are allocated.

S5.61—In Region 2, the establishment and operation of stations in the maritime radionavigation service in the bands 70–90 kHz and 110–130 kHz shall be subject to agreement obtained under No. S9.21 with administrations whose services, operating in accordance with the Table, may be affected. However, stations of the fixed, maritime mobile, and radiolocation services shall not cause harmful interference to stations in the maritime radionavigation service established under such agreements.

S5.62—Administrations that operate stations in the radionavigation service in the band 90–110 kHz are urged to coordinate technical and operating characteristics in such a way as to avoid harmful interference to the services provided by these stations.

S5.63—(SUP – WRC-97).

S5.64—Only classes A1A, F1B, A2C, A3C, F1C, or F3C emissions are authorized for stations of the fixed service in the bands allocated to this service between 90 kHz and 160 kHz (148.5 kHz in Region 1) and for stations of the maritime mobile service in the bands allocated to this service between 110 kHz and 160 kHz (148.5 kHz in Region 1). Exceptionally, class J2B or J7B emissions are also authorized in the bands between 110 kHz and 160 kHz (148.5 kHz in Region 1) for stations of the maritime mobile service.

S5.65—*Different category of service:* in Bangladesh, the Islamic Republic of Iran, and Pakistan, the allocation of the

bands 112–117.6 kHz and 126–129 kHz to the fixed and maritime mobile services is on a primary basis (see No. S5.33).

S5.66—*Different category of service:* in Germany, the allocation of the band 115–117.6 kHz to the fixed and maritime mobile services is on a primary basis (see No. S5.33) and to the radionavigation service on a secondary basis (see No. S5.32).

S5.67—*Additional allocation:* in Azerbaijan, Bulgaria, Mongolia, Kyrgyzstan, Romania, Turkmenistan, and Ukraine, the band 130–148.5 kHz is also allocated to the radionavigation service on a secondary basis. Within and between these countries this service shall have an equal right to operate (WRC-97).

S5.68—*Alternative allocation:* in Angola, Botswana, Burundi, the Congo, Malawi, Dem. Rep. of the Congo, Rwanda, and South Africa, the band 160–200 kHz is allocated to the fixed service on a primary basis.

S5.69—*Additional allocation:* in Somalia, the band 200–255 kHz is also allocated to the aeronautical radionavigation service on a primary basis.

S5.70—*Alternative allocation:* in Angola, Botswana, Burundi, Cameroon, the Central African Republic, the Congo, Ethiopia, Kenya, Lesotho, Madagascar, Malawi, Mozambique, Namibia, Nigeria, Oman, Democratic Republic of the Congo, Rwanda, South Africa, Swaziland, Tanzania, Chad, Zambia, and Zimbabwe, the band 200–283.5 kHz is allocated to the aeronautical radionavigation service on a primary basis.

S5.71—*Alternative allocation:* in Tunisia, the band 255–283.5 kHz is allocated to the broadcasting service on a primary basis.

S5.72—Norwegian stations of the fixed service situated in northern areas (north of 60° N) subject to auroral disturbances are allowed to continue operation on four frequencies in the bands 283.5–490 kHz and 510–526.5 kHz.

S5.73—The band 285–325 kHz (283.5–325 kHz in Region 1) in the maritime radionavigation service may be used to transmit supplementary navigational information using narrow-band techniques, on condition that no harmful interference is caused to radiobeacon stations operating in the radionavigation service (WRC-97).

S5.74—*Additional allocation:* in Region 1, the frequency band 285.3–285.7 kHz is also allocated to the maritime radionavigation service (other than radiobeacons) on a primary basis.

S5.75—*Different category of service:* in Armenia, Azerbaijan, Belarus, Georgia, Kazakstan, Moldova, Kyrgyzstan, Russian Federation, Tajikistan, Turkmenistan, Ukraine, and the Black Sea areas of Bulgaria and Romania, the allocation of the band 315–325 kHz to the maritime radionavigation service is on a primary basis under the condition that in the Baltic Sea area, the assignment of frequencies in this band to new stations in the maritime or aeronautical radionavigation services shall be subject to prior consultation between the administrations concerned.

S5.76—The frequency 410 kHz is designated for radio direction-finding in the maritime radionavigation service. The other radionavigation services to which the band 405–415 kHz is allocated shall not cause harmful interference to radio direction-finding in the band 406.5–413.5 kHz.

S5.77—*Different category of service:* in Australia, China, the French Overseas Territories of Region 3, India, Indonesia, the Islamic Republic of Iran, Japan, Pakistan, Papua New Guinea, and Sri Lanka, the allocation of the band 415–495 kHz to the aeronautical radionavigation service is on a primary basis. Administrations in these countries shall take all practical steps necessary to ensure that aeronautical radionavigation stations in the band 435–495 kHz do not cause interference to reception by coast stations of ship stations transmitting on frequencies designated for ship stations on a worldwide basis (see No. **S52.39**).

S5.78—*Different category of service:* in Cuba, the United States of America, and Mexico, the allocation of the band 415–435 kHz to the aeronautical radionavigation service is on a primary basis.

S5.79—The use of the bands 415–495 kHz and 505–526.5 kHz (505–510 kHz in Region 2) by the maritime mobile service is limited to radiotelegraphy.

S5.79A—When establishing coast stations in the NAVTEX service on the frequencies 490 kHz, 518 kHz, and 4209.5 kHz, administrations are strongly recommended to coordinate the operating characteristics in accordance with the procedures of the International Maritime Organization (IMO) (see Resolution **339 (Rev.WRC-97)**) (WRC-97).

S5.80—In Region 2, the use of the band 435–495 kHz by the aeronautical radionavigation service is limited to nondirectional beacons not employing voice transmission.

S5.81—The bands 490–495 kHz and 505–510 kHz shall be subject to the provisions of Appendix **S13**, § 15-1, Part A2 (WRC-97).

S5.82—In the maritime mobile service, the frequency 490 kHz is, from the date of full implementation of the GMDSS (see Resolution **331 (Rev.WRC-97)**), to be used exclusively for the transmission by coast stations of navigational and meteorological warnings and urgent information to ships, by means of narrow-band direct-printing telegraphy. The conditions for use of the frequency 490 kHz are prescribed in Articles **S31** and **S52**. In using the band 415–495 kHz for the aeronautical radionavigation service, administrations are requested to ensure that no harmful interference is caused to the frequency 490 kHz (WRC-97).

S5.83—The frequency 500 kHz is an international distress and calling frequency for Morse radiotelegraphy. The conditions for its use are prescribed in Articles **S31** and **S52** and in Appendix **S13**.

S5.84—The conditions for the use of the frequency 518 kHz by the maritime mobile service are prescribed in Articles **S31** and **S52** and in Appendix **S13** (WRC-97).

S5.85—Not used.

S5.86—In Region 2, in the band 525–535 kHz the carrier power of broadcasting stations shall not exceed 1 kW during the day and 250 W at night.

S5.87—*Additional allocation:* in Angola, Botswana, Lesotho, Malawi, Mozambique, Namibia, South Africa, Swaziland, Zambia, and Zimbabwe, the band 526.5–535 kHz is also allocated to the mobile service on a secondary basis.

S5.87A—*Additional allocation:* in Uzbekistan, the band 526.5–1606.5 kHz is also allocated to the radionavigation service on a primary basis. Such use is subject to agreement obtained under No. **S9.21** with administrations concerned and

limited to ground-based radiobeacons in operation on 27 October 1997 until the end of their lifetime (WRC-97).

S5.88—*Additional allocation:* in China, the band 526.5–535 kHz is also allocated to the aeronautical radionavigation service on a secondary basis.

S5.89—In Region 2, the use of the band 1605–1705 kHz by stations of the broadcasting service is subject to the Plan established by the Regional Administrative Radio Conference (Rio de Janeiro, 1988).

The examination of frequency assignments to stations of the fixed and mobile services in the band 1625–1705 kHz shall take account of the allotments appearing in the Plan established by the Regional Administrative Radio Conference (Rio de Janeiro, 1988).

S5.90—In the band 1605–1705 kHz, in cases where a broadcasting station of Region 2 is concerned, the service area of the maritime mobile stations in Region 1 shall be limited to that provided by ground-wave propagation.

S5.91—*Additional allocation:* in the Philippines and Sri Lanka, the band 1606.5–1705 kHz is also allocated to the broadcasting service on a secondary basis (WRC-97).

S5.92—Some countries of Region 1 use radiodetermination systems in the bands 1606.5–1625 kHz, 1635–1800 kHz, 1850–2160 kHz, 2194–2300 kHz, 2502–2850 kHz, and 3500–3800 kHz, subject to agreement obtained under No. **S9.21**. The radiated mean power of these stations shall not exceed 50 W.

S5.93—*Additional allocation:* in Angola, Armenia, Azerbaijan, Belarus, Bulgaria, Georgia, Hungary, Kazakstan, Latvia, Lithuania, Moldova, Mongolia, Nigeria, Uzbekistan, Poland, Kyrgyzstan, Slovakia, the Czech Republic, Russian Federation, Tajikistan, Chad, Turkmenistan, and Ukraine, the bands 1625–1635 kHz, 1800–1810 kHz, and 2160–2170 kHz are also allocated to the fixed and land mobile services on a primary basis, subject to agreement obtained under No. **S9.21**.

S5.94 and **S5.95**—Not used.

S5.96—In Germany, Armenia, Azerbaijan, Belarus, Denmark, Estonia, Finland, Georgia, Hungary, Ireland, Israel, Jordan, Kazakstan, Latvia, Lithuania, Malta, Moldova, Norway, Uzbekistan, Poland, Kyrgyzstan, Slovakia, the Czech Republic, the United Kingdom, Russian Federation, Sweden, Tajikistan, Turkmenistan, and Ukraine, administrations may allocate up to 200 kHz to their amateur service in the bands 1715–1800 kHz and 1850–2000 kHz. However, when allocating the bands within this range to their amateur service, administrations shall, after prior consultation with administrations of neighbouring countries, take such steps as may be necessary to prevent harmful interference from their amateur service to the fixed and mobile services of other countries. The mean power of any amateur station shall not exceed 10 W.

S5.97—In Region 3, the Loran system operates either on 1850 kHz or 1950 kHz, the bands occupied being 1825–1875 kHz and 1925–1975 kHz respectively. Other services to which the band 1800–2000 kHz is allocated may use any frequency therein on condition that no harmful interference is caused to the Loran system operating on 1850 kHz or 1950 kHz.

S5.98—*Alternative allocation:* in Angola, Armenia, Austria, Azerbaijan, Belarus, Belgium, Bulgaria, Cameroon, the Congo, Denmark, Egypt, Eritrea, Spain, Ethiopia, Georgia, Greece, Italy, Kazakstan, Lebanon, Lithuania, Moldova, the

Netherlands, Syria, Kyrgyzstan, Russian Federation, Somalia, Tajikistan, Tunisia, Turkmenistan, Turkey, and Ukraine, the band 1810–1830 kHz is allocated to the fixed and mobile, except aeronautical mobile, services on a primary basis (WRC-97).

S5.99—*Additional allocation*: in Saudi Arabia, Bosnia and Herzegovina, Iraq, Libya, Uzbekistan, Slovakia, the Czech Republic, Romania, Slovenia, Chad, Togo, and Yugoslavia, the band 1810–1830 kHz is also allocated to the fixed and mobile, except aeronautical mobile, services on a primary basis (WRC-97).

S5.100—In Region 1, the authorization to use the band 1810–1830 kHz by the amateur service in countries situated totally or partially north of 40° N shall be given only after consultation with the countries mentioned in Nos. **S5.98** and **S5.99** to define the necessary steps to be taken to prevent harmful interference between amateur stations and stations of other services operating in accordance with Nos. **S5.98** and **S5.99**.

S5.101—*Alternative allocation*: in Burundi and Lesotho, the band 1810–1850 kHz is allocated to the fixed and mobile, except aeronautical mobile, services on a primary basis.

S5.102—*Alternative allocation*: in Argentina, Bolivia, Chile, Mexico, Paraguay, Peru, Uruguay and Venezuela, the band 1850–2000 kHz is allocated to the fixed, mobile, except aeronautical mobile, radiolocation, and radionavigation services on a primary basis.

S5.103—In Region 1, in making assignments to stations in the fixed and mobile services in the bands 1850–2045 kHz, 2194–2498 kHz, 2502–2625 kHz, and 2650–2850 kHz, administrations should bear in mind the special requirements of the maritime mobile service.

S5.104—In Region 1, the use of the band 2025–2045 kHz by the meteorological aids service is limited to oceanographic buoy stations.

S5.105—In Region 2, except in Greenland, coast stations and ship stations using radiotelephony in the band 2065–2107 kHz shall be limited to class J3E emissions and to a peak envelope power not exceeding 1 kW. Preferably, the following carrier frequencies should be used: 2065.0 kHz, 2079.0 kHz, 2082.5 kHz, 2086.0 kHz, 2093.0 kHz, 2096.5 kHz, 2100.0 kHz, and 2103.5 kHz. In Argentina and Uruguay, the carrier frequencies 2068.5 kHz and 2075.5 kHz are also used for this purpose, while the frequencies within the band 2072–2075.5 kHz are used as provided in No. **S5.2.165**.

S5.106—In Regions 2 and 3, provided no harmful interference is caused to the maritime mobile service, the frequencies between 2065 kHz and 2107 kHz may be used by stations of the fixed service communicating only within national borders and whose mean power does not exceed 50 W. In notifying the frequencies, the attention of the Bureau should be drawn to these provisions.

S5.107—*Additional allocation*: in Saudi Arabia, Botswana, Eritrea, Ethiopia, Iraq, Lesotho, Libya, Somalia, Swaziland and Zambia, the band 2160–2170 kHz is also allocated to the fixed and mobile, except aeronautical mobile (R), services on a primary basis. The mean power of stations in these services shall not exceed 50 W (WRC-97).

S5.108—The carrier frequency 2182 kHz is an international distress and calling frequency for radiotelephony. The

conditions for the use of the band 2173.5–2190.5 kHz are prescribed in Articles **S31** and **S52** and in Appendix **S13**.

S5.109—The frequencies 2187.5 kHz, 4207.5 kHz, 6312 kHz, 8414.5 kHz, 12577 kHz, and 16804.5 kHz are international distress frequencies for digital selective calling. The conditions for the use of these frequencies are prescribed in Article **S31**.

S5.110—The frequencies 2174.5 kHz, 4177.5 kHz, 6268 kHz, 8376.5 kHz, 12520 kHz, and 16695 kHz are international distress frequencies for narrow-band direct-printing telegraphy. The conditions for the use of these frequencies are prescribed in Article **S31**.

S5.111—The carrier frequencies 2182 kHz, 3023 kHz, 5680 kHz, and 8364 kHz and the frequencies 121.5 MHz, 156.8 MHz, and 243 MHz may also be used, in accordance with the procedures in force for terrestrial radiocommunication services, for search and rescue operations concerning manned space vehicles. The conditions for the use of the frequencies are prescribed in Article **S31** and in Appendix **S13**.

The same applies to the frequencies 10003 kHz, 14993 kHz, and 19993 kHz, but in each of these cases emissions must be confined in a band of 3 kHz about the frequency.

S5.112—*Alternative allocation*: in Bosnia and Herzegovina, Cyprus, Denmark, France, Greece, Iceland, Italy, Malta, Norway, Sri Lanka, Turkey, and Yugoslavia, the band 2194–2300 kHz is allocated to the fixed and mobile, except aeronautical mobile, services on a primary basis (WRC-97).

S5.113—For the conditions for the use of the bands 2300–2495 kHz (2498 kHz in Region 1), 3200–3400 kHz, 4750–4995 kHz, and 5005–5060 kHz by the broadcasting service, see Nos. **S5.16** to **S5.20**, **S5.21** and **S23.3** to **S23.10**.

S5.114—*Alternative allocation*: in Bosnia and Herzegovina, Cyprus, Denmark, France, Greece, Iraq, Italy, Malta, Norway, Turkey, and Yugoslavia, the band 2502–2625 kHz is allocated to the fixed and mobile, except aeronautical mobile, services on a primary basis (WRC-97).

S5.115—The carrier (reference) frequencies 3023 kHz and 5680 kHz may also be used, in accordance with Article **S31** and Appendix **S13** by stations of the maritime mobile service engaged in coordinated search and rescue operations.

S5.116—Administrations are urged to authorize the use of the band 3155–3195 kHz to provide a common worldwide channel for low-power wireless hearing aids. Additional channels for these devices may be assigned by administrations in the bands between 3155 kHz and 3400 kHz to suit local needs. It should be noted that frequencies in the range 3000 kHz to 4000 kHz are suitable for hearing aid devices, which are designed to operate over short distances within the induction field.

S5.117—*Alternative allocation*: in Bosnia and Herzegovina, Cyprus, Côte d'Ivoire, Denmark, Egypt, France, Greece, Iceland, Italy, Liberia, Malta, Norway, Sri Lanka, Togo, Turkey, and Yugoslavia, the band 3155–3200 kHz is allocated to the fixed and mobile, except aeronautical mobile, services on a primary basis (WRC-97).

S5.118—*Additional allocation*: in the United States, Japan, Mexico, Peru, and Uruguay, the band 3230–3400 kHz is also allocated to the radiolocation service on a secondary basis.

S5.119—*Additional allocation:* in Honduras, Mexico, Peru, and Venezuela, the band 3500–3750 kHz is also allocated to the fixed and mobile services on a primary basis.

S5.120—For the use of the bands allocated to the amateur service at 3.5 MHz, 7.0 MHz, 10.1 MHz, 14.0 MHz, 18.068 MHz, 21.0 MHz, 24.89 MHz, and 144 MHz in the event of natural disasters, see Resolution 640.

S5.121—Not used.

S5.122—*Alternative allocation:* in Argentina, Bolivia, Chile, Ecuador, Paraguay, Peru, and Uruguay, the band 3750–4000 kHz is allocated to the fixed and mobile, except aeronautical mobile, services on a primary basis.

S5.123—*Additional allocation:* in Botswana, Lesotho, Malawi, Mozambique, Namibia, South Africa, Swaziland, Zambia, and Zimbabwe, the band 3900–3950 kHz is also allocated to the broadcasting service on a primary basis, subject to agreement obtained under No. S9.21.

S5.124—*Additional allocation:* in Canada, the band 3950–4000 kHz is also allocated to the broadcasting service on a primary basis. The power of broadcasting stations operating in this band shall not exceed that necessary for a national service within the frontier of this country and shall not cause harmful interference to other services operating in accordance with the Table.

S5.125—*Additional allocation:* in Greenland, the band 3950–4000 kHz is also allocated to the broadcasting service on a primary basis. The power of the broadcasting stations operating in this band shall not exceed that necessary for a national service and shall in no case exceed 5 kW.

S5.126—In Region 3, the stations of those services to which the band 3995–4005 kHz is allocated may transmit standard frequency and time signals.

S5.127—The use of the band 4000–4063 kHz by the maritime mobile service is limited to ship stations using radiotelephony (see No. S52.220 and Appendix S17).

S5.128—In Afghanistan, Argentina, Armenia, Azerbaijan, Belarus, Botswana, Burkina Faso, Central African Republic, China, Georgia, India, Kazakhstan, Mali, Niger, Kyrgyzstan, Russian Federation, Tajikistan, Chad, Turkmenistan, and Ukraine, in the bands 4063–4123 kHz, 4130–4133 kHz, and 4408–4438 kHz, stations of limited power in the fixed service that are situated at least 600 km from the coast may operate on condition that harmful interference is not caused to the maritime mobile service (WRC-97).

S5.129—On condition that harmful interference is not caused to the maritime mobile service, the frequencies in the bands 4063–4123 kHz and 4130–4438 kHz may be used exceptionally by stations in the fixed service communicating only within the boundary of the country in which they are located with a mean power not exceeding 50 W.

S5.130—The conditions for the use of the carrier frequencies 4125 kHz and 6215 kHz are prescribed in Articles S31 and S52 and in Appendix S13.

S5.131—The frequency 4209.5 kHz is used exclusively for the transmission by coast stations of meteorological and navigational warnings and urgent information to ships by means of narrow-band direct-printing techniques (WRC-97).

S5.132—The frequencies 4210 kHz, 6314 kHz, 8416.5 kHz, 12579 kHz, 16806.5 kHz, 19680.5 kHz, 22376 kHz, and 26100.5 kHz are the international frequencies for the trans-

mission of maritime safety information (MSI) (see Appendix S17).

S5.133—*Different category of service:* in Armenia, Azerbaijan, Belarus, Georgia, Kazakhstan, Latvia, Lithuania, Moldova, Uzbekistan, Kyrgyzstan, Russian Federation, Tajikistan, Turkmenistan, and Ukraine, the allocation of the band 5130–5250 kHz to the mobile, except aeronautical mobile, service is on a primary basis (see No. S5.33).

S5.134—The use of the bands 5900–5950 kHz, 7300–7350 kHz, 9400–9500 kHz, 11600–11650 kHz, 12050–12100 kHz, 13570–13600 kHz, 13800–13870 kHz, 15600–15800 kHz, 17480–17550 kHz, and 18900–19020 kHz by the broadcasting service is limited to single-sideband emissions with the characteristics specified in Appendix S11 or to any other spectrum-efficient modulation techniques recommended by ITU-R. Access to these bands shall be subject to the decisions of a competent conference (WRC-97).

S5.135—SUP-WRC-97.

S5.136—The band 5900–5950 kHz is allocated, until 1 April 2007, to the fixed service on a primary basis, as well as to the following services: in Region 1 to the land mobile service on a primary basis, in Region 2 to the mobile, except aeronautical mobile (R), service on a primary basis, and in Region 3 to the mobile, except aeronautical mobile (R), service on a secondary basis, subject to application of the procedure referred to in Resolution 21 (Rev. WRC-95). After 1 April 2007, frequencies in this band may be used by stations in the above-mentioned services, communicating only within the boundary of the country in which they are located, on the condition that harmful interference is not caused to the broadcasting service. When using frequencies for these services, administrations are urged to use the minimum power required and to take account of the seasonal use of frequencies by the broadcasting service published in accordance with the Radio Regulations.

S5.137—On condition that harmful interference is not caused to the maritime mobile service, the bands 6200–6213.5 kHz and 6220.5–6525 kHz may be used exceptionally by stations in the fixed service, communicating only within the boundary of the country in which they are located, with a mean power not exceeding 50 W. At the time of notification of these frequencies, the attention of the Bureau will be drawn to the above conditions.

S5.138—The following bands:

6765–6795 kHz (centre frequency 6780 kHz),

433.05–434.79 MHz (centre frequency 433.92 MHz) in Region 1 except in the countries mentioned in No. S5.280,

61–61.5 GHz (centre frequency 61.25 GHz),

122–123 GHz (centre frequency 122.5 GHz), and

244–246 GHz (centre frequency 245 GHz)

are designated for industrial, scientific, and medical (ISM) applications. The use of these frequency bands for ISM applications shall be subject to special authorization by the administration concerned, in agreement with other administrations whose radiocommunication services might be affected. In applying this provision, administrations shall have due regard to the latest relevant ITU-R Recommendations.

S5.139—*Different category of service:* in Armenia, Azerbaijan, Belarus, Georgia, Kazakstan, Latvia, Lithuania, Moldova, Mongolia, Uzbekistan, Kyrgyzstan, Russian Federation, Tajikistan, Turkmenistan, and Ukraine, the allocation of the band 6765–7000 kHz to the land mobile service is on a primary basis (see No. **S5.33**).

S5.140—*Additional allocation:* in Angola, Iraq, Rwanda, Somalia, and Togo, the band 7000–7050 kHz is also allocated to the fixed service on a primary basis.

S5.141—*Alternative allocation:* in Egypt, Eritrea, Ethiopia, Guinea, Libya, and Madagascar, the band 7000–7050 kHz is allocated to the fixed service on a primary basis (WRC-97).

S5.142—The use of the band 7100–7300 kHz in Region 2 by the amateur service shall not impose constraints on the broadcasting service intended for use within Region 1 and Region 3.

S5.143—The band 7300–7350 kHz is allocated, until April 1, 2007, to the fixed service on a primary basis and to the land mobile service on a secondary basis, subject to application of the procedure referred to in Resolution 21 (Rev.WRC-95). After 1 April 2007, frequencies in this band may be used by stations in the above-mentioned services, communicating only within the boundary of the country in which they are located, on condition that harmful interference is not caused to the broadcasting service. When using frequencies for these services, administrations are urged to use the minimum power required and to take account of the seasonal use of frequencies by the broadcasting service published in accordance with the Radio Regulations.

S5.144—In Region 3, the stations of those services to which the band 7995–8005 kHz is allocated may transmit standard frequency and time signals.

S5.145—The conditions for the use of the carrier frequencies 8291 kHz, 12290 kHz, and 16420 kHz are prescribed in Articles **S31** and **S52** and in Appendix **S13**.

S5.146—The bands 9400–9500 kHz, 11600–11650 kHz, 12050–12100 kHz, 15600–15800 kHz, 17480–17550 kHz, and 18900–19020 kHz are allocated to the fixed service on a primary basis until 1 April 2007, subject to application of the procedure referred to in Resolution 21 (Rev.WRC-95). After 1 April 2007, frequencies in these bands may be used by stations in the fixed service, communicating only within the boundary of the country in which they are located, on condition that harmful interference is not caused to the broadcasting service. When using frequencies in the fixed service, administrations are urged to use the minimum power required and to take account of the seasonal use of frequencies by the broadcasting service published in accordance with the Radio Regulations.

S5.147—On condition that harmful interference is not caused to the broadcasting service, frequencies in the bands 9775–9900 kHz, 11650–11700 kHz, and 11975–12050 kHz may be used by stations in the fixed service communicating only within the boundary of the country in which they are located, each station using a total radiated power not exceeding 24 dBW.

S5.148—SUP - WRC-97.

S5.149—In making assignments to stations of other services to which the bands are allocated:

13360–13410 kHz,
25550–25670 kHz,

37.5–38.25 MHz.

73–74.6 MHz in Regions 1 and 3,

150.05–153 MHz in Region 1,

322–328.6 MHz*,

406.1–410 MHz,

608–614 MHz in Regions 1 and 3,

1330–1400 MHz*,

1610.6–1613.8 MHz*,

1660–1670 MHz,

1718.8–1722.2 MHz*,

2655–2690 MHz,

3260–3267 MHz*,

3332–3339 MHz*,

3345.8–3352.5 MHz*,

4825–4835 MHz*,

4950–4990 MHz,

4990–5000 MHz,

6650–6675.2 MHz*,

10.6–10.68 GHz,

14.47–14.5 GHz*,

22.01–22.21 GHz*,

22.21–22.5 GHz,

22.81–22.86 GHz*,

23.07–23.12 GHz*,

31.2–31.3 GHz,

31.5–31.8 GHz in Regions 1 and 3,

36.43–36.5 GHz*,

42.5–43.5 GHz,

42.77–42.87 GHz*,

43.07–43.17 GHz*,

43.37–43.47 GHz*,

48.94–49.04 GHz*,

72.77–72.91 GHz*,

93.07–93.27 GHz*,

97.88–98.08 GHz*,

140.69–140.98 GHz*,

144.68–144.98 GHz*,

145.45–145.75 GHz*,

146.82–147.12 GHz*,

150–151 GHz*,

174.42–175.02 GHz*,

177–177.4 GHz*,

178.2–178.6 GHz*,

181–181.46 GHz*,

186.2–186.6 GHz*,

250–251 GHz*,

257.5–258 GHz*,

261–265 GHz,

262.24–262.76 GHz*,

265–275 GHz,

265.64–266.16 GHz*,

267.34–267.86 GHz*,

271.74–272.26 GHz*

(* indicates radio astronomy use for spectral line observations), administrations are urged to take all practicable steps to protect the radio astronomy service from harmful interference. Emissions from spaceborne or airborne stations can be particularly serious sources of interference to the radio

astronomy service (see Nos. S4.5 and S4.6 and Article S29) (WRC-97).

S5.150—The following bands:

13553–13567 kHz (centre frequency 13560 kHz),
 26957–27283 kHz (centre frequency 27120 kHz),
 40.66–40.70 MHz (centre frequency 40.68 MHz),
 902–928 MHz in Region 2 (centre frequency 915 MHz),
 2400–2500 MHz (centre frequency 2450 MHz),
 5725–5875 MHz (centre frequency 5800 MHz),
 and
 24–24.25 GHz (centre frequency 24.125 GHz)

are also designated for industrial, scientific, and medical (ISM) applications. Radiocommunication services operating within these bands must accept harmful interference that may be caused by these applications. ISM equipment operating in these bands is subject to the provisions of No. S15.13.

S5.151—The bands 1570–13 600 kHz and 13 800–13 870 kHz are allocated, until 1 April 2007, to the fixed service on a primary basis and to the mobile, except aeronautical mobile (R), service on a secondary basis, subject to application of the procedure referred to in Resolution 21 (Rev.WRC-95). After 1 April 2007, frequencies in these bands may be used by stations in the above-mentioned services, communicating only within the boundary of the country in which they are located, on the condition that harmful interference is not caused to the broadcasting service. When using frequencies in these services, administrations are urged to use the minimum power required and to take account of the seasonal use of frequencies by the broadcasting service published in accordance with the Radio Regulations.

S5.152—*Additional allocation:* in Armenia, Azerbaijan, China, Côte d'Ivoire, Georgia, the Islamic Republic of Iran, Kazakhstan, Moldova, Uzbekistan, Kyrgyzstan, Russian Federation, Tajikistan, Turkmenistan, and Ukraine, the band 14 250–14 350 kHz is also allocated to the fixed service on a primary basis. Stations of the fixed service shall not use a radiated power exceeding 24 dBW (WRC-97).

S5.153—In Region 3, the stations of those services to which the band 15 995–16 005 kHz is allocated may transmit standard frequency and time signals.

S5.154—*Additional allocation:* in Armenia, Azerbaijan, Georgia, Kazakhstan, Moldova, Uzbekistan, Kyrgyzstan, Russian Federation, Tajikistan, Turkmenistan, and Ukraine, the band 18068–18168 kHz is also allocated to the fixed service on a primary basis for use within their boundaries, with a peak envelope power not exceeding 1 kW (WRC-97).

S5.155—*Additional allocation:* in Armenia, Azerbaijan, Belarus, Bulgaria, Georgia, Hungary, Kazakhstan, Moldova, Mongolia, Uzbekistan, Kyrgyzstan, Slovakia, the Czech Republic, Russian Federation, Tajikistan, Turkmenistan, and Ukraine, the band 21 850–21 870 kHz is also allocated to the aeronautical mobile (R) services on a primary basis.

S5.155A—In Armenia, Azerbaijan, Belarus, Bulgaria, Georgia, Hungary, Kazakhstan, Moldova, Mongolia, Uzbekistan, Kyrgyzstan, Slovakia, the Czech Republic, Russian Federation, Tajikistan, Turkmenistan, and Ukraine, the use of the band 21 850–21 870 kHz by the fixed service is limited to provision of services related to aircraft flight safety.

S5.155B—The band 21 870–21 924 kHz is used by the fixed service for provision of services related to aircraft flight safety.

S5.156—*Additional allocation:* in Nigeria, the band 22 720–23 200 kHz is also allocated to the meteorological aids service (radiosondes) on a primary basis.

S5.156A—The use of the band 23 200–23 350 kHz by the fixed service is limited to provision of services related to aircraft flight safety.

S5.157—The use of the band 23 350–24 000 kHz by the maritime mobile service is limited to inter-ship radiotelegraphy.

S5.158 and S5.159—Not used.

S5.160—*Additional allocation:* in Botswana, Burundi, Lesotho, Malawi, Namibia, Democratic Republic of the Congo, Rwanda, and Swaziland, the band 41–44 MHz is also allocated to the aeronautical radionavigation service on a primary basis (WRC-97).

S5.161—*Additional allocation:* in the Islamic Republic of Iran and Japan, the band 41–44 MHz is also allocated to the radiolocation service on a secondary basis.

S5.162—*Additional allocation:* in Australia and New Zealand, the band 44–47 MHz is also allocated to the broadcasting service on a primary basis.

S5.162A—*Additional allocation:* in Germany, Austria, Belgium, Bosnia and Herzegovina, China, Vatican, Denmark, Spain, Estonia, Finland, France, Ireland, Iceland, Italy, Latvia, the former Yugoslav Republic of Macedonia, Liechtenstein, Lithuania, Luxembourg, Moldova, Monaco, Norway, the Netherlands, Poland, Portugal, Slovakia, the Czech Republic, the United Kingdom, Russian Federation, Sweden, Switzerland, and Turkey, the band 46–68 MHz is also allocated to the radiolocation service on a secondary basis. This use is limited to the operation of wind profiler radars in accordance with Resolution 217 (WRC-97), (WRC-97).

S5.163—*Additional allocation:* in Armenia, Azerbaijan, Belarus, Estonia, Georgia, Hungary, Kazakhstan, Latvia, Lithuania, Moldova, Mongolia, Uzbekistan, Kyrgyzstan, Slovakia, the Czech Republic, Russian Federation, Tajikistan, Turkmenistan, and Ukraine, the bands 47–48.5 MHz and 56.5–58 MHz are also allocated to the fixed and land mobile service on a secondary basis.

S5.164—*Additional allocation:* in Albania, Germany, Austria, Belgium, Bosnia and Herzegovina, Bulgaria, Côte d'Ivoire, Denmark, Spain, Finland, France, Gabon, Greece, Ireland, Israel, Italy, Jordan, Lebanon, Libya, Liechtenstein, Luxembourg, Madagascar, Mali, Malta, Morocco, Mauritania, Monaco, Nigeria, Norway, the Netherlands, Poland, Syria, the United Kingdom, Senegal, Slovenia, Sweden, Switzerland, Swaziland, Togo, Tunisia, Turkey, and Yugoslavia the band 47–68 MHz, in Romania the band 47–58 MHz, and in the Czech Republic the band 66–68 MHz are also allocated to the land mobile service on a primary basis. However, stations of the land mobile service in the countries mentioned in connection with each band referred to in this footnote shall not cause harmful interference to, or claim protection from, existing or planned broadcasting stations of countries other than those mentioned in connection with the band (WRC-97).

S5.165—*Additional allocation:* in Angola, Cameroon, the Congo, Madagascar, Mozambique, Somalia, Sudan, Tanzania,

and Chad, the band 47–68 MHz is also allocated to the fixed and mobile, except aeronautical mobile, services on a primary basis.

S5.166—*Alternative allocation*: in New Zealand, the band 50–51 MHz is allocated to the fixed, mobile, and broadcasting services on a primary basis; the band 53–54 MHz is allocated to the fixed and mobile services on a primary basis.

S5.167—*Alternative allocation*: in Bangladesh, Brunei Darussalam, India, Indonesia, the Islamic Republic of Iran, Malaysia, Pakistan, Singapore, and Thailand, the band 50–54 MHz is allocated to the fixed, mobile, and broadcasting services on a primary basis.

S5.168—*Additional allocation*: in Australia, China, and the Democratic People's Republic of Korea, the band 50–54 MHz is also allocated to the broadcasting service on a primary basis.

S5.169—*Alternative allocation*: in Botswana, Burundi, Lesotho, Malawi, Namibia, Democratic Republic of the Congo, Rwanda, South Africa, Swaziland, Zambia, and Zimbabwe, the band 50–54 MHz is allocated to the amateur service on a primary basis.

S5.170—*Additional allocation*: in New Zealand, the band 51–53 MHz is also allocated to the fixed and mobile services on a primary basis.

S5.171—*Additional allocation*: in Botswana, Burundi, Lesotho, Malawi, Mali, Namibia, Democratic Republic of the Congo, Rwanda, South Africa, Swaziland, and Zimbabwe, the band 54–68 MHz is also allocated to the fixed and mobile, except aeronautical mobile, services on a primary basis.

S5.172—*Different category of service*: in the French Overseas Departments in Region 2, Guyana, Jamaica, and Mexico, the allocation of the band 54–68 MHz to the fixed and mobile services is on a primary basis (see No. S5.33).

S5.173—*Different category of service*: in the French Overseas Departments in Region 2, Guyana, Jamaica, and Mexico, the allocation of the band 68–72 MHz to the fixed and mobile services is on a primary basis (see No. S5.33).

S5.174—*Alternative allocation*: in Bulgaria, Hungary, Poland, and Romania, the band 68–73 MHz is allocated to the broadcasting service on a primary basis and used in accordance with the decisions in the Final Acts of the Special Regional Conference (Geneva, 1960) (WRC-97).

S5.175—*Alternative allocation*: in Armenia, Azerbaijan, Belarus, Estonia, Georgia, Kazakhstan, Latvia, Lithuania, Moldova, Mongolia, Uzbekistan, Kyrgyzstan, Russian Federation, Tajikistan, Turkmenistan, and Ukraine, the bands 68–73 MHz and 76–87.5 MHz are allocated to the broadcasting service on a primary basis. The services to which these bands are allocated in other countries and the broadcasting service in the countries listed above are subject to agreements with the neighbouring countries concerned.

S5.176—*Additional allocation*: in Australia, China, the Republic of Korea, the Philippines, the Democratic People's Republic of Korea, and Western Samoa, the band 68–74 MHz is also allocated to the broadcasting service on a primary basis.

S5.177—*Additional allocation*: in Armenia, Azerbaijan, Belarus, Bulgaria, Estonia, Georgia, Kazakhstan, Latvia, Lithuania, Moldova, Mongolia, Uzbekistan, Poland, Kyrgyzstan, Russian Federation, Tajikistan, Turkmenistan, and

Ukraine, the band 73–74 MHz is also allocated to the broadcasting service on a primary basis, subject to agreement obtained under No. S9.21 (WRC-97).

S5.178—*Additional allocation*: in Colombia, Costa Rica, Cuba, El Salvador, Guatemala, Guyana, Honduras, and Nicaragua, the band 73–74.6 MHz is also allocated to the fixed and mobile services on a secondary basis.

S5.179—*Additional allocation*: in Armenia, Azerbaijan, Belarus, Bulgaria, China, Georgia, Kazakhstan, Latvia, Lithuania, Moldova, Mongolia, Kyrgyzstan, Slovakia, the Czech Republic, Russian Federation, Tajikistan, Turkmenistan, and Ukraine, the bands 74.6–74.8 MHz and 75.2–75.4 MHz are also allocated to the aeronautical radionavigation service, on a primary basis, for ground-based transmitters only.

S5.180—The frequency 75 MHz is assigned to marker beacons. Administrations shall refrain from assigning frequencies close to the limits of the guardband to stations of other services that, because of their power or geographical position, might cause harmful interference or otherwise place a constraint on marker beacons.

Every effort should be made to improve further the characteristics of airborne receivers and to limit the power of transmitting stations close to the limits 74.8 MHz and 75.2 MHz.

S5.181—*Additional allocation*: in Germany, Austria, Cyprus, Denmark, Egypt, France, Greece, Israel, Italy, Japan, Jordan, Lebanon, Malta, Morocco, Monaco, Norway, Syria, Sweden, and Switzerland, the band 74.8–75.2 MHz is also allocated to the mobile service on a secondary basis, subject to agreement obtained under No. S9.21. To ensure that harmful interference is not caused to stations of the aeronautical radionavigation service, stations of the mobile service shall not be introduced in the band until it is no longer required for the aeronautical radionavigation service by any administration that may be identified in the application of the procedure invoked under No. S9.21 (WRC-97).

S5.182—*Additional allocation*: in Western Samoa, the band 75.4–87 MHz is also allocated to the broadcasting service on a primary basis.

S5.183—*Additional allocation*: in China, the Republic of Korea, Japan, the Philippines, and the Democratic People's Republic of Korea, the band 76–87 MHz is also allocated to the broadcasting service on a primary basis.

S5.184—*Additional allocation*: in Bulgaria and Romania, the band 76–87.5 MHz is also allocated to the broadcasting service on a primary basis and used in accordance with the decisions contained in the Final Acts of the Special Regional Conference (Geneva, 1960) (WRC-97).

S5.185—*Different category of service*: in the United States, the French Overseas Departments in Region 2, Guyana, Jamaica, Mexico, and Paraguay, the allocation of the band 76–88 MHz to the fixed and mobile services is on a primary basis (see No. S5.33).

S5.186—SUP - WRC-97.

S5.187—*Alternative allocation*: in Albania, the band 81–87.5 MHz is allocated to the broadcasting service on a primary basis and used in accordance with the decisions contained in the Final Acts of the Special Regional Conference (Geneva, 1960).

S5.188—*Additional allocation:* in Australia, the band 85–87 MHz is also allocated to the broadcasting service on a primary basis. The introduction of the broadcasting service in Australia is subject to special agreements between the administrations concerned.

S5.189—Not used.

S5.190—*Additional allocation:* in Monaco, the band 87.5–88 MHz is also allocated to the land mobile service on a primary basis, subject to agreement obtained under No. **S9.21** (WRC-97).

S5.191—Not used.

S5.192—*Additional allocation:* in China and the Republic of Korea, the band 100–108 MHz is also allocated to the fixed and mobile services on a primary basis (WRC-97).

S5.193—Not used.

S5.194—*Additional allocation:* in Azerbaijan, Lebanon, Syria, Kyrgyzstan, Somalia, and Turkmenistan, the band 104–108 MHz is also allocated to the mobile, except aeronautical mobile (R), service on a secondary basis (WRC-97).

S5.195 and **S5.196**—Not used.

S5.197—*Additional allocation:* in Germany, Austria, Cyprus, Denmark, Egypt, France, Italy, Japan, Jordan, Lebanon, Malta, Morocco, Monaco, Norway, Pakistan, Syria, and Sweden, the band 108–111.975 MHz is also allocated to the mobile service on a secondary basis, subject to agreement obtained under No. **S9.21**. To ensure that harmful interference is not caused to stations of the aeronautical radionavigation service, stations of the mobile service shall not be introduced in the band until it is no longer required for the aeronautical radionavigation service by any administration that may be identified in the application of the procedures invoked under No. **S9.21** (WRC-97).

S5.198—*Additional allocation:* The band 117.975–136 MHz is also allocated to the aeronautical mobile-satellite (R) service on a secondary basis, subject to agreement obtained under No. **S9.21** (WRC-97).

S5.199—The bands 121.45–121.55 MHz and 242.95–243.05 MHz are also allocated to the mobile-satellite service for the reception on board satellites of emissions from emergency position-indicating radiobeacons transmitting at 121.5 MHz and 243 MHz (see Appendix **S13**).

S5.200—In the band 117.975–136 MHz, the frequency 121.5 MHz is the aeronautical emergency frequency and, where required, the frequency 123.1 MHz is the aeronautical frequency auxiliary to 121.5 MHz. Mobile stations of the maritime mobile service may communicate on these frequencies under the conditions laid down in Article **S31** and Appendix **S13** for distress and safety purposes with stations of the aeronautical mobile service.

S5.201—*Additional allocation:* in Angola, Armenia, Azerbaijan, Belarus, Bulgaria, Estonia, Georgia, Hungary, the Islamic Republic of Iran, Iraq, Japan, Kazakhstan, Latvia, Moldova, Mongolia, Mozambique, Uzbekistan, Papua New Guinea, Poland, Kyrgyzstan, Slovakia, the Czech Republic, Romania, Russian Federation, Tajikistan, Turkmenistan, and Ukraine, the band 132–136 MHz is also allocated to the aeronautical mobile (OR) service on a primary basis. In assigning frequencies to stations of the aeronautical mobile (OR) service, the administration shall take account of the frequencies

assigned to stations in the aeronautical mobile (R) service (WRC-97).

S5.202—*Additional allocation:* in Saudi Arabia, Armenia, Azerbaijan, Belarus, Bulgaria, United Arab Emirates, Georgia, the Islamic Republic of Iran, Jordan, Kazakhstan, Latvia, Moldova, Oman, Uzbekistan, Poland, Syria, Kyrgyzstan, Slovakia, the Czech Republic, Romania, Russian Federation, Tajikistan, Turkmenistan, Turkey, and Ukraine, the band 136–137 MHz is also allocated to the aeronautical mobile (OR) service on a primary basis. In assigning frequencies to stations of the aeronautical mobile (OR) service, the administration shall take account of the frequencies assigned to stations in the aeronautical mobile (R) service (WRC-97).

S5.203—In the band 136–137 MHz, existing operational meteorological satellites may continue to operate, under the conditions defined in No. **S4.4** with respect to the aeronautical mobile service, until January 1, 2002. Administrations shall not authorize new frequency assignments in this band to stations in the meteorological-satellite service (WRC-97).

S5.203A—*Additional allocation:* in Israel, Mauritania, Qatar, and Zimbabwe, the band 136–137 MHz is also allocated to the fixed and mobile, except aeronautical mobile (R), services on a secondary basis until January 1, 2005, (WRC-97).

S5.203B—*Additional allocation:* in Saudi Arabia, United Arab Emirates, Jordan, Oman, and Syria, the band 136–137 MHz is also allocated to the fixed and mobile, except aeronautical mobile, services on a secondary basis until January 1, 2005 (WRC-97).

S5.204—*Different category of service:* in Afghanistan, Saudi Arabia, Bahrain, Bangladesh, Bosnia and Herzegovina, Brunei Darussalam, China, Cuba, the United Arab Emirates, India, Indonesia, the Islamic Republic of Iran, Iraq, Malaysia, Oman, Pakistan, Philippines, Qatar, Singapore, Sri Lanka, Thailand, Yemen, and Yugoslavia, the band 137–138 MHz is allocated to the fixed and mobile, except aeronautical mobile (R), services on a primary basis (see No. **S5.33**).

S5.205—*Different category of service:* in Israel and Jordan, the allocation of the band 137–138 MHz to the fixed and mobile, except aeronautical mobile, services is on a primary basis (see No. **S5.33**).

S5.206—*Different category of service:* in Armenia, Austria, Azerbaijan, Belarus, Bulgaria, Egypt, Finland, France, Georgia, Greece, Hungary, Kazakhstan, Lebanon, Moldova, Mongolia, Uzbekistan, Poland, Kyrgyzstan, Syria, Slovakia, the Czech Republic, Romania, Russian Federation, Tajikistan, Turkmenistan, and Ukraine, the allocation of the band 137–138 MHz to the aeronautical mobile (OR) service is on a primary basis (see No. **S5.33**).

S5.207—*Additional allocation:* in Australia, the band 137–144 MHz is also allocated to the broadcasting service on a primary basis until that service can be accommodated within regional broadcasting allocations.

S5.208—The use of the band 137–138 MHz by the mobile-satellite service is subject to coordination under No. **S9.11A** (WRC-97).

S5.208A—In making assignments to space stations in the mobile-satellite service in the bands 137–138 MHz, 387–390 MHz, and 400.15–401 MHz, administrations shall take all practicable steps to protect the radio astronomy service in the

bands 150.05–153 MHz, 322–328.6 MHz, 406.1–410 MHz, and 608–614 MHz from harmful interference from unwanted emissions. The threshold levels of interference detrimental to the radio astronomy service are shown in Table 1 of Recommendation ITU-R RA. 769-1 (WRC-97).

S5.209—The use of the bands 137–138 MHz, 148–150.05 MHz, 399.9–400.05 MHz, 400.15–401 MHz, 454–456 MHz, and 459–460 MHz by the mobile-satellite service is limited to nongeostationary-satellite systems (WRC-97).

S5.210—*Additional allocation:* in Austria, France, Italy, Liechtenstein, Slovakia, the Czech Republic, the United Kingdom, and Switzerland, the bands 138–143.6 MHz and 143.65–144 MHz are also allocated to the space research service (space-to-Earth) on a secondary basis (WRC-97).

S5.211—*Additional allocation:* in Germany, Saudi Arabia, Austria, Bahrain, Belgium, Bosnia and Herzegovina, Denmark, the United Arab Emirates, Spain, Finland, Greece, Ireland, Israel, Kenya, Kuwait, the former Yugoslav Republic of Macedonia, Liechtenstein, Luxembourg, Mali, Malta, Norway, the Netherlands, Qatar, the United Kingdom, Slovenia, Somalia, Sweden, Switzerland, Tanzania, Tunisia, Turkey, and Yugoslavia, the band 138–144 MHz is also allocated to the maritime mobile and land mobile services on a primary basis.

S5.212—*Alternative allocation:* in Angola, Botswana, Burundi, Cameroon, the Central African Republic, the Congo, Gabon, Gambia, Ghana, Guinea, Iraq, Jordan, Lesotho, Liberia, Libya, Malawi, Mozambique, Namibia, Nigeria, Oman, Democratic Republic of the Congo, Rwanda, Sierra Leone, South Africa, Swaziland, Chad, Togo, Zaire, Zambia, and Zimbabwe, the band 138–144 MHz is allocated to the fixed and mobile services on a primary basis.

S5.213—*Additional allocation:* in China, the band 138–144 MHz is also allocated to the radiolocation service on a primary basis.

S5.214—*Additional allocation:* in Bosnia and Herzegovina, Croatia, Eritrea, Ethiopia, Kenya, the former Yugoslav Republic of Macedonia, Malta, Slovenia, Somalia, Sudan, Tanzania, and Yugoslavia, the band 138–144 MHz is also allocated to the fixed service on a primary basis.

S5.215—Not used.

S5.216—*Additional allocation:* in China, the band 144–146 MHz is also allocated to the aeronautical mobile (OR) service on a secondary basis.

S5.217—*Alternative allocation:* in Afghanistan, Bangladesh, Cuba, Guyana, and India, the band 146–148 MHz is allocated to the fixed and mobile services on a primary basis.

S5.218—*Additional allocation:* The band 148–149.9 MHz is also allocated to the space operation service (Earth-to-space) on a primary basis, subject to agreement obtained under No. **S9.21**. The bandwidth of any individual transmission shall not exceed ± 25 kHz.

S5.219—The use of the band 148–149.9 MHz by the mobile-satellite service is subject to coordination under No. **S9.11A**. The mobile-satellite service shall not constrain the development and use of the fixed, mobile, and space operation services in the band 148–149.9 MHz.

S5.220—The use of the bands 149.9–150.05 MHz and 399.9–400.05 MHz by the mobile-satellite service is subject to coordination under No. **S9.11A**. The mobile-satellite ser-

vice shall not constrain the development and use of the radionavigation-satellite service in the bands 149.9–150.05 MHz and 399.9–400.05 MHz (WRC-97).

S5.221—Stations of the mobile-satellite service in the band 148–149.9 MHz shall not cause harmful interference to, or claim protection from, stations of the fixed or mobile services operating in accordance with the Table of Frequency Allocations in the following countries: Albania, Algeria, Germany, Saudi Arabia, Australia, Austria, Bahrain, Bangladesh, Barbados, Belarus, Belgium, Benin, Bosnia and Herzegovina, Brunei Darussalam, Bulgaria, Cameroon, China, Cyprus, Congo, the Republic of Korea, Croatia, Cuba, Denmark, Egypt, the United Arab Emirates, Eritrea, Spain, Estonia, Ethiopia, Finland, France, Gabon, Ghana, Greece, Guinea, Guinea Bissau, Hungary, India, the Islamic Republic of Iran, Ireland, Iceland, Israel, Italy, Jamaica, Japan, Jordan, Kazakhstan, Kenya, Kuwait, Latvia, The Former Yugoslav Republic of Macedonia, Lebanon, Libya, Liechtenstein, Luxembourg, Malaysia, Mali, Malta, Mauritania, Moldova, Mongolia, Mozambique, Namibia, Norway, New Zealand, Oman, Uganda, Uzbekistan, Pakistan, Panama, Papua New Guinea, Paraguay, the Netherlands, the Philippines, Poland, Portugal, Qatar, Syria, Kyrgyzstan, Slovakia, Romania, the United Kingdom, Russian Federation, Senegal, Sierra Leone, Singapore, Slovenia, Sri Lanka, South Africa, Sweden, Switzerland, Swaziland, Tanzania, Chad, Thailand, Togo, Tonga, Trinidad and Tobago, Tunisia, Turkey, Ukraine, Viet Nam, Yemen, Yugoslavia, Zambia, and Zimbabwe (WRC-97).

S5.222—Emissions of the radionavigation-satellite service in the bands 149.9–150.05 MHz and 399.9–400.05 MHz may also be used by receiving earth stations of the space research service.

S5.223—Recognizing that the use of the band 149.9–150.05 MHz by the fixed and mobile services may cause harmful interference to the radionavigation-satellite service, administrations are urged not to authorize such use in application of No. **S4.4**.

S5.224—SUP-WRC-97.

S5.224A—The use of the bands 149.9–150.05 MHz and 399.9–400.05 MHz by the mobile-satellite service (Earth-to-space) is limited to the land mobile-satellite service (Earth-to-space) until January 1, 2015 (WRC-97).

S5.224B—The allocation of the bands 149.9–150.05 MHz and 399.9–400.05 MHz to the radionavigation-satellite service shall be effective until January 1, 2015 (WRC-97).

S5.225—*Additional allocation:* in Australia and India, the band 150.05–153 MHz is also allocated to the radio astronomy service on a primary basis.

S5.226—The frequency 156.8 MHz is the international distress, safety, and calling frequency for the maritime mobile VHF radiotelephone service. The conditions for the use of this frequency are contained in Article **S31** and Appendix **S13**.

In the bands 156–156.7625 MHz, 156.8375–157.45 MHz, 160.6–160.975 MHz, and 161.475–162.05 MHz, each administration shall give priority to the maritime mobile service on only such frequencies as are assigned to stations of the maritime mobile service by the administration (see Articles **S31** and **S52** and Appendix **S13**).

Any use of frequencies in these bands by stations of other services to which they are allocated should be avoided in areas

where such use might cause harmful interference to the maritime mobile VHF radiocommunication service.

However, the frequency 156.8 MHz and the frequency bands in which priority is given to the maritime mobile service may be used for radiocommunications on inland waterways subject to agreement between interested and affected administrations and taking into account current frequency usage and existing agreements.

S5.227—In the maritime mobile VHF service the frequency 156.525 MHz is to be used exclusively for digital selective calling for distress, safety, and calling. The conditions for the use of this frequency are prescribed in Articles S31 and S52 and in Appendices S13 and S18.

S5.228—Not used.

S5.229—*Alternative allocation*: in Morocco, the band 162–174 MHz is allocated to the broadcasting service on a primary basis. The use of this band shall be subject to agreement with administrations having services, operating or planned, in accordance with the Table which are likely to be affected. Stations in existence on 1 January 1981, with their technical characteristics as of that date, are not affected by such agreement.

S5.230—*Additional allocation*: in China, the band 163–167 MHz is also allocated to the space operation service (space-to-Earth) on a primary basis, subject to agreement obtained under No. S9.21.

S5.231—*Additional allocation*: in Afghanistan, China, and Pakistan, the band 167–174 MHz is also allocated to the broadcasting service on a primary basis. The introduction of the broadcasting service into this band shall be subject to agreement with the neighbouring countries in Region 3 whose services are likely to be affected.

S5.232—*Additional allocation*: in Japan, the band 170–174 MHz is also allocated to the broadcasting service on a primary basis.

S5.233—*Additional allocation*: in China, the band 174–184 MHz is also allocated to the space research (space-to-Earth) and the space operation (space-to-Earth) services on a primary basis, subject to agreement obtained under No. S9.21. These services shall not cause harmful interference to, or claim protection from, existing or planned broadcasting stations.

S5.234—*Different category of service*: in Mexico, the allocation of the band 174–216 MHz to the fixed and mobile services is on a primary basis (see No. S5.33).

S5.235—*Additional allocation*: in Germany, Austria, Belgium, Denmark, Spain, Finland, France, Israel, Italy, Liechtenstein, Malta, Monaco, Norway, the Netherlands, the United Kingdom, Sweden, and Switzerland, the band 174–223 MHz is also allocated to the land mobile service on a primary basis. However, the stations of the land mobile service shall not cause harmful interference to, or claim protection from, broadcasting stations, existing or planned, in countries other than those listed in this footnote.

S5.236—Not used.

S5.237—*Additional allocation*: in the Congo, Eritrea, Ethiopia, Gambia, Guinea, Libya, Malawi, Mali, Senegal, Sierra Leone, Somalia, Tanzania, and Zimbabwe, the band 174–223 MHz is also allocated to the fixed and mobile services on a secondary basis, (WRC-97).

S5.238—*Additional allocation*: in Bangladesh, India, Pakistan, and the Philippines, the band 200–216 MHz is also allocated to the aeronautical radionavigation service on a primary basis.

S5.239—Not used.

S5.240—*Additional allocation*: in China and India, the band 216–223 MHz is also allocated to the aeronautical radionavigation service on a primary basis and to the radiolocation service on a secondary basis.

S5.241—In Region 2, no new stations in the radiolocation service may be authorized in the band 216–225 MHz. Stations authorized prior to 1 January 1990 may continue to operate on a secondary basis.

S5.242—*Additional allocation*: in Canada, the band 216–220 MHz is also allocated to the land mobile service on a primary basis.

S5.243—*Additional allocation*: in Somalia, the band 216–225 MHz is also allocated to the aeronautical radionavigation service on a primary basis, subject to not causing harmful interference to existing or planned broadcasting services in other countries.

S5.244—SUP-WRC-97.

S5.245—*Additional allocation*: in Japan, the band 222–223 MHz is also allocated to the aeronautical radionavigation service on a primary basis and to the radiolocation service on a secondary basis.

S5.246—*Alternative allocation*: in Spain, France, Israel, and Monaco, the band 223–230 MHz is allocated to the broadcasting and land mobile services on a primary basis (see No. S5.33) on the basis that, in the preparation of frequency plans, the broadcasting service shall have prior choice of frequencies; and allocated to the fixed and mobile, except land mobile, services on a secondary basis. However, the stations of the land mobile service shall not cause harmful interference to, or claim protection from, existing or planned broadcasting stations in Morocco and Algeria.

S5.247—*Additional allocation*: in Saudi Arabia, Bahrain, the United Arab Emirates, Jordan, Oman, Qatar, and Syria, the band 223–235 MHz is also allocated to the aeronautical radionavigation service on a primary basis.

S5.248 and **S5.249**—Not used.

S5.250—*Additional allocation*: in China, the band 225–235 MHz is also allocated to the radio astronomy service on a secondary basis.

S5.251—*Additional allocation*: in Nigeria, the band 230–235 MHz is also allocated to the aeronautical radionavigation service on a primary basis, subject to agreement obtained under No. S9.21.

S5.252—*Alternative allocation*: in Botswana, Lesotho, Malawi, Mozambique, Namibia, South Africa, Swaziland, Zambia, and Zimbabwe, the bands 230–238 MHz and 246–254 MHz are allocated to the broadcasting service on a primary basis, subject to agreement obtained under No. S9.21.

S5.253—Not used.

S5.254—The bands 235–322 MHz and 335.4–399.9 MHz may be used by the mobile-satellite service, subject to agreement obtained under No. S9.21, on condition that stations in this service do not cause harmful interference to those of other services operating or planned to be operated in accordance with the Table of Frequency Allocations.

S5.255—The bands 312–315 MHz (Earth-to-space) and 387–390 MHz (space-to-Earth) in the mobile-satellite service may also be used by nongeostationary-satellite systems. Such use is subject to coordination under No. **S9.11A**.

S5.256—The frequency 243 MHz is the frequency in this band for use by survival craft stations and equipment used for survival purposes (see Appendix **S13**).

S5.257—The band 267–272 MHz may be used by administrations for space telemetry in their countries on a primary basis, subject to agreement obtained under No. **S9.21**.

S5.258—The use of the band 328.6–335.4 MHz by the aeronautical radionavigation service is limited to Instrument Landing Systems (glide path).

S5.259—*Additional allocation:* in Germany, Austria, Cyprus, the Republic of Korea, Denmark, Egypt, Spain, France, Greece, Israel, Italy, Japan, Jordan, Malta, Morocco, Monaco, Norway, the Netherlands, Syria, and Sweden, the band 328.6–335.4 MHz is also allocated to the mobile service on a secondary basis, subject to agreement obtained under No. **S9.21**. To ensure that harmful interference is not caused to stations of the aeronautical radionavigation service, stations of the mobile service shall not be introduced in the band until it is no longer required for the aeronautical radionavigation service by any administration that may be identified in the application of the procedure invoked under No. **S9.21** (WRC-97).

S5.260—Recognizing that the use of the band 399.9–400.05 MHz by the fixed and mobile services may cause harmful interference to the radionavigation satellite service, administrations are urged not to authorize such use in application of No. **S4.4**.

S5.261—Emissions shall be confined in a band of ± 25 kHz about the standard frequency 400.1 MHz.

S5.262—*Additional allocation:* in Saudi Arabia, Armenia, Azerbaijan, Bahrain, Belarus, Bosnia and Herzegovina, Bulgaria, Colombia, Costa Rica, Cuba, Egypt, the United Arab Emirates, Ecuador, Estonia, Georgia, Hungary, Indonesia, the Islamic Republic of Iran, Iraq, Israel, Jordan, Kazakhstan, Kuwait, Liberia, Malaysia, Moldova, Nigeria, Uzbekistan, Pakistan, the Philippines, Qatar, Syria, Kyrgyzstan, Slovakia, Romania, Russian Federation, Singapore, Somalia, Sri Lanka, Tajikistan, Turkmenistan, Ukraine, and Yugoslavia, the band 400.05–401 MHz is also allocated to the fixed and mobile services on a primary basis.

S5.263—The band 400.15–401 MHz is also allocated to the space research service in the space-to-space direction for communications with manned space vehicles. In this application, the space research service will not be regarded as a safety service.

S5.264—The use of the band 400.15–401 MHz by the mobile-satellite service is subject to coordination under No. **S9.11A**. The power flux-density limit indicated in Annex I of Appendix **S5** shall apply until such time as a competent world radiocommunication conference revises it.

S5.265—Not used.

S5.266—The use of the band 406–406.1 MHz by the mobile-satellite service is limited to low-power satellite emergency position-indicating radiobeacons (see also Article **S31** and Appendix **S13**).

S5.267—Any emission capable of causing harmful interference to the authorized uses of the band 406–406.1 MHz is prohibited.

S5.268—Use of the band 410–420 MHz by the space research service is limited to communications within 5 km of an orbiting, manned space vehicle. The power flux-density at the surface of the Earth produced by emissions from extra-vehicular activities shall not exceed -153 dB(W/m²) for $0^\circ \leq \delta \leq 5^\circ$, $-153 + 0.077(\delta - 5)$ dB(W/m²) for $5^\circ \leq \delta \leq 70^\circ$ and -148 dB(W/m²) for $70^\circ \leq \delta \leq 90^\circ$, where δ is the angle of arrival of the radio-frequency wave and the reference bandwidth is 4 kHz. No. **S4.10** does not apply to extra-vehicular activities. In this frequency band the space research (space-to-space) service shall not claim protection from, nor constrain the use and development of, stations of the fixed and mobile services (WRC-97).

S5.269—*Different category of service:* in Australia, the United States, India, Japan, and the United Kingdom, the allocation of the bands 420–430 MHz and 440–450 MHz to the radiolocation service is on a primary basis (see No. **S5.33**).

S5.270—*Additional allocation:* in Australia, the United States, Jamaica, and the Philippines, the bands 420–430 MHz and 440–450 MHz are also allocated to the amateur service on a secondary basis.

S5.271—*Additional allocation:* in Azerbaijan, Belarus, China, Estonia, India, Latvia, Lithuania, Kyrgyzstan, Turkmenistan, and Ukraine, the band 420–460 MHz is also allocated to the aeronautical radionavigation service (radio altimeters) on a secondary basis (WRC-97).

S5.272—*Different category of service:* in France, the allocation of the band 430–434 MHz to the amateur service is on a secondary basis (see No. **S5.32**).

S5.273—*Different category of service:* in Denmark, Libya, and Norway, the allocation of the bands 430–432 MHz and 438–440 MHz to the radiolocation service is on a secondary basis (see No. **S5.32**).

S5.274—*Alternative allocation:* in Denmark, Norway, and Sweden, the bands 430–432 MHz and 438–440 MHz are allocated to the fixed and mobile, except aeronautical mobile, services on a primary basis.

S5.275—*Additional allocation:* in Bosnia and Herzegovina, Croatia, Estonia, Finland, Latvia, the former Yugoslav Republic of Macedonia, Libya, Slovenia, and Yugoslavia, the bands 430–432 MHz and 438–440 MHz are also allocated to the fixed and mobile, except aeronautical mobile, services on a primary basis (WRC-97).

S5.276—*Additional allocation:* in Afghanistan, Algeria, Saudi Arabia, Bahrain, Bangladesh, Brunei Darussalam, Burkina Faso, Burundi, Egypt, the United Arab Emirates, Ecuador, Eritrea, Ethiopia, Greece, Guinea, India, Indonesia, the Islamic Republic of Iran, Iraq, Israel, Italy, Jordan, Kenya, Kuwait, Lebanon, Libya, Liechtenstein, Malaysia, Malta, Nigeria, Oman, Pakistan, the Philippines, Qatar, Syria, Democratic People's Republic of Korea, Singapore, Somalia, Switzerland, Tanzania, Thailand, Togo, Turkey, and Yemen, the band 430–440 MHz is also allocated to the fixed service on a primary basis and the bands 430–435 MHz and 438–440 MHz are also allocated to the mobile, except aeronautical mobile, service on a primary basis (WRC-97).

S5.277—*Additional allocation:* in Angola, Armenia, Azerbaijan, Belarus, Cameroon, the Congo, Djibouti, Gabon, Georgia, Hungary, Kazakhstan, Latvia, Mali, Moldova, Mongolia, Uzbekistan, Pakistan, Poland, Kyrgyzstan, Slovakia, the Czech Republic, Romania, Russian Federation, Rwanda, Tajikistan, Chad, Turkmenistan, and Ukraine, the band 430–440 MHz is also allocated to the fixed service on a primary basis (WRC-97).

S5.278—*Different category of service:* in Argentina, Colombia, Costa Rica, Cuba, Guyana, Honduras, Panama, and Venezuela, the allocation of the band 430–440 MHz to the amateur service is on a primary basis (see No. S5.33).

S5.279—*Additional allocation:* in Mexico, the bands 430–435 MHz and 438–440 MHz are also allocated on a primary basis to the land mobile service, subject to agreement obtained under No. S9.21.

S5.280—In Germany, Austria, Bosnia and Herzegovina, Croatia, The Former Yugoslav Republic of Macedonia, Liechtenstein, Portugal, Slovenia, Switzerland, and Yugoslavia, the band 433.05–434.79 MHz (centre frequency 433.92 MHz) is designated for industrial, scientific, and medical (ISM) applications. Radiocommunication services of these countries operating within this band must accept harmful interference which may be caused by these applications. ISM equipment operating in this band is subject to the provisions of No. S15.13.

S5.281—*Additional allocation:* in the French Overseas Departments in Region 2 and India, the band 433.75–434.25 MHz is also allocated to the space operation service (Earth-to-space) on a primary basis. In France and in Brazil, the band is allocated to the same service on a secondary basis.

S5.282—In the bands 435–438 MHz, 1260–1270 MHz, 2400–2450 MHz, 3400–3410 MHz (in Regions 2 and 3 only), and 5650–5670 MHz, the amateur-satellite service may operate subject to not causing harmful interference to other services operating in accordance with the Table (see No. S5.43). Administrations authorizing such use shall ensure that any harmful interference caused by emissions from a station in the amateur-satellite service is immediately eliminated in accordance with the provisions of No. S25.11. The use of the bands 1260–1270 MHz and 5650–5670 MHz by the amateur-satellite service is limited to the Earth-to-space direction.

S5.283—*Additional allocation:* in Austria, the band 438–440 MHz is also allocated to the fixed and mobile, except aeronautical mobile, services on a primary basis.

S5.284—*Additional allocation:* in Canada, the band 440–450 MHz is also allocated to the amateur service on a secondary basis.

S5.285—*Different category of service:* in Canada, the allocation of the band 440–450 MHz to the radiolocation service is on a primary basis (see No. S5.33).

S5.286—The band 449.75–450.25 MHz may be used for the space operation service (Earth-to-space) and the space research service (Earth-to-space), subject to agreement obtained under No. S9.21.

S5.286A—The use of the bands 454–456 MHz and 459–460 MHz by the mobile-satellite service is subject to coordination under No. S9.11A (WRC-97).

S5.286B—The use of the band 454–455 MHz in the countries listed in No. S5.286D, 455–456 MHz and 459–460 MHz

in Region 2, and 454–456 MHz and 459–460 MHz in the countries listed in No. S5.286E, by stations in the mobile-satellite service, shall not cause harmful interference to, or claim protection from, stations of the fixed or mobile services operating in accordance with the Table of Frequency Allocations (WRC-97).

S5.286C—The use of the band 454–455 MHz in the countries listed in No. S5.286D, 455–456 MHz and 459–460 MHz in Region 2, and 454–456 MHz and 459–460 MHz in the countries listed in No. S5.286E, by stations in the mobile-satellite service, shall not constrain the development and use of the fixed and mobile services operating in accordance with the Table of Frequency Allocations (WRC-97).

S5.286D—*Additional allocation:* in Canada, the United States, Mexico, and Panama, the band 454–455 MHz is also allocated to the mobile-satellite service (Earth-to-space) on a primary basis (WRC-97).

S5.286E—*Additional allocation:* in Cape Verde, Indonesia, Nepal, Nigeria, and Papua New Guinea, the bands 454–456 MHz and 459–460 MHz are also allocated to the mobile-satellite (Earth-to-space) service on a primary basis (WRC-97).

S5.287—In the maritime mobile service, the frequencies 457.525 MHz, 457.550 MHz, 457.575 MHz, 467.525 MHz, 467.550 MHz, and 467.575 MHz may be used by on-board communication stations. Where needed, equipment designed for 12.5 kHz channel spacing using also the additional frequencies 457.5375 MHz, 457.5625 MHz, 467.5375 MHz, and 467.5625 MHz may be introduced for on-board communications. The use of these frequencies in territorial waters may be subject to the national regulations of the administration concerned. The characteristics of the equipment used shall conform to those specified in Recommendation ITU-R M.1174 (see Resolution 341 (WRC-97)) (WRC-97).

S5.288—In the territorial waters of the United States and the Philippines, the preferred frequencies for use by on-board communication stations shall be 457.525 MHz, 457.550 MHz, 457.575 MHz, and 457.600 MHz paired, respectively, with 467.750 MHz, 467.775 MHz, 467.800 MHz, and 467.825 MHz. The characteristics of the equipment used shall conform to those specified in Recommendation ITU-R M.1174.

S5.289—Earth exploration-satellite service applications, other than the meteorological-satellite service, may also be used in the bands 460–470 MHz and 1690–1710 MHz for space-to-Earth transmissions subject to not causing harmful interference to stations operating in accordance with the Table.

S5.290—*Different category of service:* in Afghanistan, Armenia, Azerbaijan, Belarus, China, Japan, Kazakhstan, Mongolia, Uzbekistan, Kyrgyzstan, Slovakia, the Czech Republic, Russian Federation, Tajikistan, Turkmenistan, and Ukraine, the allocation of the band 460–470 MHz to the meteorological-satellite service (space-to-Earth) is on a primary basis (see No. S5.33), subject to agreement obtained under No. S9.21 (WRC-97).

S5.291—*Additional allocation:* in China, the band 470–485 MHz is also allocated to the space research (space-to-Earth) and the space operation (space-to-Earth) services on a primary basis subject to agreement obtained under No. S9.21

and subject to not causing harmful interference to existing and planned broadcasting stations.

S5.291A—*Additional allocation*: in Germany, Austria, Denmark, Estonia, Finland, Liechtenstein, Norway, the Netherlands, the Czech Republic and Switzerland, the band 470–494 MHz is also allocated to the radiolocation service on a secondary basis. This use is limited to the operation of wind profiler radars in accordance with Resolution 217 (WRC-97) (WRC-97).

S5.292—*Different category of service*: in Mexico, and Venezuela, the allocation of the band 470–512 MHz to the fixed and mobile services, and in Argentina and Uruguay to the mobile service, is on a primary basis (see No. S5.33), subject to agreement obtained under No. S9.21.

S5.293—*Different category of service*: in Chile, Colombia, Cuba, the United States, Guyana, Honduras, Jamaica, Mexico, and Panama, the allocation of the bands 470–512 MHz and 614–806 MHz to the fixed and mobile services is on a primary basis (see No. S5.33), subject to agreement obtained under No. S9.21.

S5.294—*Additional allocation*: in Burundi, Cameroon, the Congo, Ethiopia, Israel, Kenya, Lebanon, Libya, Malawi, Senegal, Sudan, Syria, and Yemen, the band 470–582 MHz is also allocated to the fixed service on a secondary basis.

S5.295—Not used.

S5.296—*Additional allocation*: in Germany, Austria, Belgium, Cyprus, Denmark, Spain, Finland, France, Ireland, Israel, Italy, Libya, Malta, Morocco, Monaco, Norway, the Netherlands, Portugal, Syria, the United Kingdom, Sweden, Switzerland, Swaziland, and Tunisia, the band 470–790 MHz is also allocated on a secondary basis to the land mobile service, intended for applications ancillary to broadcasting. Stations of the land mobile service in the countries listed in this footnote shall not cause harmful interference to existing or planned stations operating in accordance with the Table of Frequency Allocations in countries other than those listed in this footnote (WRC-97).

S5.297—*Additional allocation*: in Costa Rica, Cuba, El Salvador, the United States, Guatemala, Guyana, Honduras, Jamaica, Mexico and Venezuela, the band 512–608 MHz is also allocated to the fixed and mobile services on a primary basis, subject to agreement obtained under No. S9.21.

S5.298—*Additional allocation*: in India, the band 549.75–550.25 MHz is also allocated to the space operation service (space-to-Earth) on a secondary basis.

S5.299—Not used.

S5.300—*Additional allocation*: in Israel, Libya, Syria, and Sudan, the band 582–790 MHz is also allocated to the fixed and mobile, except aeronautical mobile, services on a secondary basis.

S5.301—Not used.

S5.302—*Additional allocation*: in the United Kingdom, the band 590–598 MHz is also allocated to the aeronautical radionavigation service on a primary basis. All new assignments to stations in the aeronautical radionavigation service, including those transferred from the adjacent bands, shall be subject to coordination with the Administrations of the following countries: Germany, Belgium, Denmark, Spain, France, Ireland, Luxembourg, Morocco, Norway, and the Netherlands.

S5.303—Not used.

S5.304—*Additional allocation*: in the African Broadcasting Area (see Nos. S5.10 to S5.13), the band 606–614 MHz is also allocated to the radio astronomy service on a primary basis.

S5.305—*Additional allocation*: in China, the band 606–614 MHz is also allocated to the radio astronomy service on a primary basis.

S5.306—*Additional allocation*: in Region 1, except in the African Broadcasting Area (see Nos. S5.10 to S5.13), and in Region 3, the band 608–614 MHz is also allocated to the radio astronomy service on a secondary basis.

S5.307—*Additional allocation*: in India, the band 608–614 MHz is also allocated to the radio astronomy service on a primary basis.

S5.308—Not used.

S5.309—*Different category of service*: in Costa Rica, El Salvador, and Honduras, the allocation of the band 614–806 MHz to the fixed service is on a primary basis (see No. S5.33), subject to agreement obtained under No. S9.21.

S5.310—SUP-WRC-97.

S5.311—Within the frequency band 620–790 MHz, assignments may be made to television stations using frequency modulation in the broadcasting-satellite service subject to agreement between the administrations concerned and those having services, operating in accordance with the Table, which may be affected (see Resolutions 33 (Rev. WRC-97) and 507). Such stations shall not produce a power flux-density in excess of the value $-129 \text{ dB(W/m}^2\text{)}$ for angles of arrival less than 20° (see Recommendation 705) within the territories of other countries without the consent of the administrations of those countries.

S5.312—*Additional allocation*: in Armenia, Azerbaijan, Belarus, Bulgaria, Georgia, Hungary, Kazakstan, Latvia, Moldova, Mongolia, Uzbekistan, Poland, Kyrgyzstan, Slovakia, the Czech Republic, Romania, Russian Federation, Tajikistan, Turkmenistan, and Ukraine, the band 645–862 MHz is also allocated to the aeronautical radionavigation service on a primary basis (WRC-97).

S5.313—SUP - WRC-97.

S5.314—*Additional allocation*: in Austria, Italy, Uzbekistan, the United Kingdom, and Swaziland, the band 790–862 MHz is also allocated to the land mobile service on a secondary basis (WRC-97).

S5.315—*Alternative allocation*: in Greece, Italy, Morocco, and Tunisia, the band 790–838 MHz is allocated to the broadcasting service on a primary basis.

S5.316—*Additional allocation*: in Germany, Bosnia and Herzegovina, Burkina Faso, Cameroon, Côte d'Ivoire, Croatia, Denmark, Egypt, Finland, Israel, Kenya, the former Yugoslav Republic of Macedonia, Libya, Liechtenstein, Monaco, Norway, the Netherlands, Portugal, Syria, Sweden, Switzerland, and Yugoslavia, the band 790–830 MHz, and in these same countries and in Spain, France, Gabon, and Malta, the band 830–862 MHz, are also allocated to the mobile, except aeronautical mobile, service on a primary basis. However, stations of the mobile service in the countries mentioned in connection with each band referred to in this footnote shall not cause harmful interference to, or claim protection from, stations of services operating in accordance with the Table in

countries other than those mentioned in connection with the band (WRC-97).

S5.317—*Additional allocation*: in Region 2 (except Brazil and the United States), the band 806–890 MHz is also allocated to the mobile-satellite service on a primary basis, subject to agreement obtained under No. **S9.21**. The use of this service is intended for operation within national boundaries.

S5.318—*Additional allocation*: in Canada, the United States, and Mexico, the bands 849–851 MHz and 894–896 MHz are also allocated to the aeronautical mobile service on a primary basis, for public correspondence with aircraft. The use of the band 849–851 MHz is limited to transmissions from aeronautical stations and the use of the band 894–896 MHz is limited to transmissions from aircraft stations.

S5.319—*Additional allocation*: in Belarus, Russian Federation, and Ukraine, the bands 806–840 MHz (Earth-to-space) and 856–890 MHz (space-to-Earth) are also allocated to the mobile-satellite, except aeronautical mobile-satellite (R), service. The use of these bands by this service shall not cause harmful interference to, or claim protection from, services in other countries operating in accordance with the Table of Frequency Allocations and is subject to special agreements between the administrations concerned.

S5.320—*Additional allocation*: in Region 3, the bands 806–890 MHz and 942–960 MHz are also allocated to the mobile-satellite, except aeronautical mobile-satellite (R), service on a primary basis, subject to agreement obtained under No. **S9.21**. The use of this service is limited to operation within national boundaries. In seeking such agreement, appropriate protection shall be afforded to services operating in accordance with the Table, to ensure that no harmful interference is caused to such services.

S5.321—*Alternative allocation*: in Italy, the band 838–854 MHz is allocated to the broadcasting service on a primary basis as from January 1, 1995.

S5.322—In Region 1, in the band 862–960 MHz, stations of the broadcasting service shall be operated only in the African Broadcasting Area (see Nos. **S5.10** to **S5.13**) excluding Algeria, Egypt, Spain, Libya, Morocco, Nigeria, South Africa, Tanzania, and Zimbabwe, subject to agreement obtained under No. **S9.21** (WRC-97).

S5.323—*Additional allocation*: in Armenia, Azerbaijan, Belarus, Bulgaria, Hungary, Kazakhstan, Latvia, Moldova, Mongolia, Uzbekistan, Poland, Kyrgyzstan, Slovakia, the Czech Republic, Romania, Russian Federation, Tajikistan, Turkmenistan, and Ukraine, the band 862–960 MHz is also allocated to the aeronautical radionavigation service on a primary basis. Such use is subject to agreement obtained under No. **S9.21** with administrations concerned and limited to ground-based radiobeacons in operation on 27 October 1997 until the end of their lifetime (WRC-97).

S5.324—Not used.

S5.325—*Different category of service*: in the United States, the allocation of the band 890–942 MHz to the radiolocation service is on a primary basis (see No. **S5.33**), subject to agreement obtained under No. **S9.21**.

S5.326—*Different category of service*: in Chile, the band 903–905 MHz is allocated to the mobile, except aeronautical mobile, service on a primary basis, subject to agreement obtained under No. **S9.21**.

S5.327—*Different category of service*: in Australia, the allocation of the band 915–928 MHz to the radiolocation service is on a primary basis (see No. **S5.33**).

S5.328—The band 960–1215 MHz is reserved on a worldwide basis for the use and development of airborne electronic aids to air navigation and any directly associated ground-based facilities.

S5.329—Use of the radionavigation-satellite service in the band 1215–1260 MHz shall be subject to the condition that no harmful interference is caused to the radionavigation service authorized under No. **S5.331**.

S5.330—*Additional allocation*: in Angola, Saudi Arabia, Bahrain, Bangladesh, Cameroon, China, the United Arab Emirates, Eritrea, Ethiopia, Guyana, India, Indonesia, the Islamic Republic of Iran, Iraq, Israel, Japan, Jordan, Kuwait, Lebanon, Libya, Morocco, Mozambique, Nepal, Nigeria, Pakistan, the Philippines, Qatar, Syria, Somalia, Sudan, Sri Lanka, Chad, Togo, and Yemen, the band 1215–1300 MHz is also allocated to the fixed and mobile services on a primary basis (WRC-97).

S5.331—*Additional allocation*: in Algeria, Germany, Austria, Bahrain, Belgium, Benin, Bosnia and Herzegovina, Burundi, Cameroon, China, Croatia, Denmark, the United Arab Emirates, France, Greece, India, the Islamic Republic of Iran, Iraq, Kenya, the former Yugoslav Republic of Macedonia, Liechtenstein, Luxembourg, Mali, Mauritania, Norway, Oman, Pakistan, the Netherlands, Portugal, Qatar, Senegal, Slovenia, Somalia, Sudan, Sri Lanka, Sweden, Switzerland, Turkey, and Yugoslavia, the band 1215–1300 MHz is also allocated to the radionavigation service on a primary basis.

S5.332—In the band 1215–1300 MHz, active spaceborne sensors in the earth exploration-satellite and space research services shall not cause harmful interference to, claim protection from, or otherwise impose constraints on operation or development of the radiolocation service, the radionavigation-satellite service, and other services allocated on a primary basis (WRC-97).

S5.333—SUP - WRC-97.

S5.334—*Additional allocation*: in Canada and the United States, the bands 1240–1300 MHz and 1350–1370 MHz are also allocated to the aeronautical radionavigation service on a primary basis.

S5.335—In Canada and the United States in the band 1240–1300 MHz, active spaceborne sensors in the earth exploration-satellite and space research services shall not cause interference to, claim protection from, or otherwise impose constraints on operation or development of the aeronautical radionavigation service (WRC-97).

S5.336—Not used.

S5.337—The use of the bands 1 300–1350 MHz, 2700–2900 MHz, and 9000–9200 MHz by the aeronautical radionavigation service is restricted to ground-based radars and to associated airborne transponders which transmit only on frequencies in these bands and only when actuated by radars operating in the same band.

S5.338—In Azerbaijan, Bulgaria, Mongolia, Poland, Kyrgyzstan, Slovakia, the Czech Republic, Romania, Turkmenistan, and Ukraine, existing installations of the radionavigation

service may continue to operate in the band 1350–1400 MHz (WRC-97).

S5.339—The bands 1370–1400 MHz, 2640–2655 MHz, 4950–4990 MHz, and 15.20–15.35 GHz are also allocated to the space research (passive) and earth exploration-satellite (passive) services on a secondary basis.

S5.340—All emissions are prohibited in the following bands:

1400–1427 MHz.

2690–2700 MHz, except those provided for by Nos. **S5.421** and **S5.422**,

10.68–10.7 GHz, except those provided for by No. **S5.483**.

15.35–15.4 GHz, except those provided for by No. **S5.511**.

23.6–24 GHz.

31.3–31.5 GHz,

31.5–31.8 GHz, in Region 2,

48.94–49.04 GHz, from airborne stations, 50.2–50.4 GHz,* except those provided for by No. **S5.555A**.

52.6–54.25 GHz,

86–92 GHz,

105–116 GHz,

140.69–140.98 GHz, from airborne stations and from space stations in the space-to-Earth direction,

182–185 GHz, except those provided for by No. **S5.563**,

217–231 GHz

S5.341—In the bands 1400–1727 MHz, 101–120 GHz, and 197–220 GHz, passive research is being conducted by some countries in a program for the search for intentional emissions of extraterrestrial origin.

S5.342—*Additional allocation:* in Belarus, Russian Federation, and Ukraine, the band 1429–1535 MHz is also allocated to the aeronautical mobile service on a primary basis exclusively for the purposes of aeronautical telemetry within the national territory. As of April 1, 2007, the use of the band 1452–1492 MHz is subject to agreement between the administrations concerned.

S5.343—In Region 2, the use of the band 1435–1535 MHz by the aeronautical mobile service for telemetry has priority over other uses by the mobile service.

S5.344—*Alternative allocation:* in the United States, the band 1452–1525 MHz is allocated to the fixed and mobile services on a primary basis (see also No. **S5.343**).

S5.345—Use of the band 1452–1492 MHz by the broadcasting-satellite service, and by the broadcasting service, is limited to digital audio broadcasting and is subject to the provisions of Resolution 528 (WARC-92).

S5.346—Not used.

S5.347—*Different category of service:* in Bangladesh, Bosnia and Herzegovina, Botswana, Bulgaria, Burkina Faso, Cuba, Denmark, Egypt, Greece, Ireland, Italy, Jordan, Kenya,

Mozambique, Portugal, Sri Lanka, Swaziland, Yemen, Yugoslavia, and Zimbabwe, the allocation of the band 1452–1492 MHz to the broadcasting-satellite service and the broadcasting service is on a secondary basis until April 1, 2007 (WRC-97).

S5.348—The use of the band 1492–1525 MHz by the mobile-satellite service is subject to coordination under No. **S9.11A**. However, no coordination threshold in Article **S21** for space stations of the mobile-satellite service with respect to terrestrial services shall apply to the situation referred to in No. **S5.343**. With respect to the situation referred to in No. **S5.343**, the requirement for coordination in the band 1492–1525 MHz will be determined by band overlap.

S5.348A—In the band 1492–1525 MHz, the coordination threshold in terms of the power flux-density levels at the surface of the Earth in application of No. **S9.11A** for space stations in the mobile-satellite (space-to-Earth) service, with respect to the land mobile service use for specialized mobile radios or used in conjunction with public switched telecommunication networks (PSTN) operating within the territory of Japan, shall be -150 dB(W/m²) in any 4 kHz band for all angles of arrival, instead of those given in Table S5-2 of Appendix S5. The above threshold level of the power flux-density shall apply until it is changed by a competent world radiocommunication conference.

S5.349—*Different category of service:* in Saudi Arabia, Azerbaijan, Bahrain, Bosnia, and Herzegovina, Cameroon, Egypt, the United Arab Emirates, France, the Islamic Republic of Iran, Iraq, Israel, Kazakstan, Kuwait, the former Yugoslav Republic of Macedonia, Lebanon, Morocco, Mongolia, Oman, Qatar, Syria, Kyrgyzstan, Romania, Turkmenistan, Ukraine, Yemen, and Yugoslavia, the allocation of the band 1525–1530 MHz to the mobile, except aeronautical mobile, service is on a primary basis (see No. **S5.33**) (WRC-97).

S5.350—*Additional allocation:* in Azerbaijan, Kyrgyzstan, Turkmenistan, and Ukraine, the band 1525–1530 MHz is also allocated to the aeronautical mobile service on a primary basis (WRC-97).

S5.351—The bands 1525–1544 MHz, 1545–1559 MHz, 1626.5–1645.5 MHz, and 1646.5–1660.5 MHz shall not be used for feeder links of any service. In exceptional circumstances, however, an earth station at a specified fixed point in any of the mobile-satellite services may be authorized by an administration to communicate via space stations using these bands.

S5.352—SUP - WRC-97.

S5.352A—In the band 1525–1530 MHz, stations in the mobile-satellite service, except stations in the maritime mobile-satellite service, shall not cause harmful interference to, or claim protection from, stations of the fixed service in France and French Overseas Territories in Region 3, Algeria, Saudi Arabia, Egypt, Guinea, India, Israel, Italy, Jordan, Kuwait, Mali, Malta, Morocco, Mauritania, Nigeria, Oman, Pakistan, Philippines, Qatar, Syria, Tanzania, Viet Nam, and Yemen notified prior to 1 April 1998 (WRC-97).

S5.353—SUP - WRC-97.

S5.353A—In applying the procedures of No. **S9.11A** to the mobile-satellite service in the bands 1530–1544 MHz and 1626.5–1645.5 MHz, priority shall be given to accommodating the spectrum requirements for distress, urgency, and safety communications of the Global Maritime Distress and

* The allocation to the earth exploration-satellite service (passive) and the space research service (passive) in the band 50.2–50.4 GHz should not impose undue constraints on the use of the adjacent bands by the primary allocated services in those bands (WRC-97).

Safety System (GMDSS). Maritime mobile-satellite distress, urgency, and safety communications shall have priority access and immediate availability over all other mobile satellite communications operating within a network. Mobile-satellite systems shall not cause unacceptable interference to, or claim protection from, distress, urgency, and safety communications of the GMDSS. Account shall be taken of the priority of safety-related communications in the other mobile-satellite services (WRC-97). (See Resolution 218 (WRC-97).)

S5.354—The use of the bands 1525–1559 MHz and 1626.5–1660.5 MHz by the mobile-satellite services is subject to coordination under No. **S9.11A**.

S5.355—*Additional allocation:* in Bahrain, Bangladesh, the Congo, Egypt, the United Arab Emirates, Eritrea, Ethiopia, the Islamic Republic of Iran, Iraq, Israel, Jordan, Kuwait, Lebanon, Malta, Morocco, Oman, Qatar, Syria, Somalia, Sudan, Sri Lanka, Chad, Togo, Yemen, and Zambia, the bands 1540–1645.5 MHz and 1646.5–1660 MHz are also allocated to the fixed service on a secondary basis (WRC-97).

S5.356—The use of the band 1544–1545 MHz by the mobile-satellite service (space-to-Earth) is limited to distress and safety communications (see Article **S31**).

S5.357—Transmissions in the band 1545–1555 MHz from terrestrial aeronautical stations directly to aircraft stations, or between aircraft stations, in the aeronautical mobile (R) service are also authorized when such transmissions are used to extend or supplement the satellite-to-aircraft links.

S5.357A—In applying the procedures of No. **S9.11A** to the mobile-satellite service in the bands 1545–1555 MHz and 1646.5–1656.5 MHz, priority shall be given to accommodating the spectrum requirements of the aeronautical mobile-satellite (R) service providing transmission of messages with priority 1 to 6 in Article **S44**. Aeronautical mobile-satellite (R) service communications with priority 1 to 6 in Article **S44** shall have priority access and immediate availability, by preemption if necessary, over all other mobile-satellite communications operating within a network. Mobile-satellite systems shall not cause unacceptable interference to, or claim protection from, aeronautical mobile-satellite (R) service communications with priority 1 to 6 in Article **S44**. Account shall be taken of the priority of safety-related communications in the other mobile-satellite services (WRC-97) (See Resolution 218 (WRC-97)); (WRC-97).

S5.358—SUP-WRC-97.

S5.359—*Additional allocation:* in Germany, Saudi Arabia, Armenia, Austria, Azerbaijan, Belarus, Benin, Bulgaria, Cameroon, Spain, France, Gabon, Georgia, Greece, Guinea, Guinea-Bissau, Hungary, Jordan, Kazakhstan, Kuwait, Latvia, Libya, Mali, Mauritania, Moldova, Mongolia, Nigeria, Uganda, Uzbekistan, Pakistan, Poland, Syria, Kyrgyzstan, the Democratic People's Republic of Korea, Romania, Russian Federation, Senegal, Swaziland, Tajikistan, Tanzania, Turkmenistan, Ukraine, Zambia, and Zimbabwe the bands 1550–1645.5 MHz and 1646.5–1660 MHz are also allocated to the fixed service on a primary basis. Administrations are urged to make all practicable efforts to avoid the implementation of new fixed-service stations in the bands 1550–1555 MHz, 1610–1645.5 MHz, and 1646.5–1660 MHz.

S5.360 to **S5.362**—SUP-WRC-97.

S5.362A—In the United States, in the bands 1555–1559 MHz and 1656.5–1660.5 MHz, the aeronautical mobile-satellite (R) service shall have priority access and immediate availability, by preemption if necessary, over all other mobile-satellite communications operating within a network. Mobile-satellite systems shall not cause unacceptable interference to, or claim protection from, aeronautical mobile-satellite (R) service communications with priority 1 to 6 in Article **S44**. Account shall be taken of the priority of safety-related communications in the other mobile-satellite services (WRC-97).

S5.363—*Alternative allocation:* in Sweden, the band 1590–1626.5 MHz is allocated to the aeronautical radionavigation service on a primary basis.

S5.364—The use of the band 1610–1626.5 MHz by the mobile-satellite service (Earth-to-space) and by the radiodetermination-satellite service (Earth-to-space) is subject to coordination under No. **S9.11A**. A mobile earth station operating in either of the services in this band shall not produce a peak eirp density in excess of –15 dB(W/4 kHz) in the part of the band used by systems operating in accordance with the provisions of No. **S5.366** (to which No. **S4.10** applies), unless otherwise agreed by the affected administrations. In the part of the band where such systems are not operating, the mean eirp density of a mobile earth station shall not exceed –3 dB(W/4 kHz). Stations of the mobile-satellite service shall not claim protection from stations in the aeronautical radionavigation service, stations operating in accordance with the provisions of No. **S5.366**, and stations in the fixed service operating in accordance with the provisions of No. **S5.359**. Administrations responsible for the coordination of mobile-satellite networks shall make all practicable efforts to ensure protection of stations operating in accordance with the provisions of No. **S5.366**.

S5.365—The use of the band 1613.8–1626.5 MHz by the mobile-satellite service (space-to-Earth) is subject to coordination under No. **S9.11A**.

S5.366—The band 1610–1626.5 MHz is reserved on a worldwide basis for the use and development of airborne electronic aids to air navigation and any directly associated ground-based or satellite-borne facilities. Such satellite use is subject to agreement obtained under No. **S9.21**.

S5.367—*Additional allocation:* The bands 1610–1626.5 MHz and 5000–5150 MHz are also allocated to the aeronautical mobile-satellite (R) service on a primary basis, subject to agreement obtained under No. **S9.21**.

S5.368—With respect to the radiodetermination-satellite and mobile-satellite services the provisions of No. **S4.10** do not apply in the band 1610–1626.5 MHz, with the exception of the aeronautical radionavigation-satellite service.

S5.369—*Different category of service:* in Angola, Australia, Burundi, China, Côte d'Ivoire, Eritrea, Ethiopia, India, the Islamic Republic of Iran, Israel, Jordan, Lebanon, Liberia, Libya, Madagascar, Mali, Pakistan, Papua New Guinea, Dem. Rep. of the Congo, Syria, Senegal, Sudan, Swaziland, Togo, and Zambia, the allocation of the band 1610–1626.5 MHz to the radiodetermination-satellite service (Earth-to-space) is on a primary basis (see No. **S5.33**), subject to agreement obtained under No. **S9.21** from countries not listed in this provision (WRC-97).

S5.370—*Different category of service:* in Venezuela, the allocation to the radiodetermination-satellite service in the band 1610–1626.5 MHz (Earth-to-space) is on a secondary basis.

S5.371—*Additional allocation:* in Region 1, the bands 1610–1626.5 MHz (Earth-to-space) and 2483.5–2500 MHz (space-to-Earth) are also allocated to the radiodetermination-satellite service on a secondary basis, subject to agreement obtained under No. **S9.21**.

S5.372—Harmful interference shall not be caused to stations of the radio astronomy service using the band 1610.6–1613.8 MHz by stations of the radiodetermination-satellite and mobile-satellite services (No. **S29.13** applies).

S5.373—Not used.

S5.373A—SUP-WRC-97.

S5.374—Mobile earth stations in the mobile-satellite service operating in the bands 1631.5–1634.5 MHz and 1656.5–1660 MHz shall not cause harmful interference to stations in the fixed service operating in the countries listed in No. **S5.359** (WRC-97).

S5.375—The use of the band 1645.5–1646.5 MHz by the mobile-satellite service (Earth-to-space) and for inter-satellite links is limited to distress and safety communications (see Article **S31**).

S5.376—Transmissions in the band 1646.5–1656.5 MHz from aircraft stations in the aeronautical mobile (R) service directly to terrestrial aeronautical stations, or between aircraft stations, are also authorized when such transmissions are used to extend or supplement the aircraft-to-satellite links.

S5.376A—Mobile earth stations operating in the band 1660–1660.5 MHz shall not cause harmful interference to stations in the radio astronomy service (WRC-97).

S5.377—In the band 1675–1710 MHz, stations in the mobile-satellite service shall not cause harmful interference to, nor constrain the development of, the meteorological-satellite and meteorological aids services (see Resolution 213 (Rev. WRC-95)) and the use of this band shall be subject to coordination under No. **S9.11A**.

S5.378—Not used.

S5.379—*Additional allocation:* in Bangladesh, India, Indonesia, Nigeria, and Pakistan, the band 1660.5–1668.4 MHz is also allocated to the meteorological aids service on a secondary basis.

S5.379A—Administrations are urged to give all practicable protection in the band 1660.5–1668.4 MHz for future research in radio astronomy, particularly by eliminating air-to-ground transmissions in the meteorological aids service in the band 1664.4–1668.4 MHz as soon as practicable.

S5.380—The bands 1670–1675 MHz and 1800–1805 MHz are intended for use, on a worldwide basis, by administrations wishing to implement aeronautical public correspondence. The use of the band 1670–1675 MHz by stations in the systems for public correspondence with aircraft is limited to transmissions from aeronautical stations and the use of the band 1800–1805 MHz is limited to transmissions from aircraft stations.

S5.381—*Additional allocation:* in Afghanistan, Costa Rica, Cuba, India, the Islamic Republic of Iran, Malaysia, Pakistan, and Sri Lanka, the band 1690–1700 MHz is also

allocated to the fixed and mobile, except aeronautical mobile, services on a primary basis (WRC-97).

S5.382—*Different category of service:* in Saudi Arabia, Armenia, Austria, Azerbaijan, Bahrain, Belarus, Bosnia and Herzegovina, Bulgaria, the Congo, Egypt, the United Arab Emirates, Eritrea, Ethiopia, Guinea, Hungary, Iraq, Israel, Jordan, Kazakhstan, Kuwait, the former Yugoslav Republic of Macedonia, Lebanon, Mauritania, Moldova, Mongolia, Oman, Uzbekistan, Poland, Qatar, Syria, Kyrgyzstan, Romania, Russian Federation, Somalia, Tajikistan, Tanzania, Turkmenistan, Ukraine, Yemen and Yugoslavia, the allocation of the band 1690–1700 MHz to the fixed and mobile, except aeronautical mobile, services is on a primary basis (see No. **S5.33**), and in the Democratic People's Republic of Korea, the allocation of the band 1690–1700 MHz to the fixed service is on a primary basis (see No. **S5.33**) and to the mobile, except aeronautical mobile, service on a secondary basis (WRC-97).

S5.383—Not used.

S5.384—*Additional allocation:* in India, Indonesia, and Japan, the band 1700–1710 MHz is also allocated to the space research service (space-to-Earth) on a primary basis (WRC-97).

S5.385—*Additional allocation:* The bands 1718.8–1722.2 MHz, 150–151 GHz, 174.42–175.02 GHz, 177–177.4 GHz, 178.2–178.6 GHz, 181–181.46 GHz, 186.2–186.6 GHz, and 257.5–258 GHz are also allocated to the radio astronomy service on a secondary basis for spectral line observations.

S5.386—*Additional allocation:* The band 1750–1850 MHz is also allocated to the space operation (Earth-to-space) and space research (Earth-to-space) services in Region 2, in Australia, India, Indonesia, and Japan on a primary basis, subject to agreement obtained under No. **S9.21**, having particular regard to troposcatter systems.

S5.387—*Additional allocation:* in Armenia, Azerbaijan, Belarus, Georgia, Kazakhstan, Mali, Mongolia, Uzbekistan, Kyrgyzstan, Slovakia, the Czech Republic, Romania, Russian Federation, Tajikistan, Turkmenistan, and Ukraine, the band 1770–1790 MHz is also allocated to the meteorological-satellite service on a primary basis, subject to agreement obtained under No. **S9.21** (WRC-97).

S5.388—The bands 1885–2025 MHz and 2110–2200 MHz are intended for use, on a worldwide basis, by administrations wishing to implement International Mobile Telecommunications-2000 (IMT-2000). Such use does not preclude the use of these bands by other services to which they are allocated. The bands should be made available for IMT-2000 in accordance with Resolution 212 (Rev. WRC-97) (WRC-97).

S5.389—Not used.

S5.389A—The use of the bands 1980–2010 MHz and 2170–2200 MHz by the mobile-satellite service is subject to coordination under No. **S9.11A** and to the provisions of Resolution 716 (WRC-95). The use of these bands shall not commence before January 1, 2000; however, the use of the band 1980–1990 MHz in Region 2 shall not commence before January 1, 2005.

S5.389B—The use of the band 1980–1990 MHz by the mobile-satellite service shall not cause harmful interference to or constrain the development of the fixed and mobile services in Argentina, Brazil, Canada, Chile, Ecuador, the United

States, Honduras, Jamaica, Mexico, Peru, Suriname, Trinidad and Tobago, Uruguay, and Venezuela.

S5.389C—The use of the bands 2010–2025 MHz and 2160–2170 MHz in Region 2 by the mobile-satellite service shall not commence before January 1, 2002 and is subject to coordination under No. **S9.11A** and to the provisions of Resolution 716 (WRC-95) (WRC-97).

S5.389D—In Canada and the United States the use of the bands 2010–2025 MHz and 2160–2170 MHz by the mobile-satellite service shall not commence before January 1, 2000.

S5.389E—The use of the bands 2010–2025 MHz and 2160–2170 MHz by the mobile-satellite service in Region 2 shall not cause harmful interference to or constrain the development of the fixed and mobile services in Regions 1 and 3.

S5.389F—In Algeria, Benin, Cape Verde, Egypt, Mali, Syria, and Tunisia, the use of the bands 1980–2010 MHz and 2170–2200 MHz by the mobile-satellite service shall neither cause harmful interference to the fixed and mobile services nor hamper the development of those services prior to January 1, 2005, nor shall the former service request protection from the latter services.

S5.390—In Argentina, Brazil, Chile, Colombia, Cuba, Ecuador, and Suriname, the use of the bands 2010–2025 MHz and 2160–2170 MHz by the mobile-satellite services shall not cause harmful interference to stations in the fixed and mobile services before January 1, 2005. After this date, the use of these bands is subject to coordination under No. **S9.11A** and to the provisions of Resolution 716 (WRC-95) (WRC-97).

S5.391—In making assignments to the mobile service in the bands 2025–2110 MHz and 2200–2290 MHz, administrations shall not introduce high-density mobile systems, as described in Recommendation ITU-R SA.1154, and shall take that Recommendation into account for the introduction of any other type of mobile system (WRC-97).

S5.392—Administrations are urged to take all practicable measures to ensure that space-to-space transmissions between two or more nongeostationary satellites, in the space research, space operations, and Earth exploration-satellite services in the bands 2025–2110 MHz and 2200–2290 MHz, shall not impose any constraints on Earth-to-space, space-to-Earth, and other space-to-space transmissions of those services and in those bands between geostationary and nongeostationary satellites.

S5.392A—*Additional allocation*: in Russian Federation, the band 2160–2200 MHz is also allocated to the space research service (space-to-Earth) on a primary basis until 1 January 2005. Stations in the space research service shall not cause harmful interference to, or claim protection from, stations in the fixed and mobile services operating in this frequency band.

S5.393—*Additional allocation*: in the United States, India, and Mexico, the band 2310–2360 MHz is also allocated to the broadcasting-satellite service (sound) and complementary terrestrial sound broadcasting service on a primary basis. Such use is limited to digital audio broadcasting and is subject to the provisions of Resolution 528(WARC-92) (WRC-97).

S5.394—In the United States, the use of the band 2300–2390 MHz by the aeronautical mobile service for telemetry has priority over other uses by the mobile services. In Canada, the use of the band 2300–2483.5 MHz by the aeronautical

mobile service for telemetry has priority over other uses by the mobile services.

S5.395—In France, the use of the band 2310–2360 MHz by the aeronautical mobile service for telemetry has priority over other uses by the mobile service.

S5.396—Space stations of the broadcasting-satellite service in the band 2310–2360 MHz operating in accordance with No. **S5.393** that may affect the services to which this band is allocated in other countries shall be coordinated and notified in accordance with Resolution 33 (Rev. WRC-97). Complementary terrestrial broadcasting stations shall be subject to bilateral coordination with neighbouring countries prior to their being brought into use.

S5.397—*Different category of service*: in France, the band 2450–2500 MHz is allocated on a primary basis to the radiolocation service (see No. **S5.33**). Such use is subject to agreement with administrations having services operating or planned to operate in accordance with the Table of Frequency Allocations which may be affected.

S5.398—In respect of the radiodetermination-satellite service in the band 2483.5–2500 MHz, the provisions of No. **S4.10** do not apply.

S5.399—In Region 1, in countries other than those listed in No. **S5.400**, harmful interference shall not be caused to, or protection shall not be claimed from, stations of the radiolocation service by stations of the radiodetermination-satellite service.

S5.400—*Different category of service*: in Angola, Australia, Bangladesh, Burundi, China, Eritrea, Ethiopia, India, the Islamic Republic of Iran, Jordan, Lebanon, Liberia, Libya, Madagascar, Mali, Pakistan, Papua New Guinea, Democratic Republic of the Congo, Syria, Sudan, Swaziland, Togo, and Zambia, the allocation of the band 2483.5–2500 MHz to the radiodetermination-satellite service (space-to-Earth) is on a primary basis (see No. **S5.33**), subject to agreement obtained under No. **S9.21** from countries not listed in this provision (WRC-97).

S5.401—Not used.

S5.402—The use of the band 2483.5–2500 MHz by the mobile-satellite and the radiodetermination-satellite services is subject to the coordination under No. **S9.11A**. Administrations are urged to take all practicable steps to prevent harmful interference to the radio astronomy service from emissions in the 2483.5–2500 MHz band, especially those caused by second-harmonic radiation that would fall into the 4990–5000 MHz band allocated to the radio astronomy service worldwide.

S5.403—Subject to agreement obtained under No. **S9.21**, the band 2520–2535 MHz (until 1 January 2005 the band 2500–2535 MHz) may also be used for the mobile-satellite (space-to-Earth), except aeronautical mobile-satellite, service for operation limited to within national boundaries. The provisions of No. **S9.11A** apply.

S5.404—*Additional allocation*: in India and the Islamic Republic of Iran, the band 2500–2516.5 MHz may also be used for the radiodetermination-satellite service (space-to-Earth) for operation limited to within national boundaries, subject to agreement obtained under No. **S9.21**.

S5.405—*Additional allocation*: in France, the band 2500–2550 MHz is also allocated to the radiolocation service on a

primary basis. Such use is subject to agreement with the administrations having services operating or planned to operate in accordance with the Table which may be affected.

S5.406—Not used.

S5.407—In the band 2500–2520 MHz, the power flux-density at the surface of the Earth from space stations operating in the mobile-satellite (space-to-Earth) service shall not exceed $-152 \text{ dB(W/m}^2/4 \text{ kHz)}$ in Argentina, unless otherwise agreed by the administrations concerned.

S5.408—*Additional allocation:* in the United Kingdom, the band 2500–2600 MHz is also allocated to the radiolocation service on a secondary basis.

S5.409—Administrations shall make all practicable efforts to avoid developing new tropospheric scatter systems in the band 2500–2690 MHz.

S5.410—The band 2500–2690 MHz may be used for tropospheric scatter systems in Region 1, subject to agreement obtained under No. **S9.21**.

S5.411—When planning new tropospheric scatter radio-relay links in the band 2500–2690 MHz, all possible measures shall be taken to avoid directing the antennas of these links toward the geostationary-satellite orbit.

S5.412—*Alternative allocation:* in Azerbaijan, Bulgaria, Kyrgyzstan, Turkmenistan, and Ukraine, the band 2500–2690 MHz is allocated to the fixed and mobile, except aeronautical mobile, services on a primary basis (WRC-97).

S5.413—In the design of systems in the broadcasting-satellite service in the bands between 2500 MHz and 2690 MHz, administrations are urged to take all necessary steps to protect the radio astronomy service in the band 2690–2700 MHz.

S5.414—The allocation of the frequency band 2500–2520 MHz to the mobile-satellite service (space-to-Earth) shall be effective on January 1, 2005 and is subject to coordination under No. **S9.11A**.

S5.415—The use of the bands 2500–2690 MHz in Region 2 and 2500–2535 MHz and 2655–2690 MHz in Region 3 by the fixed-satellite service is limited to national and regional systems, subject to agreement obtained under No. **S9.21**, giving particular attention to the broadcasting-satellite service in Region 1. In the direction space-to-Earth, the power flux-density at the Earth's surface shall not exceed the values given in Article **S21**, Table **S21-4**.

S5.415A—*Additional allocation:* in Japan, subject to agreement obtained under No. **S9.21**, the band 2515–2535 MHz may also be used for the aeronautical mobile-satellite service (space-to-Earth) for operation limited to within its national boundary from January 1, 2000 (WRC-97).

S5.416—The use of the band 2520–2670 MHz by the broadcasting-satellite service is limited to national and regional systems for community reception, subject to agreement obtained under No. **S9.21**. The power flux-density at the Earth's surface shall not exceed the values given in Article **S21**, Table **S21-4**.

S5.417—*Alternative allocation:* in Germany and Greece, the band 2520–2670 MHz is allocated to the fixed service on a primary basis.

S5.418—*Additional allocation:* in Bangladesh, Belarus, China, Rep. of Korea, India, Japan, Pakistan, Russian Federation, Singapore, Sri Lanka, Thailand, and Ukraine the band 2535–2655 MHz is also allocated to the broadcasting-satellite

service (sound) and complementary terrestrial broadcasting service on a primary basis. Such use is limited to digital audio broadcasting and is subject to provisions of Resolution 528 (WARC-92). The provisions of No. **S5.416** and Article **S21**, Table **S21-4**, do not apply to this additional allocation.

S5.419—The allocation of the frequency band 2670–2690 MHz to the mobile-satellite service shall be effective from January 1, 2005. When introducing systems of the mobile-satellite service in this band, administrations shall take all necessary steps to protect the satellite systems operating in this band prior to 3 March 1992. The coordination of mobile-satellite systems in the band shall be in accordance with No. **S9.11A**.

S5.420—The band 2655–2670 MHz (until January 1, 2005 the band 2655–2690 MHz) may also be used for the mobile-satellite (Earth-to-space), except aeronautical mobile-satellite, service for operation limited to within national boundaries, subject to agreement obtained under No. **S9.21**. The coordination under No. **S9.11A** applies.

S5.420A—*Additional allocation:* in Japan, subject to agreement obtained under No. **S9.21**, the band 2670–2690 MHz may also be used for the aeronautical mobile-satellite service (Earth-to-space) for operation limited to within its national boundary from January 1, 2000 (WRC-97).

S5.421—*Additional allocation:* in Germany and Austria, the band 2690–2695 MHz is also allocated to the fixed service on a primary basis. Such use is limited to equipment in operation by 1 January 1985.

S5.422—*Additional allocation:* in Saudi Arabia, Armenia, Azerbaijan, Bahrain, Belarus, Bosnia and Herzegovina, Brunei Darussalam, the Central African Republic, the Congo, Côte d'Ivoire, Cuba, Egypt, the United Arab Emirates, Eritrea, Ethiopia, Gabon, Georgia, Guinea, Guinea-Bissau, the Islamic Republic of Iran, Iraq, Israel, Jordan, Kazakhstan, Lebanon, Malaysia, Mali, Morocco, Mauritania, Moldova, Mongolia, Nigeria, Oman, Uzbekistan, Pakistan, the Philippines, Qatar, Syria, Kyrgyzstan, Dem. Rep. of the Congo, Romania, Russian Federation, Somalia, Tajikistan, Tunisia, Turkmenistan, Ukraine, Yemen, Yugoslavia, and Zambia, the band 2690–2700 MHz is also allocated to the fixed and mobile, except aeronautical mobile, services on a primary basis. Such use is limited to equipment in operation by 1 January 1985 (WRC-97).

S5.423—In the band 2700–2900 MHz, ground-based radars used for meteorological purposes are authorized to operate on a basis of equality with stations of the aeronautical radionavigation service.

S5.424—*Additional allocation:* in Canada, the band 2850–2900 MHz is also allocated to the maritime radionavigation service, on a primary basis, for use by shore-based radars.

S5.425—In the band 2900–3100 MHz, the use of the shipborne interrogator-transponder system (SIT) shall be confined to the sub-band 2930–2950 MHz.

S5.426—The use of the band 2900–3100 MHz by the aeronautical radionavigation service is limited to ground-based radars.

S5.427—In the bands 2900–3100 MHz and 9300–9500 MHz, the response from radar transponders shall not be capable of being confused with the response from radar beacons (racons) and shall not cause interference to ship or aeronauti-

cal radars in the radionavigation service, having regard, however, to No. S4.9.

S5.428—*Additional allocation:* in Azerbaijan, Bulgaria, Cuba, Kazakhstan, Mongolia, Poland, Kyrgyzstan, Romania, Turkmenistan, and Ukraine, the band 3100–3300 MHz is also allocated to the radionavigation service on a primary basis (WRC-97).

S5.429—*Additional allocation:* in Saudi Arabia, Bahrain, Bangladesh, Brunei Darussalam, China, the Congo, the Republic of Korea, the United Arab Emirates, India, Indonesia, the Islamic Republic of Iran, Iraq, Israel, Japan, Jordan, Kuwait, Lebanon, Libya, Malaysia, Oman, Pakistan, Qatar, Syria, Democratic People's Republic of Korea and Yemen, the band 3300–3400 MHz is also allocated to the fixed and mobile services on a primary basis. The countries bordering the Mediterranean shall not claim protection for their fixed and mobile services from the radiolocation service (WRC-97).

S5.430—*Additional allocation:* in Azerbaijan, Bulgaria, Cuba, Mongolia, Poland, Kyrgyzstan, Romania, Turkmenistan, and Ukraine, the band 3300–3400 MHz is also allocated to the radionavigation service on a primary basis (WRC-97).

S5.431—*Additional allocation:* in Germany, Israel, Nigeria, and the United Kingdom, the band 3400–3475 MHz is also allocated to the amateur service on a secondary basis.

S5.432—*Different category of service:* in the Republic of Korea, Indonesia, Japan and Pakistan, the allocation of the band 3400–3500 MHz to the mobile, except aeronautical mobile, service is on a primary basis (see No. S5.33) (WRC-97).

S5.433—In Regions 2 and 3, in the band 3400–3600 MHz the radiolocation service is allocated on a primary basis. However, all administrations operating radiolocation systems in this band are urged to cease operations by 1985. Thereafter, administrations shall take all practicable steps to protect the fixed-satellite service and coordination requirements shall not be imposed on the fixed-satellite service.

S5.434—SUP-WRC-97.

S5.435—In Japan, in the band 3620–3700 MHz, the radiolocation service is excluded.

S5.436—Not used.

S5.437—*Additional allocation:* in Germany and Norway, the band 4200–4210 MHz is also allocated to the fixed service on a secondary basis (WRC-97).

S5.438—Use of the band 4200–4400 MHz by the aeronautical radionavigation service is reserved exclusively for radio altimeters installed on board aircraft and for the associated transponders on the ground. However, passive sensing in the earth exploration-satellite and space research services may be authorized in this band on a secondary basis (no protection is provided by the radio altimeters).

S5.439—*Additional allocation:* in China, the Islamic Republic of Iran, and Libya, the band 4200–4400 MHz is also allocated to the fixed service on a secondary basis (WRC-97).

S5.440—The standard frequency and time signal-satellite service may be authorized to use the frequency 4202 MHz for space-to-Earth transmissions and the frequency 6427 MHz for Earth-to-space transmissions. Such transmissions shall be confined within the limits of 2 MHz of these frequencies, subject to agreement obtained under No. S9.21.

S5.441—The use of the bands 4500–4800 MHz (space-to-Earth), 6725–7025 MHz (Earth-to-space) by the fixed-satellite service shall be in accordance with the provisions of Appendix S30B. The use of the bands 10.7–10.95 GHz (space-to-Earth), 11.2–11.45 GHz (space-to-Earth), and 12.75–13.25 GHz (Earth-to-space) by geostationary-satellite systems in the fixed-satellite service shall be in accordance with the provisions of Appendix S30B. The use of the bands 10.7–10.95 GHz (space-to-Earth), 11.2–11.45 GHz (space-to-Earth), and 12.75–13.25 GHz (Earth-to-space) by nongeostationary-satellite systems in the fixed-satellite service shall be in accordance with the provisions of Resolution 130 (WRC-97) (WRC-97).

S5.442—In the bands 4825–4835 MHz and 4950–4990 MHz, the allocation to the mobile service is restricted to the mobile, except aeronautical mobile, service.

S5.443—*Different category of service:* in Argentina, Australia, and Canada, the allocation of the bands 4825–4835 MHz and 4950–4990 MHz to the radio astronomy service is on a primary basis (see No. S5.33).

S5.444—The band 5000–5150 MHz is to be used for the operation of the international standard system (microwave landing system) for precision approach and landing. The requirements of this system shall take precedence over other uses of this band. For the use of this band, No. S5.444A and Resolution 114 (WRC-95) apply.

S5.444A—*Additional allocation:* The band 5091–5150 MHz is also allocated to the fixed-satellite service (Earth-to-space) on a primary basis. This allocation is limited to feeder links of nongeostationary mobile-satellite systems and is subject to coordination under No. S9.11A. In the band 5091–5150 MHz, the following conditions also apply:

Prior to January 1, 2010, the use of the band 5091–5150 MHz by feeder links of non-geostationary-satellite systems in the mobile-satellite service shall be made in accordance with Resolution 114 (WRC-95).

Prior to January 1, 2010, the requirements of existing and planned international standard systems for the aeronautical radionavigation service which cannot be met in the 5000–5091 MHz band shall take precedence over other uses of this band.

After January 1, 2008, no new assignments shall be made to stations providing feeder links of nongeostationary mobile-satellite systems.

After January 1, 2010, the fixed-satellite service will become secondary to the aeronautical radionavigation service.

S5.445—Not used.

S5.446—*Additional allocation:* in the countries listed in Nos. S5.369 and S5.400, the band 5150–5216 MHz is also allocated to the radiodetermination-satellite service (space-to-Earth) on a primary basis, subject to agreement obtained under No. S9.21. In Region 2, the band is also allocated to the radiodetermination-satellite service (space-to-Earth) on a primary basis. In Regions 1 and 3, except those countries listed in Nos. S5.369 and S5.400, the band is also allocated to the radiodetermination-satellite service (space-to-Earth) on a secondary basis. The use by the radiodetermination-satellite service is limited to feeder links in conjunction with the radiodetermina-

tion-satellite service operating in the bands 1610–1626.5 MHz and/or 2483.5–2500 MHz. The total power flux-density at the Earth's surface shall in no case exceed -159 dBW/m² in any 4 kHz band for all angles of arrival.

S5.447—*Additional allocation*: in Germany, Austria, Belgium, Denmark, Spain, Finland, France, Greece, Israel, Italy, Japan, Jordan, Lebanon, Liechtenstein, Luxembourg, Malta, Morocco, Norway, Pakistan, the Netherlands, Portugal, Syria, the United Kingdom, Sweden, Switzerland, and Tunisia, the band 5150–5250 MHz is also allocated to the mobile service, on a primary basis, subject to agreement obtained under No. **S9.21**.

S5.447A—The allocation to the fixed-satellite service (Earth-to-space) is limited to feeder links of nongeostationary-satellite systems in the mobile-satellite service and is subject to coordination under No. **S9.11A**.

S5.447B—*Additional allocation*: The band 5150–5216 MHz is also allocated to the fixed-satellite service (space-to-Earth) on a primary basis. This allocation is limited to feeder links of nongeostationary-satellite systems in the mobile-satellite service and is subject to provisions of No. **S9.11A**. The power flux-density at the Earth's surface produced by space stations of the fixed-satellite service operating in the space-to-Earth direction in the band 5150–5216 MHz shall in no case exceed -164 dB(W/m²) in any 4 kHz band for all angles of arrival.

S5.447C—Administrations responsible for fixed-satellite service networks in the band 5150–5250 MHz operated under Nos. **S5.447A** and **S5.447B** shall coordinate on an equal basis in accordance with No. **S9.11A** with administrations responsible for nongeostationary-satellite networks operated under No. **S5.446** and brought into use prior to 17 November 1995. Satellite networks operated under No. **S5.446** brought into use after 17 November 1995 shall not claim protection from, and shall not cause harmful interference to, stations of the fixed-satellite service operated under Nos. **S5.447A** and **S5.447B**.

S5.447D—The allocation of the band 5250–5255 MHz to the space research service on a primary basis is limited to active spaceborne sensors. Other uses of the band by the space research service are on a secondary basis (WRC-97).

S5.448—*Additional allocation*: in Austria, Azerbaijan, Bulgaria, Libya, Mongolia, Kyrgyzstan, Slovakia, the Czech Republic, Romania, Turkmenistan and Ukraine, the band 5250–5350 MHz is also allocated to the radionavigation service on a primary basis (WRC-97).

S5.448A—The use of the frequency band 5250–5350 MHz by the earth exploration-satellite (active) and space research (active) services shall not constrain the future development and deployment of the radiolocation service (WRC-97).

S5.448B—The earth exploration-satellite (active) service operating in the band 5350–5460 MHz shall not cause harmful interference to, or constrain the use and development of, the aeronautical radionavigation service (WRC-97).

S5.449—The use of the band 5350–5470 MHz by the aeronautical radionavigation service is limited to airborne radars and associated airborne beacons.

S5.450—*Additional allocation*: in Austria, Azerbaijan, Bulgaria, the Islamic Republic of Iran, Mongolia, Kyrgyzstan, Slovakia, the Czech Republic, Romania, Turkmenistan, and Ukraine, the band 5470–5650 MHz is also allocated to the

aeronautical radionavigation service on a primary basis (WRC-97).

S5.451—*Additional allocation*: in the United Kingdom, the band 5470–5850 MHz is also allocated to the land mobile service on a secondary basis. The power limits specified in Nos. **S21.2**, **S21.3**, **S21.4**, and **S21.5** shall apply in the band 5725–5850 MHz.

S5.452—Between 5600 MHz and 5650 MHz, ground-based radars used for meteorological purposes are authorized to operate on a basis of equality with stations of the maritime radionavigation service.

S5.453—*Additional allocation*: in Saudi Arabia, Bahrain, Bangladesh, Brunei Darussalam, Cameroon, the Central African Republic, China, the Congo, the Republic of Korea, Egypt, the United Arab Emirates, Gabon, Guinea, India, Indonesia, the Islamic Republic of Iran, Iraq, Israel, Japan, Jordan, Kuwait, Lebanon, Libya, Madagascar, Malaysia, Nigeria, Oman, Pakistan, the Philippines, Qatar, Syria, Democratic People's Republic of Korea, Singapore, Swaziland, Tanzania, Chad, and Yemen, the band 5650–5850 MHz is also allocated to the fixed and mobile services on a primary basis (WRC-97).

S5.454—*Different category of service*: in Armenia, Azerbaijan, Belarus, Bulgaria, Georgia, Kazakstan, Mongolia, Uzbekistan, Kyrgyzstan, Russian Federation, Tajikistan, Turkmenistan, and Ukraine, the allocation of the band 5670–5725 MHz to the space research service is on a primary basis (see No. **S5.33**) (WRC-97).

S5.455—*Additional allocation*: in Armenia, Azerbaijan, Belarus, Bulgaria, Cuba, Georgia, Hungary, Kazakstan, Latvia, Moldova, Mongolia, Uzbekistan, Poland, Kyrgyzstan, Slovakia, Russian Federation, Tajikistan, Turkmenistan and Ukraine, the band 5670–5850 MHz is also allocated to the fixed service on a primary basis.

S5.456—*Additional allocation*: in Germany and in Cameroon, the band 5755–5850 MHz is also allocated to the fixed service on a primary basis.

S5.457—Not used.

S5.458—In the band 6425–7075 MHz, passive microwave sensor measurements are carried out over the oceans. In the band 7075–7250 MHz, passive microwave sensor measurements are carried out. Administrations should bear in mind the needs of the Earth exploration-satellite (passive) and space research (passive) services in their future planning of the bands 6425–7025 MHz and 7075–7250 MHz.

S5.458A—In making assignments in the band 6700–7075 MHz to space stations of the fixed-satellite service, administrations are urged to take all practicable steps to protect spectral line observations of the radio astronomy service in the band 6650–6675.2 MHz from harmful interference from unwanted emissions.

S5.458B—The space-to-Earth allocation to the fixed-satellite service in the band 6700–7075 MHz is limited to feeder links for nongeostationary satellite systems of the mobile-satellite service and is subject to coordination under No. **S9.11A**. The use of the band 6700–7075 MHz (space-to-Earth) by feeder links for nongeostationary satellite systems in the mobile-satellite service is not subject to No. **S22.2**.

S5.458C—Administrations making submissions in the band 7025–7075 MHz (Earth-to-space) for geostationary-satellite systems in the fixed-satellite service after 17 November

1995 shall consult on the basis of relevant ITU-R Recommendations with the administrations that have notified and brought into use nongeostationary-satellite systems in this frequency band before 18 November 1995 upon request of the latter administrations. This consultation shall be with a view to facilitating shared operation of both geostationary-satellite systems in the fixed-satellite service and nongeostationary-satellite systems in this band.

S5.459—*Additional allocation*: in Russian Federation, the frequency bands 7100–7155 MHz and 7190–7235 MHz are also allocated to the space operation service (Earth-to-space) on a primary basis, subject to agreement obtained under No. **S9.21** (WRC-97).

S5.460—*Additional allocation*: The band 7145–7235 MHz is also allocated to the space research (Earth-to-space) service on a primary basis, subject to agreement obtained under No. **S9.21**. The use of the band 7145–7190 MHz is restricted to deep space; no emissions to deep space shall be effected in the band 7190–7235 MHz.

S5.461—*Additional allocation*: The bands 7250–7375 MHz (space-to-Earth) and 7900–8025 MHz (Earth-to-space) are also allocated to the mobile-satellite service on a primary basis, subject to agreement obtained under No. **S9.21**.

S5.461A—The use of the band 7450–7550 MHz by the meteorological-satellite service (space-to-Earth) is limited to geostationary-satellite systems. Nongeostationary meteorological-satellite systems in this band notified before 30 November 1997 may continue to operate on a primary basis until the end of their lifetime (WRC-97).

S5.461B—The use of the band 7750–7850 MHz by the meteorological-satellite service (space-to-Earth) is limited to nongeostationary satellite systems (WRC-97).

S5.462—SUP - WRC-97.

S5.462A—In Regions 1 and 3 (except for Japan), in the band 8025–8400 MHz, the earth exploration-satellite service using geostationary satellites shall not produce a power flux-density in excess of the following provisional values for angles of arrival (q), without the consent of the affected administration:

- 174 dB(W/m²) in a 4 kHz band for $0^\circ \leq \theta < 5^\circ$
- 174 + 0.5 ($\theta - 5$) dB(W/m²) in a 4 kHz band for $5^\circ \leq \theta < 25^\circ$
- 164 dB(W/m²) in a 4 kHz band for $25^\circ \leq \theta \leq 90^\circ$

These values are subject to study under Resolution 124 (WRC-97) (WRC-97).

S5.463—Aircraft stations are not permitted to transmit in the band 8025–8400 MHz (WRC-97).

S5.464—SUP - WRC-97.

S5.465—In the space research service, the use of the band 8400–8450 MHz is limited to deep space.

S5.466—*Different category of service*: in Israel, Malaysia, Singapore, and Sri Lanka, the allocation of the band 8400–8500 MHz to the space research service is on a secondary basis (see No. **S5.32**) (WRC-97).

S5.467—*Alternative allocation*: in the United Kingdom, the band 8400–8500 MHz is allocated to the radiolocation and space research services on a primary basis.

S5.468—*Additional allocation*: in Saudi Arabia, Bahrain, Bangladesh, Brunei Darussalam, Burundi, Cameroon, China,

the Congo, Costa Rica, Egypt, the United Arab Emirates, Gabon, Guyana, Indonesia, the Islamic Republic of Iran, Iraq, Jamaica, Jordan, Kuwait, Lebanon, Libya, Malaysia, Mali, Morocco, Mauritania, Nepal, Nigeria, Oman, Pakistan, Qatar, Syria, Democratic People's Republic of Korea, Senegal, Singapore, Somalia, Swaziland, Tanzania, Chad, Togo, Tunisia, and Yemen, the band 8500–8750 MHz is also allocated to the fixed and mobile services on a primary basis (WRC-97).

S5.469—*Additional allocation*: in Armenia, Azerbaijan, Belarus, Bulgaria, Georgia, Hungary, Kazakhstan, Lithuania, Moldova, Mongolia, Uzbekistan, Poland, Kyrgyzstan, Slovakia, the Czech Republic, Romania, Russian Federation, Tajikistan, Turkmenistan, and Ukraine, the band 8500–8750 MHz is also allocated to the land mobile and radionavigation services on a primary basis.

S5.469A—In the band 8550–8650 MHz, stations in the earth exploration-satellite service (active) and space research service (active) shall not cause harmful interference to, or constrain the use and development of, stations of the radiolocation service (WRC-97).

S5.461B—The use of the band 7750–7850 MHz by the meteorological-satellite service (space-to-Earth) is limited to nongeostationary satellite systems (WRC-97).

S5.462—SUP - WRC-97.

S5.462A—In Regions 1 and 3 (except for Japan), in the band 8025–8400 MHz, the earth exploration-satellite service using geostationary satellites shall not produce a power flux-density in excess of the following provisional values for angles of arrival (q), without the consent of the affected administration:

- 174 dB(W/m²) in a 4 kHz band for $0^\circ \leq \theta < 5^\circ$
- 174 + 0.5 ($\theta - 5$) dB(W/m²) in a 4 kHz band for $5^\circ \leq \theta < 25^\circ$
- 164 dB(W/m²) in a 4 kHz band for $25^\circ \leq \theta \leq 90^\circ$

These values are subject to study under Resolution 124 (WRC-97) (WRC-97).

S5.463—Aircraft stations are not permitted to transmit in the band 8025–8400 MHz (WRC-97).

S5.464—SUP - WRC-97.

S5.465—In the space research service, the use of the band 8400–8450 MHz is limited to deep space.

S5.466—*Different category of service*: in Israel, Malaysia, Singapore and Sri Lanka, the allocation of the band 8400–8500 MHz to the space research service is on a secondary basis (see No. **S5.32**) (WRC-97).

S5.467—*Alternative allocation*: in the United Kingdom, the band 8400–8500 MHz is allocated to the radiolocation and space research services on a primary basis.

S5.468—*Additional allocation*: in Saudi Arabia, Bahrain, Bangladesh, Brunei Darussalam, Burundi, Cameroon, China, the Congo, Costa Rica, Egypt, the United Arab Emirates, Gabon, Guyana, Indonesia, the Islamic Republic of Iran, Iraq, Jamaica, Jordan, Kuwait, Lebanon, Libya, Malaysia, Mali, Morocco, Mauritania, Nepal, Nigeria, Oman, Pakistan, Qatar, Syria, Democratic People's Republic of Korea, Senegal, Singapore, Somalia, Swaziland, Tanzania, Chad, Togo, Tunisia, and Yemen, the band 8500–8750 MHz is also allocated to the fixed and mobile services on a primary basis (WRC-97).

S5.469—*Additional allocation:* in Armenia, Azerbaijan, Belarus, Bulgaria, Georgia, Hungary, Kazakhstan, Lithuania, Moldova, Mongolia, Uzbekistan, Poland, Kyrgyzstan, Slovakia, the Czech Republic, Romania, Russian Federation, Tajikistan, Turkmenistan, and Ukraine, the band 8500–8750 MHz is also allocated to the land mobile and radionavigation services on a primary basis.

S5.469A—In the band 8550–8650 MHz, stations in the earth exploration-satellite service (active) and space research service (active) shall not cause harmful interference to, or constrain the use and development of, stations of the radiolocation service (WRC-97).

S5.470—The use of the band 8750–8850 MHz by the aeronautical radionavigation service is limited to airborne Doppler navigation aids on a centre frequency of 8800 MHz.

S5.471—*Additional allocation:* in Algeria, Germany, Bahrain, Belgium, China, the United Arab Emirates, France, Greece, Indonesia, the Islamic Republic of Iran, Libya, the Netherlands, Qatar, and Sudan, the bands 8825–8850 MHz and 9000–9200 MHz are also allocated to the maritime radionavigation service, on a primary basis, for use by shore-based radars only.

S5.472—In the bands 8850–9000 MHz and 9200–9225 MHz, the maritime radionavigation service is limited to shore-based radars.

S5.473—*Additional allocation:* in Armenia, Austria, Azerbaijan, Belarus, Bulgaria, Cuba, Georgia, Hungary, Kazakhstan, Moldova, Mongolia, Uzbekistan, Poland, Kyrgyzstan, Slovakia, the Czech Republic, Romania, Russian Federation, Tajikistan, Turkmenistan, and Ukraine, the bands 8850–9000 MHz and 9200–9300 MHz are also allocated to the radionavigation service on a primary basis.

S5.474—In the band 9200–9500 MHz, search and rescue transponders (SART) may be used, having due regard to the appropriate ITU-R Recommendation (see also Article S31).

S5.475—The use of the band 9300–9500 MHz by the aeronautical radionavigation service is limited to airborne weather radars and ground-based radars. In addition, ground-based radar beacons in the aeronautical radionavigation service are permitted in the band 9300–9320 MHz on condition that harmful interference is not caused to the maritime radionavigation service. In the band 9300–9500 MHz, ground-based radars used for meteorological purposes have priority over other radiolocation devices.

S5.476—In the band 9300–9320 MHz in the radionavigation service, the use of shipborne radars, other than those existing on 1 January 1976, is not permitted until 1 January 2001.

S5.476A—In the band 9500–9800 MHz, stations in the earth exploration-satellite service (active) and space research service (active) shall not cause harmful interference to, or constrain the use and development of, stations of the radionavigation and radiolocation services (WRC-97).

S5.477—*Different category of service:* in Algeria, Saudi Arabia, Austria, Bahrain, Bangladesh, Brunei Darussalam, Cameroon, the Republic of Korea, Egypt, the United Arab Emirates, Eritrea, Ethiopia, Guyana, India, Indonesia, the Islamic Republic of Iran, Iraq, Jamaica, Japan, Jordan, Kuwait, Lebanon, Liberia, Malaysia, Nigeria, Oman, Pakistan, Qatar, Democratic People's Republic of Korea, Sin-

gapore, Somalia, Sudan, Sweden, Trinidad and Tobago, and Yemen, the allocation of the band 9800–10000 MHz to the fixed service is on a primary basis (see No. S5.33) (WRC-97).

S5.478—*Additional allocation:* in Azerbaijan, Bulgaria, Kazakhstan, Mongolia, Kyrgyzstan, Azerbaijan, Slovakia, the Czech Republic, Romania, Turkmenistan, and Ukraine, the band 9800–10000 MHz is also allocated to the radionavigation service on a primary basis (WRC-97).

S5.479—The band 9975–10025 MHz is also allocated to the meteorological-satellite service on a secondary basis for use by weather radars.

S5.480—*Additional allocation:* in Brazil, Costa Rica, Ecuador, Guatemala, Honduras, and Mexico, the band 10–10.45 GHz is also allocated to the fixed and mobile services on a primary basis (WRC-97).

S5.481—*Additional allocation:* in Germany, Angola, China, Ecuador, Spain, Japan, Morocco, Nigeria, Oman, Democratic People's Republic of Korea, Sweden, Tanzania, and Thailand, the band 10.45–10.5 GHz is also allocated to the fixed and mobile services on a primary basis.

S5.482—In the band 10.6–10.68 GHz, stations of the fixed and mobile, except aeronautical mobile, services shall be limited to a maximum equivalent isotropically radiated power of 40 dBW and the power delivered to the antenna shall not exceed –3 dBW. These limits may be exceeded subject to agreement obtained under No. S9.21. However, in Saudi Arabia, Armenia, Azerbaijan, Bahrain, Bangladesh, Belarus, China, the United Arab Emirates, Georgia, India, Indonesia, the Islamic Republic of Iran, Iraq, Japan, Kazakhstan, Kuwait, Latvia, Lebanon, Moldova, Nigeria, Uzbekistan, Pakistan, the Philippines, Qatar, Syria, Kyrgyzstan, Russian Federation, Tajikistan, Turkmenistan, and Ukraine, the restrictions on the fixed and mobile, except aeronautical mobile, services are not applicable.

S5.483—*Additional allocation:* in Saudi Arabia, Armenia, Azerbaijan, Bahrain, Belarus, Bosnia and Herzegovina, China, Colombia, the Republic of Korea, Costa Rica, Egypt, the United Arab Emirates, Georgia, the Islamic Republic of Iran, Iraq, Israel, Japan, Jordan, Kazakhstan, Kuwait, Latvia, Lebanon, Moldova, Mongolia, Uzbekistan, Pakistan, Qatar, Kyrgyzstan, Democratic People's Republic of Korea, Romania, Russian Federation, Tajikistan, Turkmenistan, Ukraine, Yemen, and Yugoslavia, the band 10.68–10.7 GHz is also allocated to the fixed and mobile, except aeronautical mobile, services on a primary basis. Such use is limited to equipment in operation by 1 January 1985 (WRC-97).

S5.484—In Region 1, the use of the band 10.7–11.7 GHz by the fixed-satellite service (Earth-to-space) is limited to feeder links for the broadcasting-satellite service.

S5.484A—The use of the bands 10.95–11.2 GHz (space-to-Earth), 11.45–11.7 GHz (space-to-Earth), 11.7–12.2 GHz (space-to-Earth) in Region 2, 12.2–12.75 GHz (space-to-Earth) in Region 3, 12.5–12.75 GHz (space-to-Earth) in Region 1, 13.75–14.5 GHz (Earth-to-space), 17.8–18.6 GHz (space-to-Earth), 19.7–20.2 GHz (space-to-Earth), 27.5–28.6 GHz (Earth-to-space), and 29.5–30 GHz (Earth-to-space) by nongeostationary and geostationary-satellite systems in the fixed-satellite service is subject to the provisions of Resolution 130 (WRC-97). The use of the band 17.8–18.1 GHz (space-to-Earth) by nongeostationary fixed-satellite service

systems is also subject to the provisions of Resolution 538 (WRC-97).

S5.485—In Region 2, in the band 11.7–12.2 GHz, transponders on space stations in the fixed-satellite service may be used additionally for transmissions in the broadcasting-satellite service, provided that such transmissions do not have a maximum eirp greater than 53 dBW per television channel and do not cause greater interference or require more protection from interference than the coordinated fixed-satellite service frequency assignments. With respect to the space services, this band shall be used principally for the fixed-satellite service.

S5.486—*Different category of service:* in Mexico and the United States, the allocation of the band 11.7–12.1 GHz to the fixed service is on a secondary basis (see No. S5.32).

S5.487—In the band 11.7–12.5 GHz in Regions 1 and 3, the fixed, fixed-satellite, mobile, except aeronautical mobile, and broadcasting services, in accordance with their respective allocations, shall not cause harmful interference to broadcasting-satellite stations operating in accordance with the provisions of Appendix S30.

S5.487A—*Additional allocation:* in Region 1, the band 11.7–12.5 GHz, in Region 2, the band 12.2–12.7 GHz and, in Region 3, the band 11.7–12.2 GHz are also allocated to the fixed-satellite service (space-to-Earth) on a primary basis, limited to nongeostationary systems and subject to the provisions of Resolution 538 (WRC-97) (WRC-97).

S5.488—The use of the bands 11.7–12.2 GHz by the fixed-satellite service in Region 2 and 12.2–12.7 GHz by the broadcasting-satellite service in Region 2 is limited to national and subregional systems. The use of the band 11.7–12.2 GHz by the fixed-satellite service in Region 2 is subject to previous agreement between the administrations concerned and those having services, operating or planned to operate in accordance with the Table, that may be affected (see Articles S9 and S11). For the use of the band 12.2–12.7 GHz by the broadcasting-satellite service in Region 2, see Appendix S30.

S5.489—*Additional allocation:* in Peru, the band 12.1–12.2 GHz is also allocated to the fixed service on a primary basis.

S5.490—In Region 2, in the band 12.2–12.7 GHz, existing and future terrestrial radiocommunication services shall not cause harmful interference to the space services operating in conformity with the broadcasting-satellite Plan for Region 2 contained in Appendix S30.

S5.491—*Additional allocation:* in Region 3, the band 12.2–12.5 GHz is also allocated to the fixed-satellite (space-to-Earth) service on a primary basis, limited to national and subregional systems. The power flux-density limits in Article S21, Table S21-4 shall apply to this frequency band. The introduction of the service in relation to the broadcasting-satellite service in Region 1 shall follow the procedures specified in Article 7 of Appendix S30, with the applicable frequency band extended to cover 12.2–12.5 GHz.

S5.492—Assignments to stations of the broadcasting-satellite service in conformity with the appropriate regional Plan in Appendix S30 may also be used for transmissions in the fixed-satellite service (space-to-Earth), provided that such transmissions do not cause more interference or require more protection from interference than the broadcasting-satellite

service transmissions operating in conformity with this Plan. With respect to the space services, this band shall be used principally for the broadcasting-satellite service (WRC-97).

S5.493—The broadcasting-satellite service in the band 12.5–12.75 GHz in Region 3 is limited to a power flux-density not exceeding $-111 \text{ dB(W/m}^2\text{)}/27 \text{ MHz}$ for all conditions and for all methods of modulation at the edge of the service area (WRC-97).

S5.494—*Additional allocation:* in Algeria, Angola, Saudi Arabia, Bahrain, Cameroon, the Central African Republic, the Congo, Côte d'Ivoire, Egypt, the United Arab Emirates, Eritrea, Ethiopia, Gabon, Ghana, Guinea, Iraq, Israel, Jordan, Kuwait, Lebanon, Libya, Madagascar, Mali, Morocco, Mongolia, Nigeria, Qatar, Democratic Republic of the Congo, Syria, Senegal, Somalia, Sudan, Chad, Togo, and Yemen, the band 12.5–12.75 GHz is also allocated to the fixed and mobile, except aeronautical mobile, services on a primary basis (WRC-97).

S5.495—*Additional allocation:* in Bosnia and Herzegovina, Croatia, Denmark, France, Greece, Liechtenstein, Monaco, Norway, Uganda, Portugal, Romania, Slovenia, Switzerland, Tanzania, Tunisia, and Yugoslavia, the band 12.5–12.75 GHz is also allocated to the fixed and mobile, except aeronautical mobile, services on a secondary basis (WRC-97).

S5.496—*Additional allocation:* in Austria, Azerbaijan, Kyrgyzstan, Turkmenistan, and Ukraine, the band 12.5–12.75 GHz is also allocated to the fixed service and the mobile, except aeronautical mobile, service on a primary basis. However, stations in these services shall not cause harmful interference to fixed-satellite service earth stations of countries in Region 1 other than those listed in this footnote. Coordination of these earth stations is not required with stations of the fixed and mobile services of the countries listed in this footnote. The power flux-density limit at the Earth's surface given in Article S21, Table S21-4, for the fixed-satellite service shall apply on the territory of the countries listed in this footnote (WRC-97).

S5.497—The use of the band 13.25–13.4 GHz by the aeronautical radionavigation service is limited to Doppler navigation aids.

S5.498—SUP - WRC-97.

S5.498A—The Earth exploration-satellite (active) and space research (active) services operating in the band 13.25–13.4 GHz shall not cause harmful interference to, or constrain the use and development of, the aeronautical radionavigation service (WRC-97).

S5.499—*Additional allocation:* in Bangladesh, India, and Pakistan, the band 13.25–14 GHz is also allocated to the fixed service on a primary basis.

S5.500—*Additional allocation:* in Algeria, Angola, Saudi Arabia, Bahrain, Brunei Darussalam, Cameroon, the Republic of Korea, Egypt, the United Arab Emirates, Gabon, Indonesia, the Islamic Republic of Iran, Iraq, Israel, Jordan, Kuwait, Lebanon, Madagascar, Malaysia, Mali, Malta, Morocco, Mauritania, Nigeria, Pakistan, Qatar, Syria, Senegal, Singapore, Sudan, Chad, and Tunisia, the band 13.4–14 GHz is also allocated to the fixed and mobile services on a primary basis (WRC-97).

S5.501—*Additional allocation*: in Austria, Azerbaijan, Bulgaria, Hungary, Japan, Mongolia, Kyrgyzstan, Romania, the United Kingdom, Turkmenistan, and Ukraine, the band 13.4–14 GHz is also allocated to the radionavigation service on a primary basis (WRC-97).

S5.501A—The allocation of the band 13.4–13.75 GHz to the space research service on a primary basis is limited to active spaceborne sensors. Other uses of the band by the space research service are on a secondary basis (WRC-97).

S5.501B—In the band 13.4–13.75 GHz, the Earth exploration-satellite (active) and space research (active) services shall not cause harmful interference to, or constrain the use and development of, the radiolocation service (WRC-97).

S5.502—In the band 13.75–14 GHz, the eirp of any emission from an earth station in the fixed-satellite service shall be at least 68 dBW, and should not exceed 85 dBW, with a minimum antenna diameter of 4.5 m. In addition the eirp, averaged over one second, radiated by a station in the radiolocation or radionavigation services toward the geostationary-satellite orbit shall not exceed 59 dBW.

S5.503—In the band 13.75–14 GHz, geostationary space stations in the space research service for which information for advance publication has been received by the Bureau prior to 31 January 1992 shall operate on an equal basis with stations in the fixed-satellite service; after that date, new geostationary space stations in the space research service will operate on a secondary basis. The eirp density of emissions from any earth station in the fixed-satellite service shall not exceed 71 dBW in any 6 MHz band in the frequency range 13.772–13.778 GHz until those geostationary space stations in the space research service for which information for advance publication has been received by the Bureau prior to 31 January 1992 cease to operate in this band. Automatic power control may be used to increase the eirp density above 71 dBW in any 6 MHz band in this frequency range to compensate for rain attenuation, to the extent that the power-flux density at the fixed-satellite service space station does not exceed the value resulting from use of an eirp of 71 dBW in any 6 MHz band in clear sky conditions.

S5.503A—Until January 1, 2000, stations in the fixed-satellite service shall not cause harmful interference to nongeostationary space stations in the space research and Earth exploration-satellite services. After that date, these nongeostationary space stations will operate on a secondary basis in relation to the fixed-satellite service. Additionally, when planning earth stations in the fixed-satellite service to be brought into service between January 1, 2000 and January 1, 2001, in order to accommodate the needs of spaceborne precipitation radars operating in the band 13.793–13.805 GHz, advantage should be taken of the consultation process and the information given in Recommendation ITU-R SA.1071.

S5.504—The use of the band 14–14.3 GHz by the radionavigation service shall be such as to provide sufficient protection to space stations of the fixed-satellite service.

S5.505—*Additional allocation*: in Algeria, Angola, Saudi Arabia, Australia, Bahrain, Bangladesh, Botswana, Brunei Darussalam, Cameroon, China, the Congo, the Republic of Korea, Egypt, the United Arab Emirates, Gabon, Guatemala, Guinea, India, Indonesia, the Islamic Republic of Iran, Iraq, Israel, Japan, Jordan, Kuwait, Lesotho, Lebanon, Malaysia,

Mali, Morocco, Mauritania, Oman, Pakistan, the Philippines, Qatar, Syria, the Democratic People's Republic of Korea, Senegal, Singapore, Somalia, Sudan, Swaziland, Tanzania, Chad, and Yemen, the band 14–14.3 GHz is also allocated to the fixed service on a primary basis (WRC-97).

S5.506—The band 14–14.5 GHz may be used, within the fixed-satellite service (Earth-to-space), for feeder links for the broadcasting-satellite service, subject to coordination with other networks in the fixed-satellite service. Such use of feeder links is reserved for countries outside Europe.

S5.507—Not used.

S5.508—*Additional allocation*: in Germany, Austria, Bosnia and Herzegovina, France, Greece, Ireland, Iceland, Italy, the former Yugoslav Republic of Macedonia, Libya, Liechtenstein, Portugal, the United Kingdom, Slovenia, Switzerland, Turkey, and Yugoslavia, the band 14.25–14.3 GHz is also allocated to the fixed service on a primary basis (WRC-97).

S5.509—*Additional allocation*: in Japan and Pakistan the band 14.25–14.3 GHz is also allocated to the mobile, except aeronautical mobile, service on a primary basis.

S5.510—The use of the band 14.5–14.8 GHz by the fixed-satellite service (Earth-to-space) is limited to feeder links for the broadcasting-satellite service. This use is reserved for countries outside Europe.

S5.511—*Additional allocation*: in Saudi Arabia, Bahrain, Bosnia and Herzegovina, Cameroon, Egypt, the United Arab Emirates, Guinea, the Islamic Republic of Iran, Iraq, Israel, Kuwait, Lebanon, Libya, Pakistan, Qatar, Syria, Slovenia, Somalia, and Yugoslavia, the band 15.35–15.4 GHz is also allocated to the fixed and mobile services on a secondary basis (WRC-97).

S5.511A—Use of the band 15.43–15.63 GHz by the fixed-satellite service (space-to-Earth (see Resolution 123 (WRC-97)) and Earth-to-space) is limited to feeder links of nongeostationary systems in the mobile-satellite service, subject to coordination under No. S9.11A. In the space-to-Earth direction, the minimum earth station elevation angle above and gain toward the local horizontal plane and the minimum coordination distances to protect an earth station from harmful interference shall be in accordance with Recommendation ITU-R S.1341. Also in the space-to-Earth direction, harmful interference shall not be caused to stations of the radio astronomy service using the band 15.35–15.4 GHz. The threshold levels of interference and associated power flux-density limits that are detrimental to the radio astronomy service are given in Recommendation ITU-R RA.769-1. Special measures will need to be employed to protect the radio astronomy service in the band 15.35–15.4 GHz (WRC-97).

S5.511B—SUP – WRC-97.

S5.511C—Stations operating in the aeronautical radionavigation service shall limit the effective eirp in accordance with Recommendation ITU-R S.1340. The minimum coordination distance required to protect the aeronautical radionavigation stations (No. S4.10 applies) from harmful interference from feeder-link earth stations and the maximum eirp transmitted toward the local horizontal plane by a feeder-link earth station shall be in accordance with Recommendation ITU-R S.1340 (WRC-97).

S5.511D—Fixed-satellite service systems for which complete information for advance publication has been received by the Bureau by 21 November 1997 may operate in the bands 15.4–15.43 GHz and 15.63–15.7 GHz in the space-to-Earth direction and 15.63–15.65 GHz in the Earth-to-space direction. In the bands 15.4–15.43 GHz and 15.65–15.7 GHz, emissions from a nongeostationary space station shall not exceed the power flux-density limits at the Earth's surface of -146 dB(W/m²/MHz) for any angle of arrival. In the band 15.63–15.65 GHz, where an administration plans emissions from a nongeostationary space station that exceed -146 dB(W/m²/MHz) for any angle of arrival, it shall coordinate under No. **S9.11A** with the affected administrations. Stations in the fixed-satellite service operating in the band 15.63–15.65 GHz in the Earth-to-space direction shall not cause harmful interference to stations in the aeronautical radionavigation service (No. **S4.10** applies) (WRC-97).

S5.512—*Additional allocation*: in Algeria, Angola, Saudi Arabia, Austria, Bahrain, Bangladesh, Bosnia and Herzegovina, Brunei Darussalam, Cameroon, the Congo, Costa Rica, Egypt, El Salvador, the United Arab Emirates, Finland, Guatemala, India, Indonesia, the Islamic Republic of Iran, Jordan, Kuwait, Libya, Malaysia, Morocco, Mozambique, Nepal, Nicaragua, Oman, Pakistan, Qatar, Singapore, Slovenia, Somalia, Sudan, Swaziland, Tanzania, Chad, Yemen, and Yugoslavia, the band 15.7–17.3 GHz is also allocated to the fixed and mobile services on a primary basis (WRC-97).

S5.513—*Additional allocation*: in Israel, the band 15.7–17.3 GHz is also allocated to the fixed and mobile services on a primary basis. These services shall not claim protection from, or cause harmful interference to, services operating in accordance with the Table in countries other than those included in No. **S5.512**.

S5.513A—Spaceborne active sensors operating in the band 17.2–17.3 GHz shall not cause harmful interference to, or constrain the development of, the radiolocation and other services allocated on a primary basis (WRC-97).

S5.514—*Additional allocation*: in Algeria, Germany, Angola, Saudi Arabia, Austria, Bahrain, Bangladesh, Bosnia and Herzegovina, Cameroon, Costa Rica, El Salvador, the United Arab Emirates, Finland, Guatemala, Honduras, India, the Islamic Republic of Iran, Iraq, Israel, Japan, Jordan, Kuwait, Libya, Nepal, Nicaragua, Oman, Pakistan, Qatar, Slovenia, Sudan, Sweden, and Yugoslavia, the band 17.3–17.7 GHz is also allocated to the fixed and mobile services on a secondary basis. The power limits given in Nos. **S21.3** and **S21.5** shall apply (WRC-97).

S5.515—In the band 17.3–17.8 GHz, sharing between the fixed-satellite service (Earth-to-space) and the broadcasting-satellite service shall also be in accordance with the provisions of §1 of Annex 4 of Appendix **S30A/30A**.

S5.516—The use of the band 17.3–18.1 GHz by geostationary-satellite systems in the fixed-satellite service (Earth-to-space) is limited to feeder links for the broadcasting-satellite service. For the use of the band 17.3–17.8 GHz in Region 2 by feeder links for the broadcasting-satellite service in the band 12.2–12.7 GHz, see Article **S11**. The use of the bands 17.3–18.1 GHz (Earth-to-space) in Regions 1 and 3 and 17.8–18.1 GHz (Earth-to-space) in Region 2 by nongeostationary-

satellite systems in the fixed-satellite service is subject to the provisions of Resolution 538 (WRC-97) (WRC-97).

S5.517—In Region 2, the allocation to the broadcasting-satellite service in the band 17.3–17.8 GHz shall come into effect on April 1, 2007. After that date, use of the fixed-satellite (space-to-Earth) service in the band 17.7–17.8 GHz shall not claim protection from and shall not cause harmful interference to operating systems in the broadcasting-satellite service.

S5.518—*Different category of service*: in Region 2, the allocation of the band 17.7–17.8 GHz to the mobile service is on a primary basis until March 31, 2007.

S5.519—*Additional allocation*: The band 18.1–18.3 GHz is also allocated to the meteorological-satellite service (space-to-Earth) on a primary basis. Its use is limited to geostationary satellites and shall be in accordance with the provisions of Article **S21**, Table **S21-4**.

S5.520—The use of the band 18.1–18.4 GHz by the fixed-satellite service (Earth-to-space) is limited to feeder links for the broadcasting-satellite service.

S5.521—*Alternative allocation*: in Germany, Denmark, the United Arab Emirates, Greece, Slovakia, and the Czech Republic, the band 18.1–18.4 GHz is allocated to the fixed, fixed-satellite (space-to-Earth), and mobile services on a primary basis (see No. **S5.33**). The provisions of No. **S5.519** also apply (WRC-97).

S5.522—In making assignments to stations in the fixed and mobile services, administrations are invited to take account of passive sensors in the Earth-exploration satellite and space research services operating in the band 18.6–18.8 GHz. In this band, administrations should endeavour to limit as far as possible both the power delivered by the transmitter to the antenna and the eirp to reduce the risk of interference to passive sensors to a minimum.

S5.523—In assigning frequencies to stations in the fixed-satellite service in the direction space-to-Earth, administrations are requested to limit as far as practicable the power flux-density at the Earth's surface in the band 18.6–18.8 GHz, to reduce the risk of interference to passive sensors in the earth exploration-satellite and space research services.

S5.523A—The use of the bands 18.8–19.3 GHz (space-to-Earth) and 28.6–29.1 GHz (Earth-to-space) by geostationary and nongeostationary fixed-satellite service networks is subject to the application of the provisions of No. **S9.11A**, and No. **S22.2** does not apply. Administrations having geostationary-satellite networks under coordination prior to 18 November 1995 shall cooperate to the maximum extent possible to coordinate pursuant to No. **S9.11A** with nongeostationary-satellite networks for which notification information has been received by the Bureau prior to that date, with a view to reaching results acceptable to all the parties concerned. Nongeostationary-satellite networks shall not cause unacceptable interference to geostationary fixed-satellite service networks for which complete Appendix **S4** notification information is considered as having been received by the Bureau prior to 18 November 1995 (WRC-97).

S5.523B—The use of the band 19.3–19.6 GHz (Earth-to-space) by the fixed-satellite service is limited to feeder links for nongeostationary-satellite systems in the mobile-satellite service. Such use is subject to the application of the provisions of No. **S9.11A**, and No. **S22.2** does not apply.

S5.523C—No. **S22.2** of the Radio Regulations shall continue to apply in the bands 19.3–19.6 GHz and 29.1–29.4 GHz, between feeder links of nongeostationary mobile-satellite service networks and those fixed-satellite service networks for which complete Appendix **S4** coordination information, or notification information, is considered as having been received by the Bureau prior to 18 November 1995 (WRC-97).

S5.523D—The use of the band 19.3–19.7 GHz (space-to-Earth) by geostationary fixed-satellite service systems and by feeder links for nongeostationary-satellite systems in the mobile-satellite service is subject to the application of the provisions of No. **S9.11A**, but not subject to the provisions of No. **S22.2**. The use of this band for other nongeostationary fixed-satellite service systems, or for the cases indicated in Nos. **S5.523C** and **S5.523E**, is not subject to the provisions of No. **S9.11A** and shall continue to be subject to Articles **S9** (except No. **S9.11A**) and **S11** procedures, and to the provisions of No. **S22.2** (WRC-97).

S5.523E—No. **S22.2** of the Radio Regulations shall continue to apply in the bands 19.6–19.7 GHz and 29.4–29.5 GHz, between feeder links of nongeostationary mobile-satellite service networks and those fixed-satellite service networks for which complete Appendix **S4** coordination information, or notification information, is considered as having been received by the Bureau by 21 November 1997 (WRC-97).

S5.524—*Additional allocation:* in Afghanistan, Algeria, Angola, Saudi Arabia, Bahrain, Bangladesh, Brunei Darussalam, Cameroon, China, the Congo, the Republic of Korea, Costa Rica, Egypt, the United Arab Emirates, Gabon, Guatemala, Guinea, India, Islamic Republic of Iran, Iraq, Israel, Japan, Jordan, Kuwait, Lebanon, Malaysia, Mali, Morocco, Mauritania, Nepal, Nigeria, Oman, Pakistan, the Philippines, Qatar, Dem. Rep. of the Congo, Syria, Democratic People's Republic of Korea, Singapore, Somalia, Sudan, Tanzania, Chad, Togo, and Tunisia, the band 19.7–21.2 GHz is also allocated to the fixed and mobile services on a primary basis. This additional use shall not impose any limitation on the power flux-density of space stations in the fixed-satellite service in the band 19.7–21.2 GHz and of space stations in the mobile-satellite service in the band 19.7–20.2 GHz where the allocation to the mobile-satellite service is on a primary basis in the latter band (WRC-97).

S5.525—To facilitate interregional coordination between networks in the mobile-satellite and fixed-satellite services, carriers in the mobile-satellite service that are most susceptible to interference shall, to the extent practicable, be located in the higher parts of the bands 19.7–20.2 GHz and 29.5–30 GHz.

S5.526—In the bands 19.7–20.2 GHz and 29.5–30 GHz in Region 2, and in the bands 20.1–20.2 GHz and 29.9–30 GHz in Regions 1 and 3, networks which are both in the fixed-satellite service and in the mobile-satellite service may include links between earth stations at specified or unspecified points or while in motion, through one or more satellites for point-to-point and point-to-multipoint communications.

S5.527—In the bands 19.7–20.2 GHz and 29.5–30 GHz, the provisions of No. **S4.10** do not apply with respect to the mobile-satellite service.

S5.528—The allocation to the mobile-satellite service is intended for use by networks that use narrow spot-beam antennas and other advanced technology at the space stations. Administrations operating systems in the mobile-satellite service in the band 19.7–20.1 GHz in Region 2 and in the band 20.1–20.2 GHz shall take all practicable steps to ensure the continued availability of these bands for administrations operating fixed and mobile systems in accordance with the provisions of No. **S5.524**.

S5.529—The use of the bands 19.7–20.1 GHz and 29.5–29.9 GHz by the mobile-satellite service in Region 2 is limited to satellite networks that are both in the fixed-satellite service and in the mobile-satellite service as described in No. **S5.526**.

S5.530—In Regions 1 and 3, the allocation to the broadcasting-satellite service in the band 21.4–22 GHz shall come into effect on April 1, 2007. The use of this band by the broadcasting-satellite service after that date and on an interim basis prior to that date is subject to the provisions of Resolution 525 (WARC-92).

S5.531—*Additional allocation:* in Japan, the band 21.4–22 GHz is also allocated to the broadcasting service on a primary basis.

S5.532—The use of the band 22.21–22.5 GHz by the Earth exploration-satellite (passive) and space research (passive) services shall not impose constraints upon the fixed and mobile, except aeronautical mobile, services.

S5.533—The intersatellite service shall not claim protection from harmful interference from airport surface detection equipment stations of the radionavigation service.

S5.534—*Additional allocation:* in Japan, the band 24.65–25.25 GHz is also allocated to the radionavigation service on a primary basis until 2008.

S5.535—In the band 24.75–25.25 GHz, feeder links to stations of the broadcasting-satellite service shall have priority over other uses in the fixed-satellite service (Earth-to-space). Such other uses shall protect and shall not claim protection from existing and future operating feeder-link networks to such broadcasting satellite stations.

S5.535A—The use of the band 29.1–29.5 GHz (Earth-to-space) by the fixed-satellite service is limited to geostationary-satellite systems and feeder links to nongeostationary-satellite systems in the mobile-satellite service. Such use is subject to the application of the provisions of No. **S9.11A**, but not subject to the provisions of No. **S22.2**, except as indicated in Nos. **S5.523C** and **S5.523E** where such use is not subject to the provisions of No. **S9.11A** and shall continue to be subject to Articles **S9** (except No. **S9.11A**) and **S11** procedures, and to the provisions of No. **S22.2** (WRC-97).

S5.536—Use of the 25.25–27.5 GHz band by the intersatellite service is limited to space research and Earth exploration-satellite applications, and also transmissions of data originating from industrial and medical activities in space.

S5.536A—Administrations installing earth exploration-satellite earth stations cannot claim protection from fixed and mobile stations operated by neighbouring administrations. In addition, earth stations operating in the earth exploration-satellite service should take into account Recommendation ITU-R SA. 1278 (WRC-97).

S5.536B—In Germany, Saudi Arabia, Austria, Belgium, Brazil, Bulgaria, China, the Republic of Korea, Denmark,

Egypt, United Arab Emirates, Spain, Estonia, Finland, France, Hungary, India, Islamic Republic of Iran, Ireland, Israel, Italy, Jordan, Kenya, Kuwait, Lebanon, Libya, Liechtenstein, Lithuania, Moldova, Norway, Oman, Uganda, Pakistan, the Philippines, Poland, Portugal, Syria, Slovakia, Czech Republic, Romania, the United Kingdom, Singapore, Sweden, Switzerland, Tanzania, Turkey, Viet Nam, and Zimbabwe, earth stations operating in the Earth exploration-satellite service in the band 25.5–27 GHz shall not claim protection from, or constrain the use and deployment of, stations of the fixed and mobile services (WRC-97).

S5.537—Space services using nongeostationary satellites operating in the intersatellite service in the band 27–27.5 GHz are exempt from the provisions of No. **S22.2**.

S5.538—*Additional allocation*: The bands 27.500–27.501 GHz and 29.999–30.000 GHz are also allocated to the fixed-satellite service (space-to-Earth) on a primary basis for the beacon transmissions intended for uplink power control. Such space-to-Earth transmissions shall not exceed an equivalent isotropically radiated power (eirp) of 10 dBW in the direction of adjacent satellites on the geostationary-satellite orbit. In the band 27.500–27.501 GHz, such space-to-Earth transmissions shall not produce a power flux-density in excess of the values specified in Article **S21**, Table **S21-4** on the Earth's surface.

S5.539—The band 27.5–30 GHz may be used by the fixed-satellite service (Earth-to-space) for the provision of feeder links for the broadcasting-satellite service.

S5.540—*Additional allocation*: The band 27.501–29.999 GHz is also allocated to the fixed-satellite service (space-to-Earth) on a secondary basis for beacon transmissions intended for uplink power control.

S5.541—In the band 28.5–30 GHz, the earth exploration-satellite service is limited to the transfer of data between stations and not to the primary collection of information by means of active or passive sensors.

S5.541A—Feeder links of nongeostationary networks in the mobile-satellite service and geostationary networks in the fixed-satellite service operating in the band 29.1–29.5 GHz (Earth-to-space) shall employ uplink adaptive power control or other methods of fade compensation, such that the earth station transmissions shall be conducted at the power level required to meet the desired link performance while reducing the level of mutual interference between both networks. These methods shall apply to networks for which Appendix **S4** coordination information is considered as having been received by the Bureau after 17 May 1996 and until they are changed by a future competent world radiocommunication conference. Administrations submitting Appendix **S4** information for coordination before this date are encouraged to utilize these techniques to the extent practicable. These methods are also subject to review by ITU-R (see Resolution 121 (Rev. WRC-97)) (WRC-97).

S5.542—*Additional allocation*: in Algeria, Saudi Arabia, Bahrain, Bangladesh, Brunei Darussalam, Cameroon, China, the Congo, the Republic of Korea, Egypt, the United Arab Emirates, Eritrea, Ethiopia, Guinea, India, the Islamic Republic of Iran, Iraq, Japan, Jordan, Kuwait, Lebanon, Malaysia, Mali, Morocco, Mauritania, Nepal, Pakistan, the Philippines, Qatar, Syria, Democratic People's Republic of Korea, Somalia,

lia, Sudan, Sri Lanka, and Chad, the band 29.5–31 GHz is also allocated to the fixed and mobile services on a secondary basis. The power limits specified in Nos. **S21.3** and **S21.5** shall apply (WRC-97).

S5.543—The band 29.95–30 GHz may be used for space-to-space links in the Earth exploration-satellite service for telemetry, tracking, and control purposes, on a secondary basis.

S5.544—In the band 31–31.3 GHz the power flux-density limits specified in Article **S21**, Table **S21-4** shall apply to the space research service.

S5.545—*Different category of service*: in Armenia, Azerbaijan, Belarus, Bulgaria, Georgia, Kazakstan, Mongolia, Kyrgyzstan, Russian Federation, Tajikistan, Turkmenistan, and Ukraine, the allocation of the band 31–31.3 GHz to the space research service is on a primary basis (see No. **S5.33**) (WRC-97).

S5.546—*Different category of service*: in Saudi Arabia, Armenia, Azerbaijan, Belarus, Bulgaria, Egypt, United Arab Emirates, Spain, Estonia, Finland, Georgia, Hungary, the Islamic Republic of Iran, Israel, Jordan, Kazakstan, Latvia, Lebanon, Moldova, Mongolia, Uzbekistan, Poland, Syria, Kyrgyzstan, Romania, the United Kingdom, Russian Federation, Tajikistan, Turkmenistan, Turkey, and Ukraine, the allocation of the band 31.5–31.8 GHz to the fixed and mobile, except aeronautical mobile, services is on a primary basis (see No. **S5.33**) (WRC-97).

S5.547—The bands 31.8–33.4 GHz, 51.4–52.6 GHz, 55.78–59 GHz, and 64–66 GHz are available for high-density applications in the fixed service (see Resolution 726 (WRC-97)) (WRC-97).

S5.547A—Use of the band 31.8–33.4 GHz by the fixed service shall be in accordance with Resolution 126 (WRC-97) (WRC-97).

S5.547B—*Alternative allocation*: in the United States, the band 31.8–32 GHz is allocated to the radionavigation and space research (deep space) (space-to-Earth) services on a primary basis (WRC-97).

S5.547C—*Alternative allocation*: in the United States, the band 32–32.3 GHz is allocated to the intersatellite, radionavigation, and space research (deep space) (space-to-Earth) services on a primary basis (WRC-97).

S5.547D—*Alternative allocation*: in the United States, the band 32.3–33 GHz is allocated to the intersatellite and radionavigation services on a primary basis (WRC-97).

S5.547E—*Alternative allocation*: in the United States, the band 33–33.4 GHz is allocated to the radionavigation service on a primary basis (WRC-97).

S5.548—In designing systems for the intersatellite and radionavigation services in the band 32–33 GHz, and for the space research service (deep space) in the band 31.8–32.3 GHz, administrations shall take all necessary measures to prevent harmful interference between these services, bearing in mind the safety aspects of the radionavigation service (see Recommendation 707).

S5.549—*Additional allocation*: in Saudi Arabia, Bahrain, Bangladesh, Egypt, the United Arab Emirates, Gabon, Indonesia, the Islamic Republic of Iran, Iraq, Israel, Jordan, Kuwait, Lebanon, Libya, Malaysia, Mali, Malta, Morocco, Mauritania, Nepal, Nigeria, Oman, Pakistan, the Philippines,

Qatar, Democratic Republic of the Congo, Syria, Senegal, Singapore, Somalia, Sudan, Sri Lanka, Togo, Tunisia, and Yemen. the band 33.4–36 GHz is also allocated to the fixed and mobile services on a primary basis (WRC-97).

S5.550—*Different category of service:* in Armenia, Azerbaijan, Belarus, Bulgaria, Georgia, Kazakhstan, Mongolia, Uzbekistan, Kyrgyzstan, Russian Federation, Tajikistan, Turkmenistan, and Ukraine, the allocation of the band 34.7–35.2 GHz to the space research service is on a primary basis (see No. **S5.33**) (WRC-97).

S5.551—SUP - WRC-97.

S5.551A—In the band 35.5–36.0 GHz, active spaceborne sensors in the earth exploration-satellite and space research services shall not cause harmful interference to, claim protection from, or otherwise impose constraints on operation or development of the radiolocation service, the meteorological aids service and other services allocated on a primary basis (WRC-97).

S5.551B—The use of the band 41.5–42.5 GHz by the fixed-satellite service (space-to-Earth) is subject to Resolution 128 (WRC-97) (WRC-97).

S5.551C—*Alternative allocation:* in the French overseas territories in Regions 2 and 3, the Republic of Korea and India, the band 40.5–42.5 GHz is allocated to the broadcasting, broadcasting-satellite, and fixed services on a primary basis (WRC-97).

S5.551D—*Additional allocation:* in Algeria, Saudi Arabia, Bahrain, Benin, Cameroon, Egypt, United Arab Emirates, Israel, Jordan, Kuwait, Lebanon, Libya, Mali, Morocco, Mauritania, Nigeria, Oman, Qatar, Syria, Tunisia, and Yemcut, the band 40.5–42.5 GHz is also allocated to the fixed-satellite service (space-to-Earth) on a primary basis. The use of this band by the fixed-satellite service shall be in accordance with Resolution 134 (WRC-97) (WRC-97).

S5.551E—Use of the band 40.5–42.5 GHz by the fixed-satellite service shall be in accordance with Resolution 134 (WRC-97) (WRC-97).

S5.551F—*Different category of service:* in Japan, the allocation of the band 41.5–42.5 GHz to the mobile service is on a primary basis (see No. **S5.33**) (WRC-97).

S5.552—The allocation of the spectrum for the fixed-satellite service in the bands 42.5–43.5 GHz and 47.2–50.2 GHz for Earth-to-space transmission is greater than that in the band 37.5–39.5 GHz for space-to-Earth transmission in order to accommodate feeder links to broadcasting satellites. Administrations are urged to take all practicable steps to reserve the band 47.2–49.2 GHz for feeder links for the broadcasting-satellite service operating in the band 40.5–42.5 GHz.

S5.552A—The allocation to the fixed service in the bands 47.2–47.5 GHz and 47.9–48.2 GHz is designated for use by high-altitude platform stations. The use of the bands 47.2–47.5 GHz and 47.9–48.2 GHz is subject to the provisions of Resolution 122 (WRC-97) (WRC-97).

S5.553—In the bands 43.5–47 GHz, 66–71 GHz, 95–100 GHz, 134–142 GHz, 190–200 GHz, and 252–265 GHz, stations in the land mobile service may be operated subject to not causing harmful interference to the space radiocommunication services to which these bands are allocated (see No. **S5.43**).

S5.554—In the bands 43.5–47 GHz, 66–71 GHz, 95–100 GHz, 134–142 GHz, 190–200 GHz, and 252–265 GHz, satellite links connecting land stations at specified fixed points are also authorized when used in conjunction with the mobile-satellite service or the radionavigation-satellite service.

S5.555—*Additional allocation:* The bands 48.94–49.04 GHz, 97.88–98.08 GHz, 140.69–140.98 GHz, 144.68–144.98 GHz, 145.45–145.75 GHz, 146.82–147.12 GHz, 250–251 GHz, and 262.24–262.76 GHz are also allocated to the radio astronomy service on a primary basis.

S5.555A—The band 50.2–50.4 GHz is also allocated, on a primary basis, to the fixed and mobile services until July 1, 2000 (WRC-97).

S5.556—In the bands 51.4–54.25 GHz, 58.2–59 GHz, 64–65 GHz, 72.77–72.91 GHz, and 93.07–93.27 GHz, radio astronomy observations may be carried out under national arrangements.

S5.556A—Use of the bands 54.25–56.9 GHz, 57–58.2 GHz, and 59–59.3 GHz by the intersatellite service is limited to satellites in the geostationary-satellite orbit. The single-entry power flux-density at all altitudes from 0 km to 1000 km above the Earth's surface produced by a station in the intersatellite service, for all conditions and for all methods of modulation, shall not exceed -147 dB(W/m²/100 MHz) for all angles of arrival (WRC-97).

S5.556B—*Additional allocation:* in Japan, the band 54.25–55.78 GHz is also allocated to the mobile service on a primary basis for low-density use (WRC-97).

S5.557—*Additional allocation:* in Japan, the band 55.78–58.2 GHz is also allocated to the radiolocation service on a primary basis (WRC-97).

S5.558—In the bands 55.78–58.2 GHz, 59–64 GHz, 66–71 GHz, 116–134 GHz, 170–182 GHz, and 185–190 GHz, stations in the aeronautical mobile service may be operated subject to not causing harmful interference to the intersatellite service (see No. **S5.43**) (WRC-97).

S5.558A—Use of the band 56.9–57 GHz by intersatellite systems is limited to links between satellites in geostationary-satellite orbit and to transmissions from nongeostationary satellites in high-Earth orbit to those in low-Earth orbit. For links between satellites in the geostationary-satellite orbit, the single-entry power flux-density at all altitudes from 0 km to 1000 km above the Earth's surface, for all conditions and for all methods of modulation, shall not exceed -147 dB(W/m²/100 MHz) for all angles of arrival (WRC-97).

S5.559—In the bands 59–64 GHz and 126–134 GHz, airborne radars in the radiolocation service may be operated subject to not causing harmful interference to the intersatellite service (see No. **S5.43**).

S5.560—In the band 78–79 GHz radars located on space stations may be operated on a primary basis in the Earth exploration-satellite service and in the space research service.

S5.561—In the band 84–86 GHz, stations in the fixed, mobile and broadcasting services shall not cause harmful interference to broadcasting-satellite stations operating in accordance with the decisions of the appropriate frequency assignment planning conference for the broadcasting-satellite service.

SS.562—The use of the band 94–94.1 GHz by the Earth exploration-satellite (active) and space research (active) services is limited to spaceborne cloud radars (WRC-97).

SS.563—*Additional allocation:* in the United Kingdom, the band 182–185 GHz is also allocated to the fixed and mobile services on a primary basis.

SS.564—*Additional allocation:* in Germany, Argentina, Spain, Finland, France, India, Italy, and the Netherlands, the band 261–265 GHz is also allocated to the radio astronomy service on a primary basis (WRC-97).

SS.565—The frequency band 275–400 GHz may be used by administrations for experimentation with, and development of, various active and passive services. In this band a need has been identified for the following spectral line measurements for passive services:

radio astronomy service: 278–280 GHz and 343–348 GHz;

Earth exploration-satellite service (passive) and space research service (passive): 275–277 GHz, 300–302 GHz, 324–326 GHz, 345–347 GHz, 363–365 GHz, and 379–381 GHz.

Future research in this largely unexplored spectral region may yield additional spectral lines and continuum bands of interest to the passive services. Administrations are urged to take all practicable steps to protect these passive services from harmful interference until the next competent world radio-communication conference.

Old Numbering Scheme

459—In the Region 2 polar areas (north of 60°N and south of 60°S), which are subject to auroral disturbances, the aeronautical fixed service is the primary service in the band 160–190 kHz.

471—The bands 490–495 kHz and 505–510 kHz shall be subject to the provisions of No. 3018 until the entry into force of the reduced guardband in accordance with Resolution 210 (Mob-87).

472—The frequency 500 kHz is the international distress and calling frequency for Morse radiotelegraphy. The conditions for its use are prescribed in Articles 37, 38, N 38 and 60.

472A—In the maritime mobile service, the frequency 490 kHz is, from the date of full implementation of the GMDSS (see Resolution 331 (Mob-87)), to be used exclusively for the transmission by coast stations of navigational and meteorological warnings and urgent information to ships, by means of narrow-band direct-printing telegraphy. The conditions for use of the frequency 490 kHz are prescribed in Articles N 38 and 60 and Resolution 329 (Mob-87). In using the band 415–495 kHz for the aeronautical radionavigation service, administrations are requested to ensure that no harmful interference is caused to the frequency 490 kHz.

474—The conditions for the use of frequency 518 kHz by the maritime mobile service are prescribed in Articles 38, N38 and 60 (see Resolution 324 (Mob-87) and Article 14A).

480—In Region 2, the use of the band 1605–1705 kHz by stations of the broadcasting service is subject to the plan established by the Regional Administrative Radio Conference (Rio de Janeiro, 1988.)

In Region 2, in the band 1625–1705 kHz, the relationship between the broadcasting, fixed, and mobile services is shown in No. 419. However, the examination of frequency assignments to stations of the fixed and mobile services in the band 1625–1705 kHz under No. 1241 shall take account of the allotments appearing in the Plan established by the Regional Administrative Radio Conference (Rio de Janeiro, 1988).

591—Subject to agreement obtained under the procedure set forth in Article 14, the band 117.975–137 MHz is also allocated to the aeronautical mobile-satellite (R) service on a secondary basis and on the condition that harmful interference is not caused to the aeronautical mobile (R) service.

599A—The use of the band 137–138 MHz by the mobile-satellite service is subject to the application of the coordination and notification procedures set forth in Resolution 46. However, coordination of a space station of the mobile-satellite service with respect to terrestrial services is required only if the power flux-density produced by the station exceeds $-125 \text{ dB(W/m}^2/4 \text{ kHz)}$ at the Earth's surface. The above power flux-density limit shall apply until such time as a competent world administrative radio conference revises it. In making assignments to the space stations in the mobile-satellite service in the above band, administrations shall take all practicable steps to protect the radio astronomy service in the 150.05–153 MHz band from harmful interference from unwanted emissions.

599B—The use of the bands 137–138 MHz, 148–149.9 MHz, and 400.15–401 MHz by the mobile-satellite service and the band 149.9–150.05 MHz by the land mobile-satellite service is limited to nongeostationary-satellite systems.

608A—The use of the band 148–149.9 MHz by the mobile-satellite service is subject to the application of the coordination and notification procedures set forth in Resolution 46 (WARC-92). The mobile-satellite service shall not constrain the development and use of fixed, mobile, and space operation services in the band 148–149.9 MHz. Mobile earth stations in the mobile-satellite service shall not produce a power flux-density in excess of $-150 \text{ dB(W/m}^2/4 \text{ kHz)}$ outside national boundaries.

608B—The use of the band 149.9–150.05 MHz by the land mobile-satellite service is subject to the application of the coordination and notification procedures set forth in Resolution 46 (WARC-92). The land mobile-satellite service shall not constrain the development and use of the radionavigation-satellite service in the band 149.9–150.05 MHz. Land mobile earth stations of the land mobile-satellite service shall not produce power flux-density in excess of $-150 \text{ dB(W/m}^2/4 \text{ kHz)}$ outside national boundaries.

647B—The use of the band 400.15–401 MHz by the mobile-satellite service is subject to the application of the coordination and notification procedures set forth in Resolution 46. However, coordination of a space station of the mobile-satellite service with respect to terrestrial services is required only if the power flux-density produced by the station exceeds $-125 \text{ dB(W/m}^2/4 \text{ kHz)}$ at the Earth's surface. The above power flux-density limit shall apply until such time as a competent world administrative radio conference revises it. In making assignments to the space stations in the mobile-satellite service in the above band, administrations shall take all practicable steps to protect the radio astronomy service in the

band 406.1–410 MHz from harmful interference from unwanted emissions.

669—In the maritime mobile service, the frequencies 457.525 MHz, 457.550 MHz, 457.575 MHz, 467.525 MHz, 467.550 MHz, and 467.575 MHz may be used by on-board communication stations. The use of these frequencies in territorial waters may be subject to the national regulations of the administration concerned. The characteristics of the equipment used shall conform to those specified in Appendix 20.

733—The bands 1610–1626.5 MHz, 5000–5250 MHz, and 15.4–15.7 GHz are also allocated to the aeronautical mobile-satellite (R) service on a primary basis. Such use is subject to agreement obtained under the procedure set forth in Article 14.

753F—The use of the band 2483.5–2500 MHz by the mobile-satellite and the radiodetermination-satellite services is subject to the application of the coordination and notification procedures set forth in Resolution 46 (WARC-92). Coordination of space stations of the mobile-satellite and radiodetermination-satellite services with respect to terrestrial services is required only if the power flux-density produced at the Earth's surface exceeds the limits in No. 2566. In respect of assignments operating in this band, the provisions of Section II, paragraph 2.2 of Resolution 46 (WARC-92) shall also be applied to geostationary transmitting space stations with respect to terrestrial stations.

792A—The use of the bands 4500–4800 MHz, 6725–7025 MHz, 10.7–10.95 GHz, 11.2–11.45 GHz, and 12.75–13.25 GHz by the fixed-satellite service shall be in accordance with the provisions of Appendix 30B.

796—The band 5000–5250 MHz is to be used for the operation of the international standard system (microwave landing system) for precision approach and landing. The requirements of this system shall take precedence over other uses of this band.

797—The bands 5000–5250 MHz and 15.4–15.7 GHz are also allocated to the fixed-satellite service and the intersatellite service, for connection between one or more earth stations at specified fixed points on the Earth and space stations, when these services are used in conjunction with the aeronautical radionavigation and/or aeronautical mobile (R) service. Such use shall be subject to agreement obtained under the procedure set forth in Article 14.

909—In the bands 54.25–58.2 GHz, 59–64 GHz, 116–134 GHz, 170–182 GHz, and 185–190 GHz, stations in the aeronautical mobile service may be operated subject to not causing harmful interference to the intersatellite service (see No. 435).

917—In the bands 140.69–140.98 GHz all emissions from airborne stations, and from space stations in the space-to-Earth direction, are prohibited.

Non-Government (NG) Footnotes

NG2—Facsimile broadcasting stations may be authorized in the band 88–108 MHz.

NG3—Control stations in the domestic public mobile radio service may be authorized frequencies in the bands 72–73 and 75.4–76 MHz on the condition that harmful interference will not be caused to operational fixed stations.

NG4—The use of the frequencies in the band 152.84–153.38 MHz may be authorized, in any area, to remote pickup

broadcast base and mobile stations on the condition that harmful interference will not be caused to stations operating in accordance with the Table of Frequency Allocations.

NG6—Stations in the public safety radio services authorized as of June 30, 1958, to use frequencies in the band 159.51–161.79 MHz in areas other than Puerto Rico and the Virgin Islands may continue such operation, including expansion of existing systems, on the condition that harmful interference will not be caused to stations in the services to which these bands are allocated. In Puerto Rico and the Virgin Islands this authority is limited to frequencies in the band 160.05–161.37 MHz. No new public radio service system will be authorized to operate on these frequencies.

NG12—Frequencies in the bands 454.40–455 MHz and 459.40–460 MHz may be assigned to domestic public land and mobile stations to provide a two-way air-ground public radiotelephone service.

NG17—Stations in the land transportation radio services authorized as of May 15, 1958 to operate on the frequency 161.61 MHz may, upon proper application, continue to be authorized for such operation, including expansion of existing systems, on the condition that harmful interference will not be caused to the operation of any authorized station in the maritime mobile service. No new land transportation radio service system will be authorized to operate on 161.61 MHz.

NG19—Fixed stations associated with the maritime mobile service may be authorized, for purposes of communication with coast stations, to use frequencies assignable to ship stations in this band on the condition that harmful interference will not be caused to services operating in accordance with the Table of Frequency Allocations.

NG23—Frequencies in the band 2100–2200 MHz may also be assigned to stations in the international fixed public radio service located south of latitude 25°30' N in the State of Florida and in US Possessions in the Caribbean area, provided, however, no new assignments in the band 2150–2162 MHz will be made to such stations after February 25, 1974.

NG28—The frequency band 160.86–161.40 MHz is available for assignment to remote pickup base and remote pickup mobile stations in Puerto Rico and the Virgin Islands only on a shared basis with the land transportation radio service.

NG31—Stations in the Rural Radio Service licensed for Basic Exchange Telecommunications Radio Service may be authorized to use some frequencies in the bands 816–820 MHz (fixed subscriber) and 861–865 MHz (central office or base), on a co-primary basis with private land mobile radio licensees, pursuant to part 22 subpart H.

NG41—Frequencies in the bands 3700–4200 MHz, 5925–6425 MHz, and 10.7–11.7 GHz may also be assigned to stations in the international fixed public and international control services located in US Possessions in the Caribbean area.

NG42—Non-Government stations in the radiolocation service shall not cause harmful interference to the amateur service.

NG47—In Alaska frequencies within the band 2655–2690 MHz are not available for assignment to terrestrial stations.

NG49—The following frequencies may be authorized for mobile operations in the Manufacturers Radio Service subject

to the condition that no interference is caused to the reception of television stations operating on Channels 4 and 5 and that their use is limited to a manufacturing facility:

72.02, 72.04, 72.06, 72.08, 72.10, 72.12, 72.14, 72.16, 72.18, 72.20, 72.22, 72.24, 72.26, 72.28, 72.30, 72.32, 72.34, 72.36, 72.38, and 72.40 MHz.

Further, the following frequencies may be authorized for mobile operations in the Special Industrial Radio Service, Manufacturers Radio Service, Railroad Radio Service, and Forest Products Radio Service subject to the condition that no interference is caused to the reception of television stations operating on Channels 4 and 5 and that their use is limited to a railroad yard, manufacturing plant, logging site, mill, or similar industrial facility.

72.44, 72.48, 72.52, 72.56, 72.60, 75.44, 75.48, 75.52, 75.56, and 75.60 MHz.

NG51—In Puerto Rico and the Virgin Islands only, the bands 150.8–150.98 MHz and 150.98–151.49 MHz are allocated exclusively to the business radio service.

NG53—In the band 12.7–13.15 GHz, television pickup stations and CARS pickup stations shall be assigned channels on a co-equal basis and shall operate on a secondary basis to fixed stations operating in accordance with the Table of Frequency Allocations. In the 13.15–13.20 GHz band television pickup stations and CARS pickup stations shall be assigned on an exclusive basis in the top 100 markets, as set out in Section 76.51.

NG56—In the bands 72.0–73.0 and 75.4–76.0 MHz, the use of mobile radio remote control of models is on a secondary basis to all other fixed and mobile operations. Such operations are subject to the condition that interference will not be caused to common carrier domestic public stations, to remote control of industrial equipment operating in the 72–76 MHz band, or to the reception of television signal on Channels 4 (66–72 MHz) or 5 (76–82 MHz). Television interference shall be considered to occur whenever reception of regularly used television signals is impaired or destroyed, regardless of the strength of the television signal or the distance to the television station.

NG59—The frequencies 37.60 and 37.85 MHz may be authorized only for use by base, mobile, and operational fixed stations participating in an interconnected or coordinated power service utility system.

NG63—Television Broadcast translator stations holding valid licenses on November 15, 1971, to operate in the frequency band 806–890 MHz (Channels 70–83), may continue to operate in this band, pursuant to periodic license renewals, on a secondary basis to the land mobile radio service.

NG66—The frequency band 470–512 MHz is allocated for use in the broadcasting and land mobile radio services. In the land mobile services, it is available for assignment in the domestic public, public safety, industrial, and land transportation radio services at, or in the vicinity of, 11 urbanized areas of the United States, as set forth in the following table. Additionally, in the land mobile services, TV Channel 16 is available for assignment in the public safety radio services at, or in the vicinity of, Los Angeles. Such use in the land mobile ser-

vices is subject to the conditions set forth in Parts 22 and 90 of this chapter.

Urbanized Area	TV Channel
New York, NY—Northeastern New Jersey	14, 15
Los Angeles, CA	14, 20
Chicago, IL—Northwestern Indiana	14, 15
Philadelphia, PA—New Jersey	19, 20
San Francisco—Oakland, CA	16, 17
Boston, MA	14, 16
Washington, DC—Maryland-Virginia	17, 18
Pittsburgh, PA	14, 18
Miami, FL	14
Houston, TX	17
Dallas, TX	16

NG70—In Puerto Rico and the Virgin Islands only, the bands 159.240–159.435 and 160.410–160.620 MHz are also available for assignment to base stations and mobile stations in the special industrial radio service.

NG101—The use of the band 2500–2690 MHz by the broadcasting-satellite service is limited to domestic and regional systems for community reception of educational television programming and public service information. Such use is subject to agreement among administrations concerned and those having services operating in accordance with the table, which may be affected. Unless such agreement includes the use of higher values, the power flux-density at the earth's surface produced by emissions from a space station in this service shall not exceed those values set forth in Part 73 of the rules for this frequency band.

NG102—Use of the fixed-satellite service in the bands 2500–2655 MHz (space-to-earth) and 2655–2690 MHz (Earth-to-space) is limited as follows:

- (a) For common carrier use in Alaska, for intra-Alaska service only, and in the mid- and western Pacific areas including American Samoa, Guam, the Northern Mariana Islands, and Hawaii;
- (b) for educational use in the contiguous United States, Alaska, and the mid- and western Pacific areas including American Samoa, Guam, the Northern Mariana Islands, and Hawaii.

Such use is subject to agreement with administrations having services operating in accordance with the Table, which may be affected. In the band 2500–2655 MHz, unless such agreement includes the use of higher values, the power flux-density at the Earth's surface produced by emissions from a space station in this service shall not exceed the values set forth in Part 25 of the Rules for this frequency band.

NG104—The use of the bands 10.7–11.7 and 12.75–13.25 GHz in the fixed-satellite service is limited to international systems, i.e., other than domestic systems.

NG111—The band 157.4375–157.4625 MHz may be used for one-way paging operations in the special emergency radio service.

NG112—The frequencies 25.04, 25.08, 150.980, 154.585, 158.445, 159.480, 454.000, and 459.000 MHz may be authorized to stations in the petroleum radio service for use primarily in oil spill containment and cleanup operations and secondarily in regular land mobile communication.

NG114—In the Gulf of Mexico offshore from the Louisiana–Texas coast, the frequency band 476–494 MHz (TV Channels 15, 16, and 17) is allocated to the Domestic Public and Private Land Mobile Radio Services in accordance with the regulations set forth in parts 22 and 90 respectively.

NG115—In the 174 to 216 MHz band wireless microphones may be authorized to operate on a secondary, noninterfering basis, subject to terms and conditions set forth in Part 74 of these Rules and Regulations.

NG117—The frequencies 156.050 and 156.175 MHz may be assigned to stations in the maritime mobile service for commercial and port operations in the New Orleans Vessel Traffic Service (VTS) area and the frequency 156.250 MHz may be assigned to stations in the maritime mobile service for port operating in the New Orleans and Houston VTS areas.

NG118—Television translator relay stations may be authorized to use frequencies in the 2025–2130 MHz band on a secondary basis to stations operating in accordance with the Table of Frequency Allocations.

NG120—Frequencies in the band 928–960 MHz may be assigned for multiple address systems and mobile operations on a primary basis as specified in 47 C.F.R. Part 101.

NG124—Within designated segments of the bands that comprise 30.85–47.41 MHz, 150.8–159.465 MHz, and 453.0125–467.9875 MHz, police licensees are authorized to operate low-power radio transmitters on a secondary, noninterference basis in accordance with the provisions of 47 C.F.R. §§ 2.803 and 90.20(e)(5).

NG127—In Hawaii, the frequency band 488–494 MHz is allocated exclusively to the fixed service for use by common carrier control and repeater stations for point-to-point inter-island communications only.

NG128—In the band 535–1705 kHz, AM broadcast licensees or permittees may use their AM carrier on a secondary basis to transmit signals intended for both broadcast and nonbroadcast purposes. In the band 88–108 MHz, FM broadcast licensees or permittees are permitted to use subcarriers on a secondary basis to transmit signals intended for both broadcast and nonbroadcast purposes. In the bands 54–72, 76–88, 174–216, 470–608, and 614–806 MHz, TV broadcast licensees or permittees are permitted to use subcarriers on a secondary basis for both broadcast and nonbroadcast purposes.

NG129—In Alaska, the bands 76–88 MHz and 88–100 MHz are also allocated to the Fixed service on a secondary basis. Broadcast stations operating in these bands shall not cause interference to non-Government fixed operations authorized prior to January 1, 1982.

NG133—Stations authorized in the band 73–74.6 MHz as of December 1, 1961 may continue to operate until December 31, 1985. Such stations shall not be required to afford protection to radio astronomy observations within the United States and possessions; however, such stations must afford protection to the observatories of other countries.

NG134—In the band 10.45–10.5 GHz non-Government stations in the radiolocation service shall not cause harmful interference to the amateur and amateur-satellite services.

NG135—In the 420–430 MHz band the amateur service is not allocated north of line A (def. Sec. 2.1).

NG139—Pending the adoption of further specific rules concerning usage of the band 12.2–12.7 GHz by the fixed and broadcasting-satellite services, systems in these services may be authorized subject to the condition that adjustments in certain systems design or technical parameters may become necessary during the system lifetime. The necessity for such adjustments, and their extent, will be dependent upon the Final Acts of the 1983 Regional Administrative Radio Conference and subsequent Commission decisions.

NG140—Pending the adoption of further specific rules concerning usage of the band 17.3–17.8 GHz by the fixed-satellite service for the purpose of providing feeder links to the broadcasting-satellite service, systems may be authorized for this purpose subject to the condition that adjustments in certain systems design or technical parameters may become necessary during the system lifetime. The necessity for such adjustments, and their extent, will be dependent upon the Final Acts of the 1983 Regional Administrative Radio Conference and subsequent Commission decisions.

NG141—The frequencies 42.40 MHz and 44.10 MHz are authorized on a primary basis in the State of Alaska for meteor burst communications by fixed stations in the Rural Radio Service operating under the provisions of Part 22 of this chapter. The frequencies 44.20 MHz and 45.90 MHz are authorized on a primary basis in Alaska for meteor burst communications by fixed private radio stations operating under the provisions of Part 90 of the chapter. The private radio station frequencies may be used by Common Carrier stations on a secondary, noninterference basis and the Common Carrier frequencies may be used by private radio stations for meteor burst communications on a secondary, noninterference basis. Users shall cooperate to the extent practical to minimize potential interference. Stations utilizing meteor burst communications shall not cause harmful interference to stations of other radio services operating in accordance with the Table of Frequency Allocations.

NG142—TV broadcast stations authorized to operate in the bands 54–72, 76–88, 174–216, 470–512, and 512–806 MHz may use a portion of the television vertical blanking interval for the transmission of telecommunications signals, on the conditions that harmful interference will not be caused to the reception of primary services and that such telecommunications services must accept any interference caused by primary services operating in these bands.

NG143—In the band 11.7–12.2 GHz protection from harmful interference shall be afforded to transmissions from space stations not in conformance with International Footnote 839 only if the operations of such space stations impose no unacceptable constraints on operations or orbit locations of space stations in conformance with 839.

NG144—Stations authorized as of September 9, 1983 to use frequencies in the band 17.7–19.7 GHz may, upon proper application, continue to be authorized for such operation.

NG145—In the band 11.7–12.2 GHz, transponders on space stations in the fixed-satellite service may be used addi-

tionally for transmissions in the broadcasting-satellite service, provided that such transmissions do not have a maximum eirp greater than 53 dBW per television channel and do not cause greater interference or require more protection from interference than the coordinated fixed-satellite service frequency assignments. With respect to the space services, this band shall be used principally for the fixed-satellite service.

NG147—Stations in the broadcast auxiliary service and private radio services licensed as of July 25, 1985, or on a subsequent date following as a result of submitting an application for license on or before July 25, 1985, may continue to operate on a primary basis with the in-orbit-satellite service and the radio-determination satellite service.

NG148—The frequencies 154.585 MHz, 159.480 MHz, 160.725 MHz, 160.785 MHz, 454.000 MHz, and 459.000 MHz may be authorized to maritime mobile stations for off-shore radiolocation and associated telecommand operations.

NG149—The frequency bands 54–72 MHz, 76–88 MHz, 174–216 MHz, 470–512 MHz, 512–608 MHz, and 614–746 MHz are also allocated to the fixed service to permit subscription television operations in accordance with Part 73 of the rules.

NG151—In the frequency bands 824–849 MHz and 869–894 MHz, cellular land mobile licensees are permitted to offer auxiliary services on a secondary basis subject to the provisions of Part 22.

NG152—The band 219–220 MHz is also allocated to the amateur service on a secondary basis for stations participating, as forwarding stations, in point-to-point fixed digital message forwarding systems, including intercity packet backbone networks.

NG153—The 2145–2150 MHz and 2160–2165 MHz bands are reserved for future emerging technologies on a co-primary basis with the fixed and mobile services. Allocations to specific services will be made in future proceedings.

NG154—The 157.1875–157.45 MHz and 161.775–162.0125 MHz bands are also allocated to the land mobile service for assignment to stations as described in 47 CFR Part 90.

NG155—The bands 159.500–159.675 MHz and 161.375–161.550 MHz are allocated to the maritime service as described in Part 80 of this chapter. Additionally, the frequencies 159.550, 159.575, and 159.600 MHz are available for low-power intership communications.

NG158—The frequency bands 764–776 MHz and 794–806 MHz are available for assignment exclusively to the public safety services, to be defined in Docket No. WT 96–86.

NG159—Full-power analog television stations licensed pursuant to applications filed before January 2, 2001, and new digital television (DTV) broadcasting operations in the 746–806 MHz band, will be entitled to protection from harmful interference until the end of the DTV transition period. After the end of the DTV transition period, the Commission may assign licenses in the 746–806 MHz band without regard to existing television and DTV operations.

Low-power television and television translators in the 746–806 MHz band must cease operations in the band at the end of the DTV transition period.

Frequency Allotments

Allotment of 27 575 and 27 585 kHz for Short-Distance Low-Power Service

These allotments are to provide for intermittent miscellaneous US Government short-distance low-power radiocommunications, radio signaling, and the control of remote objects or devices by means of radio (where the radiated power exceeds the limit established under Part 7.9).

The designated frequencies are allotted for use by US Government agencies and may be authorized for use by agencies as required upon application. All stations operating on these frequencies shall meet the conditions and standards established for this service.

The designated frequencies are available on a shared basis only and will not be authorized for exclusive use of any one agency. No protection from interference can be assured to any station operating in this service. Services involving safety of life and property should not employ these frequencies in view of their unprotected status. All transmissions are to be restricted to official US Government business that requires the use of radio.

Stations in this service shall utilize FCC type-accepted or type-approved Citizens Radio Band equipment or the equivalent. The maximum transmitter output power shall be 5 W.

Stations shall be identified in accordance with the regulations of each agency.

The only class of station authorized is Mobile (including portable-type operation).

Frequencies 27 575 and 27 585 kHz with 6KA2A, 6KA2D, and 6KA3E emission are designated for the US Government short-distance low-power radio service.

All applications for the use of these frequencies must bear the note S159, which reads, "US Government short-distance low-power service."

Allotments in the Band 1 710–1 850 MHz for Fixed Security Surveillance Systems

The frequencies 1 720, 1 740, 1 760, 1 780, and 1 800 MHz are allotted for use in fixed security surveillance systems, on a secondary basis to other stations operating in accordance with the Government Table of Frequency allocations.

Allotment of 163.100, 418.050, and 418.575 MHz for Wide Area Use

(Assignments made under these procedures shall no longer be authorized after December 31, 2004.)

1. The frequencies 163.100, 418.050, and 418.575 MHz are allotted for use by all US Government agencies and are to satisfy intermittent wide area requirements of a transient nature. Coordination in accordance with Section 8.3.18 is not required.
2. All operations shall be authorized in accordance with Chapter 9 of the Manual. The frequencies are available on a shared nonpriority basis only

and will not be authorized nor are they intended for the exclusive use of any one agency. No protection from interference will be provided to any station operating on these frequencies from other stations operating on the same frequency. The use of equipment with coded squelch is strongly encouraged to reduce nuisance interference from other users.

3. The intent of these allotments is for use by Government stations in the Land or Maritime Mobile Services, which are unlikely to cause harmful interference to other stations operating in these Services (Table of Services, Station Classes, and Stations, Section 6.1.4), and the following restrictions apply:
 - (a) The minimum ERP (Effective Radiated Power) necessary to support the intended use shall be employed.
 - (b) The maximum base or mobile station transmitter output power shall not exceed 30 W.
 - (c) The gain of the base station antenna shall not exceed 6 dBi.
 - (d) The height of the base station antenna shall not exceed 6 m above the height of the structure supporting the antenna.
 - (e) Station classes are limited to FB, FC, ML, MLP, MS, MSP. However, these station classes may be suffixed with the letter "R" as applicable and restricted by 3.(e)(1) below.
 - (1) Mobile and/or transportable repeater transmitting stations (as defined in Section 6.1.2, Paragraph 3 and Section 9.8.2, Paragraph 15c) may be authorized only on 163.1 or 418.05 MHz but are restricted to a maximum placement period of 45 days and a maximum transmitter output power of 30 W. (For transportable repeater stations, Note S362 applies.) The repeater receive frequencies are 168.35 or 408.4 MHz.
 - (f) Stations in the Fixed Service (FX station class) may be authorized but are restricted to transportable stations (Note S362 applies).
 - (g) All equipment shall conform to Part 5.6 of the Manual.
4. All applications utilizing these allotted frequencies must be affixed with Record Note S352 "This assignment is for intermittent wide area requirements of transient, itinerant nature pursuant to Section 4.2.3 of the Manual."

Allotment of 168.350, 408.400, and 418.075 MHz for Common Use Frequencies

(Assignments made under these procedures shall no longer be authorized after December 31, 2004.)

1. The frequencies 168.350, 408.400, and 418.075 MHz are allotted for use by all US Government agencies and are to provide for radiocommunications that do not justify the assigning of a radio

frequency exclusively to that use, i.e., the frequency can be shared with other users. Coordination in accordance with Section 8.3.18 is not required.

2. All operations shall be authorized in accordance with Chapter 9 of the Manual. The frequencies are available on a shared nonpriority basis only and will not be authorized nor are they intended for the exclusive use of any one agency. No protection from interference will be provided to any station operating on these frequencies from other stations operating on the same frequency. The use of equipment with coded squelch is strongly encouraged to reduce nuisance interference from other users.
3. The intent of these allotments is for use by Government stations in the Land or Maritime Mobile Services, which are unlikely to cause harmful interference to other stations operating in these Services (Table of Services, Station Classes, and Stations, Section 6.1.4), and the following restrictions apply:
 - (a) The minimum ERP necessary to support the intended use shall be employed.
 - (b) The maximum base or mobile station transmitter output power shall not exceed 30 W.
 - (c) The gain of the base station antenna shall not exceed 6 dBi.

TIME SCALES

The world operates on a single international time system, Coordinated Universal Time (UTC), an average of the best clocks in the world. The time scales of the Department of Commerce National Institute of Standards and Technology (NIST) and the US Naval Observatory (USNO) are synchronized to the international time scale maintained by the Bureau International des Poids et Mesures (BIPM) in Sèvres, France. NIST and USNO also maintain atomic clocks that contribute to the UTC average.*

There is an international agreement that UTC scales within a given country are to be kept (if possible) to within one microsecond of this international scale. In the United States, there are two principal timing centers: UTC(NIST) operated by NIST in Boulder, CO, and UTC(USNO) operated by the USNO Time Services Division in Washington, DC. UTC(NIST) serves the civilian sector, and the UTC(USNO) has direct responsibilities for Department of Defense activities. The official United States time is maintained and distributed by both NIST and USNO.

Since both UTC(NIST) and UTC(USNO) are synchronized to UTC, they are also synchronized to each

* Bodson, D., Adair, R. T., and Meister, M. D. "Time and Frequency Information in Telecommunications Systems Standardized by Federal Standard 1002A." *Proceedings of the IEEE*, Vol. 79, No. 7, July 1991, pp. 1077-1078.

other, typically within 50 ns. Greater accuracies with respect to these sources can be achieved by obtaining the necessary corrections from the respective agency.

A new United States standard for time reference, Federal Standard 1002A, was developed in 1991 by the Federal Government in cooperation with the telecommunications industry.

International Definition of the Second

The second is one of seven base units in the International System of Units (SI). The base units are used to derive other units of physical quantities. Use of the SI means that physical quantities such as the second and hertz are defined and measured in the same way throughout the world.

There have been several definitions of the SI second. Until 1956, the definition was based on the *mean solar day*, or one revolution of the earth on its axis. The *mean solar second* was defined as 1/86 400 of the mean solar day and provided the basis for several astronomical time scales known as Universal Time (UT):

UT0—The original mean solar time scale, based on the rotation of the earth on its axis. UT0 was first kept with pendulum clocks. When quartz clocks became available, astronomers noticed errors in UT0 caused by polar motion and developed the UT1 time scale.

UT1—The most widely used astronomical time scale, UT1 improves upon UT0 by correcting for longitudinal shifts of the observing station caused by polar motion. Because the earth's rotational rate is not uniform the uncertainty of UT1 is about 2 to 3 ms per day.

UT2—Mostly of historical interest, UT2 is a smoothed version of UT1 that corrects for deviations in the period of the earth's rotation caused by angular momenta of the earth's core, mantle, oceans, and atmosphere.

The *ephemeris second* served as the SI second from 1956 to 1967. The ephemeris second was a fraction of the tropical year, or the interval between the annual vernal equinoxes, which occur on or about March 21. The tropical year was defined as 31 556 925.9747 ephemeris seconds. Determining the precise instant of the equinox is difficult, and this limited the uncertainty of Ephemeris Time (ET) to ± 50 ms over a 9-year interval. ET was used mainly by astronomers and was replaced by *Terrestrial Time* (TT) in 1984, equal to International Atomic Time (TAI) + 32.184 s. The uncertainty of TT is ± 10 μ s.

The era of atomic time keeping formally began in 1967, when the SI second was redefined based on the resonance frequency of the cesium atom:

The duration of 9,192,631,770 periods of the radiation corresponding to the transition between

two hyperfine levels of the ground state of the cesium-133 atom.

Owing to the atomic second, time interval and frequency can now be measured with less uncertainty and more resolution than any other physical quantity. Today, the best time and frequency standards can realize the SI second with uncertainties of $\approx 1 \times 10^{-15}$. Physical realizations of the other base SI units have much larger uncertainties (Table 4).

International Atomic Time (TAI) is an atomic time scale that attempts to realize the SI second as closely as possible. TAI is maintained by the Bureau International des Poids et Mesures in Sèvres, France. The BIPM averages data collected from more than 200 atomic time and frequency standards located at more than 40 laboratories, including the National Institute of Standards and Technology.

Coordinated Universal Time (UTC) runs at the same rate as TAI. However, it differs from TAI by an integral number of seconds. This difference increases when *leap seconds* occur. When necessary, leap seconds are added to UTC on either June 30 or December 31. The purpose of adding leap seconds is to keep atomic time (UTC) within ± 0.9 s of astronomical time (UT1). Some time codes contain a UT1 correction that can be applied to UTC to obtain UT1.

Leap seconds have been added to UTC at a rate of slightly less than once per year, beginning in 1972. UT1 is currently losing about 700 to 800 ms per year with respect to UTC. This means that atomic seconds are shorter than astronomical seconds and that UTC runs faster than UT1. There are two reasons for this. The first reason involves the definition of the atomic second, which made it slightly shorter than the astronomical second to begin with. The second reason is that the earth's rotational rate is gradually slowing down and the astronomical second is gradually getting longer. When a positive leap second is added to UTC, the sequence of events is:

23 h 59 m 59 s
23 h 59 m 60 s
0 h 0 m 0 s

TABLE 4. UNCERTAINTIES OF PHYSICAL REALIZATIONS OF THE BASE SI UNITS

SI Base Unit	Physical Quantity	Uncertainty
candela	luminous intensity	10^{-3}
mole	amount of substance	10^{-6}
kelvin	thermodynamic temperature	10^{-6}
ampere	electric current	10^{-7}
kilogram	mass	10^{-8}
meter	length	10^{-11}
second	time interval	10^{-15}

The insertion of the leap second creates a minute that is 61 s long. This “stops” UTC for 1 s, so that UT1 can catch up.

Time Zones

The world is divided into 24 time zones, each 15° of longitude, or 1 hour angle, apart. The meridian of Greenwich, England, is the center of the zero zone, which extends to 7.5° east and west. Proceeding eastward from Greenwich, the zones are numbered 1 to 12 with the prefix “plus” to indicate the hour angle to be added to Universal Time to obtain local “standard time”. Proceeding westward, the zones are numbered 1 to 12 with the prefix “minus” to indicate the hour angle to be subtracted from Universal Time to obtain local “standard time.” For example, Washington, DC, at longitude 77° West, is in time zone -5.

The actual boundaries of time zones are defined by law or custom in the various countries and generally do not coincide with the theoretical zone, even in some places at sea. In many areas, local legal “standard time” differs by 60 or 30 minutes from theoretical standard time. Fig. 3 in this chapter is a map showing the standard time zones of the world; Chart 1 in Chapter 49 relates the times in several major cities of the world.

TIME AND FREQUENCY STANDARDS

The stability of time and frequency standards is closely related to their quality factor, or Q . The Q of an oscillator is its resonance frequency divided by its resonance width. The resonance frequency is the natural frequency of the oscillator. The resonance width is the range of possible values where the oscillator will run. A high- Q resonator will not oscillate at all unless it is near its resonance frequency. Obviously a high resonance frequency and a narrow resonance width are both advantages when seeking a high Q . Generally speaking, the higher the Q , the more stable the oscillator, since a high Q means that an oscillator will stay close to its natural resonance frequency.

Quartz oscillators achieve the highest Q of any mechanical-type device. Oscillators based on the atomic resonance of rubidium, hydrogen, and cesium achieve much higher Q factors. Each oscillator type is briefly described below and Table 5 provides a summary.

Quartz Oscillators

Quartz crystal oscillators are by far the most common time and frequency standard. An estimated two billion (2×10^9) quartz oscillators are manufactured annually. Most are small devices built for wristwatches, clocks, and electronic circuits. However, they are also found inside test and measurement equipment,

such as counters, signal generators, and oscilloscopes, and interestingly enough, inside every atomic oscillator.

A typical Q for a quartz oscillator ranges from 10^4 to 10^6 . The maximum Q for a high-stability quartz oscillator can be estimated as $Q = 16 \text{ million}/f$, where f is the resonance frequency in MHz.

Environmental changes such as temperature, humidity, pressure, and vibration can change the resonance frequency of a quartz crystal, and there are several designs that reduce the environmental problems. The *oven-controlled crystal oscillator* (OCXO) encloses the crystal in a temperature-controlled chamber called an oven. An alternate solution to the temperature problem is the *temperature-compensated crystal oscillator* (TCXO). In a TCXO, the signal from a temperature sensor generates a correction voltage that is applied to a voltage-variable reactance, or varactor. The varactor then produces a frequency change equal and opposite to the frequency change produced by temperature. TCXOs are used when high stability over a wide temperature range is not required.

Quartz oscillators have excellent short-term stability. An OCXO might be stable (σ, τ , at $\tau = 1 \text{ s}$) to 1×10^{-12} . The limitations in short-term stability are mainly due to noise from electronic components in the oscillator circuits. Long term stability is limited by *aging*, or a change in frequency with time caused by internal changes in the oscillator. Aging is usually a nearly linear change in the resonance frequency that can be either positive or negative, and occasionally, a reversal in aging direction occurs. Aging has many possible causes including a buildup of foreign material on the crystal, changes in the oscillator circuitry, or changes in the quartz material or crystal structure. A high-quality OCXO might age at a rate of $< 5 \times 10^{-9}$ per year, while a TCXO might age 100 times faster.

Because of aging and environmental factors such as temperature and vibration, it is hard to keep even the best quartz oscillators within 1×10^{-10} of their nominal frequency without constant adjustment. For this reason, atomic oscillators are used for applications that require higher long-term accuracy and stability.

Rubidium Oscillators

Rubidium oscillators are the lowest priced members of the atomic oscillator family. They operate at 6 834 682 608 Hz, the resonance frequency of the rubidium atom (^{87}Rb), and use the rubidium frequency to control the frequency of a quartz oscillator. A microwave signal derived from the crystal oscillator is applied to the ^{87}Rb vapor within a cell, forcing the atoms into a particular energy state. An optical beam is then pumped into the cell and is absorbed by the atoms as it forces them into a separate energy state. A photocell detector measures how much of the beam is absorbed and tunes a quartz oscillator to a frequency that maximizes the amount of light absorption. The quartz oscillator is then locked to the resonance frequency of rubidium,

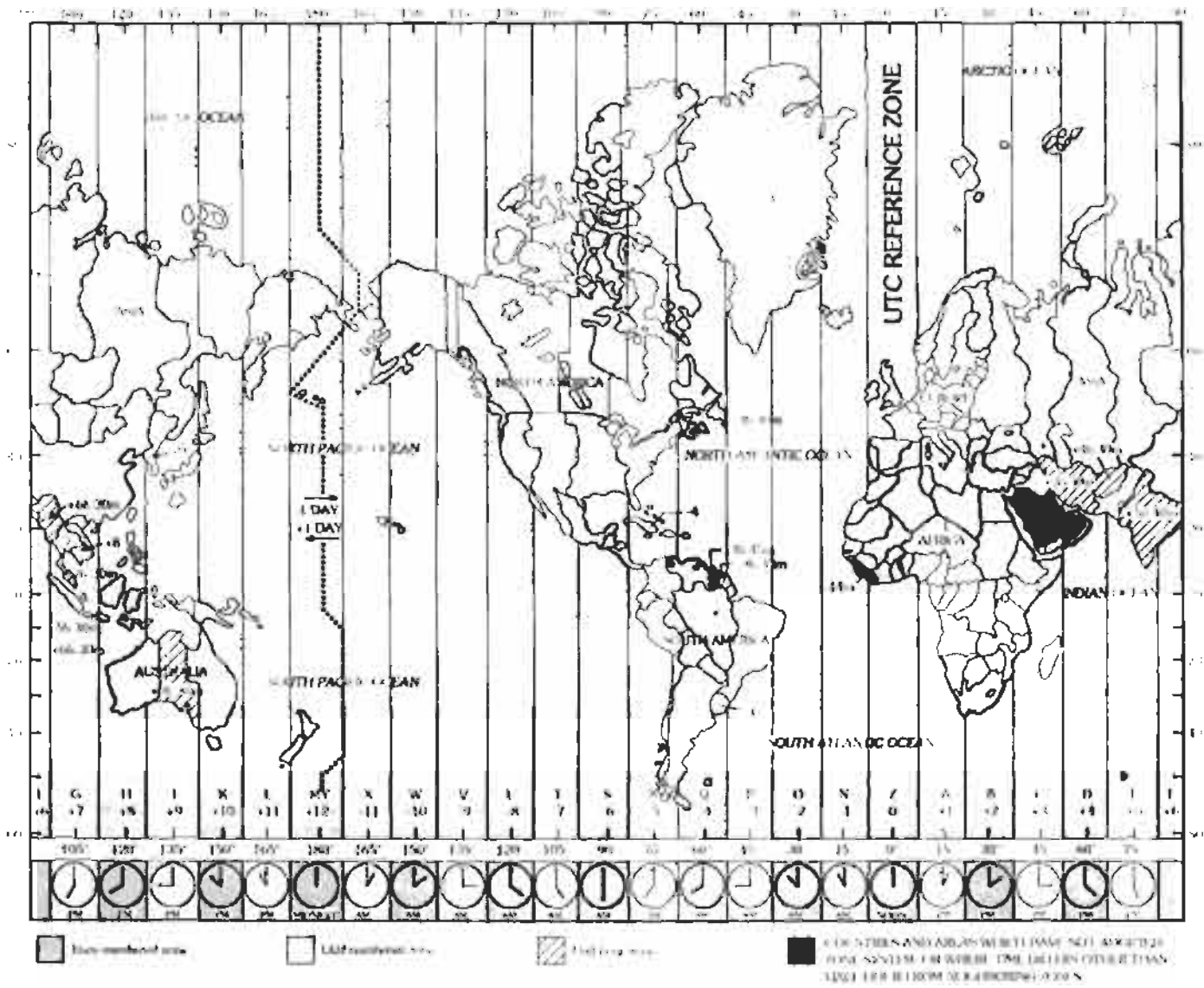


Fig. 1. Standard time zones of the world and their relationship to UTC.

TABLE 5. SUMMARY OF OSCILLATOR TYPES

Oscillator Type	Quartz (TCXO)	Quartz (OCXO)	Rubidium	Commercial Cesium Beam	Hydrogen Maser
Q	10^4 to 10^6	3.2×10^6 (5 MHz)	10^7	10^8	10^9
Resonance frequency	Various	Various	6.834 682 608 GHz	9.192 631 770 GHz	1.420 405 752 GHz
Leading cause of failure	None	None	Rubidium lamp (15 years plus)	Cesium beam tube (3 to 25 years)	Hydrogen depletion (7 years plus)
Stability, $\sigma_y(\tau)$, $\tau = 1$ s	1×10^{-8} to 1×10^{-9}	1×10^{-12}	5×10^{-11} to 5×10^{-12}	5×10^{-11} to 5×10^{-12}	1×10^{-12}
Noise floor, $\sigma_y(\tau)$	1×10^{-9} ($\tau = 1$ to 10^2 s)	1×10^{-12} ($\tau = 1$ to 10^2 s)	1×10^{-12} ($\tau = 10^3$ to 10^5 s)	1×10^{-14} ($\tau = 10^5$ to 10^7 s)	1×10^{-15} ($\tau = 10^3$ to 10^5 s)
Aging/year	5×10^{-7}	5×10^{-9}	1×10^{-10}	None	$\approx 1 \times 10^{-13}$
Frequency offset after warm-up	1×10^{-6}	1×10^{-8} to 1×10^{-10}	5×10^{-10} to 5×10^{-12}	5×10^{-12} to 1×10^{-14}	1×10^{-12} to 1×10^{-13}
Warm-up period	< 10 s to 1×10^{-6}	< 5 min to 1×10^{-8}	< 5 min to 5×10^{-10}	30 min to 5×10^{-12}	24 h to 1×10^{-12}

and standard frequencies are derived and provided as outputs.

Rubidium oscillators continue to get smaller and less expensive, and they offer perhaps the best price/performance ratio of any oscillator. Their long-term stability is much better than that of a quartz oscillator and they are also smaller, more reliable, and less expensive than cesium oscillators.

The Q of a rubidium oscillator is about 10^7 . The shifts in the resonance frequency are mainly caused by collisions of the rubidium atoms with other gas molecules. These shifts limit the long-term stability. Stability ($\sigma_y\tau$, at $\tau = 1$ s) is typically 1×10^{-11} , and about 1×10^{-12} at one day. The frequency offset of a rubidium oscillator ranges from 5×10^{-10} to 5×10^{-12} after a warm-up period of a few minutes, and so they meet the accuracy requirements of most applications without adjustment.

Cesium Oscillators

Cesium oscillators are primary frequency standards since the SI second is defined using the resonance frequency of the cesium atom (^{133}Cs), which is 9 192 631 770 Hz. A properly working cesium oscillator should be close to its nominal frequency without adjustment, and there should be no change in frequency due to aging.

Commercially available oscillators use cesium beam technology. Inside a cesium oscillator, ^{133}Cs atoms are heated to a gas in an oven. Atoms from the gas leave the oven in a high-velocity beam that travels through a vacuum tube toward a pair of magnets. The magnets serve as a gate that allows only atoms of a particular magnetic energy state to pass into a microwave cavity, where they are exposed to a microwave frequency derived from a quartz oscillator. If the microwave fre-

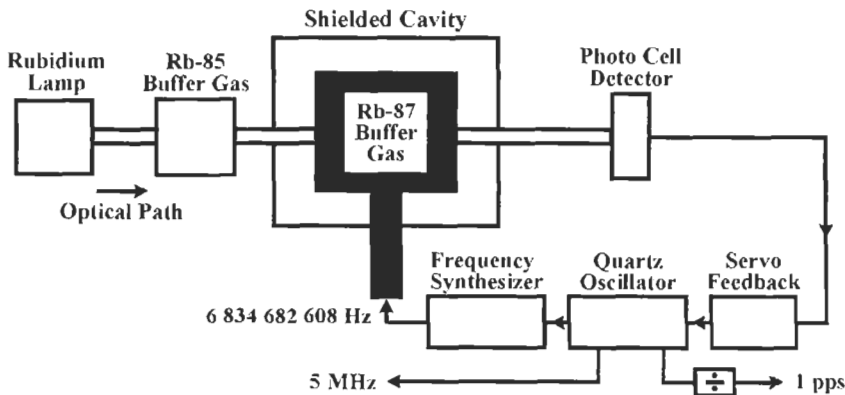


Fig. 4. Rubidium oscillator.

quency matches the resonance frequency of cesium, the cesium atoms will change their magnetic energy state.

The atomic beam then passes through another magnetic gate near the end of the tube. Those atoms that changed their energy state while passing through the microwave cavity are allowed to proceed to a detector at the end of the tube. Atoms that did not change state are deflected away from the detector. The detector produces a feedback signal that continually tunes the quartz oscillator in a way that maximizes the number of state changes so that the greatest number of atoms reaches the detector. Standard output frequencies are derived from the locked quartz oscillator (Fig. 5).

The Q of a commercial cesium standard is a few parts in 10^8 . The beam tube is typically < 0.5 m in length, and the atoms travel at velocities of > 100 ms inside the tube. This limits the observation time to a few milliseconds, and the resonance width to a few hundred hertz. Stability (σ, τ , at $\tau = 1$ s) is typically 5×10^{-12} and reaches a noise floor near 1×10^{-14} at about one day, extending out to weeks or months. The frequency offset is typically near 1×10^{-12} after a warm-up period of 30 minutes.

The current state-of-the-art in cesium technology is the *cesium fountain* oscillator, named after its fountain-like movement of cesium atoms. A cesium fountain named NIST-F1 serves as the primary standard of time and frequency for the United States.

A cesium fountain works by releasing a gas of cesium atoms into a vacuum chamber. Six infrared laser beams are directed at right angles to each other at the center of the chamber. The lasers gently push the cesium atoms together into a ball. In the process of creating this ball, the lasers slow down the movement of the atoms and cool them to temperatures a few thousandths of a degree above absolute zero. This reduces their thermal velocity to a few centimeters per second.

Two vertical lasers gently toss the ball upward and then all of the lasers are turned off. This little push is just enough to loft the ball about a meter high through a microwave-filled cavity. Under the influence of gravity, the ball then falls back down through the microwave cavity. The round trip up and down through the microwave cavity lasts for about 1 s and is limited only by the force of gravity pulling the atoms to the ground. During the trip, the atomic states of the atoms might or might not be altered as they interact with the microwave signal. When their trip is finished, another laser is pointed at the atoms. Those atoms whose states were altered by the microwave signal emit photons (a state known as *fluorescence*) that are counted by a detector. This process is repeated many times while the microwave signal in the cavity is tuned to different frequencies. Eventually, a microwave frequency is found that alters the states of most of the cesium atoms and maximizes their fluorescence. This frequency is the cesium resonance (Fig. 6).

The Q of a cesium fountain is about 10^{10} , or about 100 times higher than a traditional cesium beam. Although the resonance frequency is the same, the resonance width is much narrower (< 1 Hz), owing to the longer observation times made possible by the combination of laser cooling and the fountain design. The combined frequency uncertainty of NIST-F1 is estimated at $< 2 \times 10^{-15}$.

Hydrogen Masers

The *hydrogen maser* is the most elaborate and expensive commercially available frequency standard. The word *maser* is an acronym that stands for microwave amplification by stimulated emission of radiation. Masers operate at the resonance frequency of the hydrogen atom, which is 1 420 405 752 Hz.

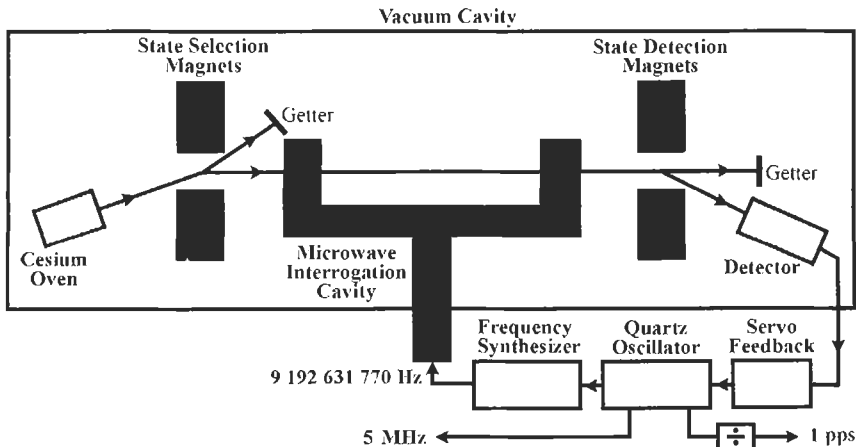


Fig. 5. Cesium beam oscillator.

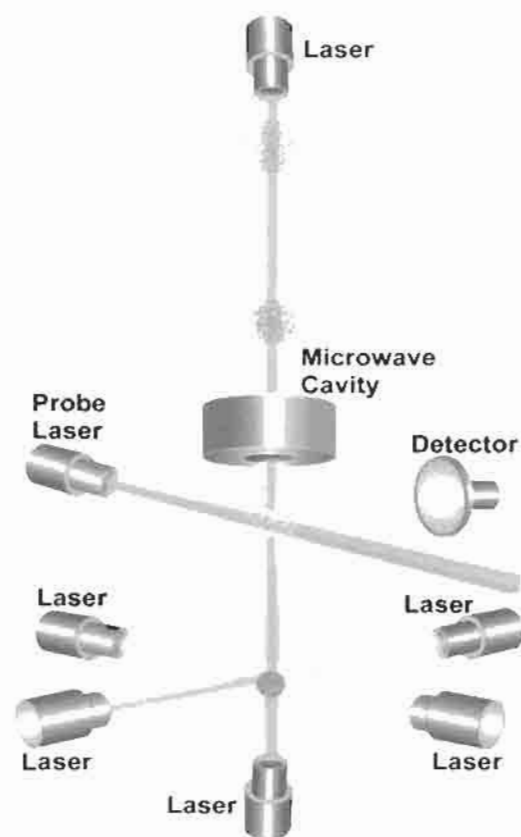


Fig. 6. Cesium fountain oscillator.

A hydrogen maser works by sending hydrogen gas through a magnetic gate that only allows atoms in certain energy states to pass through. The atoms that

make it through the gate enter a storage bulb surrounded by a tuned, resonant cavity. Once inside the bulb, some atoms drop to a lower energy level, releasing photons of microwave frequency. These photons stimulate other atoms to drop their energy level, and they in turn release additional photons. In this manner, a self-sustaining microwave field builds up in the bulb. The tuned cavity around the bulb helps to redirect photons back into the system to keep the oscillation going. The result is a microwave signal that is locked to the resonance frequency of the hydrogen atom and that is continually emitted as long as new atoms are fed into the system. This signal keeps a quartz crystal oscillator in step with the resonance frequency of hydrogen (Fig. 7).

The resonance frequency of hydrogen is much lower than that of cesium, but the resonance width of a hydrogen maser is usually just a few hertz. Therefore, the Q is about 10^9 , or at least one order of magnitude better than a commercial cesium standard. As a result, the short-term stability is better than a cesium standard for periods out to a few days—typically $< 1 \times 10^{-12}$ (σ_y, τ , at $\tau = 1$ s) and reaching a noise floor of $\approx 1 \times 10^{-15}$ after about 1 hour. However, when measured for more than a few days or weeks, a hydrogen maser might fall below a cesium oscillator's performance. The stability decreases because of changes in the cavity's resonance frequency over time.

TIME AND FREQUENCY BROADCASTS

Short-Wave Services—WWV and WWVH

The Department of Commerce National Institute of Standards and Technology operates two high-frequency (short-wave) radio stations, WWV at Fort Collins,

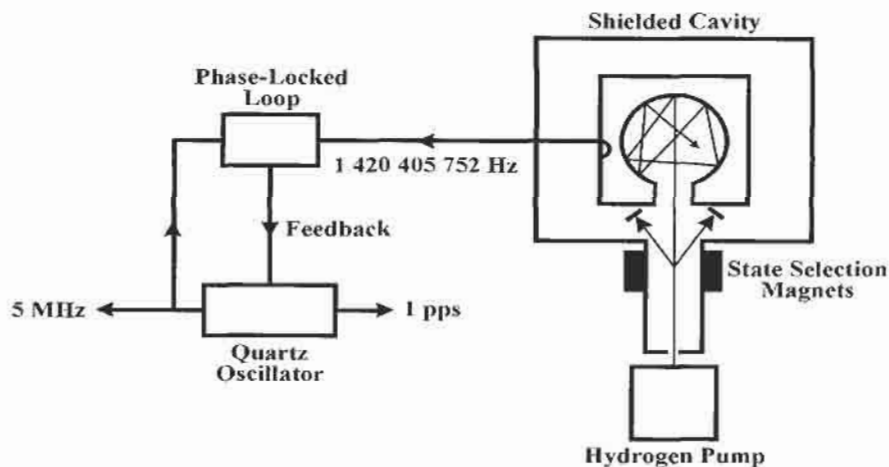


Fig. 7. Hydrogen maser oscillator.

Colorado, and WWVH at Kauai, Hawaii. Both stations broadcast continuous time and frequency signals on 2.5, 5, 10, and 15 MHz. WWV also broadcasts on 20 MHz. All frequencies provide the same information. Although radio reception conditions in the high-frequency band vary greatly with factors such as location, time of year, time of day, the particular frequency being used, atmospheric and ionospheric propagation conditions, and the type of receiving equipment used, at least one frequency should be usable at all times. As a general rule, frequencies above 10 MHz work best in the daytime, and the lower frequencies work best at night.

Services provided by WWV and WWVH include:

- Time announcements
- Standard time intervals
- Standard frequencies
- UT1 time corrections
- BCD time code
- Geophysical alerts
- Marine storm warnings
- Global Positioning System (GPS) status reports

Figs. 8 and 9 show the hourly broadcast schedules of these services.

The station addresses and telephone numbers are:

NIST Radio Station WWV
2000 East County Rd. 58
Fort Collins, CO 80524
(303) 497-3914

NIST Radio Station WWVH
P.O. Box 417
Kekaha, Kauai, HI 96752
(808) 335-4361

Accuracy and Stability—WWV and WWVH are both referenced to the primary NIST frequency standard and atomic time scales located in Boulder, Colorado. The transmitted frequency is accurate to within a few parts in 10^{11} and stable to 1×10^{-12} per day. However, the received accuracy is far less due to various propagation effects. The usable received accuracy is a few parts in 10^9 for frequency and about 1 millisecond for timing.

Facilities and Services—The facilities and services of the NIST radio stations are summarized in Table 6.

The WWV antennas are half-wave vertical antennas that radiate omnidirectional patterns. There are actually five antennas at the station site, one for each frequency.

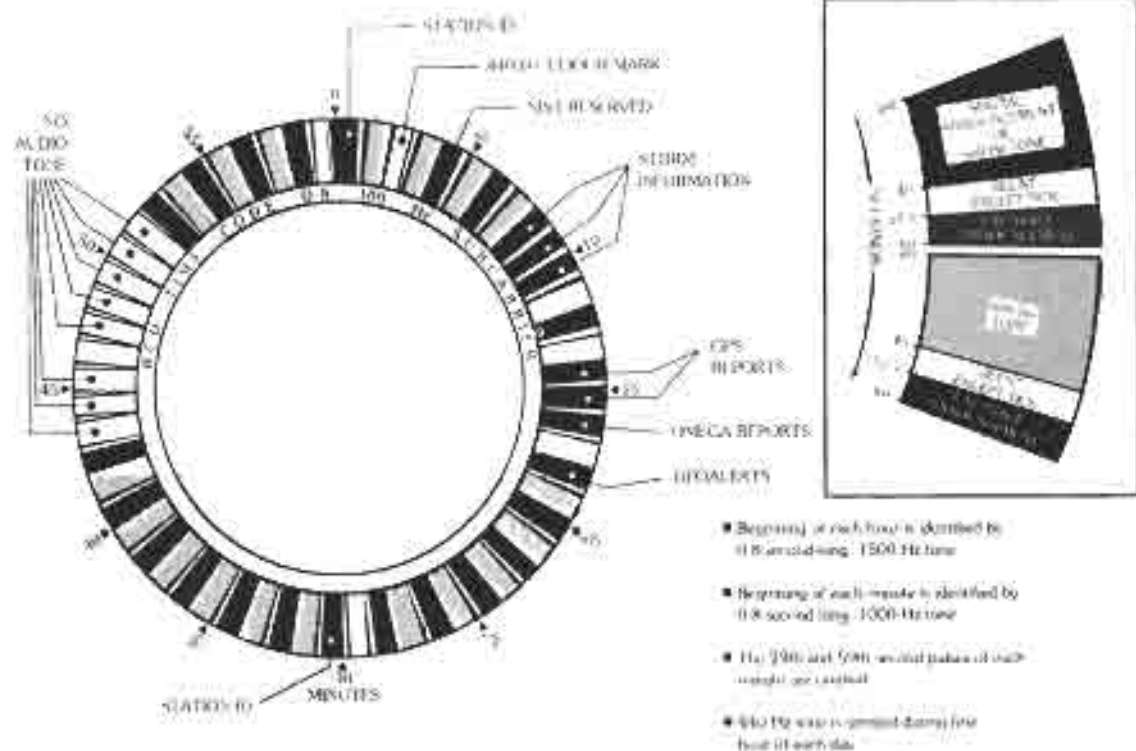


Fig. 8. The hourly broadcast schedule of WWV.

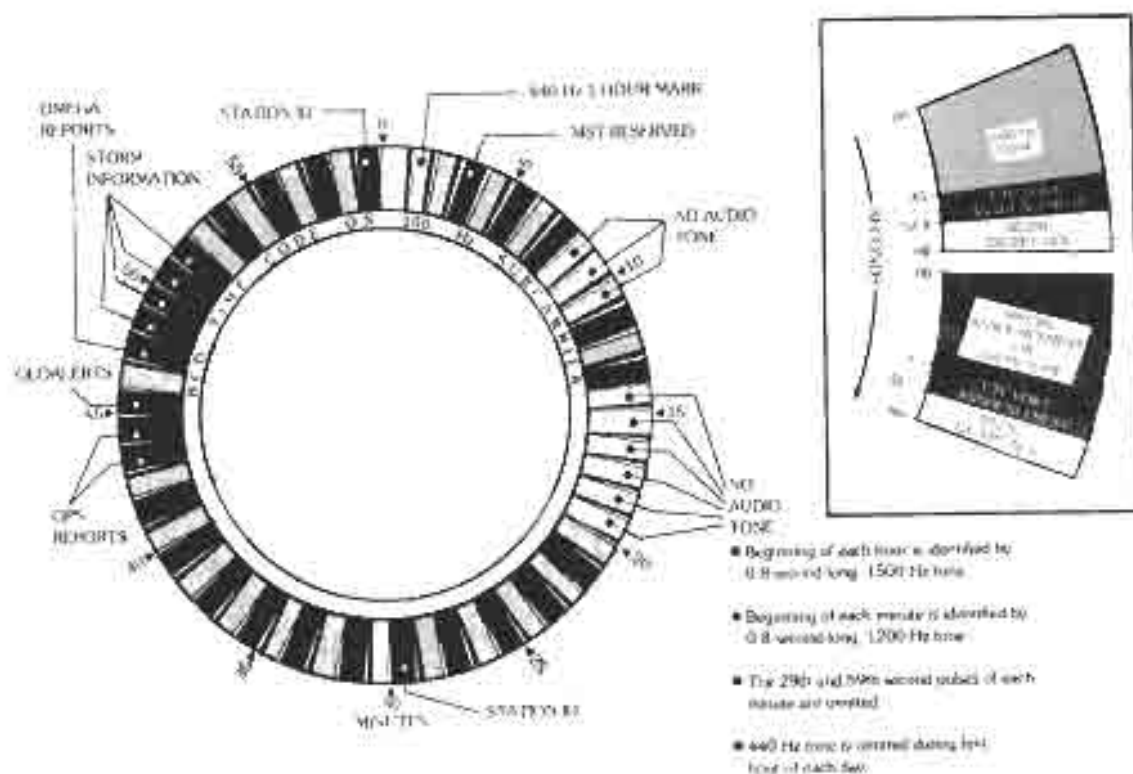


Fig. 9. The hourly broadcast schedule of WWVH.

TABLE 6. SUMMARY OF NIST RADIO BROADCAST SERVICES

Characteristics and Services	WWV	WWVH	WWVB
Date Service Began	March 1923	November 1948	July 1956
Average Geographical Coordinates	40° 40' 49.0" N 105° 02' 21.0" W	21° 59' 26.0" N 150° 46' 00.0" W	40° 40' 28.3" N 105° 02' 39.5" W
Standard Carrier Frequencies	2.5 & 20 MHz; 5, 10, & 15 MHz	2.5 MHz; 5, 10, & 15 MHz	60 kHz
Power	2500 W; 10 000 W	5000 W; 10 000 W	50 000 W
Standard Audio Frequencies	440 (A above middle C), 500, & 600 Hz		—
Time Intervals	1 pulse/s; minute mark; hour mark		n. min.
Time Signals: Voice	Once per minute		—
Time Signals: Code	BCD code on 100-Hz subcarrier; 1 pulse/s		BCD code
UT1 Corrections	UT1 corrections are broadcast with an accuracy of ± 0.1 s		
Special Announcements	Geoids, Marine Storm Warnings, Global Positioning System Status Reports		—

Each antenna is connected to a single transmitter using a rigid coaxial line, and the site is designed so that no two coaxial lines cross. Each antenna is mounted on a

tower that is approximately one-half-wavelength tall. The tallest tower, for 2.5 MHz, is about 60 m tall. The shortest tower, for 20 MHz, is about 7.5-m tall. The

top half of each antenna is a quarter-wavelength radiating element. The bottom half of each antenna consists of nine quarter-wavelength wires that connect to the center of the tower and slope downward to the ground at a 45° angle. This sloping skirt functions as the lower half of the radiating system and also guys the antenna. The WWV antenna coordinates are listed in Table 7.

The 2.5-MHz antenna at WWVH is nearly identical to its WWV counterpart. However, the 5-, 10-, and 15-MHz antennas are phased-array vertical dipoles. They consist of two half-wave vertical dipoles that are separated by a quarter-wavelength and driven 90° out of phase. These antennas radiate a cardioid pattern with the maximum gain pointed toward the west. The WWVH antenna coordinates are listed in Table 8.

Both stations use double sideband amplitude modulation. The modulation level is 50 percent for the steady tones, 25 percent for the BCD time code, 100 percent for the seconds pulses and the minute and hour markers, and 75 percent for the voice announcements.

Time Announcements—Voice announcements are made from WWV and WWVH once every minute. Since both stations can be heard in some locations, a man's voice is used on WWV and a woman's voice is used on WWVH to reduce confusion. The WWVH announcement occurs first, at about 15 s before the minute. The WWV announcement follows at about 7.5 s before the minute. Though the announcements occur at different times, the tone markers are transmitted at the exact same time from both stations. However, they may not be received at exactly the same instant owing to differences in the propagation delays from the two station sites.

The announced time is *Coordinated Universal Time* (UTC). UTC was established by international agreement in 1972, and is governed by the International Bureau of Weights and Measures (BIPM) in Sèvres,

France. Coordination with the international UTC time scale keeps NIST time signals in close agreement with signals from other time and frequency stations throughout the world.

UTC differs from local time at a given location by a specific number of hours. The number of hours depends on the number of time zones between the location in question and the zero meridian (which passes through Greenwich, England). When local time changes from Daylight Saving to Standard Time, or vice versa, UTC does not change. However, the difference between UTC and local time does change—by 1 hour. Fig. 3 may be used to find how many hours to add to or subtract from UTC to obtain local standard time. If DST is in effect, subtract 1 hour less in the United States than shown in Fig. 3. Thus, Eastern Daylight Time (EDT) is only 4 hours behind UTC, not 5 as shown in Fig. 5 for EST.

UTC is a 24-hour clock system. The hours are numbered beginning with 00 hours at midnight through 12 hours at noon to 23 hours and 59 minutes just before the next midnight.

Standard Time Intervals—On WWV and WWVH, the seconds pulses are heard every second except on the 29th and 59th seconds of each minute. The first pulse of each hour is an 800-ms pulse of 1 500 Hz. The first pulse of each minute is an 800-ms pulse of 1 000 Hz at WWV and 1 200 Hz at WWVH. The remaining seconds pulses are short audio bursts (5-ms pulses of 1 000 Hz at WWV and 1 200 Hz at WWVH) that sound like the ticking of a clock.

Each seconds pulse is preceded by 10 ms of silence and followed by 25 ms of silence. The silence makes it easier to pick out the pulse. The total 40-ms protected zone around each seconds pulse is shown in Fig. 10.

Standard Audio Frequencies—In alternate minutes during most of each hour, 500-Hz or 600-Hz audio tones are broadcast. A 440-Hz tone (the musical note A above middle C) is broadcast once each hour. In addition to being a musical standard, the 440-Hz tone

TABLE 7. WWV ANTENNA COORDINATES

Frequency (MHz)	Latitude	Longitude
2.5	40°40'55.2" N	105°02'31.3" W
5	40°40'42.1" N	105°02'24.9" W
10	40°40'47.8" N	105°02'25.1" W
15	40°40'45.0" N	105°02'24.5" W
20	40°40'53.1" N	105°02'28.5" W

TABLE 8. WWVH ANTENNA COORDINATES

Frequency (MHz)	Latitude	Longitude
2.5	21°59'20.9" N	159°45'52.4" W
5	21°59'10.8" N	159°45'44.8" W
10	21°59'18.2" N	159°45'51.3" W
15	21°59'15.3" N	159°45'50.0" W

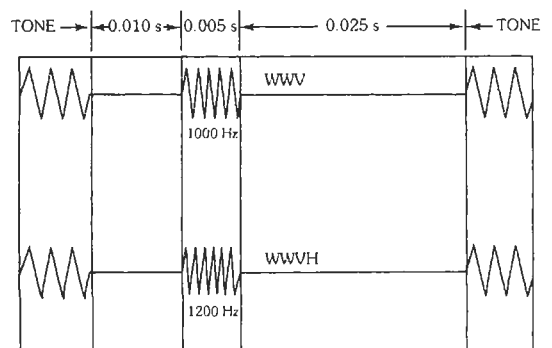


Fig. 10. Format of WWV and WWVH seconds pulses.

provides an hourly marker for chart recorders and other automated devices. The 440-Hz tone is omitted, however, during the first hour of each UTC day. (See Figs. 8 and 9 for further details.)

Silent Periods—The silent periods are without tone modulation. However, the carrier frequency, seconds pulses, time announcements, and the 100-Hz BCD time code continue during the silent periods. In general, one station will not broadcast an audio tone while the other station is broadcasting a voice message.

On WWV, the silent period extends from 43 to 46 and from 47 to 52 minutes after the hour. WWVH has two silent periods: from 8 to 11 minutes after the hour, and from 14 to 20 minutes after the hour. Minutes 29 and 59 on WWV and minutes 00 and 30 on WWVH are also silent.

BCD Time Code—WWV and WWVH continuously broadcast a binary coded decimal (BCD) time code on a 100-Hz subcarrier. The information carried by the time code includes the current minute, hour, and day of year. The time code also contains the 100-Hz frequency from the subcarrier. The 100-Hz frequency may be used as a standard with the same accuracy as the audio frequencies.

The WWV/WWVH time code is a modified version of the IRIG-H code. The code is transmitted on a 100-Hz subcarrier at a rate of 1 pulse per second. The code is in binary coded decimal format. Groups of binary digits (bits) are used to represent decimal numbers. The binary-to-decimal weighting scheme is 1–2–4–8. The least significant bit is always sent first. Table 9 shows the BCD groups and the equivalent decimal numbers.

The decimal number is obtained by multiplying each bit in the binary group by the weight of its respective column and then adding the four products together. For example, the table shows that the binary group 1010 is equal to 5. This is derived by:

$$(1 \times 1) + (0 \times 2) + (1 \times 4) + (0 \times 8) = 1 + 0 + 4 + 0 = 5$$

In the standard IRIG-H code, a "0" bit consists of exactly 20 cycles of 100-Hz amplitude modulation (200-ms duration), and a "1" bit consists of 50 cycles of 100-Hz (500-ms duration). The WWV/WWVH code differs from IRIG-H because all tones are suppressed briefly while the seconds pulses are transmitted (see "Standard Time Intervals," above).

Tone suppression also deletes the first 30 ms of each binary pulse in the time code. This makes the WWV/WWVH bits 30 ms shorter than the IRIG-H bits. Therefore, 170-ms pulses are recognized as "0" bits, and 470-ms pulses are recognized as "1" bits. The leading edge of each pulse coincides with the positive-going crossing of the 100-Hz subcarrier, but due to the tone suppression, it occurs 30 ms after the start of the second.

Refer to Fig. 11. Within 1 minute, enough bits are sent to express the minute, hour, and day of year, two digits of the current year; a leap-second warning indicator; the UT1 correction; and a Daylight Saving Time (DST) indicator. The coded time information refers to the time at the start of the one-minute frame. Seconds are determined by counting pulses within the frame. Two BCD groups each are needed to express the hour (00 to 23), minute (00 to 59), and year (00–99); and three groups are needed to express the day of year (001 to 366). Some bits in the BCD groups are unused, but may provide additional information in the future. To represent units, tens, or hundreds, the basic 1–2–4–8 weights are multiplied by 1, 10, or 100 as appropriate.

Each frame begins with a unique spacing of pulses that mark the start of a new minute. During the first second of the minute, no pulse is transmitted. This creates a 1-s (1000-ms) hole. Since the pulses are already delayed 30 ms by the tone suppression, the UTC minute actually begins 1030 ms (1.03 s) earlier than the first pulse in the frame.

For synchronization purposes, a position identifier pulse is transmitted every 10 s. The position identifier pulse lasts for 770 ms (77 cycles of 100 Hz).

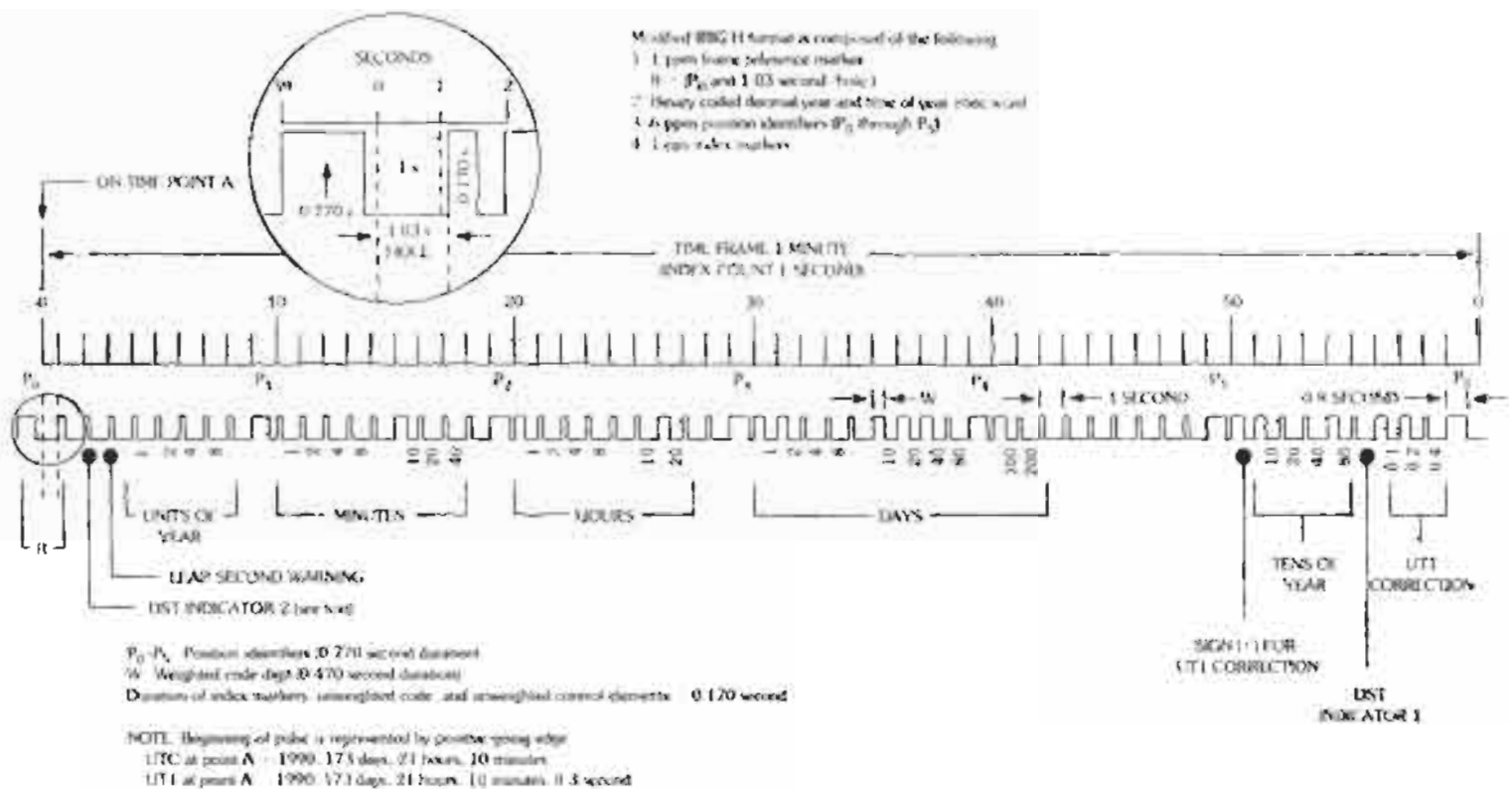
UT1 corrections are sent during the final 10 s of each frame. These corrections are to the nearest 0.1 s. The UT1 correction is expressed with bits called control functions. Control function 1 occurs at 50 s, and tells whether the UT1 correction is negative or positive. If a "0" bit is sent the correction is negative, and if a "1" bit is sent the correction is positive. Control functions 7, 8, and 9 tell the amount of the UT1 correction. They occur at 56, 57, and 58 s, respectively. Since the UT1 corrections are in tenths of seconds, the binary-to-decimal weights, are multiplied by 0.1.

During the first half of 1991, an improved procedure for alerting users of impending DST changes was implemented. Two bits, number 2 and number 55, are used to indicate whether "Standard Time" or "Daylight Saving Time" is in effect at any particular time. Bit number 55 is changed from "0" to "1" at 0000 UTC of the day the time is to be changed from "Standard" time

TABLE 9. 1–2–4–8 BCD-DECIMAL EQUIVALENTS

Binary Group	Decimal Equivalent
Weight: 1 2 4 8	
0 0 0 0	0
1 0 0 0	1
0 1 0 0	2
1 1 0 0	3
0 0 1 0	4
1 0 1 0	5
0 1 1 0	6
1 1 1 0	7
0 0 0 1	8
1 0 0 1	9

Fig. 11. WWV and WWVH time code format.



to DST and from "1" to "0" for the change back to "Standard" time in the fall. The second DST bit, number 2, is changed in the same way as bit 55 but 24 hours later, at 0000 UTC on the day *following* the time change. The use of two separate bits in this manner allows timing receivers that are powered on during the day of a time change to be able to determine whether it is a day of change and, if so, how to adjust their current time correctly.

Also beginning during 1991, the last two digits of the current year (for example, "90" for 1990) are encoded using bits 4 through 7 (for units of year) in a 1-2-4-8 weighting sequence and bits 51 through 54 (for tens of year) in a 10-20-40-80 weighting sequence.

Bit 3 is used to indicate that a leap second is to be inserted at the end of the current month. A "1" is transmitted beginning early in the month at the end of which a leap second is to be inserted. Bit 3 is returned to "0" at 0000 UTC of the first day of the following month.

Fig. 11 shows one frame of the time code. The six position identifiers are labeled P_1 , P_2 , P_3 , P_4 , P_5 , and P_6 . The minutes, hours, days, year, and UT1 sets are marked, with the weighting factors printed below the bits. Wide pulses represent "1" bits and narrow pulses represent "0" bits. Unused bits are set to 0.

In Fig. 11, the decoded UTC at the start of the frame is 1990, 173 days, 21 hours, and 10 minutes. Since the UT1 correction is +0.3 s, the decoded UT1 is 1990, 173 days, 21 hours, 10 minutes, and 0.3 s.

UT1 Time Corrections—The UTC time scale broadcast by WWV and WWVH meets the needs of most users. UTC runs at an almost perfectly constant rate, since its rate is based on cesium atomic frequency standards. Somewhat surprisingly, some users need time *less* stable than UTC but related to the rotation of the Earth. Applications such as celestial navigation, satellite observations of the Earth, and some types of surveying require time referenced to the rotational position of the Earth. These users rely on the UT1 time scale. UT1 is derived by astronomers who monitor the speed of the Earth's rotation.

UT1 time can be obtained by applying a correction to the UTC time signals broadcast from WWV and WWVH. UT1 time corrections are included in the WWV and WWVH broadcasts at two levels of accuracy. First, for those users only needing UT1 to within 1 s, occasional corrections of exactly 1 s are inserted into the UTC time scale. These corrections, called leap seconds, keep UTC within ± 0.9 s of UT1. Leap seconds are coordinated under international agreement by the International Earth Rotation Service in Paris, France. Leap seconds can be either positive or negative, but so far, only positive leap seconds have been needed. A positive leap second is normally added every 1 or 2 years, usually on June 30 or December 31. More information about leap seconds is given below under "Dating of Events in the Vicinity of Leap Seconds."

The second level of correction is for the small number of users needing UT1 accurate to within 0.1 s. These corrections are encoded into the broadcasts by using doubled ticks during the first 16 s of each minute. The amount of correction (in tenths of 1 s) is determined by counting the number of successive doubled ticks heard each minute. The sign of the correction depends on whether the doubled ticks are in the first 8 s of the minute or in the second 8 s. If the doubled ticks are in the first 8 s (1-8) the sign is positive, and if they are in the second 8 s (9-16) the sign is negative. For example, if ticks 1, 2, and 3 are doubled, the correction is "plus" 0.3 s. This means that UT1 equals UTC + 0.3 s. If UTC is 8:45:17, then UT1 is 8:45:17.3. If ticks 9, 10, 11, and 12 are doubled, the correction is "minus" 0.4 s. If UTC is 8:45:17, then UT1 is 8:45:16.6. An absence of doubled ticks indicates that the current correction is 0.

Dating of Events in the Vicinity of Leap Seconds—Leap seconds are sometimes needed to keep UTC within ± 0.9 s of UT1. The addition or deletion of a leap second always occurs at the end of a month. By international agreement, first preference is given to December 31 or June 30. Second preference is given to March 31 or September 30, and third preference is given to any other month.

When UT1 is slow relative to UTC, a positive leap second is needed. The second is inserted beginning at 23^h 59^m 60^s of the last day of the month and ending at 0^h 0^m 0^s of the first day of the following month. The minute containing the leap second is 61-s long. Fig. 12 shows how to assign dates to events occurring near the leap second.

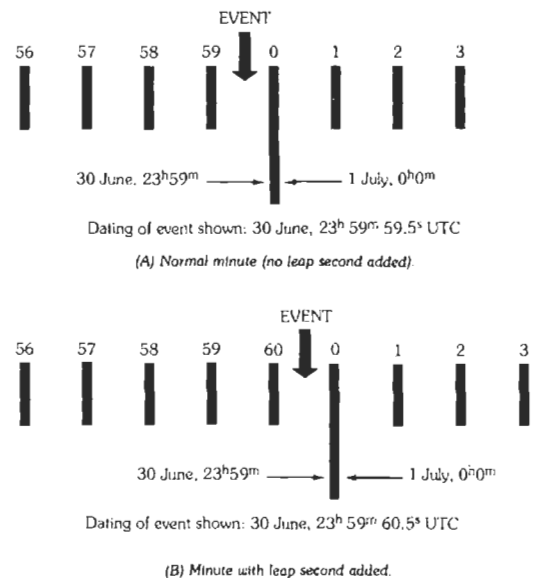


Fig. 12. Dating of events near a leap second.

Official Announcements—Announcement segments 45-s long are available to other Federal agencies by subscription (see Figs. 8 and 9). These segments are used for public service messages. The accuracy and content of these messages is the responsibility of the originating agency.

Geophysical Alerts: The National Oceanic and Atmospheric Administration (NOAA) uses WWV and WWVH to broadcast geophysical alert messages that provide information about solar–terrestrial conditions. Geophysical alerts are broadcast from WWV at 18 minutes after the hour and from WWVH at 45 minutes after the hour. The messages are less than 45 s in length and are updated every 3 hours (typically at 0000, 0300, 0600, 0900, 1200, 1500, 1800, and 2100 UTC). More frequent updates are made when necessary.

The geophysical alerts provide information about the current conditions for long-distance HF radio communications. The alerts use a standardized format and terminology that requires some explanation. Before looking at a sample message, let us define some of the terminology:

Solar flux is a measurement of the intensity of solar radio emissions with a wavelength of 10.7 cm (a frequency of about 2800 MHz). The daily solar flux measurement is recorded at 2000 UTC by the Dominion Radio Astrophysical Observatory of the Canadian National Research Council located at Penticton, British Columbia, Canada. The value broadcast is in solar flux units that range from a theoretical minimum of about 50 to numbers larger than 300. During the early part of the 11-year *sunspot cycle*, the flux numbers are low; but they rise and fall as the cycle proceeds. The numbers will remain high for extended periods around sunspot maximum.

The A and K indices are a measurement of the behavior of the *magnetic field* in and around the earth. The *K index* uses a scale from 0 to 9 to measure the change in the horizontal component of the geomagnetic field. A new K index is determined and added to the broadcast every 3 hours based on magnetometer measurements made at the Table Mountain Observatory, north of Boulder, Colorado, or an alternate middle latitude observatory. The *A index* is a daily value on a scale from 0 to 400 to express the range of disturbance of the geomagnetic field. It is obtained by converting and averaging the eight 3-hour K index values. An estimate of the A index is first announced at 2100 UTC, based on seven measurements and one estimated value. At 0000 UTC, the announced A index consists entirely of known measurements, and the word “estimated” is dropped from the announcement.

Space weather describes the conditions in space that affect earth and its technological systems. Space weather is a consequence of the behavior of the sun, the nature of earth’s magnetic field and atmosphere, and our location in the solar system.

Space weather storms observed and expected are characterized using the NOAA Space Weather Scales.

TABLE 10. NOAA SPACE WEATHER SCALES

Geomagnetic Storms	Solar Radiation Storms	Radio Blackouts	Descriptor
G5	S5	R5	Extreme
G4	S4	R4	Severe
G3	S3	R3	Strong
G2	S2	R2	Moderate
G1	S1	R1	Minor

The abbreviated Table 10 shows the levels of activity that are included in the announcements and the associated terminology. The descriptor used to identify observed or expected conditions is the maximum level reached or predicted. The NOAA Space Weather Scales are further described at <http://www.sec.noaa.gov> on the Space Environment Center’s Website.

Geomagnetic storm levels are determined by the estimated 3-hour planetary K indices derived in real time from a network of Western Hemisphere ground-based magnetometers (see Table 11).

Solar radiation storms levels are determined by the proton flux measurements made by NOAA’s primary *Geostationary Operational Environmental Satellite (GOES)* (see Table 12).

Radio blackouts are determined by the x-ray level measured by the primary *GOES* satellite (see Table 13).

Every geophysical alert consists of three parts as shown in Tables 14 and 15. Table 14 describes the

TABLE 11. GEOMAGNETIC STORM LEVELS

Planetary K Index	Geomagnetic Storm Level
K = 5	G1
K = 6	G2
K = 7	G3
K = 8	G4
K = 9	G5

TABLE 12. SOLAR RADIATION STORM LEVELS

Flux Level of ≥ 10 MeV Particles	Solar Radiation Storm Level
10	S1
10^2	S2
10^3	S3
10^4	S4
10^5	S5

TABLE 13. RADIO BLACKOUTS

Peak X-Ray Level and Flux	Radio Blackout Level
M1 and (10^{-5})	R1
M5 and (5×10^{-5})	R2
X1 and (10^{-1})	R3
X10 and (10^{-3})	R4
X20 and (2×10^{-3})	R5

TABLE 14. INFORMATION IN GEOPHYSICAL ALERT VOICE MESSAGE

Section	Information in Voice Message
1	The solar-terrestrial indices for the day: specifically the solar flux, the A index, and the K index.
2	Space weather storms observed during the previous 24 hours. Includes all observed geomagnetic storms, solar radiation storms (proton events), and radio blackouts (class M1 and greater flares).
3	Space weather expected during the following 24 hours.

TABLE 15. EXAMPLE TEXT FROM ACTUAL GEOPHYSICAL ALERT MESSAGE

Section	Example of Actual Geophysical Alert Message
1	Solar-terrestrial indices for 08 November follow. Solar flux 173 and midlatitude A index 14 The midlatitude K index at 1500 UTC on 08 November was 3.
2	Space weather for the past 24 hours has been severe. Solar radiation storm(s) reaching the S4 level is in progress. Radio blackouts(s) reaching the R2 level occurred.
Alternate Section 2	No space weather storms have been observed during the past 24 hours.
3	Space weather for the next 24 hours is expected to be severe. Solar radiation storms reaching the S4 level are expected to continue. Radio blackouts reaching the R2 level are expected.
Alternate Section 3	No space weather storms are expected during the next 24 hours.

information contained in the geophysical alert. Table 15 provides example text from an actual message.

The announcements include the descriptor of the largest space weather event observed (2) or expected (3) in the first line of each section. The remaining lines give the type of events and the level observed for each one. In the example above, no geomagnetic storm information is included because no geomagnetic storms were observed or were expected during the period. In the case where none of the three types of events are observed or expected, the announcement would contain Section 1, plus alternate Section 2 and alternate Section 3.

To hear the current geophysical alert message by telephone, dial (303) 497-3235. Inquiries regarding these messages should be addressed to: Space Weather Operations, NOAA R/SEC, 325 Broadway, Boulder, CO 80305-3328 (E-mail: swo@sec.noaa.gov; Voice: (303) 497-3171).

Marine Storm Warnings: Marine storm warnings are broadcast for the marine areas that the United States has warning responsibility for under international agreement. The storm warning information is provided by the National Weather Service. Storm warnings for the Atlantic and eastern North Pacific are broadcast by voice on WWV at 8, 9, and 10 minutes after the hour. Storm warnings for the western, eastern, southern, and north Pacific are broadcast by WWVH at 48, 49, 50, and 51 minutes after the hour. An additional segment (at 11 minutes after the hour on WWV and at 52 minutes after the hour on WWVH) is used occasionally if there are unusually widespread storm conditions. The brief voice messages warn mariners of storm threats present in their areas.

The storm warnings are based on the most recent forecasts. Updated forecasts are issued by the National Weather Service at 0500, 1100, 1700, and 2300 UTC for WWV; and at 0000, 0600, 1200, and 1800 UTC for WWVH.

A typical storm warning announcement text is as follows:

North Atlantic weather West of 35 West at 1700 UTC; Hurricane Donna, intensifying, 24 North, 60 West, moving northwest, 20 knots, winds 75 knots; storm, 65 North, 35 West, moving east, 10 knots; winds 50 knots, seas 15 feet.

For more information about marine storm warnings, write to: National Weather Service, NOAA, 1325 East West Highway, Silver Spring, MD 20910. Or see the National Weather Service Website at <http://www.nws.noaa.gov>.

Global Positioning System (GPS) Status Announcements: The United States Coast Guard sponsors two voice announcements per hour on WWV and WWVH, giving current status information about the GPS satellites and related operations. The 40-s announcements begin at 10 minutes after each hour on WWV and at 40

minutes after each hour on WWVH. For further information, contact the US Coast Guard Navigation Center, 7323 Telegraph Road, Alexandria, VA 22310, or call (703) 313-5900.

WWV and WWVH Audio Signals by Telephone—The audio portions of the WWV and WWVH broadcasts can also be heard by telephone. The accuracy of the telephone time signals is normally 30 ms or better in the continental United States. In rare instances when the telephone connection is made by satellite, there is an additional delay of 0.25 to 0.5 s.

To hear these broadcasts, dial (303) 499-7111 for WWV and (808) 335-4363 for WWVH. Callers are disconnected after 3 minutes. *These are not toll-free numbers*; callers outside the local calling area are charged for the call at regular long-distance rates.

Low-Frequency Services—WWVB

NIST radio station WWVB is located on the same site as WWV near Fort Collins, Colorado. WWVB continuously broadcasts time and frequency signals at 60 kHz, in the part of the radio spectrum known as low frequency (LF). The WWVB signal includes a time code that is used to automatically synchronize millions of radio clocks every day in the continental United States and the surrounding areas. These radio clocks include wristwatches, desk clocks, alarm clocks, wall clocks, and clocks integrated into other devices such as televisions and radios. In addition, calibration and testing laboratories use the 60-kHz carrier frequency from WWVB to calibrate electronic equipment and frequency standards.

Accuracy and Stability—WWVB is referenced to the primary NIST frequency standard and atomic time scales located in Boulder, Colorado. The transmitted frequency is accurate to within a few parts in 10^{13} and stable to 1×10^{-13} per day. Propagation effects are minor compared to those of WWV and WWVH. If it is possible to receive the groundwave signal, and proper averaging techniques are used, the received accuracy should be nearly as good as the transmitted accuracy.

Station Identification—WWVB identifies itself by advancing its carrier phase 45° at 10 minutes after the hour and returning to normal phase at 15 minutes after the hour. WWVB is also identified by its unique time code.

Radiated Power, Antenna, and Coverage—The effective radiated power from WWVB is 50 000 watts. WWVB uses two identical antennas spaced 860 m apart. Each antenna is a top-loaded dipole consisting of four 122-m towers arranged in a diamond shape. A system of cables, called a capacitance hat or top hat, is suspended between the four towers. This top hat is electrically isolated from the towers and is electrically connected to a downlead suspended from the center of

the top hat. The downlead serves as the radiating element.

The WWVB signal provides coverage of the continental United States 24 hours per day. During the nighttime hours, the coverage area extends to cover nearly all of North America and parts of South America.

WWVB Time Code—The WWVB time code is synchronized with the 60-kHz carrier and is broadcast continuously at a rate of 1 pulse per second using pulse-width modulation. Each pulse is generated by reducing the carrier power 10 dB at the start of the second, so that the leading edge of every negative-going pulse is on time. Full power is restored 0.2, 0.5, or 0.8 s later as described below.

The WWVB code contains information on the current year, day of year, hour, minute, second, status of Daylight Saving Time, leap year, and a leap-second warning.

The WWVB time code is also sent in BCD format, but the weighting is different from the WWV/WWVH weighting. Bits are sent by shifting the power of the 60-kHz carrier. The carrier power is reduced 10 dB at the start of each second. If full power is restored 200 ms later, it represents a “0” bit. If full power is restored 500 ms later, it represents a “1” bit. Reference markers and position identifiers are sent by restoring full power 800 ms later.

The binary-to-decimal weighting scheme is 8-4-2-1. The *most significant bit* is sent first. This is the reverse of the WWV/WWVH time code. The BCD groups and the equivalent decimal numbers are shown in Table 16.

The decimal number is obtained by multiplying each bit in the binary group by the weight of its respective column and then adding the four products

TABLE 16. 8-4-2-1 BCD-DECIMAL EQUIVALENTS

Binary group	Decimal Equivalent
Weight: 8 4 2 1	
0 0 0 0	0
0 0 0 1	1
0 0 1 0	2
0 0 1 1	3
0 1 0 0	4
0 1 0 1	5
0 1 1 0	6
0 1 1 1	7
1 0 0 0	8
1 0 0 1	9

together. For example, the table shows that the binary group 0101 is equal to 5. This is derived by

$$(0 \times 8) + (1 \times 4) + (0 \times 2) + (1 \times 1) = 0 + 4 + 0 + 1 = 5$$

Every minute, the WWVB time code contains the current minute, hour, day of year, two digits of the current year, a UT1 correction, a leap-second warning bit, and Daylight Saving Time (DST) and leap year indicators. Two BCD groups each are needed to express the hour (00 to 23), minute (00 to 59), and year (00–99); and three groups are needed to express the day of year (001 to 366). Some bits in the BCD groups are unused but may provide additional information in the future. To represent units, tens, or hundreds, the basic 8–4–2–1 weights are simply multiplied by 1, 10, or 100 as appropriate. The coded information refers to the time at the start of the one-minute frame. Seconds are determined by counting pulses within the frame.

Each minute begins with a frame reference pulse lasting for 0.8 s (Fig. 13). A position identifier pulse lasting for 0.8 s is transmitted every 10 s.

UT1 corrections are broadcast at seconds 36 through 43 of each frame. These corrections are to the nearest 0.1 s. The bits transmitted at seconds 36, 37, and 38 show if UT1 is positive or negative with respect to UTC. If “1” bits are sent at seconds 36 and 38, the UT1 correction is positive. If a “1” bit is sent at second 37, the UT1 correction is negative. The amount of the UT1 correction is sent in a four-bit BCD group at seconds 40, 41, 42, and 43. The binary-to-decimal weights are multiplied by 0.1, because the UT1 corrections are expressed in tenths of seconds.

The WWVB time code also contains information about leap years, DST, and leap seconds. The leap year bit is transmitted at second 55. If it is set to “1,” the current year is a leap year. The bit is set to “1” during each leap year sometime after January 1 but before February 29. It is set back to “0” shortly after January 1 of the year following the leap year. Receivers that read this bit can automatically adjust themselves during leap years.

The two DST bits are sent at 57 and 58 s after each minute. If standard time is in effect, both bits (57 and 58) are set to “0.” If DST is in effect, both bits are set to 1. On the day of change from standard to DST, bit 57 is changed from “0” to “1” at 0000 UTC. Exactly 24 hours later, bit 58 also changes from “0” to “1” at 0000 UTC. On the day of change from DST back to standard time, bit 57 goes from “1” to “0” at 0000 UTC, followed 24 hours later by bit 58. Receivers displaying local time can read the DST bits and make the one-hour adjustment automatically when time changes occur locally.

Bit 56 is used to warn users that a leap second will be inserted into the UTC(NIST) time scale at the end of the current month. The bit is set to “1” near the beginning of the month in which a leap second is added. It is reset to “0” immediately following the leap second insertion.

Fig. 13 shows one frame of the time code. The six position identifiers are labeled as P_1 , P_2 , P_3 , P_4 , P_5 , and P_0 . The minutes, hours, days, year, and UT1 sets are marked, with the weighting factors printed below the bits. Wide pulses represent “1” bits, and narrow pulses represent “0” bits. Unused bits are set to “0.”

In Fig. 13, the decoded UTC at the start of the frame is 1990, 258 days, 18 hours, and 42 minutes. Since the UT1 correction is 0.7 s, the decoded UT1 is 1990, 258 days, 18 hours, 41 minutes, 59.3 s.

Internet Time Signals—One common time and frequency application is to synchronize a computer clock to the correct date and time-of-day. This is usually done with a time code received through an Internet or telephone connection.

Internet time servers use standard timing protocols defined in a series of RFC (Request for Comments) documents. The three most common protocols are the Time Protocol, the Daytime Protocol, and the Network Time Protocol (NTP). An Internet time server waits for timing requests sent using any of these protocols, and when a request is received it sends a time code in the correct format.

Client software is available for all major operating systems, and most client software is compatible with either the Daytime Protocol or NTP. Client software that uses the Simple Network Time Protocol (SNTP) makes the same timing request as an NTP client but does less processing and provides less accuracy. Table 17 summarizes the various protocols and their port assignments, or the port where the server “listens” for a client request. NIST operates an Internet time service using multiple servers distributed around the United States. A list of IP addresses for the NIST servers and sample client software can be obtained from the NIST Time and Frequency Division Website:

<http://www.boulder.nist.gov/timefreq>

The uncertainty of Internet time signals is usually < 100 ms, but results vary with different computers, operating systems, and client software.

Telephone Time Signals—Telephone time services allow computers with analog modems to synchronize their clocks using ordinary telephone lines. These services are useful for synchronizing computers that are not on the Internet, or that reside behind an Internet firewall. One example of a telephone service is NIST’s Automated Computer Time Service (ACTS). The telephone number is (303) 494-4774.

ACTS requires a computer, a modem, and client software. When a computer connects to ACTS it receives a time code containing the month, day, year, hour, minute, second, leap second and daylight savings time indicators, and a UT1 correction. The last character in the ACTS time code is the on-time marker (OTM). To compensate for the path delay between NIST and the user, the server sends the OTM 45 ms early. If the client returns the OTM, the server can cal-

Fig. 13. WWVB time code format.

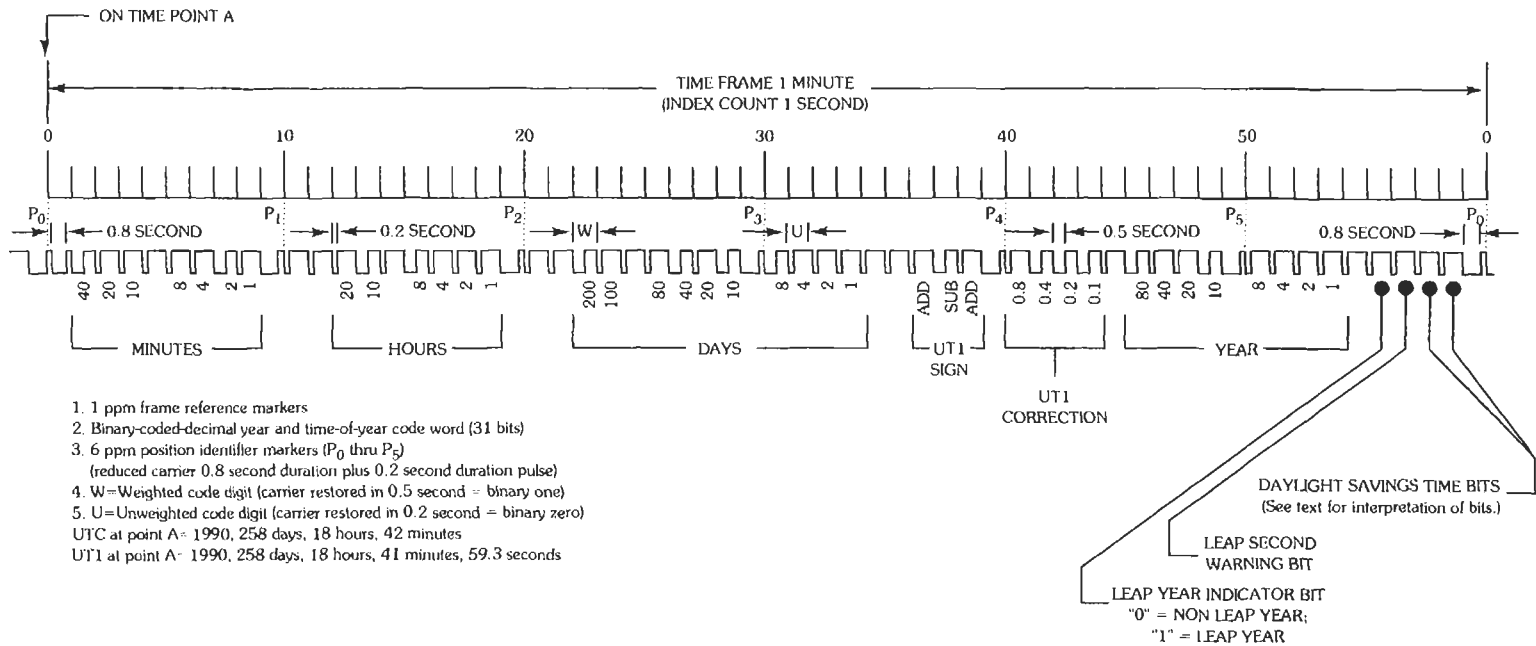


TABLE 17. INTERNET TIME PROTOCOLS

Protocol Name	Document	Format	Port Assignments
Time Protocol	RFC-868	Unformatted 32-bit binary number contains time in UTC seconds since January 1, 1900.	Port 37 tcp/ip, udp/ip
Daytime Protocol	RFC-867	Exact format not specified in standard. Only requirement is that the time code is sent as ASCII characters.	Port 13 tcp/ip, udp/ip
Network Time Protocol (NTP)	RFC-1305	The server provides a data packet with a 64-bit timestamp containing the time in UTC seconds since January 1, 1900 with a resolution of 200 ps. NTP provides accuracy of 1 to 50 ms. The client software runs continuously and gets periodic updates from the server.	Port 123 udp/ip
Simple Network Time Protocol (SNTP)	RFC-1769	The data packet sent by the server is the same as NTP, but the client software does less processing and provides less accuracy.	Port 123 udp/ip

ibrate the path using the *loop-back* method. Each time the OTM is returned, the server measures the round-trip path delay and divides this quantity by 2 to estimate the one-way path delay. This path calibration reduces the uncertainty to < 15 ms.

Global Positioning System (GPS)—GPS, well known as a versatile, global tool for positioning, has also become a primary system for distributing time and frequency. GPS receivers are fixtures in telecommunication networks and in calibration and testing laboratories. They make it possible to synchronize clocks and calibrate and control oscillators in any facility that can place an antenna outdoors for line-of-sight reception of the GPS satellites.

The GPS satellites are controlled and operated by the United States Department of Defense (USDOD). The constellation includes at least 24 satellites that orbit the earth at a height of 20 200 km in six fixed planes inclined 55° from the equator. The orbital period is 11 h 58 m, which means that a satellite will orbit the earth twice per day. By processing signals received from the satellites, a GPS receiver can determine its position with an uncertainty of < 10 m.

The GPS satellites broadcast on two carrier frequencies: L1 at 1 575.42 MHz and L2 at 1 227.6 MHz. Each satellite broadcasts a spread-spectrum waveform, called a pseudo-random noise (PRN) code on L1 and L2, and each satellite is identified by the PRN code it transmits. There are two types of PRN codes. The first type is a coarse acquisition (C/A) code with a chip rate of 1 023 chips per millisecond. The second type is a precision (P) code with a chip rate of 10 230 chips per millisecond. The C/A code is broadcast on L1, and the P code is broadcast on both L1 and L2. GPS reception is line-of-sight, which means that the antenna must have a clear view of the sky. If a clear sky view is available, the signals can be received nearly anywhere on earth.

Each satellite carries either rubidium or cesium oscillators, or a combination of both. The on-board

oscillators provide the reference for both the carrier and code broadcasts. They are steered from USDOD ground stations and are referenced to Coordinated Universal Time (UTC) maintained by the United States Naval Observatory (USNO). By mutual agreement UTC(USNO) and UTC(NIST) are maintained within 100 ns of each other, and the frequency difference between the two time scales is < 1×10^{-15} .

GPS Measurement Techniques: GPS measurements can be divided into two categories: one-way and common-view. Most GPS measurements are one-way measurements, based on direct reception of the GPS signals. Common-view measurements require more effort, including a data exchange between two receivers and postprocessing of the measurement data. For this reason, they are usually reserved for international comparisons between metrology laboratories when the measurement uncertainties must be as small as possible. Table 18 compares the GPS measurement techniques.

One-Way GPS: The one-way GPS technique uses the signals obtained from a GPS receiver as a time and frequency reference. The GPS signals are used in real time, and no postprocessing of the measurement results is required. The signal is normally used either

TABLE 18. TYPICAL UNCERTAINTY OF GPS MEASUREMENT TECHNIQUES

Technique	Timing Uncertainty, 24 h, 2σ	Frequency Uncertainty, 24 h, 2σ
One-Way	< 20 ns	< 2×10^{-13}
Single-Channel Common-View	≈ 10 ns	$\approx 1 \times 10^{-13}$
Multichannel Common-View	< 5 ns	< 5×10^{-14}

to synchronize an on-time pulse or to calibrate a frequency source.

Since the GPS satellites transmit signals that are steered to agree with UTC, the long-term accuracy of a GPS receiver has always been excellent. The performance of C/A code receivers became even better on May 2, 2000 (MJD 51666) when the USDOD set selective availability (SA) to zero. SA is a USDOD directive that can be used to intentionally introduce noise on the GPS signal to reduce its positioning and timing accuracy. Fig. 14 is a phase plot showing data from a typical GPS receiver, recorded in the days immediately before and after SA was set to zero.

The data points shown in Fig. 14 are 10-minute averages of the received signal obtained by comparing the 1 pps output from a typical GPS timing receiver to UTC(NIST) using a time interval counter. The phase plot shows that the GPS broadcasts are tightly controlled. The peak-to-peak variation in the received signal is < 30 ns at one day after SA was set to 0. This leads to excellent long-term accuracy and stability. The Allan deviation, $\sigma_y(\tau)$, plot in Fig. 15 shows that the stability of the receiver is near 1×10^{-13} at 1 day.

Common-View GPS Measurements: The common-view method is a simple but elegant way to compare two clocks or oscillators located in different places. Unlike one-way measurements that compare a clock or oscillator to GPS, a common-view measurement compares two clocks or oscillators to each other.

Fig. 16 shows how the common-view technique works. A GPS satellite (S) serves as a single reference transmitter. The two clocks or oscillators being compared (A and B) are measured against two GPS receivers. The satellite is in common view of both receivers, and its signals are simultaneously received by both.

Frequency Stability of GPS Receiver

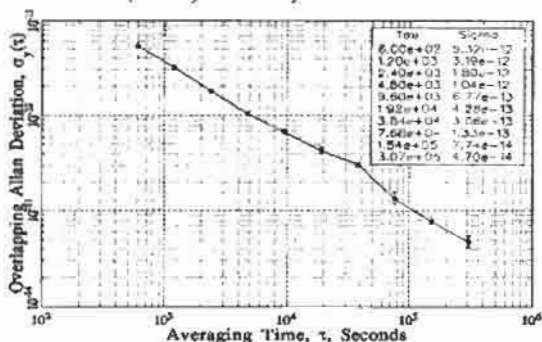


Fig. 15. Frequency stability (Allan deviation) of GPS receiver after SA was set to 0.

Each receiver compares the received signal to its local clock and records the data. Receiver A receives the signal over the path τ_{SA} and compares the reference to the local clock (S-Clock A). Receiver B receives the signal over the path τ_{SB} and records (S-Clock B). The two receivers then exchange and difference the data.

Common-view directly compares two time and frequency standards. Errors from the two paths (τ_{SA} and τ_{SB}) that are common to the reference cancel out, including the performance of the satellite clock. The final measurement result is (Clock A-Clock B) - ($\tau_{SA} - \tau_{SB}$).

The common-view technique has long been used for international comparisons of time and frequency standards. The international time scales—International Atomic Time (TAI) and Coordinated Universal Time (UTC)—are created by averaging data collected from

GPS before and after SA Deactivation

10 minute averages of all satellites in view

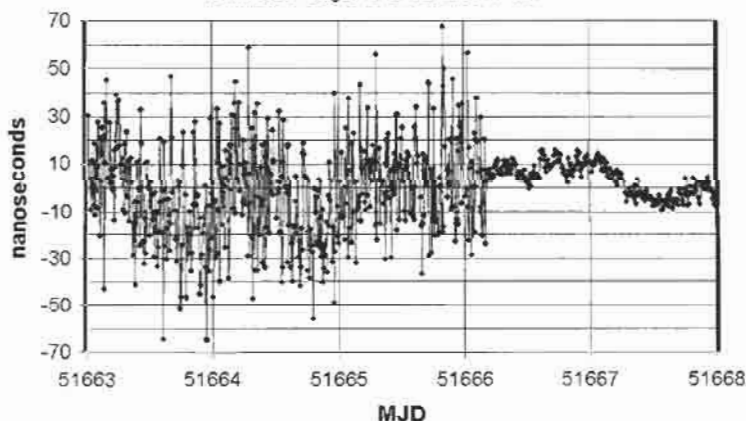


Fig. 14. Phase plot showing GPS performance before and after SA was set to 0.

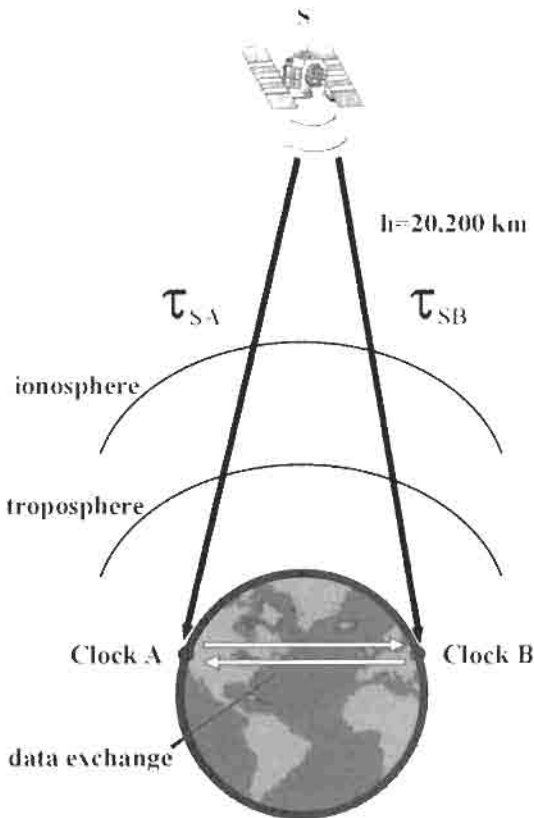


Fig. 16. The common-view measurement technique.

more than 200 atomic clocks located in more than 40 laboratories. Most of these data are collected using common-view GPS measurements. Once the data are sent to the BIPM, the TAI and UTC time scales are computed using a weighted average of all the oscillators. The stability of TAI and UTC is currently about 1 part in 10^{15} over a period of a few weeks.

Common-view measurements can be single channel or multichannel. The BIPM has conventions for recording the data used in both types of measurements. However, it is not necessary to use the BIPM format to make common-view measurements. Any data format can be used if all participants in a common-view comparison record and process their data in the same way.

BIBLIOGRAPHY

References

- Allan, D. W., Ashby, N., and Hodge, C. C. (1997). *The Science of Timekeeping*. Hewlett-Packard Application Note 1289, United States.
- Hackman, C., and Sullivan, D. B., Eds. (1996). *Time and Frequency Measurement*. American Association of Physics Teachers, College Park, Maryland.

- IEEE Standards Coordinating Committee 27. (1999). *IEEE Standard Definitions of Physical Quantities for Fundamental Frequency and Time Metrology—Random Instabilities*. Institute of Electrical and Electronics Engineers, New York.
- ITU Radiocommunication Study Group 7. (1997). *Selection and Use of Precise Frequency and Time Systems*. International Telecommunication Union, Geneva, Switzerland.
- Jespersen, J., and Fitz-Randolph, J. (1999). *From Sundials to Atomic Clocks: Understanding Time and Frequency*. NIST Monograph 155, US Government Printing Office, Washington, DC.
- Levine, J. (1999). Introduction to time and frequency metrology. *Review of Scientific Instruments* **70**, 2567–2596.
- Seidelmann, P. K., Ed. (1992). *Explanatory Supplement to the Astronomical Almanac*. University Science Books, Mill Valley, California.
- Sullivan, D. B., Allan, D. W., Howe, D. A., and Walls, F. L., Eds. (1990). *Characterization of Clocks and Oscillators*. NIST Technical Note 1337, US Government Printing Office, Washington, DC.
- Vig, J. R. (1992). *Introduction to Quartz Frequency Standards*. United States Army Research and Development Technical Report SLCET-TR-92-1.

Spectrum Management

For international technical standards, see the *Recommendations of the International Telecommunication Union*, ITU-R Publications, latest edition, available from the International Telecommunication Union, Geneva, Switzerland.

All ITU-R Recommendations in force at the end of the Radiocommunication Assembly (RA-97) (24 October 1997) are published in the 1997 Volumes (except for the few "old" Recommendations published in the 1992 RF Series Fascicule and in Volumes III and IX of the CCIR XVIIth Plenary Assembly (Düsseldorf, 1990)). ITU-R Recommendations are divided into Series according to the subject areas covered, as follows:

- SM SERIES** (SG 1) Spectrum Management
- IS SERIES** (SG 1) Interservice Sharing and Compatibility
- P SERIES—Part 1** (SG 3) Radiowave Propagation
- P SERIES—Part 2** (SG 3) Radiowave Propagation
- S SERIES** (SG 4) Fixed-Satellite Service
- SNG SERIES** (SG 4) Satellite News Gathering
- TF SERIES** (SG 7) Time Signals and Frequency Standards Emissions
- SA SERIES** (SG 7) Space Applications and Meteorology
- RA SERIES** (SG 7) Radioastronomy
- M SERIES—Part 1** (SG 8) Land Mobile Service Excluding IMT-2000
- M SERIES—Part 2** (SG 8) International Mobile Telecommunications-2000 (IMT-2000)

M SERIES—Part 3 (SG 8) Maritime Mobile Service and Aeronautical Mobile Service

M SERIES—Part 4 (SG 8) Radiodetermination Service

M SERIES—Part 5 (SG 8) Mobile-Satellite Services and Radiodetermination-Satellite Service

M SERIES—Part 6 (SG 8) Amateur Service and Amateur-Satellite Service

F SERIES—Part 1 (SG 9) Fixed Service—Radio-Relay Systems

F SERIES—Part 2 (SG 9) Fixed Service—Frequency Sharing Aspects

F SERIES—Part 3 (SG 9) Fixed Service—HF Systems

SF SERIES (SG 4–9) Frequency Sharing between the Fixed-Satellite Service and the Fixed Service

BS SERIES (SG 10) Broadcasting Service (Sound) **Rec. ITU-R BS.705-1** and **Rec. ITU-R BS.1195HF**

Transmitting and Receiving Antennas Characteristics and Diagrams—Transmitting Antenna Characteristics at VHF and UHF

BT SERIES (SG 11) Broadcasting Service (Television)

BO SERIES (SG 10-11S) Broadcasting-Satellite Service (Sound and Television)

BR SERIES (SG 10-11R) Sound and Television Recording

V SERIES (CCV) Vocabulary and Related Subjects

Communications Act of 1934, as amended. US Congress.

“Manual of Regulations and Procedures for Federal Radio Frequency Management.” Washington, DC: National Telecommunications and Information Administration, updated annually.

Title 47, Code of Federal Regulations, “Federal Communications Commission,” four volumes. Washington, DC: Office of the Federal Register, updated annually.

“Spectrum Management.” Appendix A, *NTIA TELECOM 2000: Charting the Course for a New Century*, NTIA Special Publication 88–21. Washington, DC: National Telecommunications and Information Administration, October 1988.

Matos, F. *Spectrum Management and Engineering*. New York: IEEE Press, 1985. (Reprints of 51 papers on various international and domestic spectrum management topics plus additional material on legal and regulatory matters, and spectrum management tools and methods.)

“U.S. Spectrum Management Policy: Agenda for the Future,” NTIA Special Publication 91–23. Washington, DC: National Telecommunications and Information Administration, February 1991.

Nunno, Richard M. “Radiofrequency Spectrum Management: Background, Status, and Current Issues,” CRS Report for Congress RL30829. Washington,

DC: Congressional Research Service, January 30, 2001.

Information from the Internet

Radio frequency spectrum management information is available on the Internet. The telecommunications regulatory administrations of many countries have active websites that should be consulted when dealing with spectrum issues involving those nations. Three extensively used sites are:

1. National Telecommunications and Information Administration at www.ntia.doc.gov
2. Federal Communications Commission at www.fcc.gov
3. International Telecommunication Union at www.itu.int

Spectrum Management Training

Radio frequency spectrum management training is available from two sources:

1. George Washington University, Center for Professional Development, Washington, DC. www.cpd.gwu.edu
2. For students from developing nations: United States Telecommunications Training Institute, Washington, DC. www.ussti.org

Time and Frequency

Blair, B. E., Ed. “Time and Frequency: Theory and Fundamentals.” NBS Monograph 140. Boulder, CO: National Bureau of Standards, Dept. of Commerce, 1974.

Technical standards published as reports and recommendations of the International Radio Consultative Committee (CCIR), International Telecommunication Union (ITU), Geneva, Switzerland. Reports and recommendations are revised periodically, and it is best to refer to the latest version. In particular:

“Performance of Standard-Frequency Generators,” Report 346–6 or later version, Reports of the CCIR, 1990 or later, Annex to Volume VII, Standard Frequencies and Time Signals.

“Standard-Frequency and Time-Signal Emissions,” Recommendation 460–4 or later version, Recommendations of the CCIR, 1990 or later, Volume VII, Standard Frequencies and Time Signals.

Beehler, R. E., and Lombardi, M. A. “NIST Time and Frequency Services,” NIST Special Publication 432. Boulder, CO: Dept. of Commerce, National Institute of Standards and Technology, June 1991.

2 International Telephony Recommendations

Revised by

Douglass D. Crombie and Dorothy M. Cerni

International Standards	2-3
Recommendations of the CCITT	
Recommendations of the CCIR	
Zero-Relative-Level Points and Relative Levels	2-3
Psophometric Noise and Power	2-4
Psophometric Weighting for Commercial Telephone Circuits	
Psophometric Weighting Factor	
Psophometric Power	
Conventional Telephone Signal	2-5
Telephone Circuit Loading	2-5
Nominal Mean Power During Busy Hour	
Conventional Load	
Power Levels	2-6
Maximum Power Level for Signaling Pulses	
Private Telegraph Transmission on a Rented International Circuit, With Alternative Private Telephone Service	
Simultaneous Communication by Telephony and Telegraphy on a Telephone Circuit	
Phototelegraphy Transmissions Over Telephone Circuits That Are Entirely 4-Wire Between Phototelegraph Stations	
Power Levels for Data Transmission Over Telephone Circuits	

Hypothetical Reference Circuits	2-7
General Definitions (G.212)	
Important Hypothetical Reference Circuits Defined by CCITT and CCIR	
Telephone Circuit Characteristics	2-9
Impedance of International and National Trunks	
One-Way Propagation Time	
Group Delay Distortion	
Attenuation Distortion in Worldwide Chain	
Variation of Transmission Loss With Time	
Linear Cross Talk	
Frequency Accuracy of Virtual Carrier Frequencies on an International Circuit	
Circuit Noise Objectives	2-11
Noise Characteristics for 252-Kilometer Hypothetical Reference Circuits	
Noise Characteristics for Long Circuits Not More Than 252 Kilometers in Length	
Noise Characteristics for International Circuits More Than 252 Kilometers in Length	
Noise in Actual Circuits	
Designed Objectives for Noise Produced by Modulating Equipments	
CCITT and Telegraphy	2-13
Numbering	
CCITT and Telephony	2-13
International Country Codes	
Telephone Signaling	2-13
Supervisory or Line Signals	
Ringdown Signals	
Tones	
Alternative Routing	

INTERNATIONAL STANDARDS

Administrations and operating companies throughout the world carry on studies of technical and other problems related to the interworking of their respective national telecommunication systems to provide a worldwide telecommunications network. Two international committees exist for this purpose: The International Telegraph and Telephone Consultative Committee (CCITT), and the International Radio Consultative Committee (CCIR). They operate under the auspices of the International Telecommunication Union (ITU). They promulgate their decisions in the form of Recommendations, which are published by the ITU. Generally, these Recommendations cover features of international circuits, but where essential, they deal with relevant characteristics of the national systems which may form part of international connections. This compendium collects, in condensed form, major Recommendations dealing with telephone, telegraph, and data-transmission circuits and equipment.

Recommendations of the CCITT

The CCITT develops new Recommendations, and updates existing ones, through the activities of Study Groups, whose reports are acted on at Plenary Assemblies, which meet at intervals of three or four years. The resulting Recommendations of the Second Plenary Assembly, New Delhi, 1960, were published by the ITU in a number of volumes, called collectively the Red Book. The subsequent study periods culminated in the Third Plenary Assembly, Geneva, 1964 (Blue Book); the Fourth Plenary Assembly, Mar del Plata, 1968 (White Book); the fifth Plenary Assembly, Geneva, 1972 (Green Book); the Sixth Plenary Assembly, Geneva, 1976 (Orange Book); and the Seventh Plenary Assembly, Geneva, 1980 (Yellow Book). This compendium refers to Yellow Book Recommendations, designated thus: (G.101), (H.31), (V.2), etc.

Recommendations of the CCIR

The CCIR also functions with Study Groups and Plenary Assemblies. The Eleventh Plenary Assembly was held at Oslo in 1966, the Twelfth Plenary Assembly at New Delhi in 1970, the Thirteenth Plenary Assembly at Geneva in 1974, the Fourteenth Plenary Assembly at Kyoto in 1978, and the Fifteenth Plenary Assembly at Geneva in 1982. After each Plenary Assembly, the ITU publishes volumes which contain the currently accepted Recommendations, including such Recommendations of the Plenary Assemblies at London (1953), Warsaw (1956), Los Angeles (1959), and Geneva (1963) which are still in effect. No color coding is used. This compendium deals with those Recommendations which treat point-to-point radio relay systems. A purpose of

those Recommendations is to make the performance of such systems compatible with metallic line systems which follow the CCITT Recommendations. References to the CCIR Recommendations are made thus: (CCIR,391).

NOTE: This chapter is primarily concerned with telephony, but additional material is included in chapter 38, Common Carrier Transmission. For more information on data communications, the reader is referred to chapter 26, Computer Communications Networks.

ZERO-RELATIVE-LEVEL POINTS AND RELATIVE LEVELS

Many CCITT and CCIR Recommendations specify signal or noise levels at "a point of zero relative level," or in dBm₀ or pWp₀, etc., where "0" (zero) stands for "measured at or referred to a point of zero relative level."

(A) In two-wire switching systems, the sending-end terminals of a long-distance circuit have long been considered to be at a point of zero relative level. The relative levels of all other points are calculated from this reference point, as the algebraic sum of all transmission losses and gains from it to the point in question. Any point in a circuit with the same relative level as the sending terminals is a point of zero relative level, which may be written 0 dBr (*dB* relative level). The American term for *relative level* is *transmission level*. Thus: "Zero-transmission-level point" (OTLP).

For convenience in comparing circuit noise performance, it is customary to convert absolute noise measurements made at the receiving ends of circuits having various relative levels, to absolute power levels at a zero-relative-level point. For example, -50 dBmp of noise at a -7-dBr point would be reported as -43 dBm_{0p}. Signaling-tone levels are similarly expressed. For example, a tone introduced at a -3.5-dBr point with an absolute power level of -18.5 dBm may be referred to as a -15-dBm₀ signal. The latter designation would apply to such a tone no matter where it appeared; the "0" denotes that its level is referred to a point of zero relative level. (Refer to Table 1.)

Statistics of speech power, requirements for linearity and limiting, system loading factors, cross talk, and noise have become well known in terms of their values at points of zero relative level. The proper performance of voice repeaters, carrier terminal and line equipment, radio relay systems, etc., depends on adherence to the relative levels for which they were designed. Many relative levels associated with such equipment have been standardized.

(B) In four-wire switching systems, it is often considered desirable to handle speech and signaling at lower values of absolute power through the switching

TABLE 1. RELATIVE LEVELS, PLUS LEVELS OF ABSOLUTE AND REFERRED POWER, FOR A 2-WIRE CIRCUIT

	Sending Point ○	Circuit Loss: 5 dB Direction→	Receiving Point ○
○—————○ 2-Wire Circuit			
Relative levels, or transmission levels	0 dBr		-5 dBr
"Milliwatt" test tone At any point	0 dBm 0 dBm0		-5 dBm absolute power 0 dBm0 referred power
Signaling tone At any point	-10 dBm -10 dBm0		-15 dBm absolute power -10 dBm0 referred power
Circuit noise, picked up along circuit	—		-65 dBmp absolute power
Circuit noise, referred to 0-dBr point	—		-60 dBm0p referred power

equipment than is customary in two-wire systems. In 1964, the CCITT adopted a relative level of -3.5 dBr for the sending end of a four-wire circuit, at the "virtual" switching points. These are theoretical points; their exact location depends on national practice, and the CCITT considers it unnecessary to define them. (In the American commercial system, -2 dBr is widely used.) (G.101, Section B)

Therefore, to ensure that carrier and other transmission equipment will be subjected to the same absolute speech and signaling power levels as in two-wire systems, determination of relative levels in four-wire circuits must take into account the relative level of the virtual switching points. In a four-wire circuit, there may be no actual point of zero relative level. Nevertheless, standards will continue to refer many requirements to a zero relative point. (Refer to Table 2.)

Currently, many transmission measurements are made with a standard 800- or 1000-hertz test tone, with an absolute power of 1 milliwatt at a zero-relative-level

point: a power of 0 dBm0. The actual level applied is adjusted to the relative level of the sending point. The test-tone level in dBm will be numerically equal to the relative level in dBr at any point in the circuit, but it is not proper to express *relative* levels in dBm, since dBm represents absolute power levels. If the standard-test-tone power is ever changed to another value, such as -10 dBm0, as has been tentatively proposed, the distinction between relative levels and test-tone levels will be more apparent.

PSOPHOMETRIC NOISE AND POWER

The CCITT calls a noise measuring set a "psophometer." A psophometer includes a device for measuring power through a weighting network. For measurements on commercial telephone circuits, a weighting characteristic is used which results in the objective instrument

TABLE 2. RELATIVE LEVELS, PLUS LEVELS OF ABSOLUTE AND REFERRED POWER, FOR ONE DIRECTION OF A 4-WIRE CIRCUIT

	Sending Point ○	Circuit Loss: 5 dB Direction→	Receiving Point ○
○—————○ One side of 4-Wire Circuit			
Relative levels, or transmission levels	- 3.5 dBr		- 8.5 dBr
"Milliwatt" test tone At any point	- 3.5 dBm 0 dBm0		- 8.5 dBm absolute power 0 dBm0 referred power
Signaling tone At any point	-13.5 dBm -10 dBm0		-18.5 dBm absolute power -10 dBm0 referred power
Circuit noise, picked up along circuit	—		-68.5 dBmp absolute power
Circuit noise, referred to 0-dBr point	—		-60 dBm0p referred power

measurements approximately paralleling the results of subjective tests with human observers using modern telephone sets. The CCITT weighting characteristic for commercial circuits is nominally identical with the American FIA line weighting. Psophometric noise power may be expressed in dBm0 "psophometrically weighted," or dBm0p. The conventional conversion equation used between dBm0p and dBa0 (FIA) is

$$\text{dBm0p} = \text{dBa0} - 84$$

Psophometric Weighting for Commercial Telephone Circuits

Frequency (hertz)	Level (dB)	Frequency (hertz)	Level (dB)
100	-41.0	1350	-0.65
150	-29.0	1500	-1.30
200	-21.0	1750	-2.22
250	-15.0	2000	-3.00
300	-10.6	2250	-3.60
400	-6.3	2500	-4.20
500	-3.6	2750	-4.87
600	-2.0	3000	-5.60
700	-0.9	3500	-8.5
800	0.0	4000	-15.0
900	+0.6	4200	-18.7
1000	+1.0	4500	-25.0
1100	+0.6	4700	-29.4
1200	0.0	5000	-36.0

(Extracted from G.223, which gives weighting for 84 frequencies.)

Psophometric Weighting Factor

If uniform-spectrum random noise is measured in a 3.1-kilohertz band with a flat attenuation/frequency characteristic, the noise level must be reduced by 2.5 decibels to obtain the psophometric power level. For another bandwidth *B*, the weighting factor will be equal to

$$2.5 + 10 \log_{10}(B/3.1) \text{ decibels}$$

When *B* = 4 kilohertz, for example, this gives a weighting factor of 3.6 decibels. (G.223)

Psophometric Power

Where power addition of noise can be assumed, it has been found convenient for calculations and design of international circuits to use the concept of "psophometric power."

psophometric power

$$= (\text{psophometric voltage})^2/600$$

$$= (\text{psophometric emf})^2/(4 \times 600)$$

A convenient unit is the picowatt (pW) = 10⁻¹² watt, so that

psophometric power in pW

$$= (\text{psophometric emf in mV})^2/0.0024 \quad (\text{G.212})$$

CONVENTIONAL TELEPHONE SIGNAL

For the calculation or measurement of cross-talk noise between adjacent channels, of the balance return loss for echo, and generally speaking, when it is desired to simulate the speech currents transmitted by a telephone channel, the CCITT recommends the use of a conventional telephone signal. This signal may be produced by passing the output of a generator of a uniform-spectrum random noise signal ("white noise") through a weighting network with a characteristic as shown in Fig. 1. The amount of this signal that appears in another circuit because of cross talk, etc., is measured with a psophometer or weighted-noise measuring set, with standard psophometric weighting for commercial telephone circuits. (G.227)

TELEPHONE CIRCUIT LOADING

Nominal Mean Power During Busy Hour

To simplify calculations when designing carrier systems on cables or radio links, the CCITT has adopted a *conventional* value to represent the *mean absolute power level*, at a point of zero relative level, of the speech-plus-signaling currents, etc., transmitted over a telephone channel in one direction of transmission during the busy hour, which is -15 dBm (-1.73 nepers) (mean power = 31.6 microwatts); this is the mean with time and the mean for a large batch of circuits. This total mean power of about 32 microwatts is conventionally distributed as follows (nominal mean power): 10 microwatts, all signaling and tones; 22 microwatts, to include speech currents (including echoes), carrier leak, and telegraph signals, based on a speech activity factor of 0.25 for one direction of a telephone channel. No account is taken of pilot signals, which are assumed to be an integral part of the carrier system, not affecting telephone channel power. (G.223, Section 1)

Conventional Load

It will be assumed for the calculation of intermodulation noise below the overload point that the multiplex signal during the busy hour can be represented by a

uniform-spectrum random noise signal, the mean absolute power level of which, at a zero-relative-level point, $n(\bar{P})$, is given by

$$n(\bar{P}) = -15 + 10 \log_{10} N \text{ dB}$$

for $N = 240$ or more

and

$$n(\bar{P}) = -1 + 4 \log_{10} N \text{ dB}$$

for values of N between 12 and 240

where N is the total number of telephone channels in the system.

Typical values so calculated are as follows

N	$n(\bar{P})$, dB
12	3.3
24	4.5
36	5.2
48	5.7
60	6.1
120	7.3
⋮	⋮
240	8.8
300	9.8
600	12.8
960	14.8
1800	17.0
2700	19.3

Assumed: No pre-emphasis, and use of independent amplifiers for each direction. (G.223, Section 2)

POWER LEVELS

Maximum Power Level for Signaling Pulses

For cross-talk reasons, each component of a short-duration signal should not exceed the following absolute power levels, at a zero-relative-level point.

Signaling Frequency (hertz)	Absolute Power Level at Zero-Relative-Level Point (dBm0)
800	-1
1200	-3
1600	-4
2000	-5
2400	-6
2800, 3200	-8

(G.224)

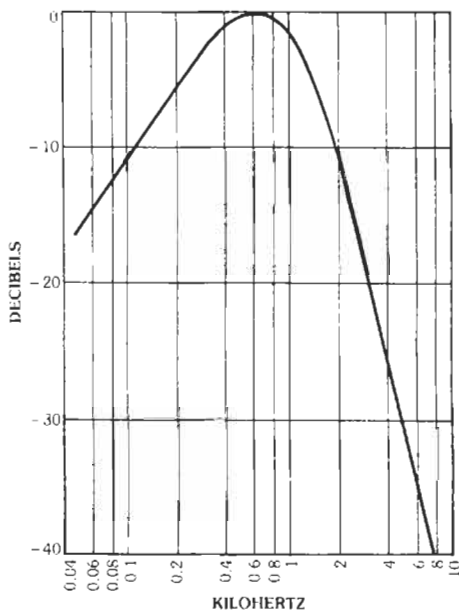


Fig. 1. Relative response curve for the weighting network of the conventional telephone signal generator (G.227).

Private Telegraph Transmission on a Rented International Circuit, With Alternative Private Telephone Service

The frequency of 1500 hertz is recommended for private telegraph transmission between subscribers permanently connected via a rented international circuit. The permissible power for a continuously transmitted telegraph marking signal is 0.3 milliwatt at a zero-relative-level point (-5 dBm0). (H.31)

Simultaneous Communication by Telephony and Telegraphy on a Telephone Circuit

A continuously transmitted telegraph signal should not exceed a level of -13 dBm0. There should not be more than three telephone circuits of this type per group, nor more than the number of supergroups in a wideband system. (H.32)

Phototelegraphy Transmissions Over Telephone Circuits That Are Entirely 4-Wire Between Phototelegraph Stations

The sent voltage for the phototelegraph signal corresponding to maximum amplitude should be so adjusted

that the absolute power level of the signal, at a zero-relative-level point, is 0 dBm0 for amplitude-modulation facsimile, and -10 dBm0 for frequency-modulation facsimile. In the former, the "black" level is about 30 decibels lower than the "white" level. (H.41)

Power Levels for Data Transmission Over Telephone Circuits

Private Wires on Carrier Systems—

(A) Maximum power output of subscriber's apparatus into line: 1 milliwatt.

(B) Continuous-tone systems (for example, frequency modulation): Maximum power level at zero-relative-level point: -10 dBm0, to be reduced to or below -20 dBm0 when data transmission is discontinued for any appreciable time.

(C) Noncontinuous-tone systems (for example, amplitude modulation): Maximum power level at zero-relative-level point: -6 dBm0, provided that busy-hour mean power in both directions of transmission added does not exceed 64 microwatts (-15 dBm0 mean level in each direction simultaneously). Also, the level of tones above 2400 hertz should conform to recommendations for signaling tones in G.224. (H.51)

Switched Telephone Network—

(A) Maximum power output of subscriber's apparatus into line: 1 milliwatt.

(B) Continuous-tone systems (for example frequency- or phase-modulation systems): The power level at the subscriber's equipment should be adjusted not to exceed -10 dBm0 (simplex operation) or -13 dBm0 (duplex operation).

(C) Noncontinuous-tone systems (for example amplitude-modulation or multifrequency systems): Higher levels may be used, if mean power at international-circuit input is limited to 64 microwatts in any hour (both directions), i.e., -15 dBm0 in each direction simultaneously. (H.51)

HYPOTHETICAL REFERENCE CIRCUITS

General Definitions (G.212)

Hypothetical Reference Circuit—A hypothetical circuit of defined length and with a specified number of terminal and intermediate equipments, this number being sufficient but not excessive. It forms the basis for the study of certain characteristics of long-distance circuits, for example, noise.

Hypothetical Reference Circuit for Telephony—This is a complete telephone circuit, between audio-frequency terminals, established on a hypothetical international telephone carrier system and having a speci-

fied length and a specified number of modulations and demodulations of the groups, supergroups, and mastergroups, these numbers being reasonably great but not having their maximum possible values.

Various hypothetical reference circuits for telephony have been defined by the CCITT and CCIR to allow the coordination of the different specifications concerning the constituent parts of the multichannel carrier telephone systems, so that the complete telephone circuits set up on these systems can meet CCITT standards.

Homogeneous Section—A section without diversion or modulation of any one of the mastergroups, supergroups, groups, or channels established on the system being considered, except for those modulations or demodulations defined at the ends of the section. All the hypothetical reference circuits defined below consist of homogeneous sections of equal length (six or nine sections). It is assumed that at the end of each homogeneous section, the channels, groups, supergroups, and mastergroups, as appropriate, are connected through at random.

Important Hypothetical Reference Circuits Defined by CCITT and CCIR

Hypothetical Reference Circuit on Symmetric Cable Pairs—This circuit is 2500 kilometers long and is set up on a symmetric-cable-pair carrier system. For each direction of transmission, it has a total of three pairs of channel modulators and demodulators, six pairs of group modulators and demodulators, and six pairs of supergroup modulators and demodulators.

Fig. 2A shows that there are 15 modulations and 15 demodulations for each direction of transmission, assuming that single-stage translations are used. There are six homogeneous sections of equal length. (G.322)

Hypothetical Reference Circuit for 4-Megahertz Systems on Coaxial Cable—This circuit is 2500 kilometers long and is set up on a 4-megahertz carrier system on 0.104/0.375-inch coaxial cable pairs. For each direction of transmission, it has a total of three pairs of channel modulators and demodulators, six pairs of group modulators and demodulators, and nine pairs of supergroup modulators and demodulators.

Fig. 2B shows that there are 18 modulations and 18 demodulations for each direction of transmission, assuming that single-stage translations are used. There are nine homogeneous sections of equal length. (G.338, Orange Book)

Hypothetical Reference Circuit for 12-Megahertz Systems on Coaxial Cable—This circuit is 2500 kilometers long and is set up on a 12-megahertz carrier system on 0.104/0.375-inch coaxial cable pairs. For each direction of transmission, it has a total of three pairs of channel modulators and demodulators, three

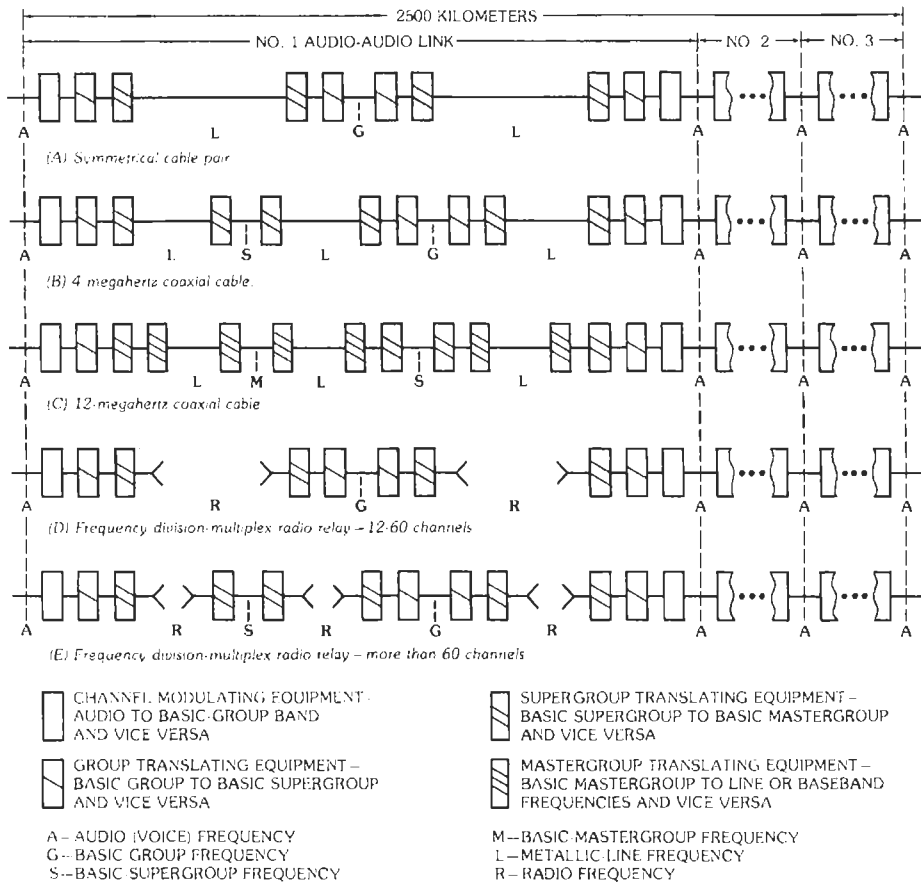


Fig. 2. Typical hypothetical reference circuits for systems using frequency-division multiplex.

pairs of group modulators and demodulators, six pairs of supergroup modulators and demodulators, and nine pairs of mastergroup modulators and demodulators.

Fig. 2C shows that there are 21 modulations and 21 demodulations for each direction of transmission, assuming that single-stage translations are used. There are nine homogeneous sections of equal length. (G.332)

Hypothetical Reference Circuit Over Radio Relay Systems With Frequency-Division Multiplex, Providing 12 to 60 Channels—This circuit is 2500 kilometers long and is set up on a carrier system providing 12 to 60 channels per radio channel over line-of-sight and near-line-of-sight radio relay systems. For each direction of transmission, it has a total of three sets of channel modulators and demodulators, six sets of group modulators and demodulators, six sets of supergroup modulators and demodulators, and six sets of radio modulators and demodulators. The circuit is divided into six homogeneous sections of equal length. (See Fig. 2D.) (CCIR, 391)

Hypothetical Reference Circuit Over Radio Relay Systems With Frequency-Division Multiplex,

Providing More than 60 Channels—This circuit is 2500 kilometers long and is set up on a carrier system providing more than 60 channels per radio channel over line-of-sight or near-line-of-sight radio relay systems. For each direction of transmission, it has a total of three sets of channel modulators and demodulators, six sets of group modulators and demodulators, nine sets of supergroup modulators and demodulators, and nine sets of radio modulators and demodulators. The circuit is divided into nine homogeneous sections of equal length. (See Fig. 2E.) (CCIR, 392)

Hypothetical Reference Circuit Over Transhorizon Radio Relay Systems—This circuit is 2500 kilometers long and is set up over a transhorizon radio relay system for telephony using frequency-division multiplex. For each direction of transmission it has a total of three sets of channel modulators and demodulators, six sets of group modulators and demodulators, and nine sets of supergroup modulators and demodulators.

It is recommended that the hypothetical reference circuit not be divided into homogeneous sections of

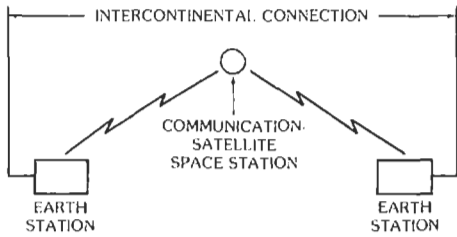


Fig. 3. Basic hypothetical reference circuit, active intercontinental communication-satellite system (CCIR, 352-3).

fixed length, because these systems, as distinct from line-of-sight systems, are usually composed of long radio sections, the lengths of which depend on local conditions and may vary considerably (between 100 and 400 kilometers). If a radio section under study is L kilometers long, the hypothetical reference circuit should be composed of $(2500/L)$ sections of this type in tandem, the value $(2500/L)$ being taken to the nearest whole number. (CCIR, 396-1)

Hypothetical Reference Circuit Over Active Intercontinental Communication-Satellite Systems—This circuit has no fixed length. For intercontinental connections, satellite links should be capable of spanning 7500 kilometers. For great-circle distances up to 25 000 kilometers, it will be necessary to connect two or three satellite links in tandem. The basic hypothetical reference circuit shall consist of one earth-satellite-earth link, as shown in Fig. 3. It shall contain one pair of modulation and demodulation equipments for translation from the baseband to the radio-frequency carrier and back. (CCIR, 352-3)

TELEPHONE CIRCUIT CHARACTERISTICS

Impedance of International and National Trunks

All circuits (whether international circuits or national two-wire or four-wire trunk circuits) terminating at the same trunk exchange should have the same nominal value of impedance as seen from the switchboard or selectors. The preferred value is 600 ohms. (G.232, Section M)

One-Way Propagation Time

The one-way propagation time of connections when echo sources exist and echo suppressors are used is considered to be: acceptable without reservation, 0-150 milliseconds; provisionally acceptable, 150-400 milliseconds; and provisionally unacceptable, 400 milliseconds or higher (should not be used except under extraordinary circumstances).

On National Extension Circuits—Probable propagation time, most distant subscriber to international center: 12 plus $(0.0064 \times \text{distance in statute miles})$ milliseconds or 12 plus $(0.004 \times \text{distance in kilometers})$ milliseconds.

International Circuits, Terrestrial (Including Submarine Cable)—100 statute miles (160 kilometers) per millisecond, including effects of terminal or intermediate multiplex equipment.

International Circuits, Communication-Satellite Systems—Moving satellite, altitude 8700 miles (14 000 kilometers): 110 milliseconds. Geostationary satellite, altitude 22 500 miles (36 000 kilometers): 260 milliseconds. (G.114)

Group Delay Distortion

The permissible differences, for a worldwide chain of 12 circuits, between the minimum group delay throughout the frequency band transmitted and the group delay at the upper and lower limits of this band in milliseconds are as follows:

Facility	At the Limits of the Frequency Band	
	Lower Limit	Upper Limit
International chain	30	15
Each of the national 4-wire extensions	15	7.5
Entire 4-wire chain	60	30

(G.113, P.13)

The CCITT does not recommend limiting values for the group delay at different frequencies, but offers information on typical values reported by four Administrations. Values are in milliseconds measured on a pair of equipments.

Administration	Frequency (hertz)				
	300	400	2000	3000	3400
Belgium (b)	4.0	2.7	1.0	1.3	2.6
France (b)	4.2	2.9	1.0	1.4	2.8
Federal Republic of Germany (b)	3.9	2.7	1.2	1.6	3.0
United Kingdom (a)	2.6	2.2	1	1.4	2.6
(b)	4.2	2.7	1.2	1.8	3.4
For 12 pairs (c)	50	35	14	22	41

(a) With in-band signaling
 (b) With out-band signaling
 (c) Typical maximum values in milliseconds for 12 pairs of equipments.

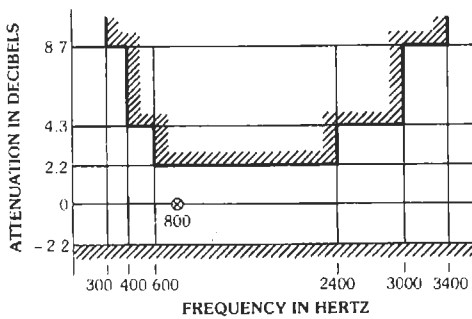


Fig. 4. Attenuation vs. frequency objective for worldwide chain of 12 circuits in terminal service (G.132).

Attenuation Distortion in Worldwide Chain

The objectives for the variation with frequency of transmission loss in the terminal condition of a worldwide chain of 12 circuits (international plus national extensions), each one routed over a single group link, are shown in Fig. 4, which assumes that no use is made of high-frequency radio circuits or 3-kilohertz channel equipment. (G.132)

Variation of Transmission Loss With Time

(A) The standard deviation for the variation in transmission loss of a circuit should not exceed 1 decibel. This can be met for circuits on single-group links with automatic regulation, and it should be met on any national circuit whether regulated or not. For other international circuits, the standard deviation should not exceed 1.5 decibels. Standard deviation is sometimes called "distribution grade."

(B) The difference between the mean value of the transmission loss and the nominal value (that is, the bias) should not exceed 0.5 dB. (G.151, Section C)

Linear Cross Talk

(A) Between circuits, the near-end or far-end cross-talk ratio (intelligible only), measured at audio frequency at a national or international center between two complete circuits in terminal service, should not be less than 58 decibels for 90 percent of all two-circuit combinations, and 52 decibels for all combinations. The planning objective for modern systems for all combinations is 58 dB.

(B) Between the go and return channels of a four-wire circuit, the near-end cross-talk ratio must be at least 43 decibels.

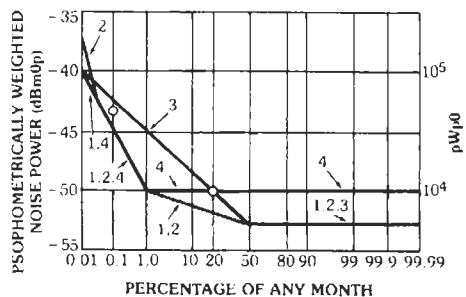
Note: The above refers to telephone circuits that are not equipped with (or used in conjunction with) modern echo suppressors designed for long propagation times. Circuits that can form part of switched connections with a long propagation time and lie between terminal half-echo suppressors of modern design should conform to high standards. (G.151, Section D)

Frequency Accuracy of Virtual Carrier Frequencies on an International Circuit

As the channel of any international telephone circuit should be suitable for voice-frequency telegraphy, the accuracy of the virtual carrier frequencies should be such that the difference between an audio frequency applied to one end of the circuit and the frequency received at the other end should not exceed 2 hertz, even when there are intermediate modulating and demodulating processes.

Considering either the hypothetical reference circuit proper to each system or the worldwide hypothetical reference connection, the accuracy of restitution of frequency that is recommended above should be assured if the channel and group carrier frequencies of the various stages have the following accuracies. (G.223)

- Virtual channel carrier frequencies $\pm 10^{-6}$
- in a group
- Group and supergroup carrier frequencies $\pm 10^{-7}$
- Mastergroup and supermastergroup carrier frequencies
- for 12-MHz system $\pm 5 \times 10^{-8}$
- for 60-MHz system $\pm 10^{-8}$
- (G.225, Section a)
- For 12-channel open wire carrier systems, channel carriers $\pm 5 \times 10^{-6}$
- (G.311)



The noise figures include 2500 pW for terminal equipment
 O - design objectives, including terminal noise.
 The numbers 1 to 4 are used to distinguish the curves

Fig. 5. Examples of distribution curves for the psophometrically weighted 1-minute mean noise power at the end of the hypothetical reference circuit (CCIR, 393-1).

CIRCUIT NOISE OBJECTIVES

Noise Characteristics for 2500-Kilometer Hypothetical Reference Circuits

The following design objectives should apply to any telephone channel provided by multichannel carrier systems on cable and radio relay links, having the same composition as the appropriate hypothetical reference circuit of 2500 kilometers. The noise is to be measured at or referred to a point of zero relative level. The objectives are intended to ensure adequate speech and signaling performance. They include an allowance of 2500 pWp0 due to frequency-division-multiplex equipment. (See Fig. 5.)

(A) The mean psophometric power during any hour shall not exceed 10 000 pWp0.

(B) The mean noise power over 1 minute shall not exceed 10 000 pWp0 for more than 20 percent of any month.

(C) The mean noise power over 1 minute shall not exceed 50 000 pWp0 for more than 0.1 percent of any month.

(D) The unweighted noise power, measured or calculated with an integrating time of 5 milliseconds, shall not exceed 1 000 000 pW0 (10^6 pW0) for more than 0.01 percent (10^{-4}) of any month. However, this is reduced to 0.001 percent (10^{-5}) of any month, or 0.1 percent of any hour, if the telephone channel is to be used for 50-baud amplitude-modulated voice-frequency telegraph, if recommended quality of telegraph operation is to be obtained. This stiffer requirement need not apply to frequency-modulated noise-frequency telegraph systems operating at 50 bauds.

In part of a hypothetical reference circuit consisting of one or more homogeneous sections, the values of noise power in (A) and (B) above, and the small percentage of time in (C) and (D) above, shall be considered to be proportional to the number of homogeneous sections involved.

These are design objectives, "and it is not intended that they be quoted in specifications for equipment or used for acceptance tests." (G.222, and CCIR, 393-3)

Noise Characteristics for Long Circuits Not More Than 2500 Kilometers in Length

Circuits on Land, in Submarine-Cable Systems, or on Radio Relay Systems—The mean psophometric power during any hour should be of the order of 4 pWp0 per kilometer, including noise due to frequency-division-multiplex equipment, except for very short circuits or those with a very complicated composition. (G.153, Section A)

Circuits on Transhorizon Relay Systems—The CCIR divides transhorizon systems into two classes, from the viewpoint of performance.

Transhorizon systems of the first class are those intended to operate between two points that might, without excessive difficulty, be connected by line-of-sight radio relay or cable systems. The hypothetical reference circuit described for transhorizon systems applies to this class. The noise power at the end of this circuit will be calculated by statistical combination of the noise power in each of its radio sections. The statistical noise distribution curve should meet requirements (A) through (D) above.

Transhorizon systems of the second class are those used between points for which other transmission systems cannot be used without excessive difficulty. If the design objectives given above cannot be met without excessive difficulty, the calculated noise-power distribution should meet the following objectives: The mean psophometric power during 1 minute must not exceed 25 000 pWp0 for more than 20 percent of any month, nor exceed 63 000 pWp0 for more than 0.5 percent of any month. Also, the unweighted noise power (integrated over 5 milliseconds) stated in (D) above should not be exceeded for more than 0.05 percent of the most unfavorable month. (CCIR, 397-3)

Noise Characteristics for International Circuits More Than 2500 Kilometers in Length

Circuits in Cable or Radio Relay Systems, With No Long Submarine-Cable Section—Such circuits, between 2500 and 25 000 kilometers in length, are generally carried in land cable or radio relay systems designed to 2500-kilometer objectives, and the number of channel demodulations seldom exceeds that in the corresponding part of the worldwide hypothetical reference chain. Automatic regulation should be used on each group link. The noise objectives are:

(A) The mean psophometric noise power during any hour, due to the line, should not exceed 3 pWp0 per kilometer; preferably, it should be 2 pWp0 per kilometer or even less. It is noted that in some countries, long overland systems (5000 kilometers or more) have the same objectives as the submarine-cable system (1 pWp0/km).

(B) For systems up to about 7500 kilometers long, the 1-minute mean power should not exceed 50 000 pWp0 for more than 0.3 percent of any month.

(C) For systems up to about 7500 kilometers long, the unweighted noise power, measured or calculated with an integrating time of 5 milliseconds, should not exceed 1 000 000 pW0 for more than 0.03 percent of any month.

Note: (B) and (C) are prorated from 2500-kilometer objectives, and proportional values should be used for

lengths between 2500 and 7500 kilometers. The CCITT does not yet recommend short-term objectives for systems longer than 7500 kilometers. (G.153, Section A)

Circuits With a Long Submarine-Cable Section—The circuit noise attributable to the submarine section, without companders, should not exceed a mean psophometric noise value of 3 pWp0 per kilometer in the worst hour on the worst channel. The mean noise power for each direction of transmission, over all channels used for the longest circuits, should not exceed 1 pWp0 per kilometer. No objectives are given for circuits equipped with companders. The other parts of the circuit should conform to normal 2500-kilometer objectives. (G.153, Section B)

Circuits on Communication-Satellite Systems—The psophometrically weighted noise power at a point of zero relative level in any telephone circuit in the basic hypothetical reference circuit (as defined in Recommendation 352) should not exceed the provisional values of 10 000 pWp0 mean in any hour, 10 000 pWp0 one-minute mean for more than 20 percent of any month, 50 000 pWp0 one-minute mean for more than 0.3 percent of any month, and 1 000 000 pW0 (5-millisecond integrating time) for more than 0.01 percent of any year. The multiplex noise is not included in these figures. (CCIR, 353-3)

Circuits on Open Wire Carrier Systems—In systems of about 10 000 kilometers in length, with lines of strict regularity of construction, accurately operating automatic line regulators, adjustment of line levels to take account of special climatic conditions, and carefully chosen repeater spacings, the following objectives may apply: The mean psophometric power, during any hour, at the end of a circuit of about 10 000 kilometers, taking into account all noise that exists, with the exception of noise due to radio transmitters, should not exceed 50 000 pWp0. This is for a reasonable distribution of wet weather in the territory crossed by the circuit. (G.153, Section D, Orange Book)

Noise in Actual Circuits

Radio Relay Circuits, 280 to 2500 Kilometers in Length—In a telephone channel of an actual radio relay system using frequency-division multiplex, whose composition does not differ appreciably from the hypothetical reference circuit and whose length L is between 280 and 2500 kilometers, the psophometrically weighted noise power, excluding multiplex noise, should not exceed (A) $3L$ pWp0 mean power in any hour, (B) $3L$ pWp0 one-minute mean power for more than 20 percent of any month, and (C) 47 500 pWp0 one-minute mean power for more than $(L/2500) \times 0.1$ percent of any month, as a planning objective. (CCIR, 359-4)

Radio Relay Circuits, Over Real Links, 50 to 2500 Kilometers in Length—If, for planning reasons, the composition of a real link differs substantially from the hypothetical reference circuit, the noise power of a circuit of length L , carried in one or more baseband sections of frequency-division-multiplex radio links, should not exceed the values in Table 3.

National Circuits on Carrier Systems Over Very Short Distances—Assuming that circuits in an international connection, making use of frequency-division-multiplex carrier systems over very short distances, can be limited in number to four, the mean psophometric power should not exceed 2000 pWp0 per circuit during any hour, including cross talk.

For CCITT Recommendations for pulse-code-modulation systems see G.123.

Design Objectives for Noise Produced by Modulating Equipments

The mean psophometric power, which corresponds to the noise produced by all modulating equipment mentioned in the definition of the hypothetical reference circuit in question, should not exceed 2500 pWp0. This value includes noise due to various causes, such as thermal noise, intermodulation, cross talk, power supplies, etc. Its allocation between the various equipments can be left somewhat to the discretion of designers, but the following maximum values are recommended for the modulating equipments

TABLE 3. MAXIMUM NOISE POWER

Length of System, (km)	Mean Power in Any Hour	1-Minute Mean Power for More Than 20% of Any Month	1-Minute Mean Power of 47 500 pWp0, for More Than Stated % of Any Month
50–280	$3L$ pWp0 + 200 pWp0	$3L$ pWp0 + 200 pWp0	$(280/2500) \times 0.1\%$
280–840	$3L$ pWp0 + 200 pWp0	$3L$ pWp0 + 200 pWp0	$(L/2500) \times 0.1\%$
840–1670	$3L$ pWp0 + 400 pWp0	$3L$ pWp0 + 400 pWp0	$(L/2500) \times 0.1\%$
1670–2500	$3L$ pWp0 + 600 pWp0	$3L$ pWp0 + 600 pWp0	$(L/2500) \times 0.1\%$

Note 1: Frequency-division-multiplex noise is not included.

Note 2: The hourly mean-noise-power objective and its subdivision are under study.

	pWp0
One pair of channel modulators	200-400
One pair of group modulators	60-100
One pair of supergroup modulators	60-100
One pair of mastergroup modulators	40-60

(G.222, Section d)

Pacific: 702 Guam, 704 Hawaii (RCA), 705 Hawaii (ITTWC), 71 Australia, 72 Japan, 75 Philippines.

Asia: 801 Korea, 802 Hong Kong, 81 India, 85 China, 88 Iran.

Africa: 901 Libya, 94 Ghana, 95 South Africa, 975 Niger, 991 Angola, 992 Mozambique.

CCITT AND TELEGRAPHY

The CCITT Yellow Book contains the CCITT Recommendations, including those adopted by the Seventh Plenary Assembly held in Geneva in 1980. The recommendations on Telegraph Technique are included in Volume VII. Those on Data Transmission are included in Volume VIII, and the Recommendations on Telegraph Operations and Tariffs are contained in Volume II.

Numbering

There is a worldwide system of Destination Indicators for the telegraph-message retransmission network. These indicators consist of two letters signifying the country and its telegraph network (if more than one) followed by two letters signifying the town on that network. Examples: Vienna AUWI, Panama City (Tropical Radio) PAPA, Balboa (ITTCACR) PZBA, Stockholm SWSM, San Francisco (ITT Worldcom) UISF.

The CCITT has approved a worldwide numbering system for telex services. The telex destination code consists of two or three numerical digits signifying the country or network within the country. The destination code is followed by the telex subscriber's national number, also consisting of numerical digits.

The telex system provides also for *designation* codes, for identifying the country and network of the originator of a communication. The designation code consists of two letters, the same two letters that compose the first half of the message-retransmission-system *destination* indicator.

Examples of destination codes are:

North and Central America: 200 Cuba, 205 Puerto Rico (RCA), 206 Puerto Rico (ITTWC), 207 Puerto Rico (C & W), 21 Canada (except TWX), 22 Mexico, 25 USA (TWX), 271 Guatemala, 290 Bermuda, 292 Virgin Islands.

South America: 304 Surinam, 305 Paraguay, 31 Venezuela, 36 Peru, 381 Brazil (Radio Brazil), 383 Brazil (PTT), 387 Argentina (ITTCM), 390 Netherlands Antilles, 391 Trinidad.

Europe: 400 Canary Islands, 403 Spain, 409 Algeria, 41 Germany, 46 Belgium, 492 Syria, 496 Kuwait, 501 Iceland, 51 United Kingdom, 57 Finland.

Eastern Europe: 601 Greece, 606 Israel, 61 Hungary, 65 Romania.

CCITT AND TELEPHONY

International Country Codes

The addressing signals of worldwide automatic telephony consist of the national telephone number, as used for long-distance dialing within a country, prefixed by a country code. Country codes are grouped by continental regions; for example, the country codes of all South American countries begin with "5." Where the national numbering system includes more than one country, the country code may also include the countries included in the national system. Thus the country code for the United States—"1"—includes Canada and some other countries. The following are examples of some country codes, grouped by world numbering regions or zones, as assigned by the Third Plenary Assembly of the CCITT in Geneva in 1964.

Zone 1—Code 1: USA, Canada, Mexico and Central America, Bahamas, Bermuda, Jamaica, French Antilles, Netherlands Antilles.

Zone 2—Africa: 51 countries, 48 country codes (Algeria, Morocco, Tunisia, Libya in one group—the Maghreb—code 21). United Arab Republic 20, South Africa 27, 45 three-digit codes.

Zones 3 and 4—Europe, Iceland, Malta, Cyprus: 17 two-digit and 13 three-digit country codes. Examples: France 33, Spain 34, Italy 39, United Kingdom 44, Germany 49, Iceland 354, Finland 401, Hungary 402.

Zone 5—South America and Cuba: 6 two-digit and 8 three-digit country codes. Examples: Cuba 53, Argentina 54, Brazil 55, Chile 56, Columbia 57, Venezuela 58, Peru 596.

Zone 6—Southwestern Pacific: 6 two-digit and 14 three-digit country codes. Examples: Malaysia 60, Australia 61, Indonesia 62, Philippines 63, New Zealand 64, Thailand 66, Guam 682.

Zone 7—Country code 7: Soviet Union

Zone 8—Northwestern Pacific: 4 two-digit and 6 three-digit country codes. Examples: Japan 81, Korea 82, Vietnam 84, China (Formosa) 85, Hong Kong 852, Mongolia 854, Laos 856.

Zone 9—East: 5 two-digit and 15 three-digit country codes. Examples: India 91, Burma 95, Iran 98, Lebanon 961, Saudi Arabia 966, Israel 972.

TELEPHONE SIGNALING

CCITT signaling systems have been standardized for international use. General descriptions are given in

TABLE 4. CCITT SIGNALING SYSTEMS

No.	Systems
1	500/20-hertz system used in the international manual service (ringdown). Used until 1980.
2	600/750-hertz two-frequency system. Never used in international service. <i>International Automatic and Semiautomatic Systems</i>
3	For unidirectional operation of circuits. Uses one in-band frequency (2280 hertz) for the transmission of both line and interregister signals; used for terminal traffic; in general not to be used for new installations. Used until late 1970s.
4	For unidirectional operation of circuits (circuits seized from one end only). Uses two in-band frequencies (2040 and 2400 hertz) for the end-to-end transmission of both line and register signals; used for international intracountry traffic; suitable for terminal and transit traffic; in the latter case two or three circuits equipped with System No. 4 may be switched in tandem. Suitable for submarine- or land-cable circuits and microwave radio circuits; not applicable to TASI-equipped systems. Capable of interworking with System No. 5
5	For both-way operation of circuits. Uses two in-band signaling frequencies (2400 and 2600 hertz) for the link-by-link transmission of line signals, and six in-band frequencies (700, 900, 1100, 1300, 1500, and 1700 hertz) in a two-out-of-six code (numerical information transmitted <i>en bloc</i>) for the link-by-link transmission of register signals; used for intercontinental traffic. Suitable for submarine- or land-cable circuits and microwave links, whether or not TASI is used; suitable for terminal or transit traffic—in the latter case, two or more circuits equipped with System No. 5 may be switched in tandem but are subject to possible undesirable delays if all are TASI-equipped. Capable of interworking with System No. 4
6	Based on principles of common channel signaling techniques in which the signaling is removed from the voice path; for both analog and digital transmission facilities; signaling link carries information by serial mode of synchronous data transmission; link-by-link operation; can be associated, quasi-associated, or non-associated mode; capable of interworking with No. 4, No. 5, and No. 5 bis. (Yellow Book, VI.3)
7	A common channel system specified in 1980; conceived for the digital environment; preferred system for integrated digital networks (IDN) for both telephony and circuit switched data; designed for interexchange signaling in integrated services digital network (ISDN); optimized for operation over 64 kb/s digital channels; employs packet-type information transfer. Flexible enough to evolve with future enhanced service requirements. (Yellow Book, VI.6)

Table 4, and some of the signaling characteristics are given in Table 5.

Signals in communications are used for passing information, for identifying the called subscriber or addressee (with resulting internal system signals concerned with the establishment of a connection), and for supervising and controlling the connection once it has been established.

Information signals may be analog (voice, telemetry, or facsimile) or digital (teleprinter or data).

Addressing signals may be dial pulse, multifrequency, or binary. They are not needed once a communication has been established.

(A) *Dial pulse* signals consist of a series of from one to ten pulses representing the corresponding numerical digits 1 to 9 and 0. The pulses are breaks in a continuous direct current on the line, usually lasting from 58 to 67 percent of the time interval between the starts of successive pulses. These breaks in direct current may have to be converted into pulses in a tone, or to frequency shifts between tones, in order to pass through some media or multiplexing systems. Dial pulse speeds are usually 10 pulses per second, although high-speed dials of nearly twice that speed are sometimes used.

(B) *Multifrequency* signals represent numerical digits by one pulse (or sometimes four pulses) of a specific frequency combination. The two-out-of-six and four-by-four multifrequency codes are shown in Tables 6 and 7. Inherent in both is the constant-ratio error-control principle (the simultaneous receipt of three or more, or one only, frequencies indicates an error). Table 8 shows the US Army two-out-of-five numerical code. The CCITT two-voice-frequency code, consisting of two frequencies sent one at a time in four pulses, is given in Table 9.

(C) *Binary* signals for addressing are usually in a numeric or alphanumeric code used also for information signals, such as in telegraph addresses or headings. Multifrequency signals are sometimes directly converted to binary signals by changing the two-frequencies-out-of-six code to a corresponding two-time-slots-out-of-six synchronous six-element binary code. Another example of a binary code used only for addressing is the CCITT one-voice-frequency code, and four-element start-stop code given in Table 10. (The CCITT two-voice-frequency code is not truly binary since it uses a third condition—no tone—in addition to the two tones.)

Supervisory or Line Signals

Supervisory or line signals cannot be generated exclusively by registers because they are required during the entire use of the connection, after the registers that established the connection have been disconnected. Supervisory signals are an extension of the original basic signals of ringing and of closing a line loop to allow direct current to flow. Supervisory signals may be classified as spurt (discontinuous) and continuous.

TABLE 5. LINE SIGNALS IN THREE CCITT SYSTEMS

Signal	Direction	CCITT No. 3 (1 VF)	CCITT No. 4 (2 VF)	CCITT No. 5 (2 VF)
Seize, terminal	→	X	PX	X
Seize, transit	→		PY	
Start pulsing, terminal	←	X	X	Y
Start pulsing, transit	←		Y	
End of pulsing (ST)	→	250 ms*	xSxSxSx*	1500+1700 hertz*
Busy	←	XX	PX	Y
Acknowledge	→			X
Answer	←	XSX	PY	X
Acknowledge	→			X
Clear back (on-hook)	←	XX	PX	Y
Acknowledge	→			X
Ring forward	→	XSX	PYY	Y(850±200 ms)
Clear forward (disconnect)	→	XXSXX	PXX	X+Y
Release guard (disconnect acknowledge)	←	XXSXX	PYY	X+Y
		X: 2280±6 hertz, 150±30 ms XX: 2280±6 hertz, 600±120 ms S: 100±20 ms silence	X: 2040±6 hertz, 100±20 ms Y: 2400±6 hertz, 100±20 ms XX: 350±70 ms YY: 350±70 ms S: 35±7 ms; x: 2040 hertz, 35 ms P: (2040 hertz, 2400 hertz), 150+30ms	X: 2400±6 hertz Y: 2600±6 hertz

* Combination No. 15 of address code.

TABLE 6. MULTIFREQUENCY NUMERICAL CODE USED BY CCITT (TWO-OUT-OF-SIX)

Digit	Frequencies	Weighting
1	700 + 900	0 + 1
2	700 + 1100	0 + 2
3	900 + 1100	1 + 2
4	700 + 1300	0 + 4
5	900 + 1300	1 + 4
6	1100 + 1300	2 + 4
7	700 + 1500	0 + 7
8	900 + 1500	1 + 7
9	1100 + 1500	2 + 7
0	1300 + 1500	4 + 7
Code 11	700 + 1700	0 + 11
Code 12	900 + 1700	1 + 11
KP	1100 + 1700	2 + 11
KP2	1300 + 1700	4 + 11
ST	1500 + 1700	7 + 11

Continuous signals are based on conditions of on-hook and off-hook, representing the condition of locked or flowing direct current on the subscriber's line, and their extension to trunk signaling is given in Table 11. Either condition is continuous and may be detected at any time. On the other hand, discontinuous signals must be recorded, and the condition represented is presumed to continue until a new signal is sent. Supervisory signaling in a backward direction (toward the calling end) is also needed in automatic working and is also described as on-hook or off-hook, although on two-wire metallic circuits the signaling condition is usually a reversal of flow of the direct current rather than an interruption.

Continuous supervisory signaling over longer distances is effected by use of signaling paths distinct from the voice path. These signaling paths may be telegraph legs of a composite telegraph system, simplexing of the voice pair, or special tones inside or outside of the voice

TABLE 7. NUMERICAL FOUR-BY-FOUR MULTIFREQUENCY CODE

Touch-Tone or Touch Calling						US Air Force 412L					
Low group (hertz)	697	1	2	3	High group (hertz)	1620	1740	1860	(1980)	High group (hertz)	
	770	4	5	6							
	852	7	8	9							
	941	spare	0	spare							
	1209	1336	1477	(1633)							

Note: Each digit is composed of one frequency from the low group and one frequency from the high group. The frequencies have been chosen to minimize voice simulation.

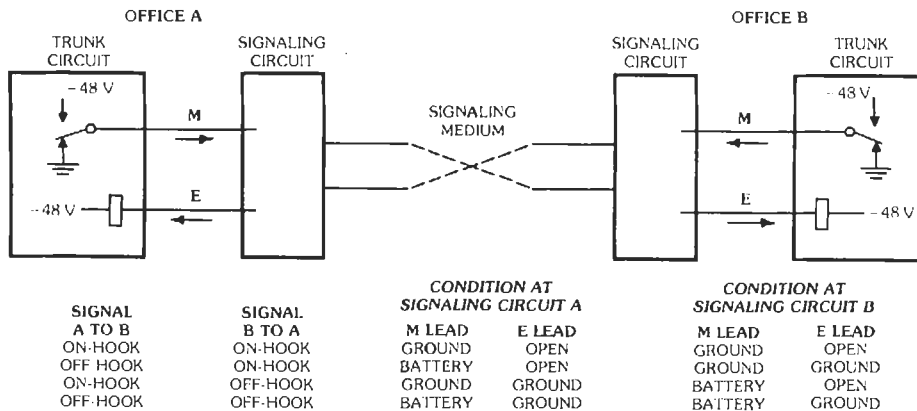


Fig. 6. E and M signaling. The E lead receives open or ground signals from the signaling circuit. The M lead sends ground or battery signals to the signaling circuit.

TABLE 8. US ARMY TA-341/PT NUMERICAL CODE

Digit	Frequencies
1	2100 + 2300
2	2300 + 2500
3	1900 + 2700
4	1900 + 2100
5	2500 + 2700
6	2300 + 2700
7	2100 + 2500
8	1900 + 2300
9	2100 + 2700
0	1900 + 2500

channel. Whatever the paths, they are extensions of separate direct-current leads from the trunk circuits, known as *E* and *M* leads. The relation between the conditions of these leads and the on-hook or off-hook signaling conditions they represent is shown in Fig. 6. The usual method of extending the *E* and *M* leads over tone channels is shown in Fig. 7. The frequency of the tone used is preferably higher than 2000 hertz and is usually 2600 hertz on four-wire circuits. A 3825-hertz signaling system is known as "out-of-band" and must be built into the carrier system using it.

Spurt signaling avoids the necessity of using distinct signaling paths. For manual operation, ringdown signaling is usually satisfactory. For semiautomatic or automatic operation, however, more elaborate systems are required. Voice-frequency signals are used, and they are distinguished from voice transmissions by filters and timing. Single-frequency systems use one or two long or short pulses of the specific signaling frequency, with the pulses carefully timed to within minimum and

maximum limits. Two-frequency systems use a gate-opening prefix pulse of the two frequencies together, followed, without a silent period, by a long or short pulse of either of the two frequencies. The submarine-cable signaling system (CCITT No. 5) uses two frequencies in the "compel" mode, wherein the signal frequencies are sent until acknowledged, and the acknowledging signal is sent until the original signal is stopped. No gate-opening prefix is used.

Ringdown Signals

Ringdown signals are spurts of ringing current (16 to 25 hertz) applied usually through the ringing key of an operator and intended to operate a bell, ringer, or drop at the called end. The current may be generated by a manually operated magneto or by a ringing machine with or without automatically inserted silent periods. Ringing to telephone subscribers in automatic central offices is stopped or "tripped" automatically by relay action resulting from the subscriber's off-hook condition. Ringing signals may be converted to 500 or 1000 hertz, usually interrupted at a 20-hertz rate, to pass through voice channels of carrier equipment. A ringing signal to a manual switchboard usually lights a switchboard lamp, which can be darkened again only by local action and not by stopping or repeating the ringing signal. This characteristic makes ringdown operation unsuitable for fully automatic operation. Ringdown signaling over carrier circuits has the advantages of simplicity and of not requiring the distinct signaling channels of *E* and *M* systems.

Tones

A special case of signaling is that of information in the form of tones to the subscriber of a telephone

TABLE 9. NUMERICAL CODE, TWO-VOICE-FREQUENCY SIGNALING SYSTEM, CCITT NO. 4

Digit	Successive Elements			
	1	2	3	4
1	y	y	y	x
2	y	y	x	y
3	y	y	x	x
4	y	x	y	y
5	y	x	y	x
6	y	x	x	y
7	y	x	x	x
8	x	y	y	y
9	x	y	y	x
0	x	y	x	y

Note: The two frequencies are sent one at a time, with a silent space between pulses. The duration of both frequency and silent periods is 35±7 milliseconds. Frequencies: x = 2040±6 hertz; y = 2400±6 hertz. Power level: -9 decibels.

TABLE 10. NUMERICAL CODE, ONE-VOICE-FREQUENCY SIGNALING SYSTEM, CCITT NO. 3.

Digit	Time Elements					Stop
	Start	1	2	3	4	
1	1				1	
2	1			1		
3	1			1	1	
4	1		1			
5	1		1		1	
6	1		1	1		
7	1		1	1	1	
8	1	1				
9	1	1			1	
0	1	1		1		

Note: "1" signifies frequency present. Length of each time element is 50 milliseconds±1 percent. Frequency: 2280±6 hertz. Power level: -6 decibels.

TABLE 11. ON-HOOK AND OFF-HOOK SIGNALS

Direct Current Telephone Line	Trunk
On-hook signifies loop is open to direct current supplied from other end.	If idle, signals on-hook to other end. Seizure at calling end signals off-hook to called end. While calling end awaits answer, called end signals on-hook to calling end.
Off-hook signifies loop is closed, allowing relay at other end to operate. Signaling in reverse direction is ring-down.	Answer results in signaling off-hook from called end. If called end is not ready to receive address signals when seized, it signals off-hook to calling end until ready.

system. The basic tones are dial tone, busy tone, and ring-back tone (representing the ringing of the called subscriber's line). The dial tone is generated at the subscriber's local switching center, but the busy tone and ring-back tone, plus special tones such as no-such-number and line-out-of-order, are generated at the called subscriber's switching center and should be standardized for universal intelligibility. A proposal by the American Telephone and Telegraph Company for standard tones is described in Table 12. On some international calls, the busy and ring-back tones are

generated locally in accordance with spurt supervisory signals from the called end.

Alternative Routing

Switching systems in which the complete called number is recorded in the first center to which the subscriber is connected permit a translation of digits from those identifying the called subscriber to those most conveniently used by the switching mechanism to

TABLE 12. TABLE OF PROPOSED STANDARD OF AUDIBLE TONES IN NORTH AMERICA (FROM CCITT DOCUMENT AP III-84)

Use	Frequencies§ (hertz)				Power per Frequency at Exchange Where Tone is Applied (dBm0)	Cadence
	350	440	480	620		
Dial tone	x	x			-13	Continuous
Busy tone			x	x	-24	0.5 second on 0.5 second off
Reorder tone*			x	x	-24	0.2 second on 0.3 second off or 0.3 second on 0.2 second off
Audible ringing tone		x	x		-16	2 seconds on 4 seconds off
High tone†			x		-16	Varies according to use
Preemption tone‡		x		x	-18	Single 200/500-ms pulse
Call-waiting tone		x			-13	Single 500-ms pulse

Notes:

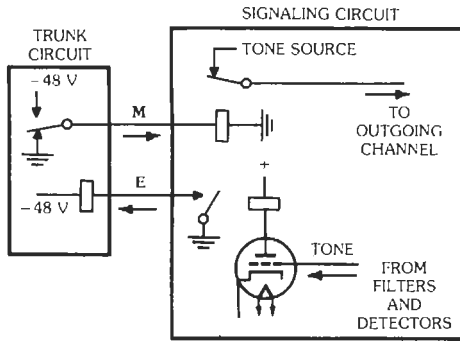
* A possible alternative is the use of a call-failure tone, which would identify the office and type of condition that prevented the successful completion of the call.

† High tone is used in many ways. For example:

- (A) Spurts of tone to indicate specific orders to operators in the manual service (order tones).
- (B) To inform operators of lines that are temporarily out of service (permanent signal tone).
- (C) To alert customers that their services are in a permanent off-hook condition.

‡ Preemption tones are used in certain private switched networks which may interconnect with the national networks.

§ Frequency limits are ±0.5 percent of nominal.



establish the desired connection. This translation of digits permits the controlling mechanism to use a variety of routes in such a way that if the first-choice route is occupied or disabled a second-choice route may be tried, etc. Safeguards are required to prevent doubling back of routes and dead-end choices. In the near future, many sophisticated plans may be expected that depend on rapid analysis of network possibilities by data-transmission means. The signals used for this purpose are interregister signals and are classed with addressing signals in that they are used in establishing a connection to the desired addressee, not being needed after a connection has been established.

SIGNAL	TONE	OPERATION	LEAD	CONDITION
ON-HOOK	ON	SENDING	M	GROUND
		RECEIVING	E	OPEN
OFF-HOOK	OFF	SENDING	M	BATTERY
		RECEIVING	E	GROUND

Fig. 7. Tone signaling.

3 Units, Constants, and Conversion Factors

*Revised by
Eugene A. Mechtly*

SI Units 3-2

SI Base Units
SI Derived Units
SI Prefixes

Additional Information on the SI 3-3

Units outside the SI
Rules and Style Conventions
Dimensionless Quantities (Quantities of Dimension 1)
Decibels and Nepers
Numerical Values
Greek Alphabet

Fundamental Physical Constants 3-7

Conversion Factors 3-13

SI UNITS*

There are two classes of units in the *International System of Units (SI)*: *base units* and *derived units*. The units of these two classes form a *coherent* set of units. This coherent set of units is designated by the name "SI Units."

SI Base Units

The SI is founded on seven *base units*, respectively, for seven *base quantities*. The base quantities are assumed to be mutually independent. These base units and base quantities are given in Table 1.

SI Derived Units

Other quantities, called *derived quantities*, are defined in terms of the seven base quantities through a system of quantity equations. *Derived units* for these derived quantities

* The information on the SI in this chapter is based on *Metric System of Measurement: The International System of Units for the United States*, published in the U.S. Federal Register, Vol. 63, No. 144, 1998, pp. 40334-40340. This Federal Register Notice can be viewed and downloaded over the World Wide Web from the website of the US Government Printing Office. Go to www.gpo.gov. Click on "Federal Register," select "1998," select "Notices," and enter "SI" as the search term.

TABLE 1. SI BASE UNITS

Base Quantity	SI Base Unit	
	Name	Symbol
Length	meter	m
Mass	kilogram	kg
Time	second	s
Electric current	ampere	A
Thermodynamic temperature	kelvin	K
Amount of substance	mole	mol
Luminous intensity	candela	cd

are obtained from the system of equations and the seven base units in a *coherent* manner. Here, the term *coherent* is used to mean a system of units where all derived units are obtained from the base units by the rules of multiplication and division, with no numerical factor other than one (1) ever occurring in the expressions for the derived units in terms of the base units.

Table 2 lists, in alphabetical order, SI derived units and reiterates the seven SI base units.

SI Prefixes

Table 3 gives the twenty SI prefixes used to form decimal multiples and submultiples of SI units.

TABLE 2. SI BASE AND DERIVED UNITS

Name	Symbol	Quantity
ampere (an SI base unit)	A	electric current
ampere per meter	A/m	magnetic field strength
ampere per square meter	A/m ²	current density
becquerel	Bq s ⁻¹	activity (of a radionuclide)
candela (an SI base unit)	cd	luminous intensity
candela per square meter	cd/m ²	luminance
coulomb	C s ⁻¹ ·A	electric charge, quantity of electricity
coulomb per cubic meter	C/m ³	electric charge density
coulomb per kilogram	C/kg	exposure (x and gamma rays)
coulomb per square meter	C/m ²	electric flux density
cubic meter	m ³	volume
cubic meter per kilogram	m ³ /kg	specific volume
degree Celsius*	°C	Celsius temperature
farad	F C/V	capacitance
farad per meter	F/m	permittivity
gray	Gy	absorbed dose, specific energy imparted, kerma, absorbed dose index
gray per second	Gy/s	absorbed dose rate
henry	H Wb/A	inductance
henry per meter	H/m	permeability
hertz	Hz s ⁻¹	frequency
joule	J N·m	energy, work, quantity of heat
joule per cubic meter	J/m ³	energy density
joule per kelvin	J/K	heat capacity, entropy
joule per kilogram	J/kg	specific energy
joule per kilogram kelvin	J/(kg·K)	specific heat capacity, specific entropy
joule per mole	J/mol	molar energy
joule per mole kelvin	J/(mol·K)	molar entropy, molar heat capacity

* The kelvin (symbol K) is the SI unit of thermodynamic temperature, T . The degree Celsius (symbol °C) is the commonly used unit of Celsius temperature, t . Celsius temperature is defined by the equation $t = T - T_0$, where $T_0 = 273.15$ kelvins.

TABLE 2 (CONT). BASE AND DERIVED UNITS

Name	Symbol	Quantity
kelvin (an SI base unit)	K	thermodynamic temperature
kilogram (an SI base unit)	kg	mass
kilogram per cubic meter	kg/m ³	density, mass density
kilogram per kilogram	kg/kg = 1	mass fraction
lumen	lm	luminous flux
lux	lx lm/m ²	illuminance
meter (an SI base unit)	m	length
meter per second	m/s	speed, velocity
meter per second squared	m/s ²	acceleration
mole (an SI base unit)	mol	amount of substance
mole per cubic meter	mol/m ³	concentration (of amount of substance)
newton	N	force
newton meter	N·m	moment of force
newton per meter	N/m	surface tension
number per meter	m ⁻¹	wave number
ohm	Ω V/A	electric resistance
pascal	Pa N/m ²	pressure, stress
pascal second	Pa·s	dynamic viscosity
radian	rad	plane angle
radian per second	rad/s	angular velocity
radian per second squared	rad/s ²	angular acceleration
reciprocal meter	m ⁻¹	wave number
second (an SI base unit)	s	time
siemens	S A/V	electric conductance
sievert	Sv	dose equivalent, equivalent dose
square meter	m ²	area
steradian	sr	solid angle
tesla	T Wb/m ²	magnetic flux density
volt	V W/A	electric potential, potential difference, electromotive force
volt per meter	V/m	electric field strength
watt	W J/s	power, radiant flux
watt per meter kelvin	W/(m·K)	thermal conductivity
watt per square meter	W/m ²	heat flux density, irradiance, power density
watt per square meter steradian	W·m ⁻² ·sr ⁻¹	radiance
watt per steradian	W/sr	radiant intensity
weber	Wb V·s	magnetic flux

TABLE 3. SI PREFIXES

Factor	Name	Symbol
10 ²⁴	yotta	Y
10 ²¹	zetta	Z
10 ¹⁸	exa	E
10 ¹⁵	peta	P
10 ¹²	tera	T
10 ⁹	giga	G
10 ⁶	mega	M
10 ³	kilo	k
10 ²	hecto	h
10	deka	da
10 ⁻¹	deci	d
10 ⁻²	centi	c
10 ⁻³	milli	m
10 ⁻⁶	micro	μ
10 ⁻⁹	nano	n
10 ⁻¹²	pico	p
10 ⁻¹⁵	femto	f
10 ⁻¹⁸	atto	a
10 ⁻²¹	zepto	z
10 ⁻²⁴	yocto	y

Example: In the expression $I = 150 \text{ mA}$, I is the symbol for a physical quantity (current), A is the symbol for the unit of current (ampere), and m is the symbol for the prefix milli. Together, m and A form the symbol for a submultiple unit of current, the milli-ampere.

Note that the kilogram, for historical reasons, is the only SI unit with a prefix (kilo) as part of its name (kilogram) and as part of its symbol (kg). In the special case of the kilogram, the prefix names are used with the name "gram" and the prefix symbols are used with the symbol "g" for gram. For example, "mg" is the symbol for milligram. With this exception, any SI prefix may be used with any SI unit.

ADDITIONAL INFORMATION ON THE SI

Additional information on the SI may be found in the NIST Special Publication 811, *Guide for the Use of the International System of Units (SI)*, by Barry N. Taylor. SP 811 is for sale by the Government Printing

Office and is also available at the NIST website: www.physics.nist.gov/cuu.

Units outside the SI

Certain units are not part of the International System of Units, that is, they are outside the SI but are important and widely used. The units in this category that are approved for use in the United States, with the SI, are given in Tables 4 and 5.

Rules and Style Conventions*

Rules and style conventions ensure that communication is not hindered by ambiguity. Some well-established rules are:

1. Symbols for units and prefixes are printed in Roman (upright) type regardless of the type used in the surrounding text.

2. Uppercase and lowercase symbols must be used for units and prefixes as defined in Tables 1 through 5 of this chapter. Failure to distinguish between upper and lower cases of the same letter is not acceptable. For example, mW (milliwatt) and MW (megawatt) must not be confused.
3. When the name of a unit is spelled out, it is always written with a lowercase initial letter unless it begins a sentence.
4. Unit symbols are unaltered in the plural.
5. Unit symbols are not followed by a period unless at the end of a sentence.
6. Symbols for units formed from other units by multiplication are indicated by means of a half-high (that is, centered) dot or space. Example: N·m or N m.
7. Symbols for units formed from other units by division are indicated by the means of a solidus (oblique stroke, /), a horizontal line, or negative exponents. For example,

$$m/s, \frac{m}{s}, \text{ or } m \cdot s^{-1}$$

* Most of these rules and style conventions are copied from the Federal Register Notice of 1998.

TABLE 4. UNITS OUTSIDE THE SI THAT ARE ACCEPTED FOR USE WITH THE SI

Name	Symbol	Value in SI Units
minute	min	1 min = 60 s
hour	h	1 h = 60 min = 3600 s
day	d	1 d = 24 h = 86 400 s
degree	°	1° = (π/180) rad
minute	'	1' = (1/60)° = (π/10 800) rad
second	"	1" = (1/60)' = (π/648 000) rad
liter*	L	1 L = 1 dm ³ = 10 ⁻³ m ³
metric ton†	t	1 t = 10 ³ kg
neper‡	Np	1 Np = 1
bel	B	1 B = (1/2) ln 10 Np

* The liter and its symbol, l, are approved for use internationally. The alternative symbol for the liter, L, was approved internationally in 1979 to avoid the risk of confusion between the letter l and the number 1. The preferred symbol for use in the United States is L. Neither a lowercase script letter l nor an uppercase script letter L are approved symbols for the liter.

† The name “metric ton” is preferred for use in the United States for the unit with symbol t and is defined as 1 t = 10³ kg. The “metric ton” is also used in some other English-speaking countries. The name “tonne” is used in many other countries.

‡ Although the neper is coherent with SI units and is accepted by the CIPM, it has not been adopted by the CGPM and is thus not an SI unit.

TABLE 5. UNITS OUTSIDE THE SI THAT ARE ACCEPTED FOR USE WITH THE SI, BUT WHOSE VALUES IN SI UNITS ARE OBTAINED EXPERIMENTALLY

Name	Symbol	Value in SI Units*
electron volt†	eV	1 eV = 1.602 177 33(49) × 10 ⁻¹⁹ J
unified atomic mass unit‡	u	1 u = 1.660 540 2(10) × 10 ⁻²⁷ kg
astronomical unit§	ua	1 ua = 1.495 978 70(30) × 10 ⁸ m

* The combined standard uncertainty (that is, estimated standard derivation) of the last two figures is shown in parentheses.

† The electronvolt is the kinetic energy acquired by an electron in passing through a potential difference of 1 V in vacuum.

‡ The unified atomic mass unit is equal to 1/12 of the mass of an unbound atom of the nuclide ¹²C at rest and in its ground state.

§ The astronomical unit is a unit of length approximately equal to the mean Earth–Sun distance. Its value is such that, when used to describe the motion of bodies in the Solar System, the heliocentric gravitation constant is (0.017 202 098 95)² ua³ · d⁻².

However, to avoid ambiguity, the solidus must not be repeated on the same line unless parentheses are used. Examples are

m/s^2 or $\text{m}\cdot\text{s}^{-2}$ *but not:* m/s/s
 $\text{m}\cdot\text{kg}/(\text{s}^3\cdot\text{A})$ or $\text{m}\cdot\text{kg}\cdot\text{s}^{-3}\cdot\text{A}^{-1}$ *but not:* $\text{m}\cdot\text{kg/s}^3/\text{A}$

Negative exponents should be used in complicated cases.

8. Prefixes symbols are printed in Roman (upright) type regardless of the type used in the surrounding text and are attached to unit symbols without a space between the prefix symbol and the unit symbol. This last rule also applies to prefix names attached to unit names. Examples include

1 ml (one milliliter)
 1 pm (one picometer)
 1 G Ω (one gigaohm)
 1 THz (one terahertz)

9. The grouping formed by a prefix symbol attached to a unit symbol constitutes a new inseparable symbol (forming a multiple or submultiple of the unit concerned), which can be raised to a positive or negative power and which can be combined with other unit symbols to form compound unit symbols. Some examples are

$$2.3 \text{ cm}^3 = 2.3 (\text{cm})^3 = 2.3 (10^{-2}\text{m})^3 = 2.3 \times 10^{-6}\text{m}^3$$

$$1 \text{ cm}^{-1} = 1 (\text{cm})^{-1} = 1 (10^{-2}\text{m})^{-1} = 10^2 \text{ m}^{-1}$$

$$5000 \mu\text{s}^{-1} = 5000 (\mu\text{s})^{-1} = 5000 (10^{-6}\text{s})^{-1} = 5000 \times 10^6\text{s}^{-1} = 5 \times 10^9 \text{ s}^{-1}$$

Prefix names are also inseparable from the unit names to which they are attached. Thus, for example, millimeter, micropascal, and meganewton are single words.

10. Compound prefix symbols, that is, prefix symbols formed by the juxtaposition of two or more prefix symbols, are not permitted. This rule also applies to compound prefix names. For example,

1 nm (one nanometer) *but not:* 1 m μm (one millimicrometer)

11. An SI prefix symbol (and name) cannot stand alone but must be attached to a unit symbol (or name). For example,

$$5 \times 10^6/\text{m}^3 \text{ *but not:* } 5 \text{ M}/\text{m}^3$$

12. In the expression for the value of a quantity, the unit symbol is placed after the numerical value and a space is left between the numerical value and the unit symbol. The only exceptions to the rule are for the unit symbols for degree, minute, and second for plane angle: $^\circ$, $'$, and $''$, respectively (see Table 4), in which case no space is left between the numerical value and the unit symbol. For example,

$$\alpha = 30^\circ 22' 8''$$

This rule means that:

- (a) The symbol $^\circ\text{C}$ for the degree Celsius is preceded by a space when one expresses the values of Celsius temperatures. Thus

$$t = 30.2 \text{ }^\circ\text{C} \text{ *but not:* } t = 30.2^\circ\text{C} \text{ or } t = 30.2^\circ \text{ C}$$

- (b) Even when the value of a quantity is used in an adjectival sense, a space is left between the numerical value and the unit symbol. (This rule recognizes that unit symbols are not like ordinary words or abbreviations but are mathematical entities and that the value of a quantity should be expressed in a way that is as independent of language as possible.) Some examples are

a 1 m end gauge *but not:* a 1-m end gauge

a 10 k Ω resistance *but not:* a 10-k Ω resistance

However, if there is any ambiguity, the words should be rearranged accordingly. For example, the statement "the samples were placed in 22 mL vials" should be replaced with the statement "the samples were placed in vials of volume 22 mL" or "the samples were placed in 22 vials of volume 1 mL," whichever was meant.

NOTE: When unit names are spelled out as is often the case in nontechnical writing, the normal rules of English apply. Thus, for example, "a roll of 35-millimeter film" is acceptable.

When there is a risk of confusion in using the standard symbol (e.g., "s" or "T") the name of the unit should be spelled out (e.g., second or tesla).

When a rarely used multiple or submultiple of an SI unit is first used in text, the full name should be stated. Subsequently, the corresponding symbols may be used alone. For example: The solar energy incident each year at the top of the Earth's atmosphere is about 5.6 yottajoules (5.6 YJ). The total energy of a solar flare is about 10 YJ.

Dimensionless Quantities (Quantities of Dimension 1)

Ratios of two quantities of the same dimension are called *dimensionless quantities* or *quantities of dimension 1*. They are expressed as pure numbers. The SI unit of all dimensionless quantities is the number 1. The two quantities may be of the same kind or they may be of a different kind.

Examples of dimensionless ratios of quantities of the same kind are: relative permittivity, relative permeability, refractive index, mass fraction, and characteristic numbers (e.g., Mach number and Reynolds number).

Examples of dimensionless ratios of quantities of different kinds are the radian, with unit symbol rad for plane angle (arc length divided by radius), and the steradian, with unit symbol sr, for solid angle (area divided by radius squared). The special names, radian

and steradian, are given to the number 1 in these two cases to aid understanding, e.g., to indicate that a particular plane angle is expressed on the basis of arc length divided by radius rather than on the basis of a full circle divided into 360 degrees.

Decibels and Nepers

The decibel, symbol dB, represents the ratio of two values of power. The number *n* of decibels is 10 times the logarithm to the base 10 of the power ratio:

$$n(\text{dB}) = 10 \log_{10} (P_2/P_1)$$

Where the power ratio is less than unity, it is usual to invert the fraction and express the answer as a loss in decibels.

In the electrical case of voltage *V*, current *I*, and impedance *Z* (when *Z* is a pure resistance):

$$P_2/P_1 = (V_2/V_1)^2(Z_1/Z_2) = (I_2/I_1)^2(Z_2/Z_1)$$

Thus,

$$n_v(\text{dB}) = 20 \log_{10} (V_2/V_1) + 10 \log_{10} (Z_1/Z_2)$$

and

$$n_i(\text{dB}) = 20 \log_{10} (I_2/I_1) + 10 \log_{10} (Z_2/Z_1)$$

Only in the special case of matched impedances (*Z*₂ = *Z*₁) does

$$n = n_{Z\text{-matched}} = 20 \log_{10} (V_2/V_1) + 20 \log_{10} (I_2/I_1)$$

The neper, symbol Np, is a dimensionless unit for expressing the ratio of two values of amplitude. The number *n* of neper is the natural logarithm of the amplitude ratio:

$$n(\text{Np}) = \ln (A_2/A_1)$$

To convert decibels to nepers, multiply by 0.1151.

Numerical Values

To facilitate the reading of numbers, the digits may be separated into groups of three, counting from the decimal sign toward the left and the right. The groups should be separated by a small space, but not by a comma or a point. In numbers of four digits, the space is usually not necessary. For example:

$$\begin{aligned} 2.141\ 596\ 73\ 772\ 7372\ 0.133\ 47 \\ \pi = 3.141\ 592\ 653\ 589 \\ e = 2.718\ 281\ 828\ 459 \end{aligned}$$

American-English documents use a period for the decimal sign. British-English documents use a raised dot (·). For all other languages, the recommended decimal sign is a comma.

If the magnitude of a number is less than unity, the decimal sign should be preceded by a zero.

The sign of multiplication of numbers is a cross (×) or a raised dot (·).

Greek Alphabet

Name	Capital	Small	Commonly Used to Designate
Alpha	Α	α	angles, coefficients, attenuation constant, absorption factor, area
Beta	Β	β	angles, coefficients, phase constant
Gamma	Γ	γ	complex propagation constant (cap), relative density, angles, electrical conductivity, propagation constant
Delta	Δ	δ	increment or decrement (cap or small), determinant (cap), density, angles
Epsilon	Ε	ε	permittivity, electric intensity
Zeta	Ζ	ζ	coordinates, coefficients
Eta	Η	η	intrinsic impedance, efficiency, surface charge density, hysteresis, coordinates
Theta	Θ	θ, θ	angular phase displacement, time constant, reluctance, angles
Iota	Ι	ι	unit vector
Kappa	Κ	κ	susceptibility, coupling coefficient, thermal conductivity
Lambda	Λ	λ	permeance (cap), wavelength, attenuation constant
Mu	Μ	μ	permeability, amplification factor, prefix micro
Nu	Ν	ν	reluctivity, frequency
Xi	Ξ	ξ	coordinates
Omicron	Ο	ο	
Pi	Π	π	3.1416 (see more precise value in text)
Rho	Ρ	ρ	resistivity, volume charge density, coordinates
Sigma	Σ	σ	summation (cap), surface charge density, complex propagation constant, electrical conductivity, leakage coefficient, deviation
Tau	Τ	τ	time constant, volume resistivity, time-phase displacement, transmission factor, density
Upsilon	Υ	υ	
Phi	Φ	φ, φ	scalar potential (cap), magnetic flux, angles
Chi	Χ	χ	electrical susceptibility, angles
Psi	Ψ	ψ	dielectric flux, phase difference, coordinates, angles
Omega	Ω	ω	resistance in ohms (cap), solid angle (cap), angular velocity

Note: Small letter is used except where capital (cap) is indicated.

**FUNDAMENTAL PHYSICAL
CONSTANTS**

The following values of the fundamental physical constants are recommended by the international *Committee on Data for Science and Technology* (CODATA). They are called "1998 CODATA Recommended Values." They are available on the World Wide Web at the URL: <http://www.physics.nist.gov/cuu/Constants/index.html>.

The 1998 values replace the 1986 set. A detailed description of the data and their analysis that led to the 1998 values is published in the *Journal of Physical and Chemical Reference Data*, Vol. 28, No. 6, 1999 and in the *Reviews of Modern Physics*, Vol. 72, No. 2, 2000. That paper was prepared under the auspices of the CODATA Task Group on Fundamental Constants.

Fundamental Physical Constants — Complete Listing

Quantity	Symbol	Value	Unit	Relative std. uncert. u_r
UNIVERSAL				
speed of light in vacuum	c, c_0	299 792 458	m s^{-1}	(exact)
magnetic constant	μ_0	$4\pi \times 10^{-7}$ $= 12.566 370 614 \dots \times 10^{-7}$	N A^{-2} N A^{-2}	(exact)
electric constant $1/\mu_0 c^2$	ϵ_0	$8.854 187 817 \dots \times 10^{-12}$	F m^{-1}	(exact)
characteristic impedance of vacuum $\sqrt{\mu_0/\epsilon_0} = \mu_0 c$	Z_0	376.730 313 461...	Ω	(exact)
Newtonian constant of gravitation	G	$6.673(10) \times 10^{-11}$	$\text{m}^3 \text{kg}^{-1} \text{s}^{-2}$	1.5×10^{-3}
	G/hc	$6.707(10) \times 10^{-39}$	$(\text{GeV}/c^2)^{-2}$	1.5×10^{-3}
Planck constant	h	$6.626 068 76(52) \times 10^{-34}$	J s	7.8×10^{-8}
in eV s		$4.135 667 27(16) \times 10^{-15}$	eV s	3.9×10^{-8}
$h/2\pi$	\hbar	$1.054 571 596(82) \times 10^{-34}$	J s	7.8×10^{-8}
in eV s		$6.582 118 89(26) \times 10^{-16}$	eV s	3.9×10^{-8}
Planck mass $(\hbar c/G)^{1/2}$	m_P	$2.1767(16) \times 10^{-8}$	kg	7.5×10^{-4}
Planck length $\hbar/m_P c = (\hbar G/c^3)^{1/2}$	l_P	$1.6160(12) \times 10^{-35}$	m	7.5×10^{-4}
Planck time $l_P/c = (\hbar G/c^5)^{1/2}$	t_P	$5.3906(40) \times 10^{-44}$	s	7.5×10^{-4}
ELECTROMAGNETIC				
elementary charge	e	$1.602 176 462(63) \times 10^{-19}$	C	3.9×10^{-8}
	e/h	$2.417 989 491(95) \times 10^{14}$	A J^{-1}	3.9×10^{-8}
magnetic flux quantum $h/2e$	Φ_0	$2.067 833 636(81) \times 10^{-15}$	Wb	3.9×10^{-8}
conductance quantum $2e^2/h$	G_0	$7.748 091 696(28) \times 10^{-5}$	S	3.7×10^{-9}
inverse of conductance quantum	G_0^{-1}	12 906.403 786(47)	Ω	3.7×10^{-9}
Josephson constant ^a $2e/h$	K_J	$483 597.898(19) \times 10^9$	Hz V^{-1}	3.9×10^{-8}
von Klitzing constant ^b $h/e^2 = \mu_0 c/2\alpha$	R_K	25 812.807 572(95)	Ω	3.7×10^{-9}
Bohr magneton $e\hbar/2m_e$	μ_B	$927.400 899(37) \times 10^{-26}$	J T^{-1}	4.0×10^{-8}
in eV T ⁻¹		$5.788 381 749(43) \times 10^{-5}$	eV T ⁻¹	7.3×10^{-9}
	μ_B/h	$13.996 246 24(56) \times 10^9$	Hz T^{-1}	4.0×10^{-8}
	μ_B/hc	46.686 4521(19)	$\text{m}^{-1} \text{T}^{-1}$	4.0×10^{-8}
	μ_B/k	0.671 7131(12)	K T^{-1}	1.7×10^{-6}
nuclear magneton $e\hbar/2m_p$	μ_N	$5.050 783 17(20) \times 10^{-27}$	J T^{-1}	4.0×10^{-8}
in eV T ⁻¹		$3.152 451 238(24) \times 10^{-8}$	eV T ⁻¹	7.6×10^{-9}
	μ_N/h	7.622 593 96(31)	MHz T^{-1}	4.0×10^{-8}
	μ_N/hc	$2.542 623 66(10) \times 10^{-2}$	$\text{m}^{-1} \text{T}^{-1}$	4.0×10^{-8}
	μ_N/k	$3.658 2638(64) \times 10^{-4}$	K T^{-1}	1.7×10^{-6}
ATOMIC AND NUCLEAR				
General				
fine-structure constant $e^2/4\pi\epsilon_0\hbar c$	α	$7.297 352 533(27) \times 10^{-3}$		3.7×10^{-9}
inverse fine-structure constant	α^{-1}	137.035 999 76(50)		3.7×10^{-9}

Fundamental Physical Constants — Complete Listing

Quantity	Symbol	Value	Unit	Relative std. uncert. μ_t
Rydberg constant $\alpha^2 m_e c / 2h$	R_∞	10 973 731.568 549(83)	m^{-1}	7.6×10^{-12}
	$R_\infty c$	$3.289 841 960 368(25) \times 10^{15}$	Hz	7.6×10^{-12}
	$R_\infty hc$	$2.179 871 90(17) \times 10^{-18}$	J	7.8×10^{-8}
$R_\infty hc$ in eV		13.605 691 72(53)	eV	3.9×10^{-8}
Bohr radius $\alpha / 4\pi R_\infty = 4\pi\epsilon_0 \hbar^2 / m_e e^2$	a_0	$0.529 177 2083(19) \times 10^{-10}$	m	3.7×10^{-9}
Hartree energy $e^2 / 4\pi\epsilon_0 a_0 = 2R_\infty hc$ $= \alpha^2 m_e c^2$	E_h	$4.359 743 81(34) \times 10^{-18}$	J	7.8×10^{-8}
in eV		27.211 3834(11)	eV	3.9×10^{-8}
quantum of circulation	$h/2m_e$	$3.636 947 516(27) \times 10^{-4}$	$m^2 s^{-1}$	7.3×10^{-9}
	h/m_e	$7.273 895 032(53) \times 10^{-4}$	$m^2 s^{-1}$	7.3×10^{-9}
Electroweak				
Fermi coupling constant ^c	$G_F / (\hbar c)^3$	$1.166 39(1) \times 10^{-5}$	GeV^{-2}	8.6×10^{-6}
weak mixing angle ^d θ_W (on-shell scheme) $\sin^2 \theta_W = s_W^2 \equiv 1 - (m_W / m_Z)^2$	$\sin^2 \theta_W$	0.2224(19)		8.7×10^{-3}
Electron, e^-				
electron mass	m_e	$9.109 381 88(72) \times 10^{-31}$	kg	7.9×10^{-8}
in u, $m_e = A_r(e) u$ (electron relative atomic mass times u)		$5.485 799 110(12) \times 10^{-4}$	u	2.1×10^{-9}
energy equivalent	$m_e c^2$	$8.187 104 14(64) \times 10^{-14}$	J	7.9×10^{-8}
in MeV		0.510 998 902(21)	MeV	4.0×10^{-8}
electron-muon mass ratio	m_e / m_μ	$4.836 332 10(15) \times 10^{-3}$		3.0×10^{-8}
electron-tau mass ratio	m_e / m_τ	$2.875 55(47) \times 10^{-4}$		1.6×10^{-4}
electron-proton mass ratio	m_e / m_p	$5.446 170 232(12) \times 10^{-4}$		2.1×10^{-9}
electron-neutron mass ratio	m_e / m_n	$5.438 673 462(12) \times 10^{-4}$		2.2×10^{-9}
electron-deuteron mass ratio	m_e / m_d	$2.724 437 1170(58) \times 10^{-4}$		2.1×10^{-9}
electron to alpha particle mass ratio	m_e / m_α	$1.370 933 5611(29) \times 10^{-4}$		2.1×10^{-9}
electron charge to mass quotient	$-e/m_e$	$-1.758 820 174(71) \times 10^{11}$	C kg^{-1}	4.0×10^{-8}
electron molar mass $N_A m_e$	$M(e), M_c$	$5.485 799 110(12) \times 10^{-7}$	kg mol^{-1}	2.1×10^{-9}
Compton wavelength $h/m_e c$	λ_C	$2.426 310 215(18) \times 10^{-12}$	m	7.3×10^{-9}
$\lambda_C / 2\pi = \alpha a_0 = \alpha^2 / 4\pi R_\infty$	$\tilde{\lambda}_C$	$386.159 2642(28) \times 10^{-15}$	m	7.3×10^{-9}
classical electron radius $\alpha^2 a_0$	r_e	$2.817 940 285(31) \times 10^{-15}$	m	1.1×10^{-8}
Thomson cross section $(8\pi/3)r_e^2$	σ_e	$0.665 245 854(15) \times 10^{-28}$	m^2	2.2×10^{-8}
electron magnetic moment to Bohr magneton ratio	μ_e	$-928.476 362(37) \times 10^{-26}$	J T^{-1}	4.0×10^{-8}
to nuclear magneton ratio	μ_e / μ_B	$-1.001 159 652 1869(41)$		4.1×10^{-12}
electron magnetic moment anomaly $ \mu_e / \mu_B - 1$	μ_e / μ_N	$-1 838.281 9660(39)$		2.1×10^{-9}
electron g -factor $-2(1 + a_e)$	a_e	$1.159 652 1869(41) \times 10^{-3}$		3.5×10^{-9}
	g_e	$-2.002 319 304 3737(82)$		4.1×10^{-12}
electron-muon magnetic moment ratio	μ_e / μ_μ	206.766 9720(63)		3.0×10^{-8}

Fundamental Physical Constants — Complete Listing

Quantity	Symbol	Value	Unit	Relative std. uncert. u_r
electron-proton magnetic moment ratio electron to shielded proton	μ_e/μ_p	- 658.2106875(66)		1.0×10^{-8}
magnetic moment ratio (H ₂ O, sphere, 25 °C)	μ_e/μ'_p	- 658.2275954(71)		1.1×10^{-8}
electron-neutron magnetic moment ratio	μ_e/μ_n	960.92050(23)		2.4×10^{-7}
electron-deuteron magnetic moment ratio	μ_e/μ_d	-2 143.923 498(23)		1.1×10^{-8}
electron to shielded helion ^c magnetic moment ratio (gas, sphere, 25 °C)	μ_e/μ'_h	864.058 255(10)		1.2×10^{-8}
electron gyromagnetic ratio $2 \mu_e /\hbar$	γ_e	$1.760\,859\,794(71) \times 10^{11}$	$s^{-1} T^{-1}$	4.0×10^{-8}
	$\gamma_e/2\pi$	28 024.9540(11)	MHz T ⁻¹	4.0×10^{-8}
Muon, μ^-				
muon mass	m_μ	$1.883\,531\,09(16) \times 10^{-28}$	kg	8.4×10^{-8}
in u, $m_\mu = A_r(\mu)$ u (muon relative atomic mass times u)		0.113 428 9168(34)	u	3.0×10^{-8}
energy equivalent	$m_\mu c^2$	$1.692\,833\,32(14) \times 10^{-11}$	J	8.4×10^{-8}
in MeV		105.658 3568(52)	MeV	4.9×10^{-8}
muon-electron mass ratio	m_μ/m_e	206.768 2657(63)		3.0×10^{-8}
muon-tau mass ratio	m_μ/m_τ	$5.945\,72(97) \times 10^{-2}$		1.6×10^{-4}
muon-proton mass ratio	m_μ/m_p	0.112 609 5173(34)		3.0×10^{-8}
muon-neutron mass ratio	m_μ/m_n	0.112 454 5079(34)		3.0×10^{-8}
muon molar mass $N_A m_\mu$	$M(\mu), M_\mu$	$0.113\,428\,9168(34) \times 10^{-3}$	kg mol ⁻¹	3.0×10^{-8}
muon Compton wavelength $h/m_\mu c$	$\lambda_{C,\mu}$	$11.734\,441\,97(35) \times 10^{-15}$	m	2.9×10^{-8}
$\lambda_{C,\mu}/2\pi$	$\tilde{\lambda}_{C,\mu}$	$1.867\,594\,444(55) \times 10^{-15}$	m	2.9×10^{-8}
muon magnetic moment	μ_μ	$-4.490\,448\,13(22) \times 10^{-26}$	J T ⁻¹	4.9×10^{-8}
to Bohr magneton ratio	μ_μ/μ_B	$-4.841\,970\,85(15) \times 10^{-3}$		3.0×10^{-8}
to nuclear magneton ratio	μ_μ/μ_N	-8.890 597 70(27)		3.0×10^{-8}
muon magnetic moment anomaly $ \mu_\mu /(e\hbar/2m_\mu) - 1$	a_μ	$1.165\,916\,02(64) \times 10^{-3}$		5.5×10^{-7}
muon g-factor $-2(1 + a_\mu)$	g_μ	-2.002 331 8320(13)		6.4×10^{-10}
muon-proton magnetic moment ratio	μ_μ/μ_p	-3.183 345 39(10)		3.2×10^{-8}
Tau, τ^-				
tau mass ^f	m_τ	$3.167\,88(52) \times 10^{-27}$	kg	1.6×10^{-4}
in u, $m_\tau = A_r(\tau)$ u (tau relative atomic mass times u)		1.907 74(31)	u	1.6×10^{-4}
energy equivalent	$m_\tau c^2$	$2.847\,15(46) \times 10^{-10}$	J	1.6×10^{-4}
in MeV		1 777.05(29)	MeV	1.6×10^{-4}

Fundamental Physical Constants — Complete Listing

Quantity	Symbol	Value	Unit	Relative std. uncert. u_r
tau-electron mass ratio	m_τ/m_e	3477.60(57)		1.6×10^{-4}
tau-muon mass ratio	m_τ/m_μ	16.8188(27)		1.6×10^{-4}
tau-proton mass ratio	m_τ/m_p	1.893 96(31)		1.6×10^{-4}
tau-neutron mass ratio	m_τ/m_n	1.891 35(31)		1.6×10^{-4}
tau molar mass $N_A m_\tau$	$M(\tau), M_\tau$	$1.907 74(31) \times 10^{-3}$	kg mol ⁻¹	1.6×10^{-4}
tau Compton wavelength $h/m_\tau c$	$\lambda_{C,\tau}$	$0.697 70(11) \times 10^{-15}$	m	1.6×10^{-4}
$\lambda_{C,\tau}/2\pi$	$\tilde{\lambda}_{C,\tau}$	$0.111 042(18) \times 10^{-15}$	m	1.6×10^{-4}
Proton, p				
proton mass	m_p	$1.672 621 58(13) \times 10^{-27}$	kg	7.9×10^{-8}
in u, $m_p = A_r(p)$ u (proton relative atomic mass times u)		1.007 276 466 88(13)	u	1.3×10^{-10}
energy equivalent	$m_p c^2$	$1.503 277 31(12) \times 10^{-10}$	J	7.9×10^{-8}
in MeV		938.271 998(38)	MeV	4.0×10^{-8}
proton-electron mass ratio	m_p/m_e	1 836.152 6675(39)		2.1×10^{-9}
proton-muon mass ratio	m_p/m_μ	8.880 244 08(27)		3.0×10^{-8}
proton-tau mass ratio	m_p/m_τ	0.527 994(86)		1.6×10^{-4}
proton-neutron mass ratio	m_p/m_n	0.998 623 478 55(58)		5.8×10^{-10}
proton charge to mass quotient	e/m_p	$9.578 834 08(38) \times 10^7$	C kg ⁻¹	4.0×10^{-8}
proton molar mass $N_A m_p$	$M(p), M_p$	$1.007 276 466 88(13) \times 10^{-3}$	kg mol ⁻¹	1.3×10^{-10}
proton Compton wavelength $h/m_p c$	$\lambda_{C,p}$	$1.321 409 847(10) \times 10^{-15}$	m	7.6×10^{-9}
$\lambda_{C,p}/2\pi$	$\tilde{\lambda}_{C,p}$	$0.210 308 9089(16) \times 10^{-15}$	m	7.6×10^{-9}
proton magnetic moment	μ_p	$1.410 606 633(58) \times 10^{-26}$	J T ⁻¹	4.1×10^{-8}
to Bohr magneton ratio	μ_p/μ_B	$1.521 032 203(15) \times 10^{-3}$		1.0×10^{-8}
to nuclear magneton ratio	μ_p/μ_N	2.792 847 337(29)		1.0×10^{-8}
proton g-factor $2\mu_p/\mu_N$	g_p	5.585 694 675(57)		1.0×10^{-8}
proton-neutron				
magnetic moment ratio	μ_p/μ_n	-1.459 898 05(34)		2.4×10^{-7}
shielded proton magnetic moment (H ₂ O, sphere, 25 °C)	μ'_p	$1.410 570 399(59) \times 10^{-26}$	J T ⁻¹	4.2×10^{-8}
to Bohr magneton ratio	μ'_p/μ_B	$1.520 993 132(16) \times 10^{-3}$		1.1×10^{-8}
to nuclear magneton ratio	μ'_p/μ_N	2.792 775 597(31)		1.1×10^{-8}
proton magnetic shielding correction $1 - \mu'_p/\mu_p$ (H ₂ O, sphere, 25 °C)	σ'_p	$25.687(15) \times 10^{-6}$		5.7×10^{-4}
proton gyromagnetic ratio $2\mu_p/\hbar$	γ_p	$2.675 222 12(11) \times 10^8$	s ⁻¹ T ⁻¹	4.1×10^{-8}
	$\gamma_p/2\pi$	42.577 4825(18)	MHz T ⁻¹	4.1×10^{-8}
shielded proton gyromagnetic ratio $2\mu'_p/\hbar$ (H ₂ O, sphere, 25 °C)	γ'_p	$2.675 153 41(11) \times 10^8$	s ⁻¹ T ⁻¹	4.2×10^{-8}
	$\gamma'_p/2\pi$	42.576 3888(18)	MHz T ⁻¹	4.2×10^{-8}
Neutron, n				

Fundamental Physical Constants — Complete Listing

Quantity	Symbol	Value	Unit	Relative std. uncert. u_r
neutron mass	m_n	$1.674\,927\,16(13) \times 10^{-27}$	kg	7.9×10^{-8}
in u, $m_n = A_r(n)$ u (neutron relative atomic mass times u)		1.008 664 915 78(55)	u	5.4×10^{-10}
energy equivalent	$m_n c^2$	$1.505\,349\,46(12) \times 10^{-10}$	J	7.9×10^{-8}
in MeV		939.565 330(38)	MeV	4.0×10^{-8}
neutron-electron mass ratio	m_n/m_e	1 838.683 6550(40)		2.2×10^{-9}
neutron-muon mass ratio	m_n/m_μ	8.892 484 78(27)		3.0×10^{-8}
neutron-tau mass ratio	m_n/m_τ	0.528 722(86)		1.6×10^{-4}
neutron-proton mass ratio	m_n/m_p	1.001 378 418 87(58)		5.8×10^{-10}
neutron molar mass $N_A m_n$	$M(n), M_n$	$1.008\,664\,915\,78(55) \times 10^{-3}$	kg mol ⁻¹	5.4×10^{-10}
neutron Compton wavelength $h/m_n c$	$\lambda_{C,n}$	$1.319\,590\,898(10) \times 10^{-15}$	m	7.6×10^{-9}
$\lambda_{C,n}/2\pi$	$\tilde{\lambda}_{C,n}$	$0.210\,019\,4142(16) \times 10^{-15}$	m	7.6×10^{-9}
neutron magnetic moment	μ_n	$-0.966\,236\,40(23) \times 10^{-26}$	J T ⁻¹	2.4×10^{-7}
to Bohr magneton ratio	μ_n/μ_B	$-1.041\,875\,63(25) \times 10^{-3}$		2.4×10^{-7}
to nuclear magneton ratio	μ_n/μ_N	-1.913 042 72(45)		2.4×10^{-7}
neutron g -factor $2\mu_n/\mu_N$	g_n	-3.826 085 45(90)		2.4×10^{-7}
neutron-electron magnetic moment ratio	μ_n/μ_e	$1.040\,668\,82(25) \times 10^{-3}$		2.4×10^{-7}
neutron-proton magnetic moment ratio	μ_n/μ_p	-0.684 979 34(16)		2.4×10^{-7}
neutron to shielded proton magnetic moment ratio (H ₂ O, sphere, 25 °C)	μ_n/μ'_p	-0.684 996 94(16)		2.4×10^{-7}
neutron gyromagnetic ratio $2 \mu_n /\hbar$	γ_n	$1.832\,471\,88(44) \times 10^8$	s ⁻¹ T ⁻¹	2.4×10^{-7}
	$\gamma_n/2\pi$	29.164 6958(70)	MHz T ⁻¹	2.4×10^{-7}
Deuteron, d				
deuteron mass	m_d	$3.343\,583\,09(26) \times 10^{-27}$	kg	7.9×10^{-8}
in u, $m_d = A_r(d)$ u (deuteron relative atomic mass times u)		2.013 553 212 71(35)	u	1.7×10^{-10}
energy equivalent	$m_d c^2$	$3.005\,062\,62(24) \times 10^{-10}$	J	7.9×10^{-8}
in MeV		1 875.612 762(75)	MeV	4.0×10^{-8}
deuteron-electron mass ratio	m_d/m_e	3 670.482 9550(78)		2.1×10^{-9}
deuteron-proton mass ratio	m_d/m_p	1.999 007 500 83(41)		2.0×10^{-10}
deuteron molar mass $N_A m_d$	$M(d), M_d$	$2.013\,553\,212\,71(35) \times 10^{-3}$	kg mol ⁻¹	1.7×10^{-10}
deuteron magnetic moment	μ_d	$0.433\,073\,457(18) \times 10^{-26}$	J T ⁻¹	4.2×10^{-8}
to Bohr magneton ratio	μ_d/μ_B	$0.466\,975\,4556(50) \times 10^{-3}$		1.1×10^{-8}
to nuclear magneton ratio	μ_d/μ_N	0.857 438 2284(94)		1.1×10^{-8}
deuteron-electron magnetic moment ratio	μ_d/μ_e	$-4.664\,345\,537(50) \times 10^{-3}$		1.1×10^{-8}
deuteron-proton magnetic moment ratio	μ_d/μ_p	0.307 012 2083(45)		1.5×10^{-8}

Fundamental Physical Constants — Complete Listing

Quantity	Symbol	Value	Unit	Relative std. uncert. u_r
deuteron-neutron				
magnetic moment ratio	μ_d/μ_n	-0.448 206 52(11)		2.4×10^{-7}
Helion, h				
helion mass ^c				
in u, $m_h = A_r(\text{h})$ u (helion relative atomic mass times u)	m_h	$5.006\,411\,74(39) \times 10^{-27}$	kg	7.9×10^{-8}
energy equivalent	$m_h c^2$	$3.014\,932\,234\,69(86)$	u	2.8×10^{-10}
in MeV		$4.499\,538\,48(35) \times 10^{-10}$	J	7.9×10^{-8}
		$2\,808.391\,32(11)$	MeV	4.0×10^{-8}
helion-electron mass ratio				
helion-proton mass ratio	m_h/m_e	5 495.885 238(12)		2.1×10^{-9}
helion molar mass $N_A m_h$	m_h/m_p	2.993 152 658 50(93)		3.1×10^{-10}
shielded helion magnetic moment (gas, sphere, 25 °C)	$M(\text{h}), M_h$	$3.014\,932\,234\,69(86) \times 10^{-3}$	kg mol ⁻¹	2.8×10^{-10}
to Bohr magneton ratio	μ'_h	$-1.074\,552\,967(45) \times 10^{-26}$	J T ⁻¹	4.2×10^{-8}
to nuclear magneton ratio	μ'_h/μ_B	$-1.158\,671\,474(14) \times 10^{-3}$		1.2×10^{-8}
shielded helion to proton magnetic moment ratio (gas, sphere, 25 °C)	μ'_h/μ_N	$-2.127\,497\,718(25)$		1.2×10^{-8}
	μ'_h/μ_p	$-0.761\,766\,563(12)$		1.5×10^{-8}
shielded helion to shielded proton magnetic moment ratio (gas/H ₂ O, spheres, 25 °C)				
	μ'_h/μ'_p	$-0.761\,786\,1313(33)$		4.3×10^{-9}
shielded helion gyromagnetic ratio $2 \mu'_h /\hbar$ (gas, sphere, 25 °C)				
	γ'_h	$2.037\,894\,764(85) \times 10^8$	s ⁻¹ T ⁻¹	4.2×10^{-8}
	$\gamma'_h/2\pi$	32.434 1025(14)	MHz T ⁻¹	4.2×10^{-8}
Alpha particle, α				
alpha particle mass				
in u, $m_\alpha = A_r(\alpha)$ u (alpha particle relative atomic mass times u)	m_α	$6.644\,655\,98(52) \times 10^{-27}$	kg	7.9×10^{-8}
energy equivalent	$m_\alpha c^2$	4.001 506 1747(10)	u	2.5×10^{-10}
in MeV		$5.971\,918\,97(47) \times 10^{-10}$	J	7.9×10^{-8}
		3 727.379 04(15)	MeV	4.0×10^{-8}
alpha particle to electron mass ratio				
alpha particle to proton mass ratio	m_α/m_e	7 294.299 508(16)		2.1×10^{-9}
alpha particle molar mass $N_A m_\alpha$	m_α/m_p	3.972 599 6846(11)		2.8×10^{-10}
	$M(\alpha), M_\alpha$	$4.001\,506\,1747(10) \times 10^{-3}$	kg mol ⁻¹	2.5×10^{-10}
PHYSICO-CHEMICAL				
Avogadro constant				
atomic mass constant	N_A, L	$6.022\,141\,99(47) \times 10^{23}$	mol ⁻¹	7.9×10^{-8}
$m_u = \frac{1}{12} m(^{12}\text{C}) = 1$ u $= 10^{-3}$ kg mol ⁻¹ / N_A	m_u	$1.660\,538\,73(13) \times 10^{-27}$	kg	7.9×10^{-8}
energy equivalent	$m_u c^2$	$1.492\,417\,78(12) \times 10^{-10}$	J	7.9×10^{-8}
in MeV		931.494 013(37)	MeV	4.0×10^{-8}
Faraday constant ^g $N_A e$	F	96 485.3415(39)	C mol ⁻¹	4.0×10^{-8}

Fundamental Physical Constants — Complete Listing

Quantity	Symbol	Value	Unit	Relative std. uncert. u_r
molar Planck constant	$N_A h$	$3.990\,312\,689(30) \times 10^{-10}$	J s mol^{-1}	7.6×10^{-9}
	$N_A h c$	$0.119\,626\,564\,92(91)$	J m mol^{-1}	7.6×10^{-9}
molar gas constant	R	$8.314\,472(15)$	$\text{J mol}^{-1} \text{K}^{-1}$	1.7×10^{-6}
Boltzmann constant R/N_A	k	$1.380\,6503(24) \times 10^{-23}$	J K^{-1}	1.7×10^{-6}
in eV K^{-1}		$8.617\,342(15) \times 10^{-5}$	eV K^{-1}	1.7×10^{-6}
	k/h	$2.083\,6644(36) \times 10^{10}$	Hz K^{-1}	1.7×10^{-6}
	k/hc	$69.503\,56(12)$	$\text{m}^{-1} \text{K}^{-1}$	1.7×10^{-6}
molar volume of ideal gas RT/p				
$T = 273.15 \text{ K}, p = 101.325 \text{ kPa}$	V_m	$22.413\,996(39) \times 10^{-3}$	$\text{m}^3 \text{mol}^{-1}$	1.7×10^{-6}
Loschmidt constant N_A/V_m	n_0	$2.686\,7775(47) \times 10^{25}$	m^{-3}	1.7×10^{-6}
$T = 273.15 \text{ K}, p = 100 \text{ kPa}$	V_m	$22.710\,981(40) \times 10^{-3}$	$\text{m}^3 \text{mol}^{-1}$	1.7×10^{-6}
Sackur-Tetrode constant (absolute entropy constant) ^h				
$\frac{5}{2} + \ln[(2\pi m_u k T_1/h^2)^{3/2} k T_1/p_0]$				
$T_1 = 1 \text{ K}, p_0 = 100 \text{ kPa}$	S_0/R	$-1.151\,7048(44)$		3.8×10^{-6}
$T_1 = 1 \text{ K}, p_0 = 101.325 \text{ kPa}$		$-1.164\,8678(44)$		3.7×10^{-6}
Stefan-Boltzmann constant $(\pi^2/60)k^4/h^3c^2$	σ	$5.670\,400(40) \times 10^{-8}$	$\text{W m}^{-2} \text{K}^{-4}$	7.0×10^{-6}
first radiation constant $2\pi hc^2$	c_1	$3.741\,771\,07(29) \times 10^{-16}$	W m^2	7.8×10^{-8}
first radiation constant for spectral radiance $2hc^2$	c_{1L}	$1.191\,042\,722(93) \times 10^{-16}$	$\text{W m}^2 \text{sr}^{-1}$	7.8×10^{-8}
second radiation constant hc/k	c_2	$1.438\,7752(25) \times 10^{-2}$	m K	1.7×10^{-6}
Wien displacement law constant $b = \lambda_{\text{max}} T = c_2/4.965\,114\,231\dots$	b	$2.897\,7686(51) \times 10^{-3}$	m K	1.7×10^{-6}

^a See the "Adopted values" table for the conventional value adopted internationally for realizing representations of the volt using the Josephson effect.

^b See the "Adopted values" table for the conventional value adopted internationally for realizing representations of the ohm using the quantum Hall effect.

^c Value recommended by the Particle Data Group, Caso et al., Eur. Phys. J. C 3(1-4), 1-794 (1998).

^d Based on the ratio of the masses of the W and Z bosons m_W/m_Z recommended by the Particle Data Group (Caso et al., 1998). The value for $\sin^2\theta_W$ they recommend, which is based on a particular variant of the modified minimal subtraction ($\overline{\text{MS}}$) scheme, is $\sin^2\hat{\theta}_W(M_Z) = 0.231\,24(24)$.

^e The helium, symbol h, is the nucleus of the ^3He atom.

^f This and all other values involving m_τ are based on the value of $m_\tau c^2$ in MeV recommended by the Particle Data Group, Caso et al., Eur. Phys. J. C 3(1-4), 1-794 (1998), but with a standard uncertainty of 0.29 MeV rather than the quoted uncertainty of $-0.26 \text{ MeV}, +0.29 \text{ MeV}$.

^g The numerical value of F to be used in coulometric chemical measurements is $96\,485.3432(76) [7.9 \times 10^{-8}]$ when the relevant current is measured in terms of representations of the volt and ohm based on the Josephson and quantum Hall effects and the internationally adopted conventional values of the Josephson and von Klitzing constants K_{J-90} and R_{K-90} given in the "Adopted values" table.

^h The entropy of an ideal monoatomic gas of relative atomic mass A_r is given by $S = S_0 + \frac{3}{2}R \ln A_r - R \ln(p/p_0) + \frac{5}{2}R \ln(T/K)$.

CONVERSION FACTORS*

Exact numerical values (based on definitions) are followed by ***

The symbol E marks a power-of-ten multiplier in the usual way for the entry of data into calculators and computers.

To Convert From	Multiply By	To Get Unit
abampere	1.*** ** E 1	A
ab coulomb	1.*** ** E 1	C
abfarad	1.*** ** E 9	F
abhenry	1.*** ** E -9	H
abmho	1.*** ** E 9	S
abohm	1.*** ** E -9	ohm

*These conversion factors are an extension of those found in Mechtly, E. A. *The International System of Units—Physical Constants and Conversion Factors*, National Aeronautics and Space Administration (1973), Publication SP-7012 (1973). Washington, D.C. 20402: U.S. Government Printing Office.

To Convert From	Multiply By	To Get Unit
abvolt	1.*** **	V
acre	4.046 9--	m ²
angstrom	1.*** **	m
are	1.*** **	m ²
astronomical unit (IAU)	1.495 978 70-	m
astronomical unit (radio)	1.495 978 9--	m
atmosphere, standard	1.013 25*	Pa
atmosphere, technical, kg _f /cm ²	9.806 65*	Pa
bar	1.*** **	Pa
barn	1.*** **	m ²
barrel (petroleum, 42 gallons)	1.589 873	m ³
barye	1.*** **	Pa
board foot (144 cubic inches)	2.359 737 216	m ³
Btu (IST before 1956)	1.055 04-	J
Btu (IST after 1956)	1.055 056	J
Btu (mean)	1.055 87-	J
Btu (thermochemical)	1.054 350	J
Btu (39 °F)	1.059 67-	J
Btu (60 °F)	1.054 68-	J
bushel (US)	3.523 907 016 688	m ³
cable	2.194 56-	m
caliber (0.01 inch)	2.54*	m
calorie (IST)	4.186 8**	J
calorie (mean)	4.190 02-	J
calorie (thermochemical)	4.184 **	J
calorie (15 °C)	4.185 80-	J
calorie (20 °C)	4.181 90-	J
calorie, kilogram (IST)	4.186 8**	J
calorie, kilogram (mean)	4.190 02-	J
calorie, kilogram (thermochemical)	4.184 **	J
carat (metric)	2.*** **	kg
centimeter of mercury (0 °C)	1.333 22-	Pa
centimeter of water (4 °C)	9.806 38-	Pa
chain (engineer or Ramden)	3.048 **	m
chain (surveyor or Gunter)	2.011 684	m
circular mil	5.067 074 8--	m ²
clo	2.003 712	K·m ² /W
cord	3.624 556 3--	m ³
cubit	4.572	m
cup	2.365 882 365	m ³
curie	3.7**	Bq
day (mean solar)	8.64**	s (mean solar)
day (sidereal)	8.616 409 0--	s (mean solar)
degree (angle)	1.745 329 251 994 3--	rad
denier (international)	1.*** **	kg/m
dram (avoirdupois)	1.771 845 195 312 5**	kg
dram (troy or apothecary)	3.887 934 6**	kg
dram (US fluid)	3.696 691 195 312 5**	m ³
dyne	1.*** **	N
electronvolt (eV)	1.602 189 2--	J
emu of xxx (see ab---)		
erg	1.*** **	J
esu of xxx (see stat---)		
faraday (based on carbon 12)	9.648 530 9--	C mol ⁻¹
faraday (chemical, before 1960)	9.649 57-	C mol ⁻¹
faraday (physical, before 1960)	9.652 19-	C mol ⁻¹
fathom	1.828 804	m
fermi (femtometer, fm)	1.*** **	m
fluid ounce (US)	2.957 352 956 25*	m ³
foot	3.048 **	m
foot (US survey)	1200/3937 **	m
	3.048 006 096	m

To Convert From	Multiply By	To Get Unit
foot of water (39.2 °F)	2.988 98- --- --- --- E 3	Pa
footcandle	1.076 391 0-- --- --- E 1	lx
footlambert	3.426 259 --- --- --- ----	cd/m ²
free fall (standard, g)	9.806 65* *** *** *** E 0	m/s ²
furlong	2.011 68- --- --- --- E 2	m
gal	1.*** *** *** *** *** E -2	m/s ²
gallon (Canadian liquid)	4.546 090 --- --- --- E -3	m ³
gallon (UK dry)	4.404 883 770 86* *** E -3	m ³
gallon (UK liquid)	4.546 092 --- --- --- E -3	m ³
gallon (US dry)	4.404 884 --- --- --- E -3	m ³
gallon (US liquid)	3.785 411 784 *** *** E -3	m ³
gamma	1.*** *** *** *** *** E -9	T
gauss	1.*** *** *** *** *** E -4	T
gilbert	7.957 747 2-- --- --- E -1	A
gill (UK)	1.420 652 --- --- --- E -4	m ³
gill (US)	1.182 941 2-- --- --- E -4	m ³
grad (angle)	9.*** *** *** *** *** E -1	degree (angle)
	1.570 796 3-- --- --- E -2	rad
grain	6.479 891 *** *** *** E -5	kg
gram	1.*** *** *** *** *** E -3	kg
hand	1.016 *** *** *** *** E -1	m
hectare	1.*** *** *** *** *** E 4	m ²
hogshead (US)	2.384 809 423 92* *** E -1	m ³
horsepower (boiler)	9.809 50- --- --- --- E 3	W
horsepower (electric)	7.46* *** *** *** *** E 2	W
horsepower (metric)	7.354 99- --- --- --- E 2	W
horsepower (UK)	7.457 --- --- --- E 2	W
horsepower (water)	7.460 43- --- --- --- E 2	W
horsepower (550 ft·lb _f /s)	7.456 998 7-- --- --- E 2	W
hour (mean solar)	3.6** *** *** *** *** E 3	s (mean solar)
hour (sidereal)	3.590 170 4-- --- --- E 3	s (mean solar)
hundredweight (long)	5.080 234 544 *** *** E 1	kg
hundredweight (short)	4.535 923 7** *** *** E 1	kg
inch	2.54* *** *** *** *** E -2	m
inch of mercury (32 °F)	3.386 389 --- --- --- E 3	Pa
inch of mercury (60 °F)	3.376 85- --- --- --- E 3	Pa
inch of water (39.2 °F)	2.490 82- --- --- --- E 2	Pa
inch of water (60 °F)	2.488 4-- --- --- --- E 2	Pa
jansky	1.*** *** *** *** *** E-26	W/(m ² Hz)
kayser	1.*** *** *** *** *** E 2	m ⁻¹
kilocalorie (IST)	4.186 8-- --- --- --- E 3	J
kilocalorie (mean)	4.190 02- --- --- --- E 3	J
kilocalorie (thermochemical)	4.184 *** *** *** *** E 3	J
kilogram force (kg _f)	9.806 65* *** *** *** E 0	N
kilopond	9.806 65* *** *** *** E 0	N
kip	4.448 221 615 260 5** E 3	N
knot (international)	5.144 444 444 --- --- E -1	m/s
lambert	1/pi* *** *** *** *** E 4	cd/m ²
	3.183 098 8-- --- --- E 3	cd/m ²
langley	4.184 --- --- --- E 4	J/m ²
lb _f (pound force, avoirdupois)	4.448 221 615 260 5** E 0	N
lb _m (pound mass, avoirdupois)	4.535 923 7** *** *** E -1	kg
league	4.828 --- --- --- E 3	m
league (international)	5.556 *** *** *** *** E 3	m
league (US nautical)	5.559 552 *** *** *** E 3	m
light year	9.460 55- --- --- --- E 15	m
link (engineer or Ramden)	3.048 *** *** *** *** E -1	m
link (surveyor or Gunter)	2.011 68- --- --- --- E -1	m
liter	1.*** *** *** *** *** E -3	m ³
maxwell	1.*** *** *** *** *** E -8	Wb
meter (based on constant c) [†]		

† On October 20, 1983, the 17th CGPM adopted the following definition: "The meter is the length of the path traveled by light in vacuum during a time interval of 1/299 792 458 of a second." This new definition has the effect of giving a fixed value to the speed of light in vacuum of $c = 299\,792\,458$ m/s exactly. (Metrologia, Vol. 19, p. 163, 1984.)

To Convert From	Multiply By	To Get Unit
meter (based on Kr 86)	1.650 763 73- --- ---	E 6 wavelengths
mho	1.*** **	E 0 S
micron	1.*** **	E -6 m
mil	2.54* **	E -5 m
mile (international)	1.609 344 **	E 3 m
mile (international nautical)	1.852 **	E 3 m
mile (UK nautical)	1.853 184 **	E 3 m
mile (US nautical)	1.852 **	E 3 m
mile (US statute)	1.609 3-- --- ---	E 3 m
millibar	1.*** **	E 2 Pa
millimeter of mercury (0 °C)	1.333 224 --- ---	E 2 Pa
minute (angle)	2.908 882 086 66- ---	E -4 rad
minute (mean solar)	6.*** **	E 1 s (mean solar)
minute (sidereal)	5.983 617 4-- ---	E 1 s (mean solar)
month (mean calendar)	2.628 **	E 6 s (mean solar)
nautical mile (international)	1.852 **	E 3 m
nautical mile (UK)	1.853 184 **	E 3 m
nautical mile (US)	1.852 **	E 3 m
ocrsted	7.957 747 2-- ---	E 1 A/m
ounce (UK fluid)	2.841 307 --- ---	E -5 m ³
ounce (US fluid)	2.957 352 956 25* **	E -5 m ³
ounce force (avoirdupois)	2.780 138 5-- ---	E -1 N
ounce mass (avoirdupois)	2.834 952 312 5** **	E -2 kg
ounce mass (troy or apothecary)	3.110 347 68* **	E -2 kg
pace	7.62- --- ---	E -1 m
parsec (IAU)	3.085 678 --- ---	E 16 m
peck (US)	8.809 767 541 72* **	E -3 m ³
pennyweight	1.555 173 84* **	E -3 kg
perch	5.029 2-- ---	E 0 m
phot	1.--- ---	E 4 lx
pica (printer's)	4.217 517 6** **	E -3 m
pint (US dry)	5.506 104 713 575 **	E -4 m ³
pint (US liquid)	4.713 764 73* **	E -4 m ³
point (printer's)	3.514 598 **	E -4 m
poise	1.*** **	E -1 N·s/m ²
pole	5.029 2-- ---	E 0 m
pound force (lb _f avoirdupois)	4.448 221 615 260 5**	**** N
pound mass (lb _m avoirdupois)	4.535 923 7** **	E -1 kg
pound mass (troy or apothecary)	3.732 417 216 **	E -1 kg
poundal	1.382 549 543 76* **	E -1 N
quad	1.055 --- ---	E 18 J
quart (US dry)	1.101 220 942 715 **	E -3 m ³
quart (US liquid)	9.463 529 46* **	E -4 m ³
rad (absorbed radiation dose)	1.*** **	E -2 Gy
rayleigh (photon emission rate)	1.*** **	E 10 m ⁻² s ⁻¹
rem (dose equivalent)	1.*** **	E -2 Sv
rhe	1.*** **	E 1 m ² /(N·s)
rod	5.029 2-- ---	E 0 m
roentgen	2.579 76* **	E -4 C/kg
rutherford	1.*** **	E 6 Bq
scruple (apothecary)	1.295 978 2** **	E -3 kg
second (angle)	4.848 136 811 --- ---	E -6 rad
second (ephemeris)	1.000 000 000 --- ---	E 0 s (SI atomic)
second (mean solar)	Consult American Ephemeris and Nautical Almanac.	s (ephemeris)
second (sidereal)	9.972 695 7-- ---	E -1 s (mean solar)
section	2.589 998 --- ---	E 6 m ²
shake	1.--- ---	E -8 s
skein	1.097 28* **	E 2 m
slug	1.459 390 29- ---	E 1 kg
span	2.286 **	E -1 m
statampere	3.335 640 --- ---	E -10 A

To Convert From	Multiply By		To Get Unit
statcoulomb	3.335 640 --- --- ---	E-10	C
statfarad	1.112 650 --- --- ---	E-12	F
stathenry	8.987 554 --- --- ---	E 11	H
statohm	8.987 554 --- --- ---	E 11	ohm
statute mile (US)	1.609 3-- --- --- ---	E 3	m
statvolt	2.997 925 --- --- ---	E 2	V
stere	1.*** *** ** ** **	E 0	m ³
stilb	1.--- --- --- --- ---	E 4	cd/m ²
stoke	1.*** *** ** ** **	E -4	m ² /s
tablespoon	1.478 676 478 125 ***	E -5	m ³
teaspoon	4.928 921 593 75* ***	E -6	m ³
tex	1.*** *** ** ** **	E -6	kg/m
therm	1.055 --- --- --- ---	E 8	J
ton (assay)	2.916 666 6-- --- ---	E -2	kg
ton (long)	1.016 046 908 8** ***	E 3	kg
ton (metric)	1.*** *** ** ** **	E 3	kg
ton (nuclear equiv of TNT)	4.184 *** *** ** **	E 9	J
ton (refrigeration)	3.516 800 --- --- ---	E 3	W
ton (register)	2.831 684 659 2** ***	E 0	m ³
ton (short, 2000 lb _m)	9.071 847 4** *** **	E 3	kg
tonne	1.*** *** ** ** **	E 3	kg
torr (0 °C)	1.333 22- --- --- ---	E 2	Pa
township	9.323 994 --- --- ---	E 7	m ²
unit pole	1.256 637 --- --- ---	E -7	Wb
yard	9.144 *** *** ** **	E -1	m
year (calendar, 365 days)	3.153 6** *** ** **	E 7	s (mean solar)
year (sidereal)	3.155 815 0-- --- ---	E 7	s (mean solar)
year (tropical)	3.155 692 6-- --- ---	E 7	s (mean solar)
year 1900 (tropical, January, day 0, hour 12)	3.155 692 597 47* ***	E 7	s (ephemeris)
	3.155 692 597 47- ---	E 7	s (SI atomic)

4 Properties of Materials

*Revised by
Eugene A. Mechtly*

General Properties of the Elements	4-3
Periodic Classification of the Elements	4-3
Physical Properties of the Elements	4-3
Galvanic Series in Sea Water	4-3
Temperature–EMF Characteristics of Thermocouples Electromotive Force and Other Properties	4-3
Conducting Materials	4-3
Semiconducting Materials	4-3
Insulating Materials	4-12
Magnetic Materials	4-17
Soft Magnetic Metals	
Permanent-Magnet Materials	
Ferrites	

Magnetostriction 4-27

Piezoelectricity 4-28

Acoustic Properties of Some Materials 4-31

GENERAL PROPERTIES OF THE ELEMENTS

Some properties of the elements are listed in Table 1.* Some of the listed quantities are defined as follows.

Atomic number Z represents the number of protons per atom.

Mass number $Z + N$ is equal to the number of protons Z plus the number of neutrons N present in the nuclei. Mass numbers of the most abundant isotopes are given in order of decreasing abundance. For example, Cadmium Cd-48 mass numbers 114–112 means that cadmium atoms of greater abundance (28.86%) have a mass number of 114; that is, the nucleus of Cd¹¹⁴ has $114 - 48 = 66$ neutrons while isotope Cd¹¹² of lower abundance (24.07%) has $112 - 48 = 64$ neutrons.

Atomic radius values listed provide a comparison of sizes (deduced from interatomic spacing of bound atoms).

Gram atomic volume in cubic centimeters gives the volume occupied in the solid state by an atom at its melting point. The gram atomic volume contains the Avogadro number of atoms.†

Electronegativity represents the relative tendency of an atom to attract shared electron pairs. The highest electronegativity is assigned to fluorine with the value 3.90.

First ionization potential is the work in electron volts required to pull 1 electron off an isolated neutral atom.

$$1 \text{ electron volt} = 1.602 \times 10^{-19} \text{ joule}^\ddagger$$

Electron work function, expressed in electronvolts, represents the energy that must be supplied to an electron to cross over the surface barrier of a metal. That energy may be supplied by heat (thermionic work function), by light (photoelectric work function), or by contact with a dissimilar metal (contact potential).

Electrochemical equivalents are expressed in joules per gram (J/g) liberated at the electrode.

PERIODIC CLASSIFICATION OF THE ELEMENTS

Fig. 1 is a periodic table of the elements.

Oxidation number is defined as the charge that an atom appears to have in a compound when electrons are counted according to certain rules:

- (A) In the free elements each atom has an oxidation number of 0.
- (B) Electrons shared between two unlike atoms are counted with the more electronegative atom.
- (C) Electrons shared between the two like atoms are divided equally between sharing atoms.

PHYSICAL PROPERTIES OF THE ELEMENTS

Some of the physical properties of the elements are listed in Table 2.

GALVANIC SERIES IN SEA WATER

In sea water, two dissimilar metals connected by a conductor form a galvanic cell. If the two metals are in different groups of Chart 1 (separated by spaces), the metal coming first in the series—starting from corroded end to protected end—will be anodic (i.e., corroded by the metal contained in the group farther from the corroded end). If the two metals are in the same group, no appreciable corrosive action will take place.

TEMPERATURE-EMF CHARACTERISTICS OF THERMOCOUPLES‡

Fig. 2 shows temperature-emf characteristics of thermocouples.

Electromotive Force and Other Properties

Electromotive force and other properties of thermocouples are listed in Tables 3 and 4.

CONDUCTING MATERIALS

Conducting materials (Tables 5–7) can be classified as follows:

Conductors: Resistivities from 10^{-6} to $10^{-4} \Omega\text{-cm}$ (1 to $100 \mu\Omega\text{-cm}$). Conductivities from 10^4 to $10^6 \text{ S}\cdot\text{cm}^{-1}$.
Semiconductors: Resistivities from 10^{-4} to $10^9 \Omega\text{-cm}$. Conductivities from 10^{-9} to $10^4 \text{ S}\cdot\text{cm}^{-1}$.

Insulators: Resistivities from 10^9 to $10^{25} \Omega\text{-cm}$. Conductivities from 10^{-25} to $10^{-9} \text{ S}\cdot\text{cm}^{-1}$.

SEMICONDUCTING MATERIALS

Some properties of semiconductor materials are listed in Table 8.

* Tables 1 and 2, Chart 1, and Fig. 1 of this chapter are partly based on data from the following sources: *Handbook of Chemistry and Physics*, 55th ed., CRC Press, Inc., 1974; *Fundamentals of Chemistry*, John Wiley & Sons; *The Encyclopedia of Electrochemistry*, Reinhold Publishing Corp., 1964; *American Institute of Physics Handbook*, 3rd ed., McGraw-Hill Book Co., 1972; *Lange's Handbook of Chemistry*, 11th ed., McGraw-Hill Book Co., 1973.

† See Chapter 3 for precise value.

‡ R. L. Weber, *Temperature Measurement and Control* (Philadelphia: Blakiston Co., 1941; pp. 68–71).

TABLE I. PROPERTIES

	Symbol	Atomic Number Z	Mass Number Z + N	Relative Atomic Mass	Atomic Radius (10^{-10} m)	Gram Atomic Volume (cm^3)
Actinium	Ac	89	227	227		
Aluminum	Al	13	27	26.98	1.25	10
Americium	Am	95	243	243		
Antimony	Sb	51	121-123	121.75	1.41	18
Argon	Ar or A	18	40	39.948	1.74	24
Arsenic	As	33	75	74.92	1.21	16
Astatine	At	85	210	210		
Barium	Ba	56	138	137.34	1.98	38
Berkelium	Bk	97	247	247		
Beryllium	Be	4	9	9.012	0.89	5
Bismuth	Bi	83	209	208.98	1.52	21
Boron	B	5	11	10.81	0.88	5
Bromine	Br	35	79-81	79.904	1.14	23
Cadmium	Cd	48	114-112	112.40	1.41	13
Calcium	Ca	20	40	40.08	1.74	26
Californium	Cf	98	251	251		
Carbon	C	6	12	12.011	0.77	5
Cerium	Ce	58	140	140.12	1.65	21
Cesium	Cs	55	133	132.905	2.35	71
Chlorine	Cl	17	35	35.453	0.99	19
Chromium	Cr	24	52	51.996	1.17	7
Cobalt	Co	27	59	58.933	1.16	7
Copper	Cu	29	63	63.546	1.17	7
Curium	Cm	96	247	247		
Dysprosium	Dy	66	164-162-163	162.50	1.59	19
Einsteinium	Es or E	99	254	254		
Erbium	Er	68	166-168-167	167.26	1.57	18
Europium	Eu	63	153-151	151.96	1.85	29
Fermium	Fm	100	257	257		
Fluorine	F	9	19	18.998	0.64	15
Francium	Fr	87	223	223		
Gadolinium	Gd	64	158-160-156	157.25	1.61	20
Gallium	Ga	31	69-71	69.72	1.25	12
Germanium	Ge	32	74-72-70	72.59	1.22	13
Gold	Au	79	197	196.967	1.34	10
Hafnium	Hf	72	180-178-177	178.49	1.44	13
Helium	He	2	4	4.003		32
Holmium	Ho	67	165	164.93	1.58	19
Hydrogen	H	1	1	1.008	0.37	13
Indium	In	49	115	114.82	1.50	16
Iodine	I	53	127	126.904	1.33	26
Iridium	Ir	77	193-191	192.22	1.26	9
Iron	Fe	26	56	55.847	1.17	7
Krypton	Kr	36	84-86	83.80	1.89	33
Lanthanum	La	57	139	138.905	1.69	22
Lawrencium	Lw	103	257	257		
Lead	Pb	82	208-206-207	207.2	1.54	18
Lithium	Li	3	7	6.940	1.23	13
Lutetium	Lu	71	175	174.97	1.56	18
Magnesium	Mg	12	24	24.305	1.36	14
Manganese	Mn	25	55	54.938	1.17	7
Mendelevium	Md or Mv	101	256	256		

OF THE ELEMENTS

Electro-negativity, Relative Scale	First Ionization Potential (eV)	Electron Work Function			Electrochemical Equiv.	
		Thermionic	Photoelectric	Contact	Valence Involved	3 600 J/g
1.1	6.9				3	0.35
1.5	5.98		4.08	3.38	3	2.98
2.05	6.05					
0	8.64		4.01	4.14	5	1.1
	5.76				n*	0.67
2.0	9.81		5.11		5	1.79
2.2						
0.9	5.21	2.11	2.48	1.73	2	0.39
1.5	9.32		3.92	3.10	2	5.94
1.9	7.29		4.25	4.17	5	0.64
2.0	8.3		4.5		3	7.43
2.85	11.81				1	0.335
1.7	8.99	4.07	4.0	2	0.477	
1.0	6.11	2.24	2.706	3.33	2	1.337
2.6	11.26	4.34	4.81		4	8.93
1.1	5.6	2.6	2.84		3	0.574
0.7	3.89	1.81	1.92	4.46	1	0.2
3.15	12.97				1	0.756
1.6	6.76	4.60	4.37	4.38	3	1.546
1.8	7.86	4.40	4.20	4.21	2	0.91
1.9	7.72	4.26	4.18	4.46	2	0.84
1.2	5.93				3	0.495
1.2	6.10				3	0.48
1.1	5.67				3	0.53
3.9	17.42				1	1.41
0.65						
1.1	6.16				3	0.513
1.6	5.99	4.12		3.80	3	1.15
1.9	7.89		4.5	4.5	4	1.48
2.4	9.22	4.32	4.82	4.46	3	0.41
1.3	7.0	3.53			4	0.600
0	24.59				n*	6.698
1.2	6.02				3	0.488
2.2	13.59				1	26.59
1.7	5.78				3	0.700
2.65	10.45		6.8		1	0.211
2.2	9.1	5.3		4.57	4	0.555
1.8	7.87	4.25	4.33	4.40	3	1.440
0	13.99				n*	0.32
1.1	5.61	3.3			3	0.579
1.8	7.42		4.05	3.94	4	0.517
1.0	5.39		2.35	2.49	1	3.862
1.2	6.15				3	0.46
1.2	7.64		3.68	3.63	2	2.204
1.5	7.43	3.83	3.76	4.14	4	1.952

TABLE 1 (CONT). PROPERTIES

	Symbol	Atomic Number Z	Mass Number Z + N	Relative Atomic Mass	Atomic Radius (10^{-10} m)	Gram Atomic Volume (cm^3)
Mercury	Hg	80	202-200-199	200.59	1.44	14
Molybdenum	Mo	42	98-96-92-95	95.94	1.29	9
Neodymium	Nd	60	142-144-146	144.24	1.64	21
Neon	Ne	10	20	20.179	1.31	17
Neptunium	Np	93	237	237.048		
Nickel	Ni	28	58	58.71	1.15	6
Niobium	Nb	41	93	92.906	1.34	11
Nitrogen	N	7	14	14.007	0.70	14
Nobelium	No	102	254	254		
Osmium	Os	76	192-190-189	190.2	1.26	9
Oxygen	O	8	16	15.999	0.66	11
Palladium	Pd	46	108-106-105	106.4	1.28	9
Phosphorus	P	15	31	30.974	1.10	17
Platinum	Pt	78	195-194-196	195.09	1.29	9
Plutonium	Pu	94	242	242		
Polonium	Po	84	209	210	1.53	
Potassium	K	19	39	39.098	2.03	46
Praseodymium	Pr	59	141	140.907	1.65	21
Promethium	Pm	61	145	145		
Protactinium	Pa	91	231	231.036		
Radium	Ra	88	226	226.025		45
Radon	Rn	86	222	222	2.14	50
Rhenium	Re	75	187-185	186.2	1.28	9
Rhodium	Rh	45	103	102.905	1.25	8
Rubidium	Rb	37	85-87	85.468	2.16	56
Ruthenium	Ru	44	102-104-101	101.07	1.24	8
Samarium	Sm	62	152-154-147	150.35	1.66	20
Scandium	Sc	21	45	44.956	1.44	15
Selenium	Se	34	80-78	78.96	1.17	16
Silicon	Si	14	28	28.086	1.17	12
Silver	Ag	47	107-109	107.868	1.34	10
Sodium	Na	11	23	22.99	1.57	24
Strontium	Sr	38	88	87.62	1.92	34
Sulfur	S	16	32	32.064	1.04	16
Tantalum	Ta	73	181	180.948	1.34	11
Technetium	Tc	43	99	98.906		
Tellurium	Te	52	130-128-126	127.60	1.37	21
Terbium	Tb	65	159	158.925	1.59	19
Thallium	Tl	81	205-203	204.37	1.55	17
Thorium	Th	90	232	232.038	1.65	20
Thulium	Tm	69	169	168.934	1.56	18
Tin	Sn	50	120-118	118.69	1.40	16
Titanium	Ti	22	48	47.90	1.32	11
Tungsten	W	74	184-186-182	183.85	1.30	10
Uranium	U	92	238	238.029	1.42	13
Vanadium	V	23	51	50.94	1.22	8
Xenon	Xe	54	132-129-131	131.30	2.09	43
Ytterbium	Yb	70	174-172-173	173.04	1.70	25
Yttrium	Y	39	89	88.906	1.62	21
Zinc	Zn	30	64-66-68	65.38	1.25	9
Zirconium	Zr	40	90-94-92	91.22	1.45	14

*n = nonvalent.

OF THE ELEMENTS

Electro-negativity, Relative Scale	First Ionization Potential (eV)	Electron Work Function			Electrochemical Equiv.	
		Thermionic	Photoelectric	Contact	Valence Involved	3 600 J/g
1.9	10.43		4.53	4.50	2	0.267
1.8	7.10	4.20	4.25	4.28	6	1.67
1.1	5.49	3.3			3	0.557
0	21.56				n*	1.33
1.3	5.8					
1.8	7.63	5.03	5.01	4.96	2	0.913
1.6	6.88	4.01	4.5			
3.05	14.53				5	9.57
2.2	8.7			4.55	4	0.56
3.5	13.62				2	3.35
2.2	8.33	4.99	4.97	4.49	4	1.005
2.15	10.48				5	4.33
2.2	9.0	5.32	5.22	5.36	4	0.549
	5.8					
2.0	8.43				6	0.766
0.8	4.34		2.24	1.60	1	0.685
1.1	5.42	2.7			3	0.571
1.1	5.55					
1.5					5	0.580
0.9	5.28				2	0.237
0	10.75				n*	0.121
1.9	7.87	5.1	5.0		7	1.007
2.2	7.46	4.80	4.57	4.52	4	1.042
0.8	4.18		2.09		1	0.314
2.2	7.37			4.52	4	1.054
1.1	5.63	3.2			3	0.535
1.3	6.54				3	1.783
2.45	9.75		4.8	4.42	6	2.037
1.9	8.15	3.59	4.52	4.2	4	3.821
1.9	7.57	3.56	4.73	4.44	1	0.248
0.9	5.14		2.28	1.9	1	1.166
1.0	5.69		2.74		2	0.612
2.6	10.36				6	5.01
1.3	7.88	4.19	4.14	4.1	5	0.741
1.9	7.28					
2.3	9.01		4.76	4.70	6	1.260
1.2	5.98				3	0.505
1.8	6.11		3.68	3.84	3	0.393
1.3	6.95	3.35	3.47	3.46	4	0.462
1.2	6.18				3	0.475
1.8	7.34		4.38	4.09	4	0.903
1.5	6.82	3.95	4.06	4.14	4	2.238
1.7	7.98	4.52	4.49	4.38	6	0.874
1.7	6.08	3.27	3.63	4.32	6	0.676
1.6	6.74	4.12	3.77	4.44	5	2.63
0	12.13				n*	0.204
1.2	6.25				3	0.465
1.3	6.38				3	0.904
1.6	9.39		3.73	3.78	2	0.820
1.6	6.84	4.21	3.82	3.60	4	1.175

TABLE 2. PHYSICAL PROPERTIES

	Symbol	Atomic Number	Density at 20 °C (g/cm ³)	Relative Hardness	Melting Point (°C)	Boiling Point (°C)
Actinium	Ac	89			1 050	3 200
Aluminum	Al	13	2.70	2.9	660	2 467
Americium	Am	95	13.67		994	2 600
Antimony	Sb	51	6.62	3	630.5	1 750
Argon	Ar or A	18	1.78*		-189.2	-185.7
Arsenic (gray)	As	33	5.73	3.5	820‡	615**
Astatine	At	85			302	337
Barium	Ba	56	3.5		725	1 640
Berkelium	Bk	97				
Beryllium	Be	4	1.82	3	1 278	2 970
Bismuth	Bi	83	9.80	2.5	271.3	1 560
Boron	B	5	2.46	9.5	2 300	2 550**
Bromine	Br	35	3.12		-7.2	58.8
Cadmium	Cd	48	8.65	2.0	320.9	765
Calcium	Ca	20	1.54		842	1 487
Californium	Cf	98				
Carbon	C	6	2.22	10†	>3 500	4 827
Cerium	Ce	58	6.9	2.5	795	3 468
Cesium	Cs	55	1.87	0.2	28.5	690
Chlorine	Cl	17	3.21*		-100.98	-34.7
Chromium	Cr	24	7.14	9	1 890	2 672
Cobalt	Co	27	8.9	5	1 495	2 900
Copper	Cu	29	8.96	3	1 083	2 595
Curium	Cm	96	13.51		1 340	
Dysprosium	Dy	66	8.54		1 407	2 600
Einsteinium	Es or E	99				
Erbium	Er	68	9.05		1 497	2 900
Europium	Eu	63	5.26		826	1 439
Fermium	Fm	100				
Fluorine	F	9	1.69*††		-220	-188
Francium	Fr	87			27	677
Gadolinium	Gd	64	7.89		1 312	3 000
Gallium	Ga	31	5.91	1.5	29.78	2 403
Germanium	Ge	32	5.36	6.2	937.4	2 830
Gold	Au	79	19.3	2.5	1 063	2 966
Hafnium	Hf	72	13.31		2 220	4 602
Helium	He	2	0.1664*		<-272‡	-268.94
Holmium	Ho	67	8.803		1 461	2 600
Hydrogen	H	1	0.08375*		-259.14	-252.8
Indium	In	49	7.31	1.2	156	2 050
Iodine	I	53	4.93		113.5	184.3
Iridium	Ir	77	22.4	6.15	2 410	4 527
Iron	Fe	26	7.87	4	1 535	3 000
Krypton	Kr	36	3.448*		-156.6	-152.3
Lanthanum	La	57	6.15		920	3 469
Lawrencium	Lw	103				
Lead	Pb	82	11.34	1.5	327.4	1 744
Lithium	Li	3	0.53	0.6	179	1 336
Lutetium	Lu	71	9.84		1 652	3 327
Magnesium	Mg	12	1.74	2	651	1 100
Manganese	Mn	25	7.44	5.0	1 244	2 097
Mendelevium	Md or Mv	101				
Mercury	Hg	80	13.55	1.5	-38.87	356.9

OF THE ELEMENTS

Latent Heat of Fusion (4.184 J/g)	Specific Heat at 20 °C (4.184 J/g°C)	Thermal Conductivity at 20 °C (W/cm°C)	Linear Expansion at 20 °C (10 ⁻⁶ /°C)	Elasticity Modulus (GN/m ²)	Tensile Strength (MN/m ²)
93	0.226	2.18	22.9	71.1	61.8
38.3	0.049	0.19	8.5-10.8	77.5	10.3
6.7	0.125	1.7 × 10 ⁻¹			
	0.082		4.7		
324	0.425	1.64	12	294	117
12.5	0.0294	0.084	13.3	31.4	
16.2	0.307		2		
13.2	0.107	0.91	29.8	53.9	70.6
	0.055		25	20.6	55.9
	0.145				
	0.165	0.24	0.6-4.3	4.9	
3.8	0.042		97		88.8
23	0.052	0.072 × 10 ⁻⁴			
	0.226				
75.6	0.11	0.69	6.2		
58.4	0.1001	0.69	12.3	206	237
50.6	0.0921	3.94	16.5	108	221
10.1					
19.2	0.079		18		
	0.073				
16.1	0.031	2.96	14.2	71.6	113
	1.25	13.9 × 10 ⁻⁴			
15	3.415	17 × 10 ⁻⁴			
	0.057	0.24	33		2.9
15.8	0.052	43.5 × 10 ⁻⁴	93		
33	0.032	1.4	6.5	515	
65	0.108	0.79	11.7	196	201
			0.89 × 10 ⁻⁴		
	0.045				
6.3	0.030	0.35	28.7	17.7	13
159	0.79	0.71	56		
88	0.249	1.55	25.2	45.1	89.7
64.8	0.107		23	157	283
2.7	0.033	0.084			

TABLE 2 (CONT). PHYSICAL

	Symbol	Atomic Number	Density at 20 °C (g/cm ³)	Relative Hardness	Melting Point (°C)	Boiling Point (°C)
Molybdenum	Mo	42	10.2	6	2 610	4 800
Nodymium	Nd	60	7.05		1 024	3 027
Neon	Ne	10	0.8387*		-248.7	-245.9
Neptunium	Np	93	20.45		640	3 902
Nickel	Ni	28	8.9	5	1 453	2 732
Niobium	Nb	41	8.57		2 468	5 127
Nitrogen	N	7	1.1649*		-209.9	-195.8
Nobelium	No	102				
Osmium	Os	76	22.48	7.0	3 000	5 000
Oxygen	O	8	1.3318*		-218.4	-183
Palladium	Pd	46	12	4.8	1 552	2 927
Phosphorus	P	15	1.82		44.1	280
Platinum	Pt	78	21.45	4.3	1 769	3 827
Plutonium	Pu	94	19.82		639.5	3 235
Polonium	Po	84	9.2		254	962
Potassium	K	19	0.86	0.5	63.65	774
Praseodymium	Pr	59	6.63		935	3 127
Promethium	Pm	61			1 027	2 027
Protactinium	Pa	91	15.37		1 227	4 027
Radium	Ra	88	5		700	1 525
Radon	Rn	86	4.40*		-71	-61.8
Rhenium	Rc	75	20		3 180	5 627
Rhodium	Rh	45	12.44	6	1 966	3 727
Rubidium	Rb	37	1.53	0.3	38.5	688
Ruthenium	Ru	44	12.2	6.5	2 250	3 900
Samarium	Sm	62	7.7		1 072	1 900
Scandium	Sc	21	2.5		1 539	2 727
Selenium	Se	34	4.81	2	217	685
Silicon	Si	14	2.4	7	1 410	2 355
Silver	Ag	47	10.49	2.7	960.8	2 212
Sodium	Na	11	0.97	0.4	97	892
Strontium	Sr	38	2.6	1.8	769	1 384
Sulfur	S	16	2.07	2.0	116	444.6
Tantalum	Ta	73	16.6	7	2 996	5 425
Technetium	Tc	43	11.49		2 200	4 700
Tellurium	Tc	52	6.24	2.3	449.5	990
Terbium	Tb	65	8.27		1 356	2 800
Thallium	Tl	81	11.85	1.2	303.5	1 457
Thorium	Th	90	11.5		1 800	4 200
Thulium	Tm	69	9.33		1 545	1 727
Tin	Sn	50	7.3	1.8	231.89	2 270
Titanium	Ti	22	4.54	4	1 675	3 260
Tungsten	W	74	19.3	7	3 410	5 660
Uranium	U	92	18.7		1 133	3 818
Vanadium	V	23	5.68		1 890	3 400
Xenon	Xe	54	5.495*		-111	-107
Ytterbium	Yb	70	6.98		824	1 427
Yttrium	Y	39	5.51		1 495	2 927
Zinc	Zn	30	7.14	2.5	419.4	907
Zirconium	Zr	40	6.4	4.7	1 852	4 377

*g/L. †diamond. ‡36 atm. §26 atm. **sublimes. ††At 15 °C.

PROPERTIES OF THE ELEMENTS

Latent Heat of Fusion (4.184 J/g)	Specific Heat at 20 °C (4.184 J/g°C)	Thermal Conductivity at 20 °C (W/cm°C)	Linear Expansion at 20 °C (10 ⁻⁶ /°C)	Elasticity Modulus (GN/m ²)	Tensile Strength (MN/m ²)
69.0	0.065 0.045	1.46 4.57 × 10 ⁻⁴	4.9	343	1 177
73.8	0.112	0.9	13.3	206	317
68.0	0.064	0.52	7.1		
6.2	0.247				
34.0	0.031	0.61	5		
3.3	0.218				
34.2	0.059	0.70	11.8	117	137
5.0	0.177		125		
27.1	0.032	0.69	8.9	147	157
3.0	0.032	0.08	54		
14.5	0.177 0.458	0.99	83		
	0.035				
50.0	0.060	1.5	8.1	29.4	
6.1	0.080 0.061		90 9.1		
16.0	0.077	0.005	37		
430.0	0.176	0.84	2.8-7.3	108	
24.3	0.056	4.08	18.9	70.6	148
27.5	0.295	1.35	71		
25					
9.3	0.175	26.4 × 10 ⁻⁴	6.4		
41.0	0.036	0.54	6.6	186	490
25.3	0.047	0.060	16.8	20.6	11.0
7.2	0.031	0.39	28		
17.0	0.028	0.41	11.1		549
14.4	0.054	0.64	23	403	13.7
100.0	0.142	0.2	8.5	83.4	
44.0	0.034	1.99	4.3	343	2 650
12.0	0.028	0.25	13.4		
98.0	0.115	0.60 5.9 × 10 ⁻⁴	8		
24.1	0.09 0.066	1.1	17.39 5.6	82.4 73.5	103 294

PERIOD	GROUP																			
	1A	2A	3B	4B	5B	6B	7B	8						1B	2B	3A	4A	5A	6A	7A
1	⁺¹ ₁ H																			⁰ ₂ He
	LIGHT METALS		HEAVY METALS						NON-METALS			INERT GAS								
2	⁺¹ ₃ Li	⁺² ₄ Be													⁺³ ₅ B	⁺² ₆ C	⁺³ ₇ N	⁻² ₈ O	⁻¹ ₉ F	⁰ ₁₀ Ne
	BRITTLE		DUCTILE						LOW-MELTING											
3	⁺¹ ₁₁ Na	⁺² ₁₂ Mg													⁺³ ₁₃ Al	⁺² ₁₄ Si	⁺³ ₁₅ P	⁻³ ₁₆ S	⁺⁴ ₁₇ Cl	⁻¹ ₁₈ Ar
4	⁺¹ ₁₉ K	⁺² ₂₀ Ca	⁺³ ₂₁ Sc	⁺² ₂₂ Ti	⁺³ ₂₃ V	⁺² ₂₄ Cr	⁺³ ₂₅ Mn	⁺² ₂₆ Fe	⁺³ ₂₇ Co	⁺² ₂₈ Ni	⁺¹ ₂₉ Cu	⁺² ₃₀ Zn	⁺³ ₃₁ Ga	⁺² ₃₂ Ge	⁺³ ₃₃ As	⁻³ ₃₄ Se	⁺⁴ ₃₅ Br	⁻¹ ₃₆ Kr		
5	⁺¹ ₃₇ Rb	⁺² ₃₈ Sr	⁺³ ₃₉ Y	⁺⁴ ₄₀ Zr	⁺³ ₄₁ Nb	⁺² ₄₂ Mo	⁺³ ₄₃ Tc	⁺⁴ ₄₄ Ru	⁺³ ₄₅ Rh	⁺² ₄₆ Pd	⁺¹ ₄₇ Ag	⁺² ₄₈ Cd	⁺³ ₄₉ In	⁺² ₅₀ Sn	⁺³ ₅₁ Sb	⁻³ ₅₂ Te	⁺⁴ ₅₃ I	⁻¹ ₅₄ Xe		
6	⁺¹ ₅₅ Cs	⁺² ₅₆ Ba	◆	⁺⁴ ₇₂ Hf	⁺⁵ ₇₃ Ta	⁺⁶ ₇₄ W	⁺⁴ ₇₅ Re	⁺³ ₇₆ Os	⁺² ₇₇ Ir	⁺¹ ₇₈ Pt	⁺² ₇₉ Au	⁺¹ ₈₀ Hg	⁺³ ₈₁ Tl	⁺² ₈₂ Pb	⁺³ ₈₃ Bi	⁻² ₈₄ Po	⁺² ₈₅ At	⁰ ₈₆ Rn		
	TRANSITION ELEMENTS (BETWEEN GROUPS 2A AND 3A)																			
7	⁺¹ ₈₇ Fr	⁺² ₈₈ Ra	★	—	—															

◆ LANTHANIDES (RARE EARTHS)	⁺³ ₅₇ La	⁺³ ₅₈ Ce	⁺³ ₅₉ Pr	⁺³ ₆₀ Nd	⁺³ ₆₁ Pm	⁺² ₆₂ Sm	⁺³ ₆₃ Eu	⁺³ ₆₄ Gd	⁺³ ₆₅ Tb	⁺³ ₆₆ Dy	⁺³ ₆₇ Ho	⁺³ ₆₈ Er	⁺³ ₆₉ Tm	⁺² ₇₀ Yb	⁺³ ₇₁ Lu
★ ACTINIDES	⁺³ ₈₉ Ac	⁺⁴ ₉₀ Th	⁺⁵ ₉₁ Pa	⁺³ ₉₂ U	⁺³ ₉₃ Np	⁺⁴ ₉₄ Pu	⁺³ ₉₅ Am	⁺³ ₉₆ Cm	⁺³ ₉₇ Bk	⁺³ ₉₈ Cf	⁺³ ₉₉ Es	⁺³ ₁₀₀ Fm	⁺³ ₁₀₁ Md	⁺² ₁₀₂ No	⁺³ ₁₀₃ Lw

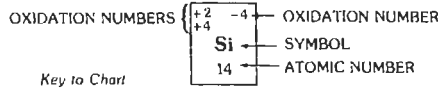


Fig. 1. Periodic classification of the elements.

INSULATING MATERIALS

The permittivity, ϵ , of an insulating material is defined by

$$\epsilon = \epsilon_0 + (P/E)$$

where,

- ϵ_0 = permittivity of free space,
- P = flux density from dipoles within the dielectric medium,
- E = electric-field intensity,
- P/E = electric susceptibility.

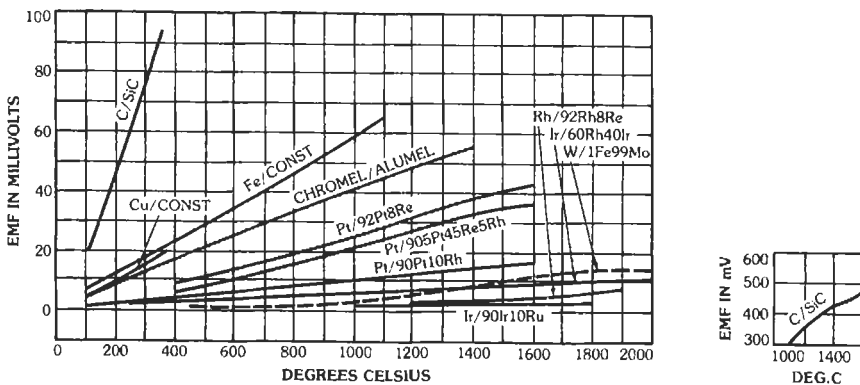


Fig. 2. Temperature-emf characteristics of thermocouples.

CHART I. GALVANIC SERIES IN SEA WATER

Corroded end (anodic)	Lead
	Tin
Magnesium	
Magnesium alloys	Muntz metal
	Manganese bronze
Zinc	Naval brass
Galvanized steel	
Galvanized wrought iron	Nickel (active)
	Inconel (active)
Aluminum:	
52SH, 4S, 3S, 2S, 53ST	
Aluminum clad	Yellow brass
	Admiralty brass
	Aluminum bronze
Cadmium	Red brass
	Copper
Aluminum:	Silicon bronze
A17ST, 17ST, 24ST	Ambrac
	70-30 copper-nickel
Mild steel	Comp. G, bronze
Wrought iron	Comp. M, bronze
Cast iron	
Ni-resist	Nickel (passive)
	Inconel (passive)
13% chromium stainless steel (type 410—active)	Monel
50-50 lead-tin solder	18-8 stainless steel type 304 (passive)
18-8 stainless steel type 304 (active)	18-8-3 stainless steel type 316 (passive)
18-8-3 stainless steel type 316 (active)	Protected end (cathodic or most noble)

It is often convenient to use the relative permittivity, ϵ_r , defined by

$$\epsilon_r = \epsilon / \epsilon_0$$

The relative permittivity, ϵ_r , is a function of temperature and frequency. From Table 9, which gives the values of ϵ_r as a function of frequencies at room temperature, it is easy to get values of ϵ from

$$\epsilon = \epsilon_r / \epsilon_0$$

In the international system of units, the permittivity of vacuum is equal to

$$\epsilon_0 = 8.854 \times 10^{-12} \text{ F/m}$$

and we have

Coulomb's Law

$$F = (1/4\pi\epsilon_0\epsilon_r)(q_1q_2/R^2)$$

Gauss's Law

$$\Phi = (\epsilon_0\epsilon_r)^{-1}\sum q_i$$

The dissipation factor of an insulating material (Table 9) is defined as the ratio of the energy dissipated to the energy stored in the dielectric per hertz, or as the tangent of the loss angle. For dissipation factors less than 0.1, the dissipation factor may be considered equal to the power factor of the dielectric, which is the cosine of the phase angle by which the current leads the voltage.

Many of the materials listed are characterized by a peak dissipation factor that occurs somewhere in the frequency range, this peak being accompanied by a rapid change in the permittivity. These effects are the result of a resonance phenomenon occurring in polar materials. The position of the dissipation-factor peak in the frequency spectrum is very sensitive to temperature. An increase in the temperature increases the frequency at which the peak occurs, as illustrated qualitatively in Fig. 3. Nonpolar materials have very low losses without a noticeable peak; the permittivity remains essentially unchanged over the frequency range.

Another effect that contributes to dielectric losses is that of ionic or electronic conduction. This loss, if present, is important usually at the lower end of the frequency range only and is distinguished by the fact that the dissipation factor varies inversely with frequency.

TABLE 3. THERMOCOUPLES AND THEIR CHARACTERISTICS

	Copper/Constantan		Iron/Constantan		Chromel/Constantan		Chromel/Alumel		Platinum/Platinum Rhodium (10)		Platinum/Platinum Rhodium (13)		Carbon/Silicon Carbide	
Composition, percent	100 Cu/60 Cu 40 Ni		100 Fe/60 Cu 40 Ni		90 Ni 10 Cr/ 55 Cu 45 Ni		90 Ni 10 Cr/ 94 Ni 2 Al 3 Mn 1 Si		Pt/90 Pt 10 Rh		Pt/87 Pt 13 Rh		C/SiC	
*Range of application, °C	-200 to + 300		-200 to + 1100		0 to + 1100		-200 to + 1200		0 to + 1450		0 to + 1450		0 to + 2000	
Resistivity, $\mu\Omega\cdot\text{cm}$	1.75	49	10	49	70	49	70	29.4	10	21				
Temperature coefficient of resistivity, per °C	0.0039	0.00001	0.005	0.00001	0.00035	0.0002	0.00035	0.000125	0.0030	0.0018				
Melting temperature, °C	1085	1190	1535	1190	1400	1190	1400	1430	1755	1700			3000	2700
emf in millivolts; reference junction at 0°C	100 °C	4.24 mV	100 °C	5.28 mV	100 °C	6.3 mV	100 °C	4.1 mV	100 °C	0.643 mV	100 °C	0.646 mV	1210 °C	353.6 mV
	200	9.06	200	10.78	200	13.3	200	8.13	200	1.436	200	1.464	1300	385.2
	300	14.42	400	21.82	400	28.5	400	16.39	400	3.251	400	3.398	1360	403.2
			600	33.16	600	44.3	600	24.90	600	5.222	600	5.561	1450	424.9
			800	45.48			800	33.31	800	7.330	800	7.927		
			1000	58.16			1000	41.31	1000	9.569	1000	10.470		
							1200	48.85	1200	11.924	1200	13.181		
							1400	55.81	1400	14.312	1400	15.940		
									1600	16.674	1600	18.680		
Influence of temperature and gas atmosphere	Subject to oxidation and alteration above 400 °C due to Cu, above 600° due to constantan wire. Ni-plating of Cu tube gives protection in acid-containing gas. Contamination of Cu affects calibration greatly. Resistance to oxide, atm. good. Resistance to reducing atm. good. Requires protection from acid fumes.		Oxidizing and reducing atmosphere have little effect on accuracy. Best used in dry atmosphere. Resistance to oxidation good to 400 °C. Resistance to reducing atmosphere good. Protect from oxygen, moisture, sulphur.		Chromel attacked by sulphurous atmosphere. Resistance to oxidation good. Resistance to reducing atmosphere poor.		Resistance to oxidizing atmosphere very good. Resistance to reducing atmosphere poor. Affected by sulphur, reducing or sulphurous gas, SO ₂ and H ₂ S.		Resistance to oxidizing atmosphere very good. Resistance to reducing atmosphere poor. Susceptible to chemical alteration by As, Si, P vapor in reducing gas (CO ₂ , H ₂ , H ₂ S, SO ₂). Pt corrodes easily above 1000 °C. Used in gas-tight protecting tube.		Resistance to oxidizing atmosphere very good. Resistance to reducing atmosphere poor. Susceptible to chemical alteration by As, Si, P vapor in reducing gas (CO ₂ , H ₂ , H ₂ S, SO ₂). Pt corrodes easily above 1000 °C. Used in gas-tight protecting tube.		Used as tube element. Carbon sheath chemically inert.	
Particular applications	Low temperature, industrial. Internal-combustion engine. Used as a tube element for measurements in steam line.		Low temperature, industrial. Steel annealing, boiler flues, tube stills. Used in reducing or neutral atmosphere.				Used in oxidizing atmosphere. Industrial. Ceramic kilns, tube stills, electric furnaces.		International Standard 630 to 1065 °C.		Similar to Pt/PtRh (10) but has higher emf.		Steel furnace and ladle temperatures. Laboratory measurements.	

* For prolonged use; can be used at higher temperature for short periods.

TABLE 4. THERMAL ELECTROMOTIVE FORCE OF PLATINUM-RHODIUM ALLOYS VERSUS PLATINUM*

Temp.(°C)	Electromotive Force (mV)							
	Percent Rhodium							
	0.5	1.0	5.0	10.0	20.0	40.0	80.0	100.0
0	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
100	+0.10	+0.18	+0.54	+0.64	+0.63	+0.65	+0.62	+0.70
200	0.20	0.37	1.16	1.43	1.44	1.52	1.49	1.61
300	0.29	0.57	1.82	2.32	2.40	2.55	2.55	2.68
400	0.39	0.76	2.49	3.25	3.47	3.70	3.77	3.91
500	0.48	0.94	3.17	4.22	4.63	4.97	5.12	5.28
600	0.58	1.12	3.86	5.22	5.87	6.36	6.60	6.77
700	0.67	1.30	4.55	6.26	7.20	7.85	8.20	8.40
800	0.76	1.48	5.25	7.33	8.59	9.45	9.92	10.16
900	0.85	1.66	5.96	8.43	10.06	11.16	11.76	12.04
1000	0.94	1.84	6.68	9.57	11.58	12.98	13.73	14.05
1100	1.03	2.02	7.42	10.74	13.17	14.90	15.81	16.18
1200	1.13	2.20	8.16	11.93	14.84	16.91	17.99	18.42

* From *Smithsonian Physical Tables*, 9th revised ed., Vol. 120 (Washington, DC: Smithsonian Institution, 1969).

TABLE 5. RESISTIVITIES OF METALS AND ALLOYS

Material	Form	Resistivity ($\times 10^{-6} \Omega \cdot \text{cm}$)	Temperature (°C)	Temperature Coefficient (°C ⁻¹)
Alumel	solid	33.3	0	0.0012
Aluminum	liquid	20.3	670	
	solid	2.62	20	0.0039
Antimony	liquid	123	800	
	solid	39.2	20	0.0036
Arsenic	solid	35	0	0.0042
Beryllium		4.57	20	
Bismuth	liquid	128.9	300	
	solid	115	20	0.004
Boron		1.8×10^{12}	0	
Brass (66 Cu 34 Zn)		3.9	20	0.002
Cadmium	liquid	34	400	
	solid	7.5	20	0.0038
Carbon	diamond	5×10^{20}	15	
	graphite	1400	20	-0.0005
Cerium		78	20	
Cesium	liquid	36.6	30	
	solid	20	20	
		18.83	0	
Chromax (15 Cr, 35 Ni, balance Fe)		100	20	0.00031
Chromel	solid	70-110	0	0.00011-0.000054
Chromium		2.6	0	
Cobalt		9.7	20	0.0033
Constantan (55 Cu, 45 Ni)		44.2	20	+0.0002
Copper (commercial annealed)	liquid	21.3	1083	
	solid	1.7241	20	0.0039
Gallium	liquid	27	30	
	solid	53	0	
Gold	liquid	30.8	1063	
	solid	2.44	20	0.0034
		2.19	0	
Hafnium		32.1	20	
Indium	liquid	29	157	
	solid	9	20	0.00498
Iridium		5.3	20	0.0039

Continued on next page

TABLE 5 (CONT). RESISTIVITIES OF METALS AND ALLOYS

Material	Form	Resistivity ($\times 10^{-6} \Omega \cdot \text{cm}$)	Temperature ($^{\circ}\text{C}$)	Temperature Coefficient ($^{\circ}\text{C}^{-1}$)
Iron		9.71	20	0.0052–0.0062
Kovar A (29 Ni, 17 Co, 0.3 Mn, balance Fe)		45–84	20	
Lead	liquid	98	400	
	solid	21.9	20	0.004
PbO ₂		92		
Lithium	liquid	45	230	0.003
	solid	9.3	20	0.005
Magnesium		4.46	20	0.004
Manganese		5	20	
MnO ₂		6 000 000	20	
Manganin (84 Cu, 12 Mn, 4 Ni)		44	20	± 0.0002
Mercury	liquid	95.8	20	0.00089
	solid	21.3	-50	
Molybdenum		5.17	0	
		4.77	20	0.0033
Monel metal (67 Ni, 30 Cu, 1.4 Fe, 1 Mn)	solid	42	20	0.002
Neodymium	solid	79	18	
Nichrome (65 Ni, 12 Cr, 23 Fe)	solid	100	20	0.00017
Nickel	solid	6.9	20	0.0047
Nickel–silver (64 Cu, 18 Zn, 18 Ni)	solid	28	20	0.00026
Niobium		12.4	20	
Osmium		9	20	0.0042
Palladium		10.8	20	0.0033
Phosphor bronze (4 Sn, 0.5 P, balance Cu)		9.4	20	0.003
Platinum		10.5	20	0.003
Plutonium		150	20	
Potassium	liquid	13	62	
	solid	7	20	0.006
Praseodymium		68	25	
Rhenium		19.8	20	
Rhodium		5.1	20	0.0046
Rubidium		12.5	20	
Ruthenium		10	20	
Selenium	solid	1.2	20	
Silver		1.62	20	0.0038
Sodium	liquid	9.7	100	
	solid	4.6	20	
Steel (0.4–0.5 C, balance Fe)		13–22	20	0.003
Steel, manganese (13 Mn, 1 C, 86 Fe)		70	20	0.001
Steel, stainless (0.1 C, 18 Cr, 8 Ni, balance Fe)		90	20	
Strontium		23	20	
Sulfur		2×10^{23}	20	
Tantalum		13.1	20	0.003
Thallium		18.1	20	0.004
Thorium		18	20	0.0021
Tin		11.4	20	0.0042
Titanium		47.8	25	
Tophet A (80 Ni, 20 Cr)		108	20	0.00014
Tungsten		5.48	20	0.0045
W ₂ O ₅		450	20	
WO ₃		2×10^{11}	20	
Uranium		29	0	0.0021
Zinc	liquid	35.3	420	
	solid	6	20	0.0037
Zirconium		40	20	0.0044

TABLE 6. ELECTRICAL RESISTIVITY OF ROCKS AND SOILS*

	Resistivity ($\Omega \cdot \text{cm}$)
Igneous Rocks	
Granite	$10^7\text{--}10^9$
Lava flow (basic)	$10^6\text{--}10^7$
Lava, fresh	$3 \times 10^5\text{--}10^6$
Quartz vein, massive	$> 10^6$
Metamorphic Rocks	
Marble	4×10^8
Marble, white	10^{10}
Marble, yellow	10^{10}
Schist, mica	10^7
Shale, bed	10^5
Shale, Nonesuch	10^4
Sedimentary Rocks	
Limestone	10^4
Limestone, Cambrian	$10^4\text{--}10^5$
Sandstone	10^5
Sandstone, eastern	$3 \times 10^3\text{--}10^4$
Unconsolidated Materials	
Clay, blue	2×10^4
Clay, fire	2×10^5
Clayey earth	$10^4\text{--}4 \times 10^4$
Gravel	10^5
Sand, dry	$10^5\text{--}10^6$
Sand, moist	$10^3\text{--}10^6$

* From *Smithsonian Physical Tables*, 9th revised ed., Vol. 120 (Washington, DC: Smithsonian Institution, 1969).

Increase in temperature increases the loss due to ionic conduction because of increased ionic mobility.

The data given on dielectric strength are accompanied by the thickness of the specimen tested because the dielectric strength, expressed in volts/mil, varies inversely with the square root of thickness, approximately.

The direct-current volume resistivity of many materials is influenced by changes in temperature or humidity. The values given in the table may be reduced several decades by raising the temperature toward the higher end of the working range of the material, or by raising the relative humidity of the air surrounding the material to above 90 percent.

MAGNETIC MATERIALS*

All materials that are magnetized by a magnetic field are called magnetic materials. Depending on the magnetic response, various kinds of magnetism are classified. Most magnetic materials of commercial

* This section was contributed by the late Gilbert Y. Chin.

TABLE 7. SUPERCONDUCTIVITY OF SOME METALS, ALLOYS, AND COMPOUNDS*

Material	Critical Temperature (K)
NbC	10.1
Niobium	9.22
TaC	9.2
Pb-As-Bi	9.0
Pb-Bi-Sb	8.9
Pb-Sn-Bi	8.5
Pb-As	8.4
MoC	7.7
Lead	7.2
N_2Pb_3	7.2
Bi_6Tl_3	6.5
Sb_2Tl_7	5.5
Lanthanum	5.2
Tantalum	4.4
Vanadium	4.3
TaSi	4.2
Mercury	4.15
PbS	4.1
Hg_5Tl_7	3.8
Tin	3.71
Indium	3.38
ZrB	2.82
WC	2.8
Rhenium	2.57
Mo_2C	2.4
Thallium	2.4
W_2C	2.05
Au_2Bi	1.84
CuS	1.6
TiN	1.4
Thorium	1.32
VN	1.3
Aluminum	1.15
Gallium	1.12
TiC	1.1
Zinc	0.95
Uranium	0.75
Osmium	0.71
Zirconium	0.54
Cadmium	0.54
Titanium	0.53
Ruthenium	0.47
Hafnium	0.35

* Information on recently discovered ceramic superconductors is posted at www.ceramics.nist.gov/srd and at www.iastate.edu.

importance today are ferromagnets or ferrimagnets. Iron, nickel, cobalt, and their alloys are examples of ferromagnets; spinel ferrites and garnets are examples of ferrimagnets. These exhibit spontaneous magnetism and develop a flux density B upon application of a magnetic field of strength H in accordance with

$$B = \mu_0 H + J$$

In this equation, $\mu_0 = 1.257 \times 10^{-6}$ H/m is the permeability of free space, and J is the magnetic polarization. The units of B and J are teslas (T) or webers

TABLE 8. PROPERTIES OF SEMICONDUCTORS

Units	IV					III-V								II-VI					
	Si	Ge	AlP*	AlAs	AlSb	GaN	GaP	GaAs _{0.5} P _{0.5}	GaAs	GaSb	InP	InAs	InSb	ZnS	ZnSe	ZnTe	CdS	CdSe	CdTe
$E_{g,ind}$ (eV)	1.11	0.651	2.451	2.141	1.621	3.191	2.261	-1.911	1.431	0.701	1.351	0.3561	0.181	3.661	2.671	2.251	2.421	1.711	1.501
$E_{g,direct}$ (eV)	1.117	1.907	0.506	0.579	0.765	0.366	0.548	0.642	0.867	1.77	0.918	3.48	6.89	0.339	0.464	0.551	0.512	0.712	0.861
$E_{g,ind} - E_{g,direct}$ (eV)	-1.91	0.811	3.61	3.141	2.221	2.811	-2.031	1.871	0.781	2.251	1.821	-0.81	0.235	3.80	2.80	2.38	2.58	1.84	1.60
$\Delta E_{spin-orbit}$ (eV)	0.038	0.29		0.29	0.75		0.08		0.34	0.80	0.11	0.41	0.82	0.07	0.45	0.91			0.31
a_0 (Å)	5.431	5.658	5.463	5.661	6.138	3.180	5.449	5.572	5.654	6.095	5.868	6.058	6.479	5.409	5.669	6.104	4.137	4.299	6.481
d (g/cm ³)	2.328	5.323	2.40	3.73	4.26	5.166	6.10	4.129	4.88	5.316	5.613	4.787	5.667	5.775	4.09	5.26	5.64	6.716	7.018
Therm. Exp. (10 ⁶ /°C)	2.56	5.92		5.20	4.88	15.59	5.8	5.9	6.8	6.7	4.5	5.19	5.04	6.2	7	8	14		5.5
Therm. Cond. (W/cm°C)	1.41	0.61	0.9	0.8	0.54		0.97	0.16	0.54	0.35	0.68	0.26	0.18	0.26	0.13	0.112	0.2	0.063	0.07
Hardness (Knoop)	1150	780	500	481	360		945		750	450	535	381	223	178	150	130	55	44	100
Melting Pt. (°C)	1415	958	>2000	1740	1080	1500	1467		1238	712	1070	943	525	1830	1520	1295	1175	1239	1098
Molecular Weight	28.09	72.60	57.95	101.90	148.74	83.73	100.70	127.06	144.64	191.48	145.79	189.74	236.58	97.45	144.34	192.99	144.48	191.37	240.02
m_e (m ₀)	1.19	1.082		0.35	0.39	0.19	0.35	0.089	0.065	0.049	0.078	0.023	0.014	0.34	0.17	0.09	0.20	0.13	0.11
m_{hh} (m ₀)	0.90	1.64																	
m_{hl} (m ₀)	0.15	0.042			0.11		0.14		0.087	0.056		0.025	0.016			0.15	0.7	0.45	0.13
μ_{300K} (cm ² /V-s)	1900	3600	80	280	90	150	190	2900	8800	6000	4700	22600	8.2-10 ⁴	140	600	340	350	650	1050
μ_{300K} (cm ² /V-s)	19000	40000	30	18	700	2700	17000	2.1-10 ⁵	10000	60000	60000	1.2-10 ⁵	1.2-10 ⁶	7000		5000	5000	15000	
μ_{300K} (cm ² /V-s)	425	2300		400	120	400	120	400	800	150	260	1700		5	28	110	45		80
μ_{300K} (cm ² /V-s)	8000	40000		4000		2000		7500	6000	1200	350	7000			900				500
Phonon LOG(TG) (meV)	51.0	28.1	62.0	49.8	42.1	49.9		36.2	29.8	42.8	30.2	24.2	36.9	31	25.3	16.8	26.2	21.2	
**Ref. Index	57.4	33.3	54.6	35.1	39.5	45.5		33.3	26.8	37.7	27.1	22.6	28.4	26	22.3	32.1	20.6	17.4	
b_s	3.5	4.1		3.1	3.4	2.1	3.3	3.4	3.4	3.8	3.1	3.5	3.9	2.4	2.8	3.1	2.5	2.6	2.8
b_c	11.94	16.0		8.16	10.24	4	9.04	-10.2	10.9	14.44	9.52	11.8	15.7	5.07	6.1	7.3	15.3	16.2	6.7
F_0	11.94	16.0		10.06	14.4		11.1	-12.4	13.18	15.69	12.35	14.55	17.72	8.3	9.2	10.4	18.9	19.3	9.4

All values at 300 K unless otherwise noted; compiled by D. E. Hill

Copyright 1971 Monsanto Company. Used by permission.

* Highest experimental Hall mobilities.

† \perp // to c axis.** Value near Eq on long λ side.† m_h/m_e .

D = "Direct" energy gap.

I = "Indirect" energy gap.

Consult Semiconductor Equipment and Materials International (SEMI) at <http://www.semi.org> for data on new semiconductor materials.

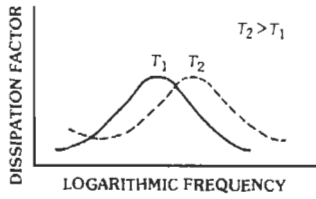


Fig. 3. Variation of dissipation factor with frequency.

per square meter (Wb/m^2), and H is in amperes per meter (A/m).

The permeability of a magnetic material is

$$\mu = B/H = \mu_0 + J/H$$

It is, however, customary to use the relative permeability, μ_r , defined by

$$\mu_r = \mu/\mu_0 = 1 + J/\mu_0 H$$

The permeability of a ferromagnet (or ferrimagnet) is a function of applied field, temperature, and frequency. For electronic transformer and inductor applications where a small field is impressed, the initial relative permeability, μ_i (often measured at $B = 4 \text{ mT}$), is the most useful quantity. For power transformer applications, values at higher fields are more useful. Sometimes the value of the maximum relative permeability, μ_m , is quoted. This is a useful quantity for materials exhibiting a square hysteresis loop ideal for amplifier-type applications. Magnetic materials are classified as soft or hard (permanent) depending on the value of coercivity H_c , which is the field strength required to reduce the flux density to zero after the material has been magnetized. Soft magnets have values of H_c less than about 1 kA/m , and hard magnets have H_c greater than about 10 kA/m . Some applications such as reed contacts make use of semihard magnets with $1 \text{ kA/m} < H_c < 10 \text{ kA/m}$.

Soft Magnetic Metals

Table 10 lists some typical commercial soft magnetic metals generally used in low-frequency transformers and inductors. They are often available as laminations, cut cores, and tape-wound cores. The nickel-iron alloys known as permalloys exhibit the highest initial relative permeability and lowest coercivity. The most commonly used materials in this category are the 48-percent nickel alloy with μ_i about 11 000 and the 80-percent nickel-5-percent molybdenum alloy with μ_i about 80 000, the latter being more expensive. The nickel-iron alloys can also be processed to exhibit a square hysteresis loop, useful for amplifier-type applications, or a skewed loop, useful for unipolar pulse transformer designs. For common 50 Hz or 60 Hz transformers used in the electrical utility industry, silicon steel is the most common material. Low-carbon

steels are very inexpensive and widely used in small motors and generators. The cobalt-iron alloys, usually known as permendur, have the highest values of saturation polarization and Curie temperature and tend to be used in high-performance, lightweight applications such as airborne motors. Appearing on the market recently is a new class of soft magnetic alloys called amorphous magnets or metallic glasses. These materials are rapidly solidified from the melt into thin tapes such that the usual crystalline structure is absent. Some amorphous magnets have been prepared to have properties similar to the silicon steels but with substantially lower core losses, while others have been prepared with properties similar to the best grades of the permalloys.

Permanent-Magnet Materials

Table 11 lists some typical permanent-magnet materials. The alnicos are brittle and hence can be used only in cast or sintered form. The hexagonal ferrites are oxides having the general formula $\text{MO} \cdot 6\text{Fe}_2\text{O}_3$, where M is barium or strontium. The magnets are prepared by ceramic techniques and are often called ceramic magnets. Large numbers of low-cost ferrite magnets are bonded in plastics and widely used in door catches, wall magnets, refrigerator door gaskets, and toys.

In recent years, a new class of high-performance permanent magnets made of cobalt rare earths, particularly those containing samarium, have been developed commercially. These magnets have the highest combination of coercivity and maximum energy product available on the market, with applications ranging from tiny wristwatch and earphone magnets to horse-power-size industrial dc and synchronous motors and generators. Some are available in plastic bonded form.

An emerging class of permanent-magnet alloys is the iron chromes, containing iron, chromium, and cobalt. These have magnetic properties similar to the alnicos, but they are sufficiently ductile that they can be cold-rolled to thin strips, drawn into fine wires, or machined or punched into intricate shapes. They also contain less cobalt than the alnicos.

Ferrites

“Ferrite” is the common term applied to a wide range of different ceramic ferromagnetic materials. Specifically, the term applies to those materials with the spinel crystal structures having the general formula XFe_2O_4 , where X is any divalent metallic ion having the proper ionic radius to fit in the spinel structure. Several ceramic ferromagnetic materials have been prepared that deviate stoichiometrically from the basic formula XFe_2O_4 , but common usage has included them in the family of ferrite materials.

The behavior of the conductivity and permittivity of ferrites can be understood by considering them as grains (crystals) of fairly low-resistance material separated by thin layers of a relatively poor conductor. Therefore, the permittivity and conductivity show a

TABLE 9. CHARACTERISTICS

Material Composition	T(°C)	Relative Permittivity at					
		(Frequency in hertz)					
		60	10 ³	10 ⁶	10 ⁸	3 × 10 ⁹	2.5 × 10 ¹⁰
Ceramics:							
Aluminum oxide	25	—	8.83	8.80	8.80	8.79	—
Barium titanate [‡]	26	1250	1200	1143	—	600	100
Calcium titanate	25	168	167.7	167.7	167.7	165	—
Magnesium oxide	25	—	9.65	9.65	9.65	—	—
Magnesium silicate	25	6.00	5.98	5.97	5.96	5.90	—
Magnesium titanate	25	—	13.9	13.9	13.9	13.8	13.7
Oxides of aluminum, silicon, magnesium, calcium, barium	24	—	6.04	6.04	—	5.90	—
Porcelain (dry process)	25	5.5	5.36	5.08	5.04	—	—
Steatite 410	25	5.77	5.77	5.77	5.77	5.7	—
Strontium titanate	25	—	233	232	232	—	—
Titanium dioxide (rutile)	26	—	100	100	100	—	—
Glasses:							
Iron-sealing glass	24	8.41	8.38	8.30	8.20	7.99	7.84
Soda-borosilicate	25	—	4.97	4.84	4.84	4.82	4.65
100% silicon dioxide (fused quartz)	25	3.78	3.78	3.78	3.78	3.78	3.78
Plastics:							
Alkyd resin	25	—	5.10	4.76	4.55	4.50	—
Cellulose acetate-butylate, plasticized	26	3.60	3.48	3.30	3.08	2.91	—
Cresylic acid-formaldehyde, 50% α-cellulose	25	5.45	4.95	4.51	3.85	3.43	3.21
Cross-linked polystyrene	25	2.59	2.59	2.58	2.58	2.58	—
Epoxy resin (Araldite CN-501)	25	—	3.67	3.62	3.35	3.09	—
Epoxy resin (Epon resin RN-48)	25	—	3.63	3.52	3.32	3.04	—
Foamed polystyrene, 0.25% filler	25	1.03	1.03	1.03	—	1.03	1.03
McLamine-formaldehyde, α-cellulose	24	—	7.57	7.00	6.0	4.93	—
McLamine-formaldehyde, 55% filler	26	—	6.00	5.75	5.5	—	—
Phenol-formaldehyde (Bakelite BM 120)	25	4.90	4.74	4.36	3.95	3.70	3.55
Phenol-formaldehyde, 50% paper laminate	26	5.25	5.15	4.60	4.04	3.57	—
Phenol-formaldehyde, 65% mica, 4% lubricants	24	5.1	5.03	4.78	4.72	4.71	—
Polycarbonate	—	3.17	3.02	2.96	—	—	—
Polychlorotrifluoroethylene	25	2.72	2.63	2.42	2.32	2.29	2.28
Polyethylene	25	2.26	2.26	2.26	2.26	2.26	2.26
Polyethylene-terephthalate	—	3.16	3.12	2.98	—	—	—
Polyethylmethacrylate	22	—	2.75	2.55	2.52	2.51	2.5
Polyhexamethylene-adipamide (nylon)	25	3.7	3.50	3.14	3.0	2.84	2.73
Polyimide	—	—	3.5	3.4	—	—	—
Polyisobutylene	25	2.23	2.23	2.23	2.23	2.23	—
Polymer of 95% vinyl-chloride, 5% vinyl-acetate	20	—	3.15	2.90	2.8	2.74	—
Polymethyl methacrylate	27	3.45	3.12	2.76	—	2.60	—
Polyphenylene oxide	—	2.55	2.55	2.55	—	2.55	—
Polypropylene	—	2.25	2.25	2.55	—	—	—
Polystyrene	25	2.56	2.56	2.56	2.55	2.55	2.54
Polytetrafluoroethylene (teflon)	22	2.1	2.1	2.1	2.1	2.1	2.08
Polyvinylcyclohexane	24	—	2.25	2.25	2.25	2.25	—
Polyvinyl formal	26	3.20	3.12	2.92	2.80	2.76	2.7
Polyvinylidene fluoride	—	8.4	8.0	6.6	—	—	—
Urea-formaldehyde, cellulose	27	6.6	6.2	5.65	5.1	4.57	—
Urethane elastomer	—	6.7-7.5	6.7-7.5	6.5-7.1	—	—	—
Vinylidene-vinyl chloride copolymer	23	5.0	4.65	3.18	2.82	2.71	—

OF INSULATING MATERIALS*†

		Dissipation Factor at				Dielectric Strength in Volts/Mil at 25 °C	DC Volume Resistivity in Ω·cm at 25 °C	Thermal Expansion (Linear) in Parts/°C	Softening Point in °C	Moisture Absorption in Percent
60	10 ³	10 ⁵	10 ⁸	3×10 ⁹	2.5×10 ¹⁰					
—	0.00057	0.00033	0.00030	0.0010	—	—	—	—	—	
0.056	0.0130	0.0105	—	0.30	0.60	75	10 ¹² –10 ¹³	—	1400–1430	0.1
0.006	0.00044	0.0002	—	0.0023	—	100	10 ¹² –10 ¹⁴	—	1510	<0.1
—	<0.0003	<0.0003	<0.0003	—	—	—	—	—	—	—
0.012	0.0034	0.0005	0.0004	0.0012	—	—	>10 ¹⁴	9.2×10 ⁻⁶	1350	0.1–1
—	0.0011	0.0004	0.0005	0.0017	0.0065	—	—	—	—	—
—	0.0019	0.0011	—	0.0024	—	—	—	7.7×10 ⁻⁵	1325	—
0.03	0.0140	0.0075	0.0078	—	—	—	—	—	—	—
—	0.0030	0.0007	0.0006	0.00089	—	—	—	—	—	—
—	0.0011	0.0002	0.0001	—	—	100	10 ¹² –10 ¹⁴	—	1510	0.1
—	0.0015	0.0003	0.00025	—	—	—	—	—	—	—
—	0.0004	0.0005	0.0009	0.00199	0.0112	—	10 ¹⁰ at250°	132×10 ⁻⁷	484	Poor
—	0.0055	0.0036	0.0030	0.0054	0.0090	—	7×10 ⁷ at250°	50×10 ⁻⁷	693	—
0.0009	0.00075	0.0001	0.0002	0.00006	0.00025	410 (0.25")	>10 ¹⁹	5.7×10 ⁻⁷	1667	—
—	0.0236	0.0149	0.0138	0.0108	—	—	—	—	—	—
0.0045	0.0097	0.018	0.017	0.028	—	250–400 (0.125")	—	11–17×10 ⁻⁵	60–121	2.3
0.098	0.033	0.036	0.055	0.051	0.038	1020 (0.033")	3×10 ¹³	3×10 ⁻⁵	>125	1.2
0.0004	0.0005	0.0016	0.0020	0.0019	—	—	—	—	—	—
—	0.0024	0.019	0.034	0.027	—	405 (0.125")	>3.8×10 ⁷	4.77×10 ⁻⁵	109 (distortion)	0.14
—	0.0038	0.0142	0.0264	0.021	—	—	—	—	—	—
<0.0002	<0.0001	<0.0002	—	0.0001	—	—	—	—	85	low
—	0.0122	0.041	0.085	0.103	—	300–400	—	—	99 (stable)	0.4–0.6
—	0.0119	0.0115	0.020	—	—	—	—	1.7×10 ⁻⁵	—	0.6
0.08	0.0220	0.0280	0.0380	0.0438	0.0390	300 (0.125")	10 ¹¹	30–40×10 ⁻⁶	<135 (distortion)	<0.6
0.025	0.0165	0.034	0.057	0.060	—	—	—	—	—	—
0.015	0.0104	0.0082	0.0115	0.0126	—	—	—	—	—	—
0.009	0.0021	0.010	—	—	—	364 (0.125")	2×10 ¹⁶	7×10 ⁻⁵	135 (deflection)	—
0.015	0.0270	0.0082	—	0.0028	0.0053	—	10 ¹⁸	—	—	—
<0.0002	<0.0002	<0.0002	0.0002	0.00031	0.0006	1200 (0.033")	10 ¹⁷	19×10 ⁻⁵ (varies)	95–105 (distortion)	0.03
0.0021	0.0047	0.016	—	—	—	4000 (0.002")	—	—	—	—
—	0.0294	0.0090	—	0.0075	0.0083	—	—	—	60 (distortion)	low
0.018	0.0186	0.0218	0.0200	0.0117	0.0105	400 (0.125")	8×10 ¹⁴	10.3×10 ⁻⁵	65 (distortion)	1.5
—	0.002	0.003	—	—	—	570	—	—	—	—
0.0004	0.0001	0.0001	0.0003	0.00047	—	600 (0.010")	—	—	25 (distortion)	low
—	0.0165	0.0150	0.0080	0.0059	—	—	—	—	—	—
0.064	0.0465	0.0140	—	0.0057	—	990 (0.030")	>5×10 ¹⁶	8–9×10 ⁻⁵	70–75 (distortion)	0.3–0.6
0.0004	0.0003	0.0007	—	0.0011	—	500 (0.125")	10 ¹⁷	5.3×10 ⁻⁵	195 (deflection)	—
<0.0005	<0.0005	<0.0005	—	—	—	650 (0.125")	6×10 ¹⁶	6–8.5×10 ⁻⁵	99–116 (deflection)	—
<0.00005	<0.00005	0.00007	<0.0001	0.00033	0.0012	500–700 (0.125")	10 ¹⁸	6–8×10 ⁻⁵	82 (distortion)	0.05
<0.0005	<0.0003	<0.0002	<0.0002	0.00015	0.0006	1000–2000 (0.005"–0.012")	10 ¹⁷	9.0×10 ⁻⁵	66 (distortion) (stable to 300)	0.00
—	0.0002	<0.0002	<0.0002	0.00018	—	—	—	—	—	—
0.003	0.0100	0.019	0.013	0.0113	0.0115	860 (0.034")	5×10 ¹⁶	7.7×10 ⁻⁵	190	1.3
0.049	0.018	0.17	—	—	—	260 (0.125")	2×10 ¹⁴	12×10 ⁻⁵	148 (deflection)	—
0.032	0.024	0.027	0.050	0.0555	—	375 (0.085")	—	2.6×10 ⁻⁵	152 (distortion)	2
0.016	0.055	—	—	—	—	450–500 (0.125")	2×10 ¹¹	10–20×10 ⁻⁵	—	—
0.042	0.063	0.057	0.0180	0.0072	—	300 (0.125")	10 ¹⁴ –10 ¹⁶	15.8×10 ⁻⁵	150	<0.1

Continued on next page

TABLE 9 (CONT). CHARACTERISTICS

Material Composition	T(°C)	Relative Permittivity at					
		(Frequency in hertz)					
		60	10 ³	10 ⁶	10 ⁸	3 × 10 ⁹	2.5 × 10 ¹⁰
Plastics (cont.):							
100% aniline-formaldehyde (Dilectene-100)	25	3.70	3.68	3.58	3.50	3.44	—
100% phenol-formaldehyde	24	8.6	7.15	5.4	4.4	3.64	—
100% polyvinyl-chloride	20	3.20	3.10	2.88	2.85	2.84	—
Organic Liquids:							
Aviation gasoline (100 octane)	25	—	—	1.94	1.94	1.92	—
Benzene (pure, dried)	25	2.28	2.28	2.28	2.28	2.28	2.28
Carbon tetrachloride	25	2.17	2.17	2.17	2.17	2.17	—
Ethyl alcohol (absolute)	25	—	—	24.5	23.7	6.5	—
Ethylene glycol	25	—	—	41	41	12	—
Jet fuel (JP-3)	25	—	—	2.08	2.08	2.04	—
Methyl alcohol (absolute analytical grade)	25	—	—	31	31.0	23.9	—
Methyl or ethyl siloxane polymer (1000 cs)	22	2.78	2.78	2.78	—	2.74	—
Monomeric styrene	22	2.40	2.40	2.40	2.40	2.40	—
Transil oil	26	2.22	2.22	2.22	2.20	2.18	—
Vaseline	25	2.16	2.16	2.16	2.16	2.16	—
Waxes:							
Beeswax, yellow	23	2.76	2.66	2.53	2.45	2.39	—
Dichloronaphthalenes	23	3.14	3.04	2.98	2.93	2.89	—
Polybutene	25	2.34	2.34	2.34	2.30	2.27	—
Vegetable and mineral waxes	25	2.3	2.3	2.3	2.3	2.25	—
Rubbers:							
Butyl rubber	25	2.39	2.38	2.35	2.35	2.35	—
GR-S rubber	25	2.96	2.96	2.90	2.82	2.75	—
Gutta-percha	25	2.61	2.60	2.53	2.47	2.40	—
Hevea rubber (pale crepe)	25	2.4	2.4	2.4	2.4	2.15	—
Hevea rubber, vulcanized (100 pts pale crepe, 6 pts sulfur)	27	2.94	2.94	2.74	2.42	2.36	—
Ncoprene rubber	24	6.7	6.60	6.26	4.5	4.00	4.0
Organic polysulfide, fillers	23	—	22.60	110	30	16	13.6
Silicone-rubber compound	25	—	3.35	3.20	3.16	3.13	—
Woods:§							
Balsa wood	26	1.4	1.4	1.37	1.30	1.22	—
Douglas fir	25	2.05	2.00	1.93	1.88	1.82	1.78
Douglas fir, plywood	25	2.1	2.1	1.90	—	—	1.6
Mahogany	25	2.42	2.40	2.25	2.07	1.88	1.6
Yellow birch	25	2.9	2.88	2.70	2.47	2.13	1.87
Yellow poplar	25	1.85	1.79	1.75	—	1.50	1.4
Miscellaneous:							
Amber (fossil resin)	25	2.7	2.7	2.65	—	2.6	—
DeKhotinsky cement	23	3.95	3.75	3.23	—	2.96	—
Gilsonite (99.9% natural bitumen)	26	2.69	2.66	2.58	2.56	—	—
Shellac (natural XL)	28	3.87	3.81	3.47	3.10	2.86	—
Mica, glass-bonded	25	—	7.45	7.39	—	—	—
Mica, glass, titanium dioxide	24	—	9.3	9.0	—	—	—
Ruby mica	26	5.4	5.4	5.4	5.4	5.4	—
Paper, royal grey	25	3.30	3.29	2.99	2.77	2.70	—
Selenium (amorphous)	25	—	6.00	6.00	6.00	6.00	6.00
Asbestos fiber-chrysotile paper	25	—	4.80	3.1	—	—	—
Sodium chloride (fresh crystals)	25	—	5.90	5.90	—	—	5.90
Soil, sandy dry	25	—	2.91	2.59	2.55	2.55	—
Soil, loamy dry	25	—	2.83	2.53	2.48	2.44	—
Ice (from pure distilled water)	-12	—	—	4.15	3.45	3.20	—
Freshly fallen snow	-20	—	3.33	1.20	1.20	1.20	—
Hard-packed snow followed by light rain	-6	—	—	1.55	—	1.5	—
Water (distilled)	25	—	—	78.2	78	76.7	34

* Mostly taken from *Tables of Dielectric Materials*, Vols. I-IV, prepared by the Laboratory for Insulation Research of the Massachusetts Institute of Technology, Cambridge, Massachusetts, January 1953; from *Dielectric Materials and Applications*, A. R. von Hippel, ed., John Wiley & Sons, New York, NY, 1954; and from *Modern Plastics Encyclopedia*, Joel Frados, ed., 1301 Avenue of the Americas, New York, NY, 1962. Materials listed are typical of a class. Further data should be sought for a particular material of interest.

OF INSULATING MATERIALS*†

Dissipation Factor at (Frequency in hertz)						Dielectric Strength in Volts/Mil at 25 °C	DC Volume Resistivity in Ω·cm at 25 °C	Thermal Expansion (Linear) in Parts/°C	Softening Point in °C	Moisture Absorption in Percent
60	10 ³	10 ⁵	10 ⁸	3×10 ⁹	2.5×10 ¹⁰					
0.0033	0.0032	0.0061	0.0033	0.0026	—	810 (0.068")	10 ¹⁶	5.4×10 ⁻⁵	125	0.06–0.08
0.15	0.082	0.060	0.077	0.052	—	277 (0.125")	—	8.3–13×10 ⁻⁵	50 (distortion)	0.42
0.0115	0.0185	0.0160	0.0081	0.0055	—	400 (0.125")	10 ¹⁴	6.9×10 ⁻⁵	54 (distortion)	0.05–0.15
—	—	—	0.0001	0.0014	—	—	—	—	—	—
<0.0001	<0.0001	<0.0001	<0.0001	<0.0001	<0.0001	—	—	—	—	—
0.007	0.0008	<0.00004	<0.0002	0.0004	—	—	—	—	—	—
—	—	0.090	0.062	0.250	—	—	—	—	—	—
—	—	-0.030	0.045	1.00	—	—	—	—	—	—
—	—	0.0001	—	0.0055	—	—	—	—	—	—
—	—	0.20	0.038	0.64	—	—	—	—	—	—
0.0001	0.00008	<0.0003	—	0.0096	—	—	—	—	—	—
0.01	0.005	<0.0003	—	0.0020	—	300 (0.100")	3×10 ¹²	—	—	0.06
0.001	<0.00001	<0.0005	0.0048	0.0028	—	300 (0.100")	—	—	-40 (pourpoint)	—
0.0004	0.0002	<0.0001	<0.0004	0.00066	—	—	—	—	—	—
—	0.0140	0.0092	0.0090	0.0075	—	—	—	—	45–64 (melts)	—
0.10	0.0110	0.0003	0.0017	0.0037	—	—	—	—	35–63 (melts)	nil
0.0002	0.0003	0.00133	0.00133	0.0009	—	—	—	—	—	—
0.0009	0.0006	0.0004	0.0004	0.00046	—	—	—	—	57	—
0.0034	0.0035	0.0010	0.0010	0.0009	—	—	—	—	—	—
0.0008	0.0024	0.0120	0.0080	0.0057	—	870 (0.040")	2×10 ¹⁵	—	—	—
0.0005	0.0004	0.0042	0.0120	0.0060	—	—	10 ¹⁵	—	—	—
0.0030	0.0018	0.0018	0.0050	0.0030	—	—	—	—	—	—
0.005	0.0024	0.0446	0.0180	0.0047	—	—	—	—	—	—
0.018	0.011	0.038	0.090	0.034	0.025	300 (0.125")	8×10 ¹²	—	—	nil
—	1.29	0.39	0.28	0.22	0.10	—	—	—	—	—
—	0.0067	0.0030	0.0032	0.0097	—	—	—	—	—	—
0.058	0.0040	0.0120	0.0135	0.100	—	—	—	—	—	—
0.004	0.0080	0.026	0.033	0.027	0.032	—	—	—	—	—
0.012	0.0105	0.0230	—	—	0.0220	—	—	—	—	—
0.008	0.0120	0.025	0.032	0.025	0.020	—	—	—	—	—
0.007	0.0090	0.029	0.040	0.033	0.026	—	—	—	—	—
0.004	0.0054	0.019	—	0.015	0.017	—	—	—	—	—
0.001	0.0018	0.0056	—	0.0090	—	2300 (0.125")	Very high	—	200	—
0.049	0.0335	0.024	—	0.021	—	—	—	9.8×10 ⁻⁵	80–85	—
0.006	0.0035	0.0016	0.0011	—	—	—	—	—	155 (melts)	—
0.006	0.0074	0.031	0.030	0.0254	—	—	10 ¹⁶	—	80	low after baking
—	0.0019	0.0013	—	—	—	—	—	—	—	—
—	0.0125	0.0026	—	0.0040	—	—	—	—	400	<0.5
0.005	0.0006	0.0003	0.0002	0.0003	—	3800–5600 (0.040")	5×10 ¹³	—	—	—
0.010	0.0077	0.038	0.066	0.056	—	202 (0.125")	—	—	—	—
—	0.0004	<0.0003	<0.0002	0.00018	0.0013	—	—	—	—	—
—	0.15	0.025	—	—	—	—	—	—	—	—
—	<0.0001	<0.0002	—	—	<0.0005	—	—	—	—	—
—	0.08	0.017	—	0.0062	—	—	—	—	—	—
—	0.05	0.018	—	0.0011	—	—	—	—	—	—
—	—	0.12	0.035	0.0009	—	—	—	—	—	—
—	0.492	0.0215	—	0.00029	—	—	—	—	—	—
—	—	0.29	—	0.0009	—	—	—	—	—	—
—	—	0.040	0.005	0.157	0.2650	—	10 ⁵	—	—	—

† Consult the American Ceramic Society at <http://www.acers.org> for recent advances in dielectric materials and optoelectronic materials for multilayer electronic devices.

‡ Relative permittivity and dissipation factor on electrical field strength.

§ Field perpendicular to grain.

TABLE 10. PROPERTIES OF SOFT MAGNETIC METALS

Name	Composition, %	Relative Permeability		Coercivity H_c (A/m)	Retentivity B_r (T)	B_{max} (T)	Resistivity ($\mu\Omega \cdot \text{cm}$)
		Initial	Maximum				
Ingot iron	99.8 Fe	150	5 000	80	0.77	2.14	10
Low carbon steel	99.5 Fe	200	4 000	100	—	2.14	12
Silicon iron, unoriented	3 Si, bal Fe	270	8 000	60	—	2.01	47
Silicon iron, grain oriented	3 Si, bal Fe	1 400	50 000	7	1.20	2.01	50
4750 alloy	48 Ni, bal Fe	11 000	80 000	2	—	1.55	48
4-79 Permalloy	4 Mo, 79 Ni, bal Fe	40 000	200 000	1	—	0.80	58
Supermalloy	5 Mo, 80 Ni, bal Fe	80 000	450 000	0.4	—	0.78	65
2V-Permendur	2V, 49 Co, bal Fe	800	8 000	160	—	2.30	40
Supermendur	2V, 49 Co, bal Fe	—	100 000	16	2.00	2.30	26
Metglas* 2605SC	$\text{Fe}_{81}\text{B}_{13.5}\text{Si}_{3.5}\text{C}_2$	—	210 000	14	1.46	1.60	125
Metglas* 2605S-3	$\text{Fe}_{70}\text{B}_{16}\text{Si}_5$	—	30 000	8	0.30	1.58	125

* Metglas is Allied Corporation's registered trademark for amorphous alloys.

TABLE 11. TYPICAL PROPERTIES OF PERMANENT-MAGNET MATERIALS

Name	Composition, %	Retentivity $B_r(T)$	Coercivity $H_c/III_c(kA/m)$	Max. Energy Product (BH) max (kJ/m ³)
Cast Alnico 2	10 Al, 19 Ni, 13 Co, 3 Cu, bal Fe	0.75	45	13.5
Cast Alnico 5	8 Al, 14 Ni, 24 Co, 3 Cu, bal Fe	1.28	51	44.0
Cast Alnico 5-7	8 Al, 14 Ni, 24 Co, 3 Cu, bal Fe	1.35	59	60.0
Cast Alnico 6	8 Al, 16 Ni, 24 Co, 3 Cu, 1 Ti, bal Fe	1.05	62	31.0
Cast Alnico 8	7 Al, 15 Ni, 35 Co, 4 Cu, 5 Ti, bal Fe	0.82	130	42.0
Cast Alnico 9	7 Al, 15 Ni, 35 Co, 4 Cu, 5 Ti, bal Fe	1.05	120	72.0
Sintered Alnico 2	10 Al, 19 Ni, 13 Co, 3 Cu, bal Fe	0.71	44	12.0
Sintered Alnico 5	8 Al, 14 Ni, 24 Co, 3 Cu, bal Fe	1.09	49	31.0
Sintered Alnico 8	7 Al, 15 Ni, 35 Co, 4 Cu, 5 Ti, bal Fe	0.74	120	32.0
Bonded ferrite	BaO·6 Fe ₂ O ₃ + organics	0.16	110/240	4.4
Bonded ferrite	BaO·6 Fe ₂ O ₃ + organics	0.24	170/215	11.0
Sintered ferrite	BaO·6 Fe ₂ O ₃	0.22	140/280	8.0
Sintered ferrite	BaO·6 Fe ₂ O ₃	0.32	240/290	20.0
Sintered ferrite	BaO·6 Fe ₂ O ₃	0.38	180/190	27.0
Sintered ferrite	SrO·6 Fe ₂ O ₃	0.36	250/290	24.0
Sintered ferrite	SrO·6 Fe ₂ O ₃	0.40	180/190	30.0
Lodex 31	16 Fe, 9 Co, 67.5 Pb, 7.5 Sb	0.63	90	27.0
Lodex 32	19.2 Fe, 10.8 Co, 63 Pb, 7 Sb	0.74	75	28.0
Lodex 41	16 Fe, 9 Co, 67.5 Pb, 7.5 Sb	0.44	79	11.0
Lodex 42	19.2 Fe, 10.8 Co, 63 Pb, 7 Sb	0.53	68	11.0
Cunife	20 Fe, 20 Ni, 60 Cu	0.54	44	12.0
Vicalloy I	39 Fe, 51 Co, 10 V	0.84	19	7.0
Iron Chromc	28 Cr, 10.5 Co, bal Fe	0.98	30	16.0
Iron Chrome	23 Cr, 15 Co, 3 V, 2 Ti, bal Fe	1.35	44	44.0
Cobalt Rare Earth	Co ₅ Sm	0.82	600/>2400	130.0
Cobalt Rare Earth	Co ₅ Sm	0.87	660/>1200	144.0
Cobalt Rare Earth	(Co, Cu, Fe) ₇ Sm	1.10	510/530	240.0

For US suppliers, contact the Magnetic Materials Producers Association, 8 South Michigan Avenue, Suite 1000, Chicago, IL 60603; (312) 456-5590. Download the PDF file "Permanent Magnet Materials" MMPA Standard No. 0100-00 from www.mmpa.org to obtain new data on neodymium magnets.

relaxation as a function of frequency with the relaxation frequency varying from 1 000 to several million hertz. Most ferrites have relatively high resistivities (~ 10⁶ Ω-cm) if they are prepared carefully so as to avoid the presence of any divalent iron ion in the material. However, if the ferrite is prepared with an appreciable amount of divalent iron, then both the conductivity and permittivity are very high. Relative permittivities as high as 100 000 and resistivities less than 1 Ω-cm have been measured in several ferrites having a small amount of divalent iron in their composition.

Since the electrical resistivity of ferrites is typically 10⁶ times that of metals, ferrite components have much lower eddy-current losses and hence are used at frequencies above about 10 kHz. Table 12 lists some of the pertinent information with respect to typical commercial ferrites used in transformers and inductors. As a general rule, manganese zinc (MnZn) ferrites are used at frequencies up to about 1 MHz, beyond which nickel zinc (NiZn) ferrites become more efficient as a result of greater electrical resistivity. Relative initial permeabilities can be as high as 20 000 in commercial

MnZn ferrites for transformer applications, while the relative loss factor, tan δ/μ, can be as low as 10⁻⁶ (at 100 kHz) for some low-loss-inductor MnZn ferrites. For filter applications, a low-temperature coefficient of permeability and a low disaccommodation factor are often critical. Table 12 also contains data for some MnZn ferrites used in the design of switched-mode power supplies where a high saturation flux density is desirable.

For microwave applications, devices make use of the nonreciprocal propagation characteristics of ferrites close to or at a gyromagnetic resonance frequency between about 1 and 100 GHz. The most important of such devices are isolators and circulators. Table 13 lists the pertinent design values of several current commercial microwave ferrites, including saturation polarization J_s , resonance line width Δ*H*, relative permeability ε_r, loss tangent tan δ, and Curie temperature T_c . Materials having a range of J_s are needed, depending on the operating frequency, ω, of the device, since for resonance, $J_s < ω/δ$, where δ is the gyromagnetic ratio. In the low-frequency range (1-5 GHz), a

TABLE 12. CHARACTERISTICS OF FERRITES

Material	MnZn Ferrites							NiZn Ferrites				
	H5A	H5B	H5C2	H5E	H6F	H6H3	H6K	H7C1	H7C2	K5	K6A	K8
Code*	H5A	H5B	H5C2	H5E	H6F	H6H3	H6K	H7C1	H7C2	K5	K6A	K8
Practical freq. (MHz)	<0.2	<0.1	<0.1	<0.01	0.2–2.0	0.01–0.8	0.01–0.3	<0.3	<0.2	<8	1–50	<200
Initial permeability	3 300	5 000	10 000	18 000	800	1 300	2 200	2 500	3 900	290	25	16
Relative loss factor, $\tan \delta/\mu_i \times 10^6$	<2.5 10 kHz	<6.5 10 kHz	<7.0 10 kHz	—	<17 1 MHz	<1.2 100 kHz	<3.5 100 kHz	—	—	<28 1 MHz	<150 10 MHz	<250 100 MHz
Temp. coef. of $\mu_i \times 10^6$, –30 to 20 °C†	–0.5 to 2.0	–0.5 to 2.0	–0.5 to 1.5	–0.5 to 2.0	—	0.3 to 2.0	0.4 to 1.2	—	—	–4.0 to 2.0	—	—
Curie temp. (°C)	>130	>130	>120	>115	>200	>200	>130	>230	>200	>280	>450	>500
Saturation flux density (T)	0.41	0.42	0.40	0.44	0.40	0.47	0.39	0.51	0.48	0.33	0.30	0.27
Disaccommodation factor $\times 10^6$, 1–10 min.‡	<3	<3	<1	<1	<12	<5	<2	—	—	<30	<20	—
Resistivity ($\Omega \cdot m$)	1	1	0.15	0.05	4	25	8	10	2	20×10^5	2.5×10^5	1.0×10^5
Applications	Transformers				Inductors			Power Supplies		Inductors		

From Kirk-Othmer, *Encyclopedia of Chemical Technology*, Vol. 14, 3rd ed. (New York: John Wiley & Sons, Inc., p. 665, Table 7).

* TDK Data Book: Ferrite Cores-2, Aug., 1978.

† $(\mu_2 - \mu_1)/\mu_1^2(T_2 - T_1)$.

‡ $(\mu_1 - \mu_2)\mu_1^2 \log(t_2/t_1)$.

TABLE 13. SELECTED MICROWAVE MATERIALS

Material	Saturation Polarization $J_s(T)$	Line Width ΔH (A/cm)	Relative Permittivity ϵ_r	Loss Tangent $\tan \delta$	Curie Temperature T_c ($^{\circ}C$)	Remarks
<i>Garnets</i>						
Y	0.180	36	15.0	2×10^{-4}	280	—
YAl	0.018	36	13.8	2×10^{-4}	105	Decreasing aluminum content
.
YAl	0.120	36	14.8	2×10^{-4}	220	.
YGd	0.073	160	15.4	2×10^{-4}	280	Decreasing Gd, low $\Delta M_s/\Delta T$, const. T_c
.
YGd	0.160	40	15.1	2×10^{-4}	280	.
YGdAl	0.040	52	14.2	2×10^{-4}	150	Similar to YAl, but lower $\Delta M_s/\Delta T$
.
YGdAl	0.140	40	15.1	2×10^{-4}	265	.
<i>Spinel</i>						
MgMnAl	0.075	96	11.3	2.5×10^{-4}	90	Decreasing aluminum content
.
MgMnAl	0.175	180	12.2	2.5×10^{-4}	225	.
MgMn	0.215	432	12.7	2.5×10^{-4}	320	—
MgMnZn	0.250	416	12.9	2.5×10^{-4}	275	Increasing zinc content
.
MgMnZn	0.280	432	13.1	2.5×10^{-4}	225	.
NiZn	0.400	272	12.3	2.5×10^{-3}	470	High $4\pi M_s$
NiZn	0.500	128	12.5	1.0×10^{-3}	375	High $4\pi M_c$
Li	0.375	520	15.0	2.5×10^{-3}	640	High T_c
LiZn	0.480	192	14.5	2.5×10^{-3}	400	High $4\pi M_c$
LiTi	0.100	240	18.0	2.5×10^{-3}	300	Decreasing titanium content
.
LiTi	0.290	440	15.2	2.5×10^{-3}	600	.

From Kirk-Othmer, *Encyclopedia of Chemical Technology*, Vol. 14, 3rd ed. (New York: John Wiley & Sons, Inc.; p. 668, Table 8). (Source: Trans-Tech, Inc.)

class of ferrites called garnets (general formula $R_3Fe_5O_{12}$ with R being yttrium or other rare earth element) is used, whereas in the high-frequency range, the spinel ferrites are used.

MAGNETOSTRICTION

The static strain $\Delta l/l$ produced by a direct-current polarizing flux density B_0 is given by

$$\Delta l/l = cB_0^2$$

where c is a material constant expressed in m^4/Wb^2 .

If a small alternating-current driving field is superimposed on a large constant polarizing field B_0 , we have

$$d(\Delta l/l) = 2cB_0B = \beta B$$

The magnetostriction stress constant Λ in newtons/weber is

$$\Lambda = \beta Y_0 = 2cB_0Y_0$$

where Y_0 is the Young's modulus for a free bar.

Nickel contracts with increasing B , so Λ is a negative. Permalloy and Alfer expand, and their Λ is positive (Fig. 4).

Table 14 gives values for the three important transducer materials in SI units: annealed nickel, 45 Permalloy (45-percent nickel, 55-percent iron), and Alfer (an alloy of 13-percent aluminum and 87-percent

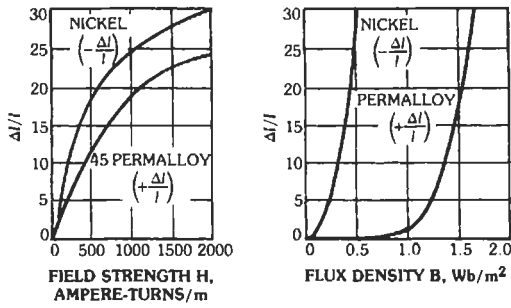


Fig. 4. Strain versus field strength (left) and versus flux density (right). From T. F. Hueter and R. H. Bolt, Eds., *Sonics* (New York: John Wiley & Sons, Inc., 1955; p. 173).

iron). For nickel the values for two different polarizing conditions are given: 160 ampere-turns per meter and 1 200 ampere-turns per meter, the latter appearing in parentheses. Table 15 compares properties of six magnetostriction materials.

PIEZOELECTRICITY

Table 16 lists piezoelectric strain coefficients d_{ij} , which are ratios of piezoelectric polarization components to components of applied stress at constant electric field (direct piezoelectric effect) and also ratios of piezoelectric strain components to applied electric field components at constant mechanical stress (con-

TABLE 14. MAGNETOMECHANICAL COEFFICIENTS OF THREE IMPORTANT MAGNETOSTRICTIVE MATERIALS AT INTERNAL POLARIZING FIELD $H_0^{*†}$

Quantity Unit	Annealed Nickel**	45 Permalloy	Alfer (13% Al, 87% Fe)
H_0 A/m	160 (1200)	600	800
B_0 V·s/m ²	0.25 (0.51)	1.43	1.15
μ	1250 (340)	1900	1150
μ_i	137 (41)	230	190
$\Delta l/l$ at H_0	-8×10^{-6} (-26×10^{-6})	14×10^{-6}	26×10^{-6}
c m ⁴ /Wb ²	-1×10^{-1}	6.9×10^{-6}	19.5×10^{-6}
Λ	-4.8×10^6 (-20×10^6)	2.7×10^6	6.7×10^6
Y_0 N/m ²	20×10^{10}	13.8×10^{10}	15×10^{10}
ρ kg/m ³	8.7×10^3	8.25×10^3	6.7×10^3
$k_c\%$ (electromechanical coupling factor)	14 (31)	12.4	27
ρ_c $\Omega \cdot m$	7×10^{-8}	7×10^{-7}	9×10^{-7}

* From T. F. Hueter and R. H. Bolt, eds., *Sonics* (New York: John Wiley & Sons, Inc., 1955; p. 175).

† The number of external ampere-turns per meter required to produce H_0 depends on the shape of the core. For closed magnetic loops H_{ext} is equal to H_0 . For rod-shaped cores external fields larger than H_0 are necessary to compensate for the demagnetizing effect of the poles at the free ends.

** The values for two different polarizing conditions are given: 160 and 1200 ampere-turns/meter, the latter in parentheses.

TABLE 15. OTHER CONSTANTS FOR SOME MAGNETOSTRICTION MATERIALS

Composition	$\Delta l/l$ at Saturation of B	Λ (N/Wb)	μ_i (H/m)	Young's modulus (N/m ²)	k_c (%)	Curie Temperature (°C)
99.9% Ni annealed	-33×10^{-6}	-20×10^6	4.3×10^{-6}	2.0×10^{11}	31	358
2 V 49 Co 49 Fe (2V Permadrur)	$+70 \times 10^{-6}$	—	—	1.7×10^{11}	20–37	980
45 Ni 55 Fe (45 Permalloy)	$+27 \times 10^{-6}$	2.7×10^6	2.9×10^{-4}	1.4×10^{11}	12	440
13 Al 87 Fe (13 Alfer)	+40	6.7×10^6	2.4×10^{-4}	1.5×10^{11}	27	500
Fe ₃ O ₄	+40	-90×10^6	1.9×10^{-2}	1.8×10^{11}	3	190
Ferrite 7 Al	-28	-28 to -44×10^6	$4-5 \times 10^{-5}$	1.68 to 1.75×10^{11}	25–30	640

verse effect). The subscripts $i = 1$ to 3 indicate electric field components, and the subscripts $j = 1$ to 6 indicate mechanical stress or strain components. These components are referred to the crystallographic principal

axes. For correlation of these to crystallographic axes, we follow Standards on Piezoelectric Crystals. For completeness, the full d matrix is listed in Fig. 5 for the various crystal classes.

TABLE 16. PIEZOELECTRIC STRAIN COEFFICIENTS FOR VARIOUS MATERIALS

(A) Cubic and Tetragonal Crystals	Composition	Class	d_{14}	d_{36}				
Sphalerite	ZnS	$\bar{4}3m$	9.7	—				
Sodium chlorate	NaClO ₃	23	5.2	—				
Sodium bromate	NaBrO ₃	23	7.3	—				
"ADP"	NH ₄ H ₂ PO	$\bar{4}2m$	-1.5	+48.0				
"KDP"	KH ₂ PO ₄	$\bar{4}2m$	+1.3	+21				
"ADA"	NH ₄ H ₂ AsO ₄	$\bar{4}2m$	+41	+31				
"KDA"	KH ₂ AsO ₄	$\bar{4}2m$	+23.5	+22				
(B) Trigonal Crystals	Class	d_{11}	d_{14}	d_{15}	d_{22}	d_{31}	d_{33}	
Quartz	32	+6.9	-2.0					
Tourmaline	3			+11.0	-0.94	+0.96	+5.4	
(C) Orthorhombic Crystals	Class	d_{14}	d_{25}	d_{36}				
Epsomite	222	-6.2	-8.2	-11.5				
Iodic acid	222	57	46	70				
Rochelle salt (30 °C)	222	+1500†	-160	+35				
NaNH ₄ tartrate	222	+56	-150	+28				
LiK tartrate	222	9.6	33.6	22.8				
LiNH ₄ tartrate	222	13.2	19.6	14.8				
(NH ₄) ₂ oxalate	222	50	11	25				
	d_{15}	d_{24}	d_{31}	d_{32}	d_{33}			
K pentaborate	9.5	1.7	-5.4	0	+5.6			
(D) Monoclinic Crystals (Class 2)	d_{14}	d_{16}	d_{21}	d_{22}	d_{23}	d_{25}	d_{34}	d_{36}
Lithium sulfate	+14.0	-12.5	+11.6	-45.0	-5.5	+16.5	-26.4	+10.0
Tartaric acid	+24.0	+15.8	-2.3	-6.5	-6.3	+1.1	-32.4	+35.0
K ₂ tartrate (DKT)	-25	+6.5	-2.2	+8.5	-10.4	-22.5	+29.4	-66.0
(NH ₄) ₂ tartrate	+9.3	-8.5	+17.6	-26.2	+1.8	-5.9	-14.0	+5.6
EDT (ethylene diarnine tartrate)	-31.1	-36.5	+30.6	+6.6	-33.8	-54.3	-51	-56.9
Cane sugar	-3.7	-7.2	+4.4	-10	+2.2	-2.6	-1.3	+1.3
(E) Polarized Polycrystalline Substance	d_{15}	d_{31}	d_{33}					
Barium titanate ceramic $\epsilon_r = 1700$	750	-235	+570					

Note: If the sign of a coefficient is not given, it is unknown (not necessarily positive).

* From *Smithsonian Physical Tables*, 9th revised ed., Vol. 120 (Washington, DC: Smithsonian Institution, 1969; p. 432).

† The coefficient d_{14} of Rochelle salt is extremely dependent on temperature and amplitude. The ratio of d_{14} to dielectric constant ϵ_r is, however, nearly constant.

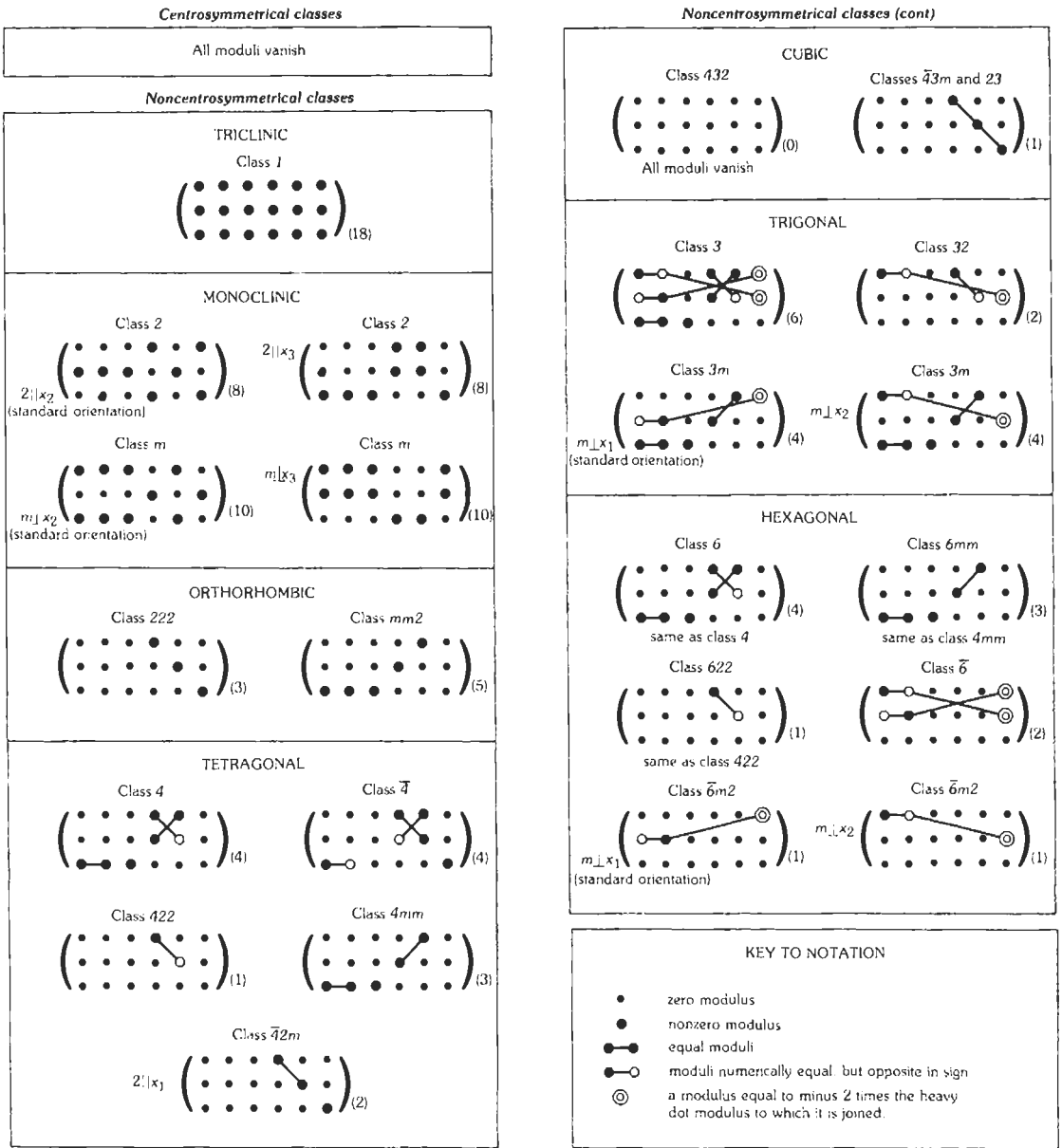


Fig. 5. Form of the (d_{ij}) matrix. From Nye, J. F., *Physical Properties of Crystals* (Oxford: Oxford University Press, 1986; pp. 123-124).

In the monoclinic system, indices 2 and 5 refer to the symmetry (b) axis, in distinction from the older convention relating indices 3 and 6 to the symmetry axis. Crystal classes are designated by international (Hermann-Mauguin) symbols. A dash in place of a coefficient indicates that it is equal by symmetry from another listed coefficient; a blank space indicates that

the coefficient is zero by symmetry. If the sign of a coefficient is not given, it is unknown, not necessarily positive.

$$\begin{aligned} \text{Unit for } d_{ij} &= 1/3 \times 10^{-12} \text{ C/N} \\ &= 1/3 \times 10^{-12} \text{ m/V} \end{aligned}$$

Coupling factor k is defined practically by

$$k^2 = \frac{\text{(mechanical energy converted into electric energy)}}{\text{(mechanical energy put into the crystal)}}$$

The converse effect is also true. The same type of relationship holds, and the coupling coefficient is numerically identical to what it was before, namely

$$k^2 = \frac{\text{(electrical energy converted into mechanical energy)}}{\text{(electrical energy put into the crystal)}}$$

d is the measure of the deflection caused by an applied voltage or the amount of charge produced by a given force (units = meters per volt or coulombs per newton)

g denotes a field produced in a piezoelectric crystal by an applied stress unit:

$$\frac{\text{V/m}}{\text{N/m}^2}$$

Equations that relate g , d , and k are:

$$g = d/\epsilon_r\epsilon_0$$

and

$$k^2 = gdE$$

where

- ϵ_r = relative permittivity of the dielectric,
- ϵ_0 = permittivity of free space = 8.854×10^{-12} F/m,
- E = Young's modulus.

Constants of some piezoelectric materials are listed in Table 17.

ACOUSTIC PROPERTIES OF SOME MATERIALS

Information regarding the acoustic properties of some materials is contained in Tables 18 and 19.

TABLE 17. CONSTANTS OF SOME PIEZOELECTRIC MATERIALS*†

Physical Property	Quartz O° X-cut	Lithium Sulfate O° Y-cut	Barium Titanate Type B	Lead Zirconate-/ Titanate		Lead Meta-/ Niobate	Units
				PZT-4	PZT-5		
Density ρ	2.65	2.06	5.6	7.6	7.7	5.8	10^3 kg/m ³
Acoustic impedance ρc	15.2	11.2	24	30.0	28.0	16	10^6 kg/m ² s
Frequency thickness constant ft	2870	2730	2740	2000	1800	1400	kHz·mm
Maximum operating temperature	550	75	70-90	250	290	500	°C
Relative permittivity	4.5	10.3	1700	1300	1700	225	—
Electromechanical coupling factor for thickness mode k_{33}	0.1	0.35	0.48	0.64	0.675	0.42	—
Electromechanical coupling factor for radial mode k_p	0.1	—	0.33	0.58	0.60	0.07	—
Elastic quality factor Q	10^6	—	400	500	75	11	—
Piezoelectric modulus for thickness mode d_{33}	2.3	16	149	285	374	85	10^{-12} m/V
Piezoelectric pressure constant g_{33}	58	175	14.0	26.1	24.8	42.5	10^{-3} (V/m)/(N/m ²)
Volume resistivity at 25 °C	$>10^{12}$	—	$>10^{11}$	$>10^{12}$	$>10^{13}$	10^9	—
Curie temperature	575	—	115	320	365	550	°C
Young's modulus E	8.0	—	11.8	8.15	6.75	2.9	10^{10} N/m ²
Rated dynamic tensile strength	—	—	—	24	27	—	MN/m ²

* From J. R. Fredrick, *Ultrasonic Engineering* (New York: John Wiley & Sons, Inc., 1965; p. 66).

† The properties of the ceramic materials can vary with slight changes in composition and processing, and hence the values that are shown should not be taken as exact.

TABLE 18. SIMPLIFIED EQUATIONS FOR SOUND INTENSITY*

Material	Mode of Vibration	Crystal Cut	Effective† Piezo Modulus H (C/m ²)	Velocity c (m/s)	Sound Density ρ (kg/m ³)	Sound Intensity Water for Airbacked Transducer \mathcal{I} (W/m ²)	Units
Quartz	Thickness	X	$H = e_{11}$ $= 0.17$	5.72×10^3	2.65×10^3	$f_0^2 V^2 \times 10^{-14}$	V (V) rms f (Hz)
	Longitudinal	X	$H = d_{11}/s'$ $= 0.18$	5.44×10^3	2.65×10^3	$0.087 \times E^2 \times 10^{-6}$	E (V/m) rms
Ammonium dihydrogen phosphate (ADP)	Longitudinal	45°Z	$H = d_{36}/2s'$ $= 0.042$	3.28×10^3	1.8×10^3	$0.6E^2 \times 10^{-6}$	E (V/m) rms
Rochelle salt (0 °C)	Longitudinal	45°X	$H = d_{14}/2s'$ $= 5.67$	3.4×10^3	1.77×10^3	$85E^2 \times 10^{-6}$	E (V/m) rms
Barium titanate (40 °C)	Thickness	Polarized normal to thickness	$H = e_{33}$ $= 10 \text{ to } 17$	5×10^3	5.5×10^3	$0.005f_0^2 V^2 \times 10^{-8}$ to $0.014f_0^2 V^2 \times 10^{-8}$	V (V) rms f (Hz)

* From T. F. Hueter and R. H. Bolt, Eds., *Sonics* (New York: John Wiley & Sons, Inc., 1955; p. 125).

† The quantity s' in the relationship $H = d_{ij}/s'$ is the effective compliance in the direction of longitudinal vibration.

TABLE 19. VELOCITIES, DENSITIES, AND CHARACTERISTIC IMPEDANCES OF VARIOUS METALS*

Metals	Velocities			Density ρ (kg/m ³)	Characteristic Impedance ρc Bulk (kg/m ² s)
	Longitudinal				
	Bulk (m/s)	Bar (m/s)	Shear (m/s)		
	$\times 10^3$	$\times 10^3$	$\times 10^3$	$\times 10^3$	$\times 10^6$
Aluminum	6.40	5.15	3.13	2.7	17.3
Beryllium	12.89	...	8.88	1.8	23.2
Brass, 70-30	4.37	3.40	2.10	8.5	37.0
Cast iron	3.50-5.60	3.0-4.7	2.2-3.2	7.2	25.0-40.0
Copper	4.80	3.65	2.33	8.9	42.5
Gold	3.24	2.03	1.20	19.3	63.0
Iron	5.96	5.18	3.22	7.9	46.8
Lead	2.40	1.25	0.79	11.3	27.2
Magnesium	5.74	4.90	3.08	1.7	9.9
Mercury	1.45	13.6	19.6
Molybdenum	6.25	...	3.35	10.2	63.7
Nickel	5.48	4.70	2.99	8.9	48.5
Platinum	3.96	2.80	1.67	21.4	85.0
Steel, mild	6.10	5.05	3.24	7.9	46.7
Silver	3.70	2.67	1.70	10.5	36.9
Tin	3.38	2.74	1.61	7.3	24.7
Titanium	5.99	...	3.12	4.50	27.0
Tungsten	5.17	...	2.88	19.3	100.0
Uranium	3.37	...	2.02	18.7	63.0
Zinc	4.17	3.81	2.48	7.1	29.6
Zirconium	4.65	...	2.30	6.4	29.8
Other Solid Materials:					
Crown glass	5.66	5.30	3.42	2.5	14.0
Granite	...	3.95	...	2.75	...
Ice	3.98	...	1.99	0.9	3.6
Nylon	1.8-2.2	1.1-1.2	2.0-2.7
Paraffin, hard	2.2	0.83	1.8
Plexiglas or Lucite	2.68	1.8	1.32	1.20	3.2
Polystyrene	2.67	1.06	2.8
Quartz, fused	5.57	5.37	3.52	2.6	14.5
Teflon	1.35	2.2	3.0
Tungsten carbide	6.66	...	3.98	10.0-15.0	66.5-98.5
Wood, oak	...	4.1	...	0.8	...
Fluids:					
Benzene	1.32	0.88	1.16
Castor oil	1.54	0.95	1.45
Glycerine	1.92	1.26	2.5
Methyl iodide	0.98	3.23	3.2
Oil, SAE 20	1.74	0.87	1.5
Water, fresh	1.48	1.00	1.48

Note: The values that are shown should not be taken as exact values because of the effects of variations in composition and processing. They are adequate for most practical purposes, however.

* From J. R. Fredrick, *Ultrasonic Engineering* (New York: John Wiley & Sons, Inc., 1965; p. 363).

5 Components or Parts

Revised by

*John J. Bohrer, Robert Blaszczyk,
Pierre J. Gerondeau, Jack Isken, Edward Mette,
Rick Price, Joseph A. Toro, and Theodore E. VanKampen*

General Standards	5-4
Color Coding	
Tolerance	
Preferred Values	
Voltage Rating	
Characteristic	
Environmental Test Methods	5-5
Standard Ambient Conditions for Measurement	5-5
Other Standard Environmental Test Conditions	5-5
Ambient Temperature	
Constant-Humidity Tests	
Cycling Humidity Tests	
High-Altitude Tests	
Vibration Tests	
Component Value Coding	5-7
Semiconductor-Diode Type-Number Coding	
Resistors—Definitions	5-7

- Resistors—Fixed Composition 5-9
- Color Code
 - Tolerance
 - Packaging Styles
 - Temperature and Voltage Coefficients
 - Noise
 - RF Effects
 - Good Design Practice
- Resistors—Fixed Wirewound 5-11
- EIA Low-Power Insulated Resistors
 - EIA Precision Resistors
 - EIA Power Resistors
- Resistors—Fixed Film 5-12
- Construction
 - Resistive Films
 - Applications
 - Technical Characteristics
- Resistors—Adjustable 5-13
- Types of Adjustable Resistors
 - Terminal Identification
 - Mounting Characteristics
- Capacitors—Definitions 5-16
- Classes of Capacitors 5-17
- Plastic Film Capacitors 5-17
- Polystyrene Film Capacitors With Foil Electrodes
 - Polyester Film Capacitors With Foil Electrodes
 - Plastic Film Capacitors With Metallized Electrodes
- Electrolytic Capacitors 5-18
- Aluminum Electrolytics
 - Tantalum-Foil Electrolytics
 - Tantalum Electrolytics With Porous Anode and Liquid Electrolyte
 - Tantalum Electrolytics With Porous Anode and Solid Electrolyte
- Ceramic Capacitors 5-22
- Class I
 - Class II
 - Class III
 - Color Code
 - Temperature Coefficient
- Paper Foil-Type Capacitors 5-23

Mica Capacitors 5-25**Construction****Applications****Type Designation****Capacitance****Temperature Coefficient****Dissipation Factor****High-Potential or Withstanding-Voltage Test****Humidity and Thermal-Shock Tests****Life****Printed Circuits 5-26****Printed-Circuit Base Materials****Conductor Materials****Manufacturing Processes****Circuit-Board Finishes****Design Considerations****Preparation of Artwork****Assembly**

GENERAL STANDARDS

Standardization of electronic components or parts is handled by several cooperating agencies.

In the US, the Electronic Industries Association (EIA)* and the American National Standards Institute (ANSI)† are active in the commercial field. Electron-tube and semiconductor-device standards are handled by the Joint Electron Device Engineering Council (JEDEC), a cooperative effort of EIA and the National Electrical Manufacturers Association (NEMA)‡. Military (MIL) standards are issued by the US Department of Defense or one of its agencies such as the Defense Electronics Supply Center (DESC).

International standardization in the electronics field is carried out by the various Technical Committees of the International Electrotechnical Commission (IEC)§. A list of the available IEC Recommendations is included in the ANSI Index (outside the US, consult the national standardization agency or the IEC). Documents from the IEC may be used directly, or their recommendations may be incorporated in whole or in part in national standards issued by the EIA or ANSI. A few broad areas may be covered by standards issued by the International Standards Organization (ISO).

These organizations establish standards for electronic components or parts (and in some cases, for equipments) to provide interchangeability among different products regarding size, performance, and identification; minimum number of sizes and designs; and uniform testing of products for acceptance. This chapter presents a brief outline of the requirements, characteristics, and designations for the major types of component parts used in electronic equipment. Such standardization offers economic advantages to both the parts user and the parts manufacturer, but is not intended to prevent the manufacture and use of other parts under special conditions.

Color Coding

The color code of Table 1 is used for marking electronic parts.

Tolerance

The maximum deviation allowed from the specified nominal value is known as the tolerance. It is usually given as a percentage of the nominal value, though for very small capacitors the tolerance may be specified in

picofarads (pF). For critical applications it is important to specify the permissible tolerance; where no tolerance is specified, components are likely to vary by ± 20 percent from the nominal value.

Do not assume that a given lot of components will have values distributed throughout the acceptable range of values. A lot ordered with a $\pm 20\%$ tolerance may include *no* parts having values within 5% of the desired nominal value; these may have been sorted out before shipment. The manufacturing process for a given lot may produce parts in a narrow range of values only, not necessarily centered in the acceptable tolerance range.

Preferred Values

To maintain an orderly progression of sizes, preferred numbers are frequently used for the nominal values. A further advantage is that all parts manufactured are salable as one or another of the preferred values. Each preferred value differs from its predecessor by a constant multiplier, and the final result is conveniently rounded to two significant figures.

ANSI Standard Z17.1-1973 covers a series of preferred numbers based on $(10)^{1/5}$ and $(10)^{1/10}$ as listed in Table 2. This series has been widely used for fixed wirewound power-type resistors and for time-delay fuses.

Because of the established practice of using ± 20 -, ± 10 -, and ± 5 -percent tolerances, a series of values based on $(10)^{1/6}$, $(10)^{1/12}$, and $(10)^{1/24}$ has been adopted by the EIA, and is now an ANSI Standard (C83.2-1971) (EIA RS-385). It is widely used for such small electronic components as fixed composition resistors and fixed ceramic, mica, and molded paper capacitors. These values are listed in Table 2. (For series with smaller steps, consult the ANSI or EIA Standard.)

Voltage Rating

Distinction must be made between the breakdown-voltage rating (test volts) and the working-voltage rating. The maximum voltage that may be applied (usually continuously) over a long period of time without causing the part to fail determines the working-voltage rating. Application of the test voltage for more than a very few minutes, or even repeated applications of short duration, may result in permanent damage or failure of the part.

Characteristic

The term "characteristic" is frequently used to include various qualities of a part such as temperature coefficient of capacitance or resistance, Q value, maximum permissible operating temperature, stability when subjected to repeated cycles of high and low temperature, and deterioration when it is subjected to moisture either as humidity or water immersion. One or two letters are assigned in EIA or MIL type designations, and the characteristic may be indicated by color coding

* EIA Engineering Dept., Washington, D.C. Index of standards is available. EIA was formerly Radio-Electronics-Television Manufacturers' Association (RETMA).

† ANSI, New York, New York. Index of standards is available. ANSI was formerly the USA Standards Institute (USASI).

‡ NEMA, New York, New York. Index of standards is available.

§ IEC, Central Office; Geneva, Switzerland. The US National Committee for the IEC operates within the ANSI.

TABLE 1. STANDARD COLOR CODE OF ELECTRONICS INDUSTRY

Color	Significant Figure	Decimal Multiplier	Tolerance in Percent*	Voltage Rating	Characteristic
Black	0	1	±20 (M)	—	A
Brown	1	10	±1 (F)	100	B
Red	2	100	±2 (G)	200	C
Orange	3	1 000	±3	300	D
Yellow	4	10 000	GMV‡	400	E
Green	5	100 000	±5(J)†, (0.50(D))§	500	F
Blue	6	1 000 000	±6(0.25(C))§	600	G
Violet	7	10 000 000	±12.5, (0.10(B))§	700	—
Gray	8	0.01†	±30, (0.05(N))§	800	I
White	9	0.1†	±10†	900	J
Gold	—	0.1	±5 (J), (0.50(E))	1 000	—
Silver	—	0.01	±10 (K)	2 000	—
No Color	—	—	±20	500	—

* Tolerance letter symbol as used in type designations has tolerance meaning as shown. ±3, ±6, ±12.5, and ±30 percent are tolerances for USA Std 40-, 20-, 10-, and 5-step series, respectively.

† Optional coding where metallic pigments are undesirable.

‡ GMV is -0 to +100-percent tolerance or Guaranteed Minimum Value.

§ For some film and other resistors only.

|| For some capacitors only.

on the part. An explanation of the characteristics applicable to a component or part will be found in the following sections covering that part.

ENVIRONMENTAL TEST METHODS

Since many component parts and equipments have the same environmental exposure, environmental test methods are becoming standardized. The principal standards follow.

EIA Standard RS-186-E (ANSI C83.58-1978): Standard Test Methods for Passive Electronic Component Parts.

IEC Publication 68: Basic Environmental Testing Procedures for Electronic Components and Electronic Equipment (published in multiple).

MIL-STD-202F: Military Standard Test Methods for Electronic and Electrical Component Parts.

MIL-STD-750B: Test Method for Semiconductor Device.

MIL-STD-810C: Military Standard Environmental Test Methods.

MIL-STD-883B: Test Methods and Procedures for Microelectronics.

MIL-STD-1344A: Test Methods for Electrical Connectors.

ASTM Standard Test Methods*—Primarily applicable to the materials used in electronic component parts.

Wherever the test methods in these standards are reasonably applicable, they should be specified in preference to other methods. This simplifies testing of a wide variety of parts, testing in widely separated locations, and comparison of data.

When selecting destructive environmental tests to determine the probable life of a part, distinguish between the environment prevailing during normal equipment operation and the environment used to accelerate deterioration. During exposure to the latter environment, the item may be out of tolerance with respect to its parameters in its normal operating-environment range. Accelerated tests are most meaningful if some relation between the degree of acceleration and component life is known. Such acceleration factors are known for many insulation systems.

STANDARD AMBIENT CONDITIONS FOR MEASUREMENT

Standard ambient conditions for measurement are listed in Table 3.

OTHER STANDARD ENVIRONMENTAL TEST CONDITIONS

Ambient Temperature

Dry heat, °C: +30, +40, (+49), +55, (+68), +70, (+71), +85, +100, +125, +155, +200 (values in parentheses not universally used).

Cold, °C: -10, -25, -40, -55, -65.

* ASTM = American Society for Testing and Materials; Philadelphia, Pa. Index of standards is available.

TABLE 2. PREFERRED VALUES*

Name of Series	USA Standard Z17.1-1973†		USA Standard C83.2-1971 (R 1977)‡		
	"5"	"10"	±20% (E6)	±10% (E12)	±5% (E24)
Percent step size	60	25	≈40	20	10
Step multiplier	$(10)^{1/5}=1.58$	$(10)^{1/10}=1.26$	$(10)^{1/6}=1.46$	$(10)^{1/12}=1.21$	$(10)^{1/24}=1.10$
Values in the series (Use decimal multipliers for smaller or larger values)	10	10 12.5 (12) }	10	10	10
	—	—	—	—	11
	—	—	—	12	12
	—	—	—	—	13
	—	—	15	15	15
	16	16	—	—	16
	—	—	—	18	18
	—	20	—	—	20
	—	—	22	22	22
	—	—	—	—	24
	25	25	—	—	—
	—	—	—	27	27
	—	31.5 (32) }	—	—	30
	—	—	—	—	—
	—	—	33	33	33
	—	—	—	—	36
	—	—	—	39	39
	40	40	—	—	—
	—	—	—	—	43
	—	—	47	47	47
	—	50	—	—	—
	—	—	—	—	51
	—	—	—	56	56
	—	—	—	—	62
	63	63	—	—	—
	—	—	68	68	68
	—	—	—	—	75
	—	80	—	—	—
	—	—	—	82	82
	—	—	—	—	91
	100	100	100	100	100

* ANSI Standard C83.2-1971 applies to most electronics components. It is the same as EIA Standard RS-385 (formerly GEN-102) and agrees with IEC Publication 63. ANSI Standard Z17.1-1973 covers preferred numbers and agrees with ISO 3 and ISO 497.

† "20" series with 12-percent steps ($(10)^{1/20}=1.122$ multiplier) and a "40" series with 6-percent steps ($(10)^{1/40}=1.059$ multiplier) are also standard.

‡ Associate the tolerance ±20%, ±10%, or ±5% only with the values listed in the corresponding column. Thus, 1200 ohms may be either ±10 or ±5, but not ±20 percent; 750 ohms may be ±5, but neither ±20 nor ±10 percent.

Constant-Humidity Tests

40 °C, 90 to 95% RH; 4, 10, 21, or 56 days.

66 °C, ≈ 100% RH: 48, 96, or 240 hours (primarily for small items).

Cycling Humidity Tests

Fig. 1 shows a number of cycling humidity tests. (See applicable chart in standard for full details.) Preconditioning is customary before starting cycle series. RH = relative humidity.

High-Altitude Tests

Information regarding high-altitude tests is given in Table 4.

Vibration Tests

The purposes of vibration tests are:

- (A) Search for resonance.
- (B) Determination of endurance (life) at resonance (or at specific frequencies).
- (C) Determination of deterioration resulting from

TABLE 3. STANDARD AMBIENT CONDITIONS FOR MEASUREMENT

	Standard	Temperature (°C)	Relative Humidity (%)	Barometric Pressure	
				mm Hg	mbar
Normal range	RS-186-D	15-35	47-75	650-800	860-1060
	IEC-68	15-35	45-75	(645-795)	860-1060
	MIL-STD-202E	15-35	45-75	650-800	(866-1066)
Closely controlled range	MIL-STD-810C	13-33	20-80	650-775	(866-1033)
	IEC-68	20±1	65±2	(645-795)	860-1060
	IEC-68	23±1	50±2	(645-795)	860-1060
	MIL-STD-202E	23±1	50±2	650-800	(866-1066)
	MIL-STD-810C	23±1.4	50±5	650-775	(866-1033)
	IEC-68	27±1	65±2	(645-795)	860-1060
	RS-186-D	25±2	50±2	650-800	860-1060

Notes:

1. Use the closely controlled range only if the properties are sensitive to temperature or humidity variations, or for referee conditions in case of a dispute. The three temperatures 20°, 23°, and 27 °C correspond to normal laboratory conditions in various parts of the world.
2. Rounded derived values are shown in parentheses ().
3. 25±2 °C, 20 to 50% relative humidity (RH) has been widely used as a closely controlled ambient for testing electronics components.

long exposure to swept frequency (or random vibration).

Recommended Frequency Ranges for Tests—Hertz: 1 to 10, 5 to 35, 10 to 55, 10 to 150, 10 to 500, 10 to 2000, 10 to 5000.

Recommended Combinations of Amplitude and Frequency—IEC Publication 68 recommends testing at constant amplitude below and constant acceleration

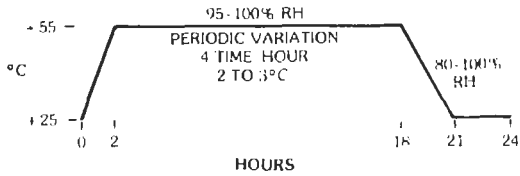
above the crossover frequency (57 to 62 hertz). MIL-STD-202 and MIL-STD-810 also follow this principle but use different crossover points and low-frequency severities. The choice of frequency range and vibration amplitude or acceleration should bear some relation to the actual service environment. Successful completion of 10⁷ vibration cycles indicates a high probability of no failures in a similar service environment. Resonances may make the equipment output unusable, although the mechanical life may be adequate.

TABLE 4. HIGH-ALTITUDE TESTS

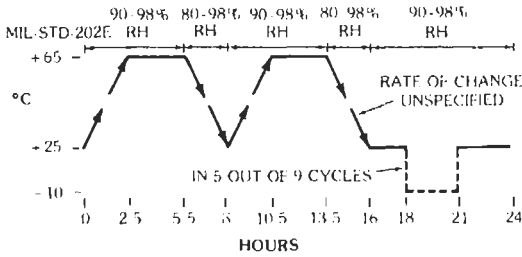
Pressure			Approximate Corresponding Altitude		Standard
mbar	mm Hg	in. Hg	feet	meters	
700	525	20.67	7 218	2 200	IEC
600	450	17.72	11 483	3 500	IEC
533	400	15.74	14 108	4 300	IEC
586	439	17.3	15 000	4 572	MIL-202
466	349	13.75	20 000	6 096	RS-186
300	225	8.86	27 900	8 500	IEC
300	226	8.88	30 000	9 144	MIL-202, RS-186
116	87.0	3.44	50 000	15 240	MIL-202, RS-186
85	63.8	2.51	52 500	16 000	IEC
44	33.0	1.30	65 600	20 000	IEC
44.4	33.0	1.31	70 000	21 336	MIL-202
20	17.2	0.677	85 300	26 000	IEC
10.6	8.00	0.315	100 000	30 480	MIL-202
1.28	1.09	0.043	150 000	45 720	MIL-202
3.18×10 ⁻⁶	2.40×10 ⁻⁶	9.44×10 ⁻⁸	656 000	200 000	MIL-202

Notes:

1. The inconsistency in the pressure-altitude relation arises from the use of different model atmospheres. For testing purposes always specify the desired pressure rather than an elevation in feet or meters.
2. Values in italics are derived from the values specified in the associated standard.



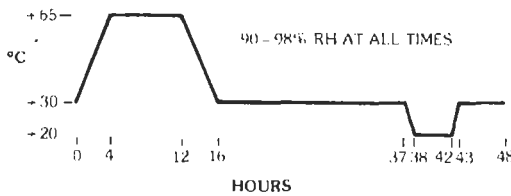
IEC-68-1. 2. or b 24-Hour Cycles



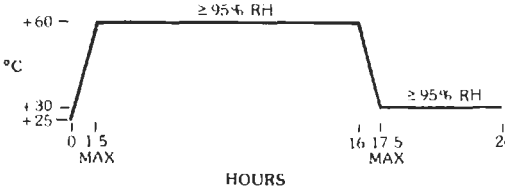
RS 186-D - 4. 10. or 28 24 Hour Cycles
IEC-68 MIL STD 202E - 10 24 Hour Cycles



MIL STD 810C - 10 24 Hour Cycles



MIL STD 810C - 5 48 Hour Cycles or 5 48 Hour Cycles Plus 480 Hours at 30°C. 90-98% RH



MIL STD 810C - 5 24-Hour Cycles

Fig. 1. Cycling humidity tests. Relative humidity for RS-186-D is 90-95% but may be uncontrolled during temperature changes.

COMPONENT VALUE CODING

Axial-lead and some other components are often color coded by circumferential bands to indicate the

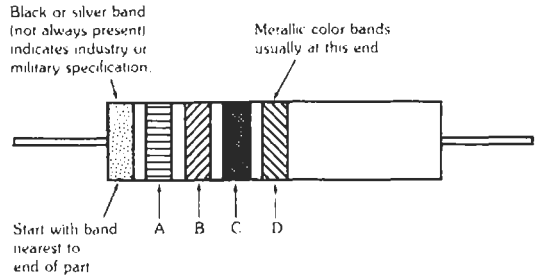


Fig. 2. Component value coding. The code of Table 1 determines values. Band A color = First significant figure of value in ohms, picofarads, or microhenries. Band B color = Second significant figure of value. Band C color = Decimal multiplier for significant figures. Band D color = Tolerance in % (if omitted, the broadest tolerance series of the part applies).

resistance, capacitance, or inductance value and its tolerance. Usually the value may be decoded as indicated in Fig. 2 and Table 1.

Sometimes, instead of circumferential bands, colored dots are used as shown in Fig. 3 and Table 5.

Semiconductor-Diode Type-Number Coding

The sequential number portion (following the "1N" of the assigned industry type number) may be indicated by color bands* as shown in Fig. 4. Colors have the numerical significance given in Table 1.

Bands J, K, L, M represent the digits in the sequential number (for two-digit numbers, band J is black). Band N is used to designate the suffix letter as shown in Table 6. Band N may be omitted in two- or three-digit number coding if not required, but it will always be present in 4-digit number coding (black if no suffix letter is required). See Table 7 for examples.

A single band indicates the cathode end of a diode or rectifier.

TABLE 5. COLOR-CODE EXAMPLES

Component Value	Band or Dot Color			
	A	B	C	D
3300 ± 20%	Orange	Orange	Red	Black or omitted
5.1 ± 10%	Green	Brown	Gold or white	Silver
1.8 megohms ± 5% (as applied to a resistor)	Brown	Gray	Green	Gold

* EIA Standard RS-236-B.

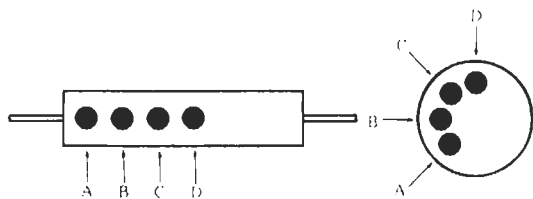


Fig. 3. Alternative methods of component value coding.

RESISTORS—DEFINITIONS

Wattage Rating: The maximum power that the resistor can dissipate, assuming (A) a specific life, (B) a standard ambient temperature, and (C) a stated long-term drift from its no-load value. Increasing the ambient temperature or reducing the allowable deviation from the initial value (more-stable resistance value) requires derating the allowable dissipation. With few exceptions, resistors are derated linearly from full wattage at rated temperature to zero wattage at the maximum temperature.

Temperature Coefficient (Resistance-Temperature Characteristic): The magnitude of change in resistance due to temperature, usually expressed in percent per degree Celsius or parts per million per degree Celsius (ppm/°C). If the changes are linear over the operating temperature range, the parameter is known as “temperature coefficient”; if nonlinear, the parameter is known as “resistance-temperature characteristic.” A large temperature coefficient and a high hot-spot temperature cause a large deviation from the nominal condition; e.g., 500 ppm/°C and 275 °C result in a resistance change of over 12 percent.

Maximum Working Voltage: The maximum voltage that may be applied across the resistor (maximum working voltage) is a function of (A) the materials used, (B) the allowable resistance deviation from the low-voltage value, and (C) the physical configuration of the resistor. Carbon composition resistors are more voltage-sensitive than other types.

Noise: An unwanted voltage fluctuation generated within the resistor. Total noise of a resistor always

TABLE 6. DIODE COLOR CODE

Color	Suffix Letter	Number
Black	—	0
Brown	A	1
Red	B	2
Orange	C	3
Yellow	D	4
Green	E	5
Blue	F	6
Violet	G	7
Gray	H	8
White	J	9

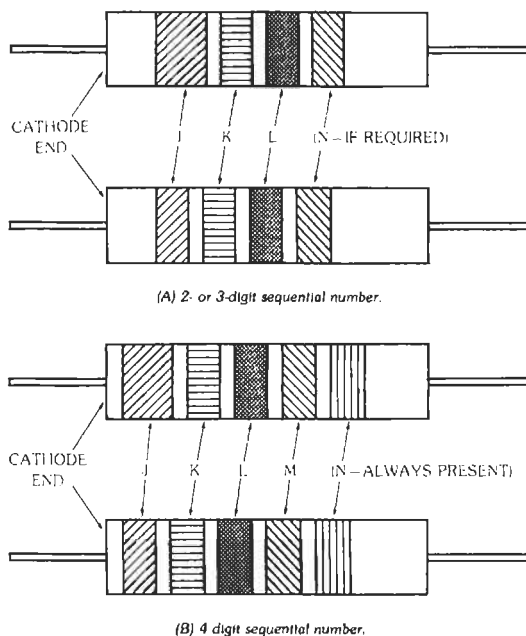


Fig. 4. Semiconductor-diode value coding.

TABLE 7. EXAMPLES OF DIODE COLOR CODING

Band	Band Color		
J	Red	Red	Orange
K	Green	Green	Blue
L	Yellow	Yellow	Violet
M	—	—	Red
N	—	Red	Black
	1N254	1N254B	1N3672

includes Johnson noise, which depends only on resistance value and the temperature of the resistance element. Depending on type of element and construction, total noise may also include noise caused by current and noise caused by cracked bodies and loose end caps or leads. For adjustable resistors, noise may also be caused by jumping of the contact over turns of wire and by imperfect electrical path between contact and resistance element.

Hot-Spot Temperature: The maximum temperature measured on the resistor due to both internal heating and the ambient operating temperature. The allowable maximum hot-spot temperature is predicated on thermal limits of the materials and the design. Since the maximum hot-spot temperature may not be exceeded under normal operating conditions, the wattage rating of the resistor must be lowered if it is operated at an ambient temperature higher than that at which the wattage rating was established. At zero dissipation, the maximum ambient around the resistor may be at its maximum hot-spot temperature. The ambient temperature for a resistor is affected by surrounding heat-

producing devices; resistors stacked together do not experience the ambient surrounding the stack except under forced cooling.

Critical Resistance Value: A resistor of specified power and voltage ratings has a critical resistance value above which the allowable voltage limits the permissible power dissipation. Below the critical resistance value, the maximum permitted voltage across the resistor is never reached at rated power.

Inductance and Other Frequency Effects: For other than wirewound resistors, the best high-frequency performance is secured if (A) the ratio of resistor length to cross section is a maximum, and (B) dielectric losses are kept low in the base material and a minimum of dielectric binder is used in composition types.

Carbon composition types exhibit little change in effective dc resistance up to frequencies of about 100 kHz. Resistance values above 0.3 megohm start to decrease in resistance at approximately 100 kHz. Above 1 MHz, all resistance values decrease.

Wirewound types have inductive and capacitive effects and are unsuited for use above 50 kHz, even when specially wound to reduce the inductance and capacitance. Wirewound resistors usually exhibit an increase in resistance at high frequencies because of "skin" effect.

Film types have the best high-frequency performance. The effective dc resistance for most resistance values remains fairly constant up to 100 MHz and decreases at higher frequencies. In general, the higher the resistance value the greater the effect of frequency.

Established-Reliability Resistors: Some resistor styles can be purchased with maximum-failure-rate guarantees. Standard-failure-rate levels are:

%/1000 hours—1.0; 0.1; 0.01; 0.001.

Resistance Value and Tolerance Choice: A calculated circuit-resistance nominal value should be checked to determine the allowable deviation in that value under the most unfavorable circuit, ambient, and life conditions. A resistor type, resistance value, and tolerance should be selected considering (A) standard resistance values (specials are uneconomical in most cases), (B) purchase tolerance, (C) resistance value changes caused by temperature, humidity, voltage, etc., and (D) long-term drift.

RESISTORS—FIXED COMPOSITION

Color Code

EIA-standard and MIL-specification requirements for color coding of fixed composition resistors are identical (see Fig. 2). The exterior body color of insulated axial-lead composition resistors is usually tan, but other colors (except black) are permitted. Nonin-

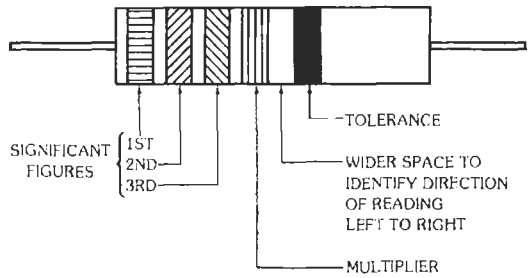


Fig. 5. Resistor value color code for three significant figures. Colors of Table 1 determine values.

sulated axial-lead composition resistors have a black body color.

If three significant figures are required, Fig. 5 shows the resistor markings (EIA Std RS-279).

Another form of resistor color coding (MIL-STD-1285A) is shown in Fig. 6. Colors have the significance shown in Table 8 for the fifth band.

Tolerance

Standard resistors are furnished in $\pm 20\%$, or $\pm 10\%$, and $\pm 5\%$ -percent tolerances, and in the preferred-value series of Table 2. "Even" values, such as 50 000 ohms, may be found in old equipment, but they are seldom used in new designs.

Packaging Styles

In addition to the familiar axial-lead configuration, resistors are also available in SIP (single in-line package), DIP (dual in-line package), flat-pack, and leadless chip configurations.

The SIP and DIP types are designed for thru-hole or socket mounting on printed circuit boards. Flat-pack and chip configurations are surface mounted.

Flat-packs, SIPs, and DIPs contain multiple resistors, which may be internally connected to form networks or to have a common terminal pin.

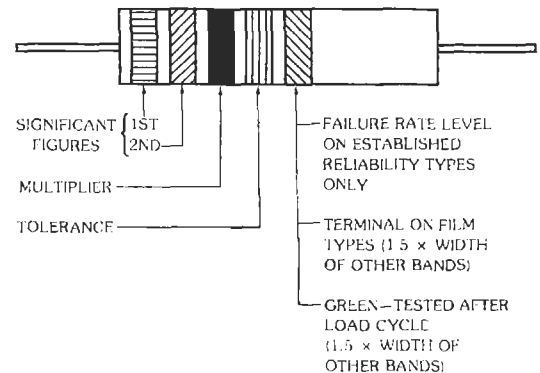


Fig. 6. Resistor color code per MIL-STD-1285A.

TABLE 8. COLOR CODE FOR FIFTH BAND (FIG. 6)

	Brown	Red	Orange	Yellow	Green	White
Failure Rate Level:						
Letter	M	P	R	S	—	—
Rate (%/1000 hrs)	1.0	0.1	0.01	0.001	—	—
Terminal	—	—	—	—	—	Weldable
Special	—	—	—	—	Fig. 6	—

The chip configuration allows increased component density and is autoinsertable. Since there are no leads, basic circuit reliability and high-frequency performance are improved.

Temperature and Voltage Coefficients

Resistors are rated for maximum wattage at an ambient temperature of 70 °C; above these temperatures up to the maximum allowable hot-spot temperature of 130 or 150 °C, it is necessary to operate at reduced wattage ratings. Resistance values are a function of voltage as well as temperature; present specifications allow a maximum voltage coefficient of resistance as given in Table 9 and permit a resistance-temperature characteristic as in Table 10.

A 1000-hour rated-load life test should not cause a change in resistance greater than 12% for 1/8-watt resistors and 10% for all other ratings. A severe cycling humidity test may cause resistance changes of 10% average and 15% maximum; 250 hours at 40 °C and 95% relative humidity may cause up to 10% change. Five temperature-change cycles, -55 °C to +85 °C, should not change the resistance value by more than 4% from the 25 °C value. Soldering the resistor in place may cause a resistance change of 3%. Always allow 1/4-inch minimum lead length; use heat-dissipating clamps when soldering confined assemblies. The preceding summary indicates that close tolerances cannot be maintained over a wide range of load and ambient conditions.

Noise

Composition resistors above 1 mcgohm have high Johnson noise levels, precluding their use in critical applications.

RF Effects

The end-to-end shunted capacitance effect may be noticeable because of the short resistor bodies and small internal distance between the ends. Operation at VHF or higher frequencies reduces the effective resistance because of dielectric losses (Boella effect).

Good Design Practice

Operate at one-half the allowable wattage dissipation for the expected ambient temperature. Provide an adequate heat sink. Mount no other heat-dissipating parts within one diameter. Use only in applications where a 15% change from the installed value is permissible or where the environment is controlled to reduce the resistance-value change.

RESISTORS—FIXED WIREWOUND

Fixed wirewound resistors are available as low-power insulated types, precision types, and power types.

EIA Low-Power Insulated Resistors*

These resistors are furnished with power ratings from 1 watt through 15 watts, in tolerances of ±5 and ±10 percent, and in resistance values from 0.1 ohm to 30 000 ohms in the preferred-value series of Table 2. They may be color coded as described in Fig. 2, but band A will be twice the width of the other bands. They may also be typographically marked in accordance with the EIA Standard.

The stability of these resistors is somewhat better than that of composition resistors, and they may be preferred except where a noninductive resistor is required.

EIA Precision Resistors†

These resistors are furnished in ±1.0-, ±0.5-, ±0.25-, ±0.1-, and ±0.05-percent tolerances and in any value from 1.0 ohm to 1.0 megohm in the preferred-value series of Table 2. Power ratings range from 0.1 watt to 0.5 watt. The maximum ambient temperature for full wattage rating is 125 °C. If the resistor is mounted in a confined area or may be required to operate in higher ambient temperatures (145 °C maximum), the allowable dissipation must be reduced in accordance with the EIA Standard.

* EIA Standard RS-344.

† EIA Standard RS-229-A.

TABLE 9. STANDARD RATINGS FOR COMPOSITION RESISTORS

Watts	Working Volts (Maximum)	Hot-Spot Temperature (°C)(Maximum)	Critical Resistance (Megohms)	Voltage Coefficient* (%/Volt) (Maximum)
1/8	150	150	0.22	0.05
1/4	250	130	0.25	0.035
1/2	350	130	0.25	0.035
1	500	130	0.25	0.02
2	500	130	0.12	0.02

* Applicable only to resistors of 1000 ohms and over.

These resistors have an inherently low noise level, approaching the thermal agitation level, and their stability is excellent—the typical change in resistance for the lifetime of the resistor will not exceed 50 percent of the initial resistance tolerance when used within the specified design limits of the EIA Standard.

The temperature coefficient of resistance over the range -55°C to $+145^{\circ}\text{C}$, referred to 25°C , may have maximums as follows:

Value	EIA Standard
Above 10 ohms	$\pm 0.002\%/^{\circ}\text{C}$
5 ohms to 10 ohms	$\pm 0.006\%/^{\circ}\text{C}$
Below 5 ohms	$\pm 0.010\%/^{\circ}\text{C}$

Where required, temperature coefficients of less than ± 20 ppm/ $^{\circ}\text{C}$ can be obtained by special selection of the resistance wire. Temperature coefficients of ± 10 ppm/ $^{\circ}\text{C}$ may be obtained by limiting the range of temperatures for testing from -40°C to $+105^{\circ}\text{C}$. The application of temperature coefficient to resistors should be limited, where possible, to the actual temperatures at which the equipment will operate.

EIA Power Resistors*

These resistors are furnished in 3 styles (strip; tubular, open end; and axial lead) and 24 power ratings

* EIA Standard RS-155-B.

ranging from 1 watt to 210 watts in tolerances of ± 1.0 percent and ± 5 percent. Resistance values range from 1.0 ohm to 182 kilohms in the preferred-value series of Table 2.

Axial-lead types are available in two general inductance classifications—inductive winding and noninductive winding. The noninductive styles have a maximum resistance value of $1/2$ the maximum resistance of inductive styles because of the special manner in which they are wound. The inductance of noninductive styles must not exceed 0.5 microhenry when measured at a test frequency of 1.0 megahertz $\pm 5\%$. However, these resistors should not be used in very-high-frequency circuits where the inductance may affect circuit operation.

The maximum ambient temperature for full wattage rating for these resistors is 25°C . When the resistors are operated at ambient temperatures above 25°C , the wattage dissipation must be reduced in accordance with the EIA Standard.

RESISTORS—FIXED FILM

Film-type resistors use a thin layer of resistive material deposited on an insulating core. The low-power types are more stable than the usual composition resistors. Except for very high-precision requirements, film-type resistors are a good alternative to accurate wirewound resistors, being both smaller and less expensive and having excellent noise characteristics.

The power types are similar in size and performance to conventional wirewound power resistors. While their 200°C maximum operating temperature limits the

TABLE 10. TEMPERATURE COEFFICIENT OF RESISTANCE FOR COMPOSITION RESISTORS

	Charac- teristic*	Percent Maximum Allowable Change from Resistance at 25 Degrees Celsius†					
		± 6.5	± 10	± 13	± 15	± 20 - ± 15	± 25 ± 15
At -55°C ambient	F						
At $+105^{\circ}\text{C}$ ambient	F	± 5	± 6	± 7.5	± 10	± 15	± 15
Nominal resistance in ohms		0 to 1 000	>1 000 to 10 000	>10 000 to 0.1 meg	>0.1 meg to 1.0 meg	>1 meg to 10 meg	>10 meg

* Resistance-temperature.

† Up to 1 megohm, data also apply to MIL Established Reliability characteristic G (= former GF).

power rating, the maximum resistance value available for a given physical size is much higher than that of the corresponding wirewound resistor.

Construction

For low resistance values, a continuous film is applied to the core, a range of values being obtained by varying the film thickness. Higher resistances are achieved by the use of a spiral pattern, a coarse spiral for intermediate values and a fine spiral for high resistance. Thus, the inductance is greater in high values, but it is likely to be far less than in wirewound resistors. Special high-frequency units having greatly reduced inductance are available.

Resistive Films

Resistive-material films presently used are micro-crystalline carbon, boron-carbon, and various metallic oxides or precious metals.

Deposited-carbon resistors have a negative temperature coefficient of 0.01 to 0.05 percent/°C for low resistance values and somewhat larger for higher values. Cumulative permanent resistance changes of 1 to 5 percent may result from soldering, overload, low-temperature exposure, and aging. Additional changes up to 5 percent are possible from moisture penetration and temperature cycling.

The introduction of a small percentage of boron into the deposited-carbon film results in a more stable unit. A negative temperature coefficient of 0.005 to 0.02 percent/°C is typical. Similarly, a metallic dispersion in the carbon film provides a negative coefficient of 0.015 to 0.03 percent/°C. In other respects, these materials are similar to standard deposited carbon. Carbon and boron-carbon resistive elements have the highest random noise of the film-type resistors.

Metallic-oxide and precious-metal-alloy films permit higher operating temperatures. Their noise characteristics are excellent. Temperature coefficients are predominantly positive, varying from 0.03 to as little as 0.0025 percent/°C.

Applications

Power ratings of film resistors are based on continuous direct-current operation or on root-mean-square operation. Power derating is necessary for operation at ambient temperatures above the rated temperature. In pulse applications, the power dissipated during each pulse and the pulse duration are more significant than average power conditions. Short high-power pulses may cause instantaneous local heating sufficient to alter or destroy the film. Excessive peak voltages may result in flashover between turns of the film element. Derating under these conditions must be determined experimentally.

Film resistors are fairly stable up to about 10 megahertz. Because of the extremely thin resistive film,

skin effect is small. At frequencies above 10 megahertz, it is advisable to use only unspiraled units if inductive effects are to be minimized (these are available in low resistance values only).

Under extreme exposure, deposited-carbon resistors deteriorate rapidly unless the element is protected. Encapsulated or hermetically sealed units are preferred for such applications. Open-circuiting in storage as the result of corrosion under the end caps has been reported in all types of film resistors. Silver-plated caps and core ends effectively overcome this problem.

Technical Characteristics

Some technical characteristics of film resistors are given in Tables 11, 12, and 13.

The MIL "RN" series of film resistors is more stable than the "RL" series and is available in a wider range of ratings. Commercial equivalents are also offered. Where stability and reliability are desired, the "RN" series is economically very competitive with the "RC" or "RL" series.

RESISTORS—ADJUSTABLE

Adjustable resistors may be divided into three separate and distinct categories, potentiometers, trimmers, and rheostats.

Potentiometers are control devices that are used where the frequency of adjustment is high. They may be operated manually by human effort or mechanically served by machine. Potentiometers are designed with a long mechanical life in view, generally from 10 000 to 100 000 cycles, with certain types having life capabilities in the millions of cycles. A cycle, or excursion, consists of wiper traverse from one limit of travel to the other limit and back.

Trimmers differ from potentiometers in that they are designed to be adjusted infrequently, sometimes only once, and normally exhibit greater setting stability once

TABLE 11. STABLE EQUIVALENTS FOR COMPOSITION RESISTORS; AXIAL LEADS; DATA FOR MIL "RL" SERIES

Watts	¼	½	1	2
Voltage rating	250	350	500	500
Critical resistance (megohms)	0.25	0.25	0.25	0.12
Maximum temp for full load—70 °C; for 0 load—150 °C				
Resistance-temperature characteristic: ±200 ppm/°C maximum				
Life-test resistance change: ±2% maximum				
Moisture resistance test: ±1.5% maximum change				
Resistance values: E24 series, same as composition resistors; tolerances 2% or 5%.				

TABLE 12. HIGH-STABILITY FILM RESISTORS; AXIAL LEADS; DATA FOR MIL "RD" SERIES*

Characteristic	Watts							
	$\frac{1}{20}$	$\frac{1}{10}$	$\frac{1}{8}$	$\frac{1}{4}$	$\frac{1}{2}$	$\frac{3}{4}$	1	2
	Voltage Rating							
B	—	—	—	—	—	—	500	750
D	—	—	200	300	350	500	—	—
C,E	200	200	250	300	350	—	500	—

Characteristic	Characteristic			
	B	C	D	E
Maximum temp (°C):				
Full load	70	125	70	125
0 load	150	175	165	175
Life-test resistance change (max)	±1%	±0.5%	±1%	±0.5%
Resistance-temperature characteristic (max ppm/°C)	±500	±50	+200 -500	±25
Moisture resistance test (max change)	±1.5%	±0.5%	±1.5%	±0.5%

Resistance Values: E96 series (E48 preferred): 1% tolerance;
E192 series (E96 preferred): 0.5%, 0.25%, 0.1% tolerances.

*Uninsulated commercial versions have lower temperature limits and greater resistance change. Color coding is the same as previously indicated for composition resistors.

set. Their employment eliminates the use of expensive precision related components and provides an easily retunable vehicle to compensate for drift or aging in related parts. Normal life designs are rated at approximately 200 excursions. Some typical trimmer characteristics are listed in Table 14.

Unlike potentiometers and trimmers, whose primary function is to control voltage, the rheostat is basically a current-controlling device. Rheostats are made in much the same manner as potentiometers and trimmers, but

with more attention paid to wiper current-carrying ability and generally higher power ratings. Some rheostats are wirewound tubular ceramic power resistors with a track of exposed wire to allow for the adjustable feature. The wiper is in the form of a dimpled clamp band that is screw-tightened at the desired setting. Due to the heat generated by these units, care should be exercised in the circuit location.

Wattage ratings of adjustable resistors apply only when all the resistance is in the circuit. To avoid

TABLE 13. POWER-TYPE FILM RESISTORS, UNINSULATED

	Axial-Lead and MIL "RD" Series			Commercial Tab-Terminal Styles				
Watts	2	4	8	7	23	25	55	115
Voltage rating	350	500	750	525	1380	2275	3675	7875
Critical resistance (kilohms)	61	62	70	39	83	208	245	540
Maximum temperature (°C):								
Full load	25	25	25	25	25	25	25	25
0 load	235	235	235	235	235	235	235	235
	275 (MIL)	275 (MIL)	275 (MIL)					

Life-test resistance change: ±5% maximum
Resistance-temperature characteristic: ±500 ppm/°C maximum
Moisture resistance test: ±3% maximum change
Resistance values: E12 series (approx. for MIL)

Tolerances: Axial lead: $\frac{1}{2}$ %, 1%, 2%, (MIL "RD" Series), 5%, 10%.

Tab lead: 1%, 2%, 5%, 10%, 20%.

TABLE 14. TYPICAL CHARACTERISTICS OF 1/2-INCH TRIMMERS

	Resistance Ranges	Typical Resistance Tolerance	Wattage Rating	Rotational Life	Temperature Coefficient of Resistance
Wirewound	10 Ω - 50 kΩ	±5%	1 W at 70 °C	200	±50 ppm/°C
Cermet	50 Ω - 2 MΩ	±10%	1 W at 70 °C	200	±100 ppm/°C
Conductive Plastic	50 Ω - 2 MΩ	±10%	0.33 W at 50 °C	1 × 10 ⁶	±500 ppm/°C
Carbon Composition	100 Ω - 2.5 MΩ	±20%	0.33 W at 50 °C	200	±10%/°C*

*Temperature characteristic.

overloading any section, never exceed the maximum rated current based on total resistance.

Types of Adjustable Resistors

Wirewound Resistors—Wirewound resistor elements are made by winding a very fine resistance wire precisely around a mandrel. Most resistance wire is made from a nickel-chrome alloy with other elemental additives to enhance its electrical characteristics.

Wirewound adjustable resistors exhibit superior independent linearity characteristics, and it is for this reason that they are frequently specified for direct motion controls. Precision wirewounds are available with independent linearity ratings as low as 0.1%. They are very stable over a wide range of operating temperatures. Panel controls and precision potentiometers are recommended for normal operating temperature ranges from -65° to + 125 °C, and most trimmers will perform satisfactorily from -55° to + 150 °C. Temperature coefficients of resistance as low as ±20 ppm/°C are available.

Wirewound elements, with rare exceptions, change value in steps as the wiper traverses each individual winding. Resolution may be improved by using multi-turn adjustable resistors. This type of construction increases the winding length and decreases incremental resistance steps from one winding to the next.

Wirewound resistive elements are usually not suitable for frequency-sensitive RF circuits because of inductive and capacitive effects. They are impractical above a resistance value determined by the winding space available and the smallest resistance wire that can be space-wound. If infinite resolution is required, then a cermet or composition element must be substituted.

Cermet-Element Resistors—Cermets are a mixture of fine metal-oxide or precious-metal particles and glass in a viscous organic vehicle. This paste is screened onto a ceramic substrate and fired at vitrifying temperatures.

Cermet adjustable resistors are designed for low to moderate adjustment life. They feature infinite resolution and are generally available having temperature coefficients of resistance of ±100 ppm/°C. Sheet re-

sistivities* are available from one ohm per square (1 Ω/□) to one megohm per square (1 MΩ/□). Ultimate resistance values are limited by substrate geometry.**

Carbon Composition Resistors—A mixture of carbon powders and a binder is molded under heat and pressure into a solid mass. In some constructions, the carbon composition is molded at the same time as the plastic substrate. This process is called comolding.

Carbon composition adjustable resistors are the least expensive and most common type of potentiometer for general electronics use. Their low noise makes them ideal for use in live audio controls.

In addition to a linear rotation-vs-resistance characteristic, composition elements can have a wide range of nonlinear output curves. Standard logarithmic curves (tapers) are used as volume controls for radio, TV, etc.

Temporary resistance changes up to ±10% can be expected when these devices are operated near the extreme limits of a temperature range of -55 °C to +120 °C. Their use is not recommended for precision controls or in varying hostile environments.

Conductive Plastic—Conductive plastic is an ink formulated from carbon, other proprietary materials, a resin, and solvent. It is applied to a substrate by screening, dipping, or comolding. The low curing temperature (150 to 300 °C) of the ink allows it to be applied to a wide variety of substrates.

Conductive-plastic potentiometers are most notable for their high rotational life, and they are used most frequently as machine-operated servo-controls. Another desirable feature is their low noise or output smoothness.

Resistance ranges in sheet resistivities of up to 50 000 Ω/□ are available with temperature coefficients of ±500 ppm/°C. As with cermets, ultimate resistance values are limited by substrate geometry.

* Sheet resistivity (Ω/□) is the resistance of a square sheet of material. It is independent of the units of length used because the resistance increases with length but decreases with width.

** For material of given volume resistivity ρ (ohm-meter), sheet resistivity given by ρ/d will increase as depth d is decreased.

Terminal Identification

Industry standards have been developed for identification of potentiometer terminals. Most potentiometer terminals are either numbered or color coded as follows:

- | | |
|----------|--|
| 1 Yellow | Always the counterclockwise element limit. |
| 2 Red | Wiper (or collector). |
| 3 Green | Always the clockwise element limit. |

For rotary potentiometers, clockwise is always defined by viewing from the specified mounting end of the potentiometer.

Mounting Characteristics

Potentiometers usually must be accessible from the outside of a product, and they are often mounted on a panel by means of a threaded bushing. Most precision potentiometers are manufactured with both bushing and servo-mount options. Servo-mount units are secured by servo clamps to assure precise shaft alignment.

Trimmers are usually not accessible from outside the instrument. Most are circuit-board mounted by their terminals and are small in size to conserve space.

CAPACITORS—DEFINITIONS

Dielectric: A dielectric is a medium that can withstand high electric stress without appreciable conduction. When such stress is applied, energy in the form of an electric charge is held by the dielectric. Most of this stored energy is recovered when the stress is removed. The only perfect dielectric in which no conduction occurs and from which the whole of the stored energy may be recovered is a perfect vacuum.

Relative Capacitivity: The relative capacitivity or relative permittivity or dielectric constant is the ratio by which the capacitance is increased when another dielectric replaces a vacuum between two electrodes.

Dielectric Absorption: Dielectric absorption is the absorption of charge by a dielectric when subjected to an electric field by other than normal polarization. This charge is not recovered instantaneously when the capacitor is short-circuited, and a decay current will continue for many minutes. If the capacitor is short-circuited momentarily, a new voltage will build up across the terminals afterward. This is the source of some danger with high-voltage dc capacitors or with ac capacitors not fitted with a discharge resistor. The phenomenon may be used as a measure of dielectric absorption.

Tangent of Loss Angle: This is a measure of the energy loss in the capacitor. It is expressed as $\tan \delta$ and is the power loss of the capacitor divided by its reactive power at a sinusoidal voltage of specified frequency. (This term also includes power factor, loss factor, and dielectric loss. The true power factor is $\cos(90 - \delta)$.)

Insulation Resistance: This is a measure of the conduction in the dielectric. Because this conduction takes a very long time to reach a stable value, it is usually measured after 2 minutes of electrification for nonelectrolytic types and 3 minutes for electrolytics. It is measured preferably at the rated working voltage or at a standardized voltage.

The insulation resistance is usually multiplied by the capacitance to give the ohm-farad value, which is the apparent discharge time constant (seconds). This is a figure of merit for the dielectric, although for small capacitances a maximum value of insulation resistance is usually also specified.

In electrolytics, the conduction is expressed as leakage current at rated working voltage. It is calculated as $\mu A/\mu FV$, which is the reciprocal of the ohm-farad value. In this case, a maximum value of leakage current is specified for small capacitances.

Leakage Current: The current flowing between two or more electrodes by any path other than the interelectrode space is termed the leakage current, and the ratio of this to the test voltage is the insulation resistance.

Impedance: Impedance is the ratio of voltage to current at a specified frequency. At high frequencies, the inductance of leads becomes a limiting factor, in which case a transfer impedance method may be employed. This then measures the impedance of the shunt path only.

DC or AC Capacitor: A dc capacitor is designed to operate on direct current only. It is normally not suitable for use above 200 volts ac because of the occurrence of discharges in internal gas bubbles (corona). An ac capacitor is designed to have freedom from internal discharges and low tangent of loss angle to minimize internal heating.

Rated Voltage and Temperature: The rated voltage is the direct operating voltage that may be applied continuously to a capacitor at the rated temperature.

Category Voltage and Temperature: The category voltage is the voltage that may be applied to the capacitor at the maximum category temperature. It differs from the rated voltage by a derating factor.

Ripple Voltage: If alternating voltages are present in addition to direct voltage, the working voltage of the capacitor is taken as the sum of the direct voltage and the peak alternating voltages. This sum must not exceed the value of the rated voltage.

In electrolytics, the permissible ripple may be expressed as a rated ripple current.

Surge Voltage: This is a voltage above the rated voltage which the capacitor will withstand for a short time.

Voltage Proof Test (Dielectric Strength): This is the highest possible voltage that may be applied without breakdown to a capacitor during qualification approval

testing to prove the dielectric. The repeated application of this voltage may cause failure.

Forming Voltage (Electrolytics): The voltage at which the anodic oxide has been formed. The thickness of the oxide layer is proportional to this voltage.

Burnout Voltage (Metallized Types): The voltage at which metallized types burn out during manufacture.

Self-Healing (Metallized Types): A momentary partial discharge of a capacitor resulting from a localized failure of the dielectric. Burning away the metallized electrodes isolates the fault and effectively restores the properties of the capacitor. The self-healing action is also called "clearing."

Equivalent Series Resistance (ESR): Equivalent series resistance (ESR) is a single resistive value that represents the sum of the ac losses (due to the leads, electrode plates, and junction terminations), the resistive losses due to leakage currents, and the resistive losses due to the inherent molecular polarization dielectric absorption factors of the base dielectric material.

Volt-Ampere Rating (VA): This is the reactive power in a capacitor when an ac voltage is applied. $VA \cos \theta$ gives the amount of heat generated in the capacitor. Since the amount of heat that can be dissipated is limited, the VA must also be limited and in some cases a VA rating is quoted. (Note that $\cos \theta = \cos(90 - \delta) \approx \tan \delta$, when δ is small.)

Scintillation: Minute and rapid fluctuations of capacitance formerly exhibited by silvered mica or silvered ceramic types but overcome by modern manufacturing techniques.

Corona Discharge: Partial discharge of a capacitor due to ionization of the gas in a bubble in the dielectric. On ac or pulse operation, this may occur in dielectric stressed above 200 volts and is a major cause of failure. On dc, such discharges are very infrequent and normally are not a cause of failure.

CLASSES OF CAPACITORS

Modern electronic circuits require the smallest possible capacitors, which are usually made with the thinnest possible dielectric material since they are for operation at low voltages. There are three broad classes of capacitors.

(A) Low-loss capacitors with good capacitance stability. These are usually of mica, glass, ceramic, or a low-loss plastic such as polypropylene or polystyrene.

(B) Capacitors of medium loss and medium stability, usually required to operate over a fairly wide range of ac and dc voltages. This need is met by paper, plastic film, or high- K ceramic types. The first two of these may have electrodes of metal foil or electrodes of evaporated metal which have a self-healing characteristic.

(C) Capacitors of the highest possible capacitance per unit volume. These are the electrolytics, which are normally made either of aluminum or tantalum. Both of these metals form extremely thin anodic oxide layers of high dielectric constant and good electrical characteristics. Contact with this oxide layer is normally by means of a liquid electrolyte that has a marked influence on the characteristics of the capacitor. In solid tantalum, the function of the electrolyte is performed by a manganese-dioxide semiconductor.

PLASTIC FILM CAPACITORS

Advances in organic chemistry have made it possible to produce materials of high molecular weight. These are formed by joining together a number of basic elements (monomers) to produce a polymer. Some of these have excellent dielectric characteristics.

Physically, they can be classified as thermoplastic or thermosetting. In the former case, the molecule consists of long chains with little or no branching, whereas in the latter the molecules are crosslinked. Thermosetting materials have no clearly defined melting point and are usually hard and brittle, making them unsuitable for the manufacture of plastic films. A cast film is usually amorphous, but by extrusion, stretching, and heat treatment, oriented crystalline films are produced with good flexibility and dielectric characteristics.

The electrical properties of the plastics depend on the structure of the molecule. If the molecule is not symmetrical, it will have a dipole moment giving increased dielectric constant. On the other hand, the dielectric constant and $\tan \delta$ are then dependent on frequency. Generally speaking, nonpolar materials have electrical characteristics that are independent of frequency, while polar materials exhibit a decrease in capacitance with increasing frequency, and $\tan \delta$ may pass through a maximum in the frequency range.

Figs. 7, 8, and 9 show some characteristics of several types of capacitors. At the present time, two classes of plastic film capacitors are recognized.

(A) *Polystyrene and Polypropylene Capacitors.* Polystyrene and polypropylene are nonpolar plastics that have excellent electrical characteristics which are independent of frequency.

(B) *Polyester Films.* Strictly speaking, these are the polyethylene terephthalates (Mylar, Melinex, Hostaphan), but the polycarbonates are now included in this group because they have similar electrical characteristics.

Plastic films for capacitor manufacture are usually of the oriented crystalline type because of their good combinations of characteristics. One important feature of some of these films is that they tend to shrink back to their original shape after being heated. This fact is sometimes exploited in manufacturing the capacitor.

Moisture usually has little effect on the dielectric properties of plastic films, and capacitors made from them require less protection than paper or mica types.

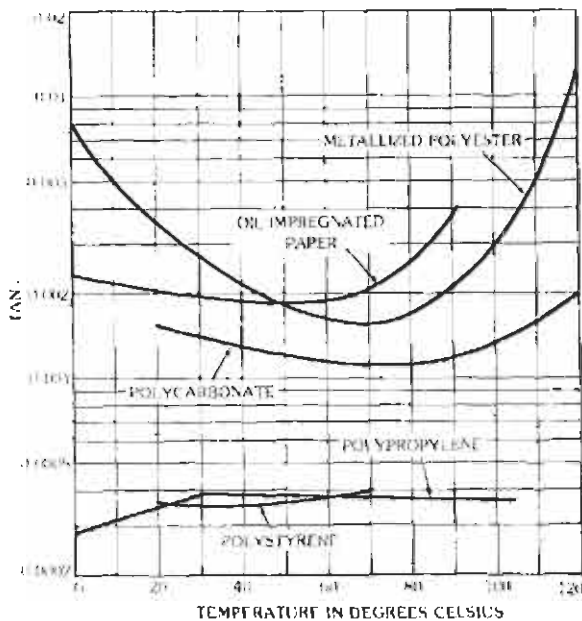


Fig. 7. Variation of $\tan \delta$ with temperature for various plastic films compared with oil-impregnated paper.

This, together with simple processing, has permitted them to be mass-produced at relatively low cost.

The electrical characteristics of capacitors made with these materials depend on the construction employed. The resistance of metallized electrodes and the shape and method of connection to the unit are particularly important.

Many plastic materials are being used in the manufacture of capacitors for which there are no internationally agreed specifications. In this case, it is necessary to obtain the relevant data from the manufacturer. There is no doubt that in the future the range of plastic film capacitors will be considerably extended.

Polystyrene Film Capacitors With Foil Electrodes

Polystyrene has excellent electrical characteristics. The film employed is of the oriented crystalline type, which makes it flexible and suitable for forming into thin films. On heat treatment, the film shrinks considerably, and this is used in the manufacturing process to obtain capacitance stability. The film is affected by greases and solvents, and care must be taken both in manufacture and in use to ensure that capacitors do not come into contact with these materials.

The power factor of polystyrene is low over the whole frequency range, but the resistance of the electrodes may result in an increase of power factor at high frequencies in the larger values, as shown in Table 15.

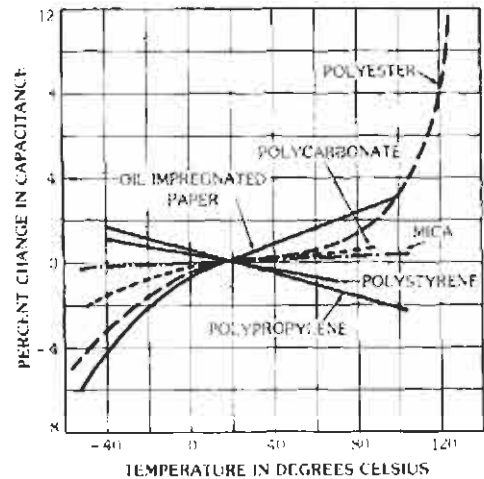


Fig. 8. Typical capacitance characteristics of various capacitors as a function of temperature. Measured on 0.1-microfarad capacitors at 1000 hertz.

Polyester Film Capacitors With Foil Electrodes

The generic title "polyester" is usually used to apply to polyethylene terephthalate. It is a slightly polar plastic film suitable for operation up to temperatures of 125 °C. Capacitors are available with foil electrodes.

Plastic Film Capacitors With Metallized Electrodes

Plastic film capacitors with metallized electrodes have superseded metallized paper capacitors in dc applications because of superior electrical characteristics, less tendency for self-healing to occur during service, higher and more stable insulation resistance, and approximately the same space factors. Three types of film are generally used, polyethylene terephthalate, polycarbonate, and polypropylene. For some purposes these are comparable, but polypropylene has a lower loss angle, and polycarbonate has a smaller change of capacitance with temperature. Polyester and polycarbonate are also available in thinner films, giving an advantage of space factor.

ELECTROLYTIC CAPACITORS

Electrolytic capacitors (Fig. 10) employ for at least one of their electrodes a "valve metal." This metal, when operated in an electrolytic cell as the anode, forms a layer of dielectric oxide. The most commonly used metals are aluminum and tantalum. The valve-metal behavior of these metals was known about 1850. Tantalum electrolytic capacitors were introduced in the

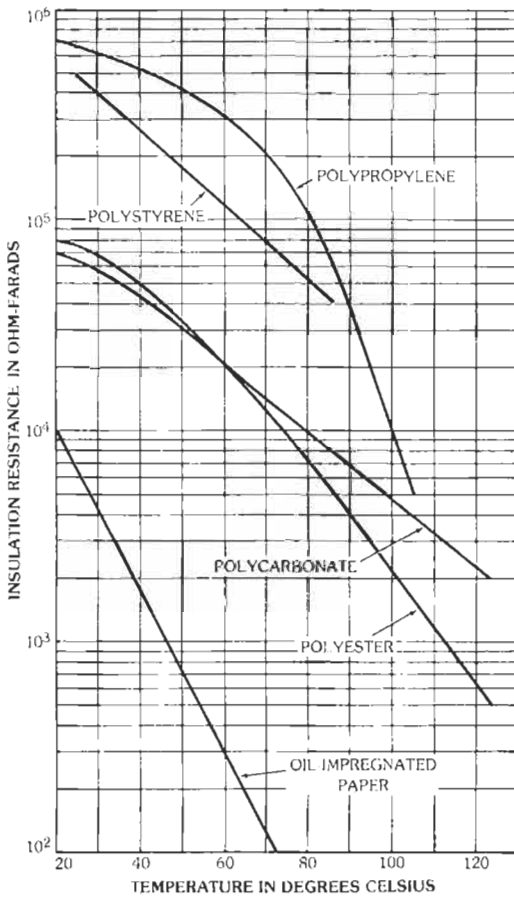


Fig. 9. Variation of insulation resistance with temperature for various plastic films compared with the resistance variation for oil-impregnated paper.

1950s because of the need for highly reliable miniature capacitors in transistor circuits over a wide temperature range. These capacitors were made possible by improved refining and powder metallurgy techniques.

The term "electrolytic capacitor" is applied to any capacitor in which the dielectric layer is formed by an electrolytic method. The capacitor does not necessarily contain an electrolyte.

The oxide layer is formed by placing the metal in a bath containing a suitable forming electrolyte, and applying voltage between the metal as anode and another electrode as cathode. The oxide grows at a rate determined by the current, but this rate of growth decreases until the oxide has reached a limiting thickness determined by the voltage. For most practical purposes, it may be assumed that the thickness of the oxide is proportional to the forming voltage.

Properties of aluminum and tantalum and their oxides are shown in Table 16.

The structure of these oxide layers plays an important part in determining their performance. Ideally they are

TABLE 15. POWER FACTOR OF POLYSTYRENE AT VARIOUS FREQUENCIES

Frequency (hertz)	Nominal Capacitance (pF)			
	Up to 1 000	1 000 to 10 000	10 000 to 100 000	Above 100 000
800	0.0003	0.0003	0.0003	0.0003
10 000	0.0003	0.0003	0.0003	0.001
100 000	0.0003	0.0005	0.001-0.003	—
1 000 000	0.001	0.002	0.005-0.02	—

amorphous, but aluminum tends to form two distinct layers, the outer one being porous. Tantalum normally forms an amorphous oxide which, under conditions of a high field strength of the oxide layer, may become crystalline. Depending on the forming electrolyte and the surface condition of the metal, there is an upper limit of voltage beyond which the oxide breaks down. The working voltage is between 25 and 90 percent (according to type) of the forming voltage at which stable operation of the oxide layer can be obtained.

To produce a capacitor, it is necessary to make contact to the oxide layer on the anode, and there are two distinct methods of doing this. The first is to use a working electrolyte that has sufficient conductivity over the temperature range to give a good power factor. There are many considerations in choosing the working electrolyte, and the choice is usually a compromise between high- and low-temperature performance. The working electrolyte also provides a rehealing feature in

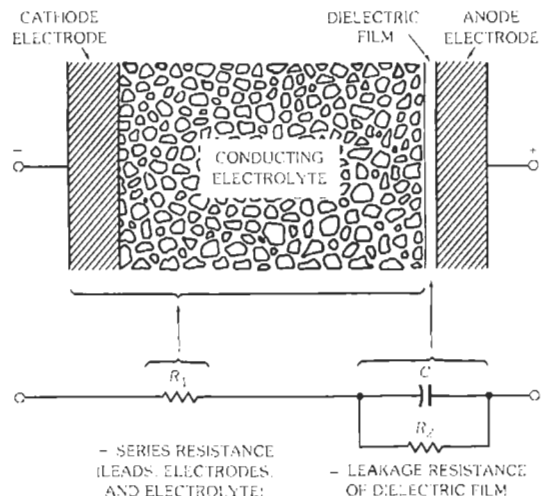


Fig. 10. Basic cell and simplified equivalent circuit for polar electrolytic capacitor.

TABLE 16. ALUMINUM AND TANTALUM PROPERTIES

Metal	Density	Principal Oxide	Dielectric Constant	Thickness (Å/V)
Aluminum	2.7	Al ₂ O ₃	8	13.5
Tantalum	16.6	Ta ₂ O ₅	27.6	17

that any faults in the oxide layer will be repaired by further anodization.

In aluminum electrolytic capacitors, the working electrolyte must be restricted to those materials in which aluminum and its oxide are inert. Corrosion can be minimized by using the highest possible purity of aluminum. This also reduces the tendency of the oxide layer to dissolve in the electrolyte, giving a better shelf life.

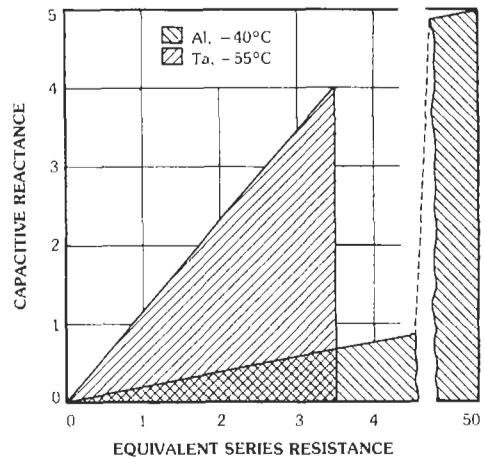
Tantalum, on the other hand, is very inert and therefore allows a wider choice of electrolyte. Since there is no gas evolution, better methods of sealing can be employed. The characteristics of aluminum and tantalum electrolytic capacitors are shown in Figs. 11 and 12.

A major problem with all electrolytic capacitors is to ensure that the electrolyte is retained within the case under all operating conditions. In the aluminum capacitor, allowance must be made for gas evolution on reforming. Even the tantalum capacitor usually must employ only organic materials for sealing, and these do not provide completely hermetic sealing. All organic materials have finite moisture transmission properties, and, therefore, at the maximum category temperature the high vapor pressure of the electrolyte results in some diffusion.

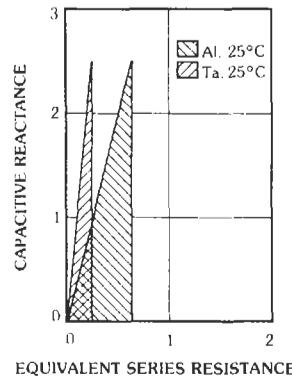
An elegant solution to this problem, using a semiconductor instead of an electrolyte, was found by Bell Telephone Laboratories. The semiconductor is manganese dioxide in a polycrystalline form and has a higher conductivity than conventional electrolyte systems. This material also provides a limited self-healing feature at a fault, resulting in oxidation of the tantalum and reduction of the manganese dioxide to a nonconducting form.

Electrolytic capacitors take many forms, and the anode may be of foil, wire, or a porous sintered body. The foil may be either plain or etched. The porous body may be made with fine or with coarse particles, and the body itself may be short and fat or long and thin. The aluminum-foil capacitor has a space factor about six times better than that of the equivalent paper capacitor, whereas tantalum capacitors are even smaller and enjoy a space factor up to 20 times better.

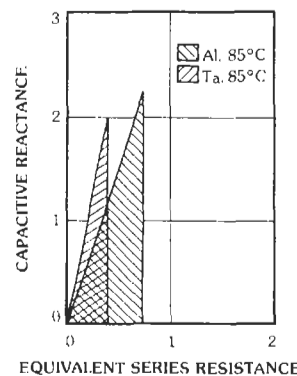
When electrolytic capacitors are operated in series-parallel, stabilizing resistors should be used to equalize the voltage distribution. It should also be noted that, even when the case is not connected to one terminal, a low resistance path exists between it and the electrodes. The case must be insulated from the chassis, particularly if the chassis and the negative terminal are not at the same potential.



(A) Low temperature.



(B) Room temperature.

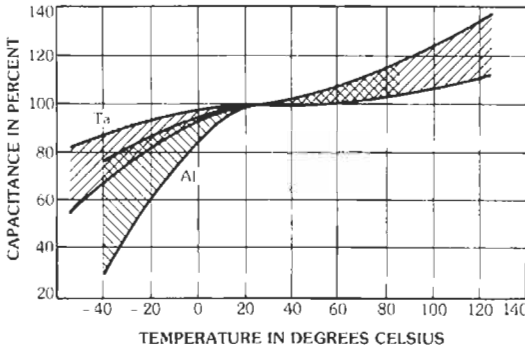


(C) High temperature.

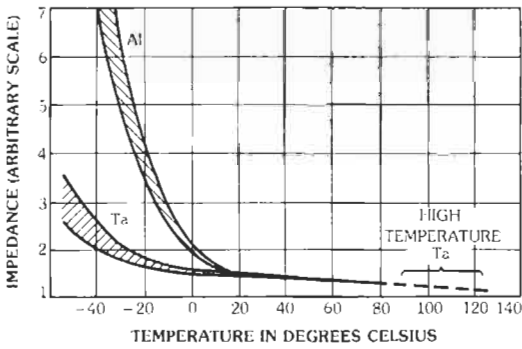
Fig. 11. Typical 120-hertz impedance diagrams for aluminum (Al) and tantalum (Ta) plain-foil polar electrolytic capacitors of 150-volt rating.

Aluminum Electrolytics

The aluminum type is the most widely known electrolytic capacitor and is used extensively in radio and television equipment. It has a space factor about six times better than the equivalent paper capacitor. Types



(A) Capacitance.



(B) Impedance.

Fig. 12. Capacitance and 120-hertz impedance as a function of temperature for aluminum (Al) and tantalum (Ta) electrolytic capacitors.

of improved reliability are now available using high-purity (better than 99.99%) aluminum.

Conventional aluminum electrolytic capacitors which have gone six months or more without voltage applied may need to be reformed. Rated voltage is applied from a dc source with an internal resistance of 1500 ohms for capacitors with a rated voltage exceeding 100 volts, or 150 ohms for capacitors with a rated voltage equal to or less than 100 volts. The voltage must be applied for one hour after reaching rated value with a tolerance of ± 3 percent. The capacitor is then discharged through a resistor of 1 ohm/volt.

Tantalum-Foil Electrolytics

The tantalum-foil type of capacitor was introduced around 1950 to provide a more reliable type of electrolytic capacitor without shelf-life limitation. It was made possible by the availability of thin, high-purity annealed tantalum foils and wires. Plain-foil types were introduced first, followed by etched types. The purity, and particularly the surface purity, of these materials plays a major part in determining the leakage current and their ability to operate at the higher working voltages.

These capacitors are smaller than their aluminum counterparts and will operate at temperatures up to about 125 °C (Figs. 13-15). The plain-foil types usually exhibit less variation of capacitance with temperature or frequency.

Tantalum Electrolytics With Porous Anode and Liquid Electrolyte

The tantalum electrolytic with porous anode and liquid electrolyte was the first type of tantalum electrolytic capacitor to be introduced and still has the best space factor. Types using sulphuric-acid electrolyte have excellent electrical characteristics up to about 70 working volts. Other types contain neutral electrolytes.

Basically, this type of capacitor consists of a sintered porous anode of tantalum powder housed in a silver or silver-plated container. The porous anode is made by pressing a high-purity tantalum powder into a cylindrical body and sintering in vacuum at about 2000 °C.

Tantalum Electrolytics With Porous Anode and Solid Electrolyte

The so-called "solid" tantalum capacitor originally developed by Bell Telephone Laboratories developed from the porous-anode type with liquid electrolyte by replacing the liquid with a semiconductor. This overcame the problem of sealing common to all other types of electrolytic capacitors. Since there is no liquid electrolyte, it is possible to use a conventional hermetic seal.

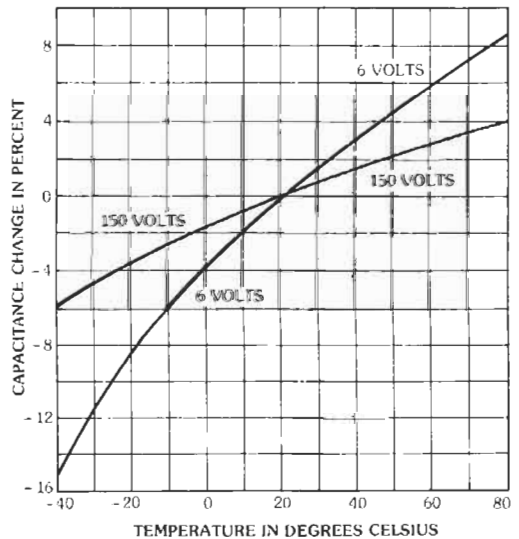


Fig. 13. Variation of capacitance with temperature for plain tantalum-foil electrolytic capacitors.

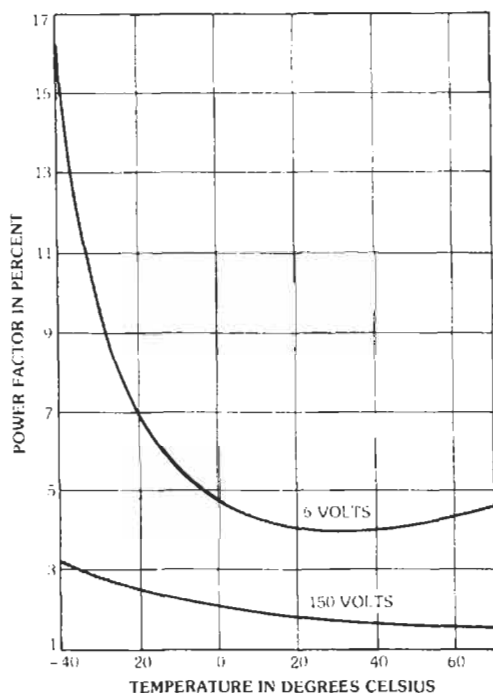


Fig. 14. Variation of power factor with temperature for plain tantalum-foil electrolytic capacitors.

CERAMIC CAPACITORS

Ceramic capacitors are defined in classes based on their distinct and inherent electrical properties.

Class I

Class I capacitors are stable temperature-compensating capacitors that have essentially linear characteristics with properties independent of frequency over the normal range. Materials are usually magnesium titanate for positive temperature coefficient of capacitance and calcium titanate for negative temperature coefficient of capacitance. Combinations of these and other materials produce a dielectric constant of 5 to 150 and temperature coefficient of capacitance of +150 to -4700 ppm/°C with tolerances of ± 15 ppm/°C.

Low- K ceramics are suitable for resonant-circuit or filter applications, particularly where temperature compensation is a requirement. Disc and tubular types are the best forms for this purpose. Stability of capacitance is good, being next to that of mica and polystyrene capacitors.

Class II

Class II includes ceramic dielectrics suitable for fixed capacitors used for bypass, coupling, and decoupling. This class is usually divided into two subgroups for which the temperature characteristics define the characteristics.

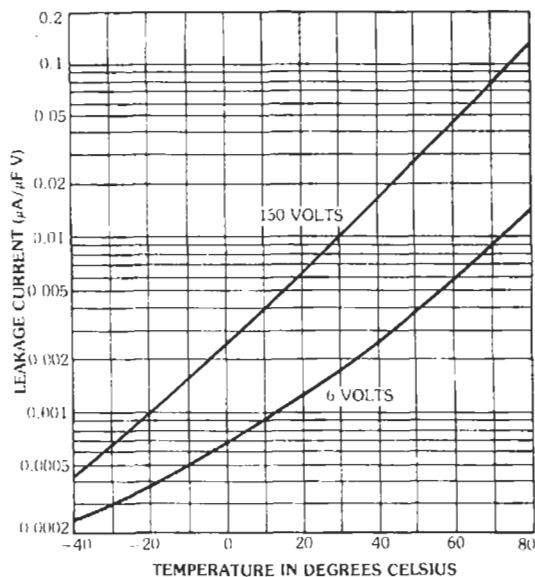


Fig. 15. Variation of leakage current with temperature for plain tantalum-foil electrolytic capacitors.

- (A) Embody stable K values of 250 to 2400 over a temperature range of -55 °C to $+125$ °C with maximum capacitance change 15% from a 25 °C ambient.
- (B) Embody combinations of materials, including titanates, with K values of 3000 to 12 000 made by keeping the Curie point near room ambient.

The high- K materials are the ferroelectrics. Because of their crystal structure, they sometimes have very high values of internal polarization, giving very high effective dielectric constants. In this way, these materials are comparable with ferromagnetic materials. Above the Curie temperature, a change of domain structure occurs that results in a change of electrical characteristics. This region is known as the paraelectric region. In common with the ferromagnetic materials, a hysteresis effect is apparent, and this makes the capacitance voltage-dependent.

The ferroelectrics are based on barium titanate, which has a peak dielectric constant of 6000 at the Curie point of 120 °C. Additions of barium stannate, barium zirconate, or magnesium titanate reduce this dielectric constant but make it more uniform over the temperature range. Thus a family of materials can be obtained with a Curie point at about room temperature and with the dielectric constant falling off on either side. The magnitude of this change increases with increasing dielectric constant. These materials exhibit a decrease of capacitance with time and, as a result of the hysteresis effect, with increasing voltage.

Inductance in the leads and element causes parallel resonance in the megahertz region. Care is necessary in

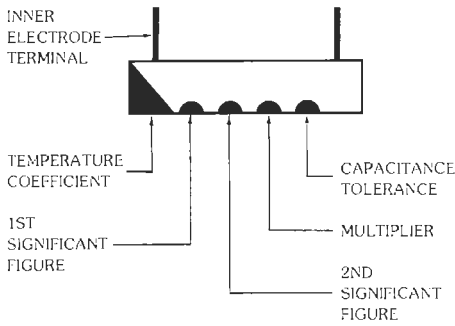
their application above about 50 megahertz for tubular styles and about 500 megahertz for disc types.

Class III

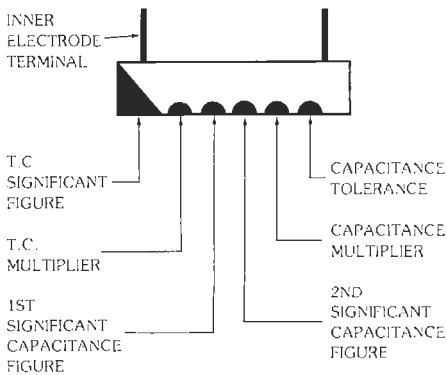
Class III includes reduced barium titanate in which a reoxidized layer or diffusion zone is the effective dielectric. Also included are internal grain boundary layer strontium titanate capacitors in which an internal insulating layer surrounds each grain. These capacitors are suitable for use in low-voltage circuits for coupling and bypass where low insulation resistance and voltage coefficient can be tolerated.

Color Code

The significance of the various colored dots for EIA Standard RS-198-B (ANSI C83.4-1972) fixed ceramic dielectric capacitors is explained by Figs. 16 and 17 and may be interpreted from Table 17.



(A) Five-dot system.



(B) Six-dot system.

Fig. 16. Color coding of EIA Class-I ceramic dielectric capacitors. See Table 17 for color code. Tubular style shown to illustrate identification of inner electrode. For disc or plate styles, color code will read from left to right as observed with lead wires downward.

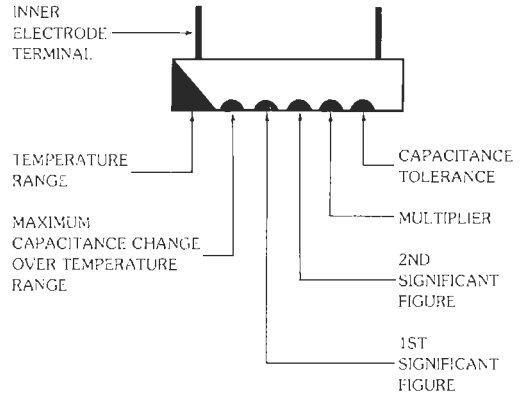


Fig. 17. Color coding of EIA Class-II ceramic dielectric capacitors. See Table 17 for color code. Tubular style shown to illustrate identification of inner electrode. For disc or plate styles, color code will read from left to right as observed with lead wires downward.

Temperature Coefficient

Standard temperature coefficients of capacitance expressed in parts per million per degree Celsius are: +150, +100, +33, 0, -33, -75, -150, -220, -330, -470, -750, -1500, -2200, -3300, and -4700.

PAPER FOIL-TYPE CAPACITORS

In general, paper capacitors have been largely replaced by plastic film capacitors in both foil and metallized constructions and in most dc electronic circuits. Impregnated paper capacitors for ac power applications are also being replaced with plastic film capacitors with both impregnated and dry constructions.

Construction usually consists of multiple layers of paper as the dielectric interleaved with aluminum or tin alloy foil electrodes. Termination is either tabs inserted during winding or leads welded to extended foils.

Paper capacitors are impregnated with stabilized waxes or oils. Most chlorinated materials, such as chlorinated diphenyls, have been banned by the EPA as impregnants. The impregnant usually has a dielectric constant of about 6 and is used for size reduction and improvement in corona start and to reduce internal discharge in ac power applications.

Included among dc applications are coupling, decoupling, bypass, smoothing filters, power-separating filters, energy-storage capacitors, etc. Included among ac applications are motor start, fluorescent lighting, interference suppression, power-factor correction, power-line coupling, distribution capacitors for high-voltage switching gear, capacitor voltage dividers for ac measurement, etc.

TABLE 17. COLOR CODE FOR CERAMIC DIELECTRIC CAPACITORS, CLASSES I AND II*

Color	Digit	Multiplier	Class I					Class II			
			Capacitance Tolerance		Temperature Coefficient ppm/°C (5-Dot System)	Temperature Coefficient Significant Figure (6-Dot System)	Temperature Coefficient Multiplier (6-Dot System)	Capacitance Tolerance (%)	Temperature Range (°C)	Maximum Capacitance Change Over Temperature Range (%)	
			10 pF or less (pF)	Over 10 pF (%)							
Black	0	1	±2.0	±20	0	0.0	-1	±20	—	±2.2	
Brown	1	10	±0.1	±1	-33	—	-10	—	+10 to +85	±3.3	
Red	2	100	—	±2	-75	1.0	-100	—	-55 to +125	±4.7	
Orange	3	1000	—	±3	-150	1.5	-1000	—	+10 to +65	±7.5	
Yellow	4	10 000	—	—	-220	2.2	-10 000	GMV	—	±10	
Green	5	—	±0.5	±5	-330	3.3	+1	±5	—	±15	
Blue	6	—	—	—	-470	4.7	+10	—	—	±22	
Violet	7	—	—	—	-750	7.5	+100	—	—	+22, -33	
Gray	8	0.01	±0.25	—	+150 to -1500	(-1000 to -5200 ppm/°C, With Black Multiplier)	+1000	+80, -20	—	+22, -56	
White	9	0.1	±1.0	±10	+100 to -750	—	+10 000	±10	—	+22, -82	
Silver	—	—	—	—	—	—	—	—	-30 to +85	±1.5	
Gold	—	—	—	—	—	—	—	—	-55 to +85	±1	

* EIA Standard RS-198-B (ANSI C83.4-1972). This standard classifies ceramic dielectric, fixed capacitors as follows:

Class I —Temperature compensating ceramics suited for resonant circuit or other applications where high Q and stability of capacitance characteristics are required.

Class II —Ceramics suited for bypass and coupling applications, or for frequency discriminating circuits where high Q and stability of capacitance characteristics are not of major importance.

Class III—Low-voltage ceramics specifically suited for transistorized or other electronic circuits for bypass, coupling, or frequency determination where dielectric losses, high insulation resistance, and capacitance stability are not of major importance.

Note: Where size permits, EIA Class-III ceramics are typographically marked as follows:

- (1) Capacitance value in microfarads. (2) Rated voltage. (3) Manufacturer's mark or EIA source code.
- (4) Capacitance value tolerance or appropriate code letter; either ±20% (Code M) or +80, -20% (Code Z).
- (5) Temperature stability code (see the EIA or ANSI Standard).

MICA CAPACITORS

Mica capacitors fall within the classification of low loss and good capacitance stability. Mica is one of the earliest dielectric materials used and has an unrivaled combination of physical and electrical characteristics. It is of mineral origin and, because of its monoclinic structure, can be readily split into thin plates. It has a dielectric constant of about 6 (largely independent of frequency) together with a very low loss.

Construction

Mica capacitor constructions of the eyelet, molded, bonded, and button styles have been largely replaced by precision polystyrene, polypropylene, and polyester plastic film capacitors and by ceramic capacitors. The main surviving style is the dipped epoxy-coated radial-leaded capacitor using a clamp-type silvered mica stack with tin electrodes.

Applications

Their low temperature coefficient of capacitance and good stability with temperature and frequency make mica capacitors a good choice for critical precision circuitry such as filter applications.

Type Designation

A comprehensive numbering system, the type designation, is used to identify mica capacitors. Type designations are of the form shown in Fig. 18. MIL specifications now require type designation marking. Color coding is now used only for EIA standard capacitors.

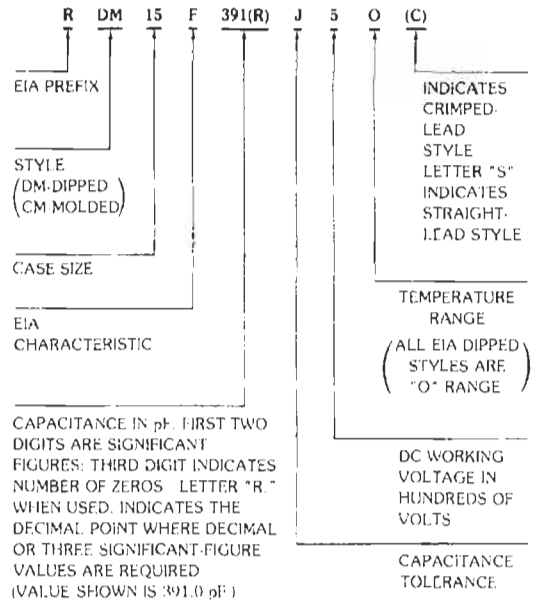
Component Designation—Fixed mica-dielectric capacitors are identified by the symbol CM. For EIA, a prefix letter R is always included, and dipped types are identified by the symbol DM.

Case Designation—The case designation is a two-digit symbol that identifies a particular case size and shape.

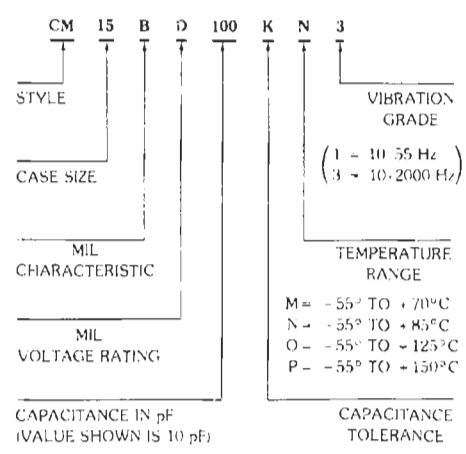
Characteristic—The EIA or MIL characteristic is indicated by a single letter in accordance with Table 18.

Capacitance Value—The nominal capacitance value in picofarads is indicated by a three-digit number. The first two digits are the first two digits of the capacitance value in picofarads. The final digit specifies the number of zeros that follow the first two digits. For EIA, if more than two significant figures are required, an additional digit is used, and the letter "R" is inserted to designate the decimal position.

Capacitance Tolerance—The symmetrical capacitance tolerance in percent is designated by a letter as shown in Table 1.



(A) EIA.



(B) MIL.

Fig. 18. Type designation for mica-dielectric capacitors.

Voltage Rating—MIL voltage ratings are designated by a single letter as follows. A = 100, B = 250, C = 300, D = 500, E = 600, F = 1000, G = 1200, H = 1500, J = 2000, K = 2500, L = 3000, M = 4000, N = 5000, P = 6000, Q = 8000, R = 10 000, S = 12 000, T = 15 000, U = 20 000, V = 25 000, W = 30 000, and X = 35 000 volts. EIA dc working voltage is a number designating hundreds of volts.

Temperature Range—MIL specifications provide for four temperature ranges, all of which have a lower limit of -55 °C; the upper limits are M = +70, N = +85, O = +125, and P = +150 °C. The EIA uses only N and O, which are identical to the MIL standard.

TABLE 18. FIXED-MICA-CAPACITOR REQUIREMENTS BY EIA AND MIL CHARACTERISTIC

EIA or MIL Characteristic	Maximum Capacitance Drift	Maximum Range of Temperature Coefficient (ppm/°C)
B	Not Specified	Not Specified
C	$\pm(0.5\% + 0.1 \text{ pF})$	± 200
D	$\pm(0.3\% + 0.1 \text{ pF})$	± 100
E	$\pm(0.1\% + 0.1 \text{ pF})$	-20 to +100
F	$\pm(0.05\% + 0.1 \text{ pF})$	0 to +70

Vibration Grade—The MIL vibration grade is a number, 1 corresponding to vibration from 10 to 55 hertz at 10g for 4.5 hours and 3 corresponding to 10 to 2000 hertz at 20 g for 12 hours.

Capacitance

Capacitance is measured at 1 megahertz for capacitors of 1000 picofarads or smaller. Larger capacitors are measured at 1 kilohertz.

Temperature Coefficient

Measurements to determine the temperature coefficient of capacitance and the capacitance drift are based on one cycle over the following temperature values (all in degrees Celsius): +25, -55, -40, -10, +25, +45, +65, +70, +85, +125, +150, +25. Measurements at +85, +125, and +150 are not made if these values are not within the applicable temperature range of the capacitor.

Dissipation Factor

The EIA and MIL specifications require that for molded and dipped capacitors the dissipation factor not exceed the values shown in Fig. 19. For potted and cast epoxy capacitors, the dissipation factor shall not exceed 0.35 percent from 1 to 1000 picofarads and 0.15 percent above 1000 picofarads.

High-Potential or Withstanding-Voltage Test

Molded or dipped mica capacitors are subjected to a test potential of twice their direct-current voltage rating.

Humidity and Thermal-Shock Tests

EIA Standard RS-153-B capacitors must withstand 5 cycles of -55, +25, +85, or +125 (as applicable), and +25 degrees Celsius thermal shock followed by a humidity test of 10 cycles (each of 24 hours) given for

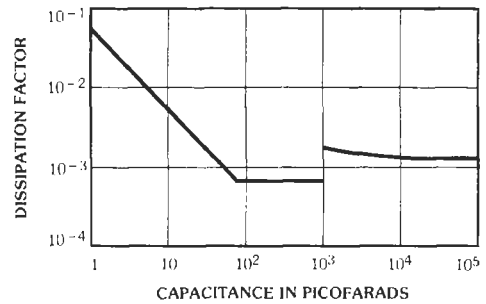


Fig. 19. EIA and MIL maximum dissipation factor at 1 megahertz for capacitance of 1000 picofarads or less and at 1 kilohertz for capacitance greater than 1000 picofarads.

EIA Standard RS-186-D in Fig. 1. Units must pass a withstanding-voltage test. Capacitance may not change by more than 1.0 percent or 1.0 picofarad, whichever is greater. Insulation resistance must meet or exceed 30 percent of the initial requirements at 25 °C (50 000 megohms for capacitances of 20 000 picofarads or less; 1000 ohm-farads for larger capacitances).

MIL Specification MIL-C-5D capacitors must withstand 5 cycles of -55; +25; +85, +125, or +150 (as applicable); and +25 degrees Celsius thermal shock followed by a humidity test of 10 cycles (each of 24 hours) given for MIL-STD-202E in Fig. 1. Units must pass a withstanding-voltage test. Capacitance may not change by more than $\pm(0.2 \text{ percent} + 0.5 \text{ picofarad})$. Insulation resistance must meet or exceed 30 percent of the initial requirements at 25 °C (100 000 megohms for 10 000 picofarads or less; 1000 megohm-microfarads for larger capacitances).

Life

Capacitors are given accelerated life tests at 85 degrees Celsius with 150 percent of rated voltage applied for 2000 hours for MIL specification or 250 hours for EIA standard. If capacitors are rated above +85 °C, the test will be at their maximum rated temperature.

PRINTED CIRCUITS

A printed circuit is a conductive circuit pattern on one or both sides of an insulating substrate. Multilayer boards with tens of levels of circuitry are manufactured with conducting thru-holes to interconnect the circuitry levels. The conductive pattern can be formed by any of several techniques after which component lead holes are drilled or punched in the substrate and components are installed and soldered in place. Printed-circuit construction is ideal for assembly of circuits that employ miniature solid-state components. Its advantages over conventional chassis and point-to-point wiring include:

(A) Considerable space savings over conventional construction methods is usually a result.

(B) A complex circuit may be modularized by using several small printed circuits instead of a single larger one. Modularization simplifies troubleshooting, circuit modification, and mechanical assembly in an enclosure.

(C) Soldering of component leads may be accomplished in an orderly sequence by hand or by dip or wave soldering.

(D) A more uniform product is produced because wiring errors are eliminated and because distributed capacitances are constant from one production unit to another.

(E) The printed-circuit method of construction lends itself to automatic assembly and testing.

(F) Using appropriate base metals, flexible cables or flexible circuits can be built.

(G) By using several layers of circuits (in proper registry) in a sandwich construction, with the conductors separated by insulating layers, relatively complex wiring can be provided.

Printed-Circuit Base Materials

Rigid printed-circuit base materials are available in thicknesses varying from $\frac{1}{64}$ to $\frac{1}{2}$ inch. The important properties of the usual materials are given in Table 19. For special applications, other rigid or flexible materials are available as follows:

(A) Glass-cloth Teflon (polytetrafluoroethylene, PTFE) laminate.

(B) Kel-F (polymonochlorotrifluoroethylene) laminate

(C) Silicone rubber (flexible)

(D) Glass-mat-polyester-resin laminate.

(E) Teflon film.

(F) Ceramic.

The most widely used base material is NEMA-XXXXP paper-base phenolic.

Conductor Materials

Copper is used almost exclusively as the conductor material, although silver, brass, and aluminum also have been used. The common thicknesses of foil are 0.0014 inch (1 oz/ft²) and 0.0028 inch (2 oz/ft²). The current-carrying capacity of a copper conductor may be determined from Fig. 20.

Manufacturing Processes

The most widely used production methods are:

(A) Etching process, wherein the desired circuit is printed on the metal-clad laminate by photographic, silk-screen, photo-offset, or other means, using an ink or lacquer resistant to the etching bath. The board is then placed in an etching bath that removes all of the unprotected metal (ferric chloride is a commonly used mordant for copper-clad laminates). After the etching is

completed, the ink or lacquer is removed to leave the conducting pattern exposed.

(B) Plating process, wherein the designed circuit pattern is printed on the unclad base material using an electrically conductive ink, and, by electroplating, the conductor is built up to the desired thickness. This method lends itself to plating through punched holes in the board for making connections from one side to the other.

(C) Other processes, including metal spraying and die stamping.

Circuit-Board Finishes

Conductor protective finishes are required on the circuit pattern to improve shelf-storage life of the circuit boards and to facilitate soldering. Some of the most widely used finishes are:

(A) Hot-solder coating (done by dip-soldering in a solder bath) is a low-cost method and gives good results where coating thickness is not critical.

(B) Silver plating used as a soldering aid but is subject to tarnishing and has a limited shelf life.

(C) Hot-rolled or plated solder coat gives good solderability and uniform coating thickness.

(D) Other finishes for special purposes are gold plate, for corrosion resistance and solderability, and electroplated rhodium over nickel, for wear resistance. Insulating coatings such as acrylic, polystyrene, epoxy, or silicone resin are sometimes applied to circuit boards to improve circuit performance under high humidity or to improve the anchorage of parts to the board. Conformal coatings are relatively thick and tend to smooth the irregular contour of the mounted items; they add less mass than encapsulation. A protective organic coating (unless excessively thick) will not improve the electrical properties of an insulating base material during long exposure to high humidity. On two-sided circuit boards, where the possibility of components shorting out the circuit patterns exists, a thin sheet of insulating material is sometimes laminated over the circuit before the parts are inserted.

Design Considerations

Before a printed-circuit layout is made, the circuit must be breadboarded and tested under the anticipated final operating conditions. This procedure will permit operating deficiencies and quirks to be detected and corrected before the time-consuming process of producing the circuit board is begun. It is important to note that certain circuits may operate differently on a printed-circuit board than on a breadboard, and appropriate corrective steps may be necessary. For example, inductive coupling between foil patterns may cause

TABLE 19. PROPERTIES OF TYPICAL PRINTED-CIRCUIT DIELECTRIC BASE MATERIALS

Material	Comparable MIL Type	Punchability	Mechanical Strength	Moisture Resistance	Insulation	Arc Resistance	Abrasive Action on Tools	Max Temperature (°C)*
NEMA type XXXP paper-base phenolic	—	Good	Good	Good	Good	Poor	No	105
NEMA type XXXPC paper-base phenolic	—	Very good	Good	Very good	Good	Poor	No	105
NEMA type FR-2 paper-base phenolic, flame resistant	—	Very good	Good	Very good	Good	Poor	No	105
NEMA type FR-3 paper-base epoxy, flame resistant	PX	Very good	Very good	Very good	Very good	Good	No	105
NEMA type FR-4 glass-fabric-base epoxy, general purpose, flame resistant	GF	Fair	Excellent	Excellent	Excellent	Very good	Yes	130 (125)
NEMA type FR-5 glass-fabric-base epoxy, temperature and flame resistant	GH	Fair	Excellent	Excellent	Excellent	Very good	Yes	155 (150)
NEMA type G-10 glass-fabric-base epoxy, general purpose	GE	Fair	Excellent	Excellent	Excellent	Very good	Yes	130 (125)
NEMA type G-11 glass-fabric-base epoxy, temperature resistant	GB	Poor	Excellent	Excellent	Excellent	Very good	Yes	155 (150)
Glass-fabric-base polytetrafluoroethylene	GT	—	Good	Excellent	Excellent	Excellent	—	(150)
Glass-fabric-base fluorinated ethylene propylene	FEP	—	Good	Excellent	Excellent	Excellent	—	(150)

* MIL-STD-275C rating shown in parentheses if different from industry rating.

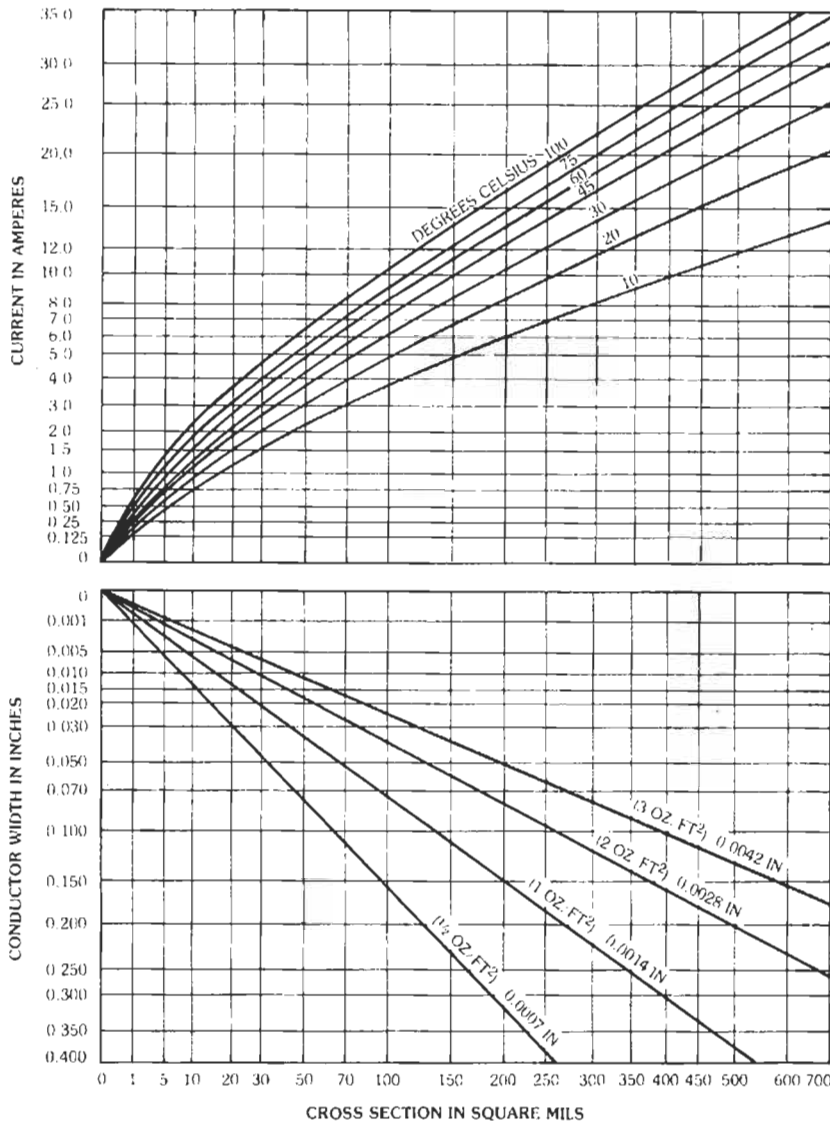


Fig. 20. Current-carrying capacity and sizes of etched copper conductors for various temperature rises above ambient. From MIL-STD-275C, 9 January 1970.

unwanted oscillation in high-frequency or amplifier circuits.

All features (terminal areas, contacts, board boundaries, holes, etc.) should be arranged to be centered at the intersections of a 0.100-, 0.050-, or 0.025-inch rectangular grid, with preference in the order stated. Many components are available with leads spaced to match the standard grids. Devices with circular lead configurations and a few other multilead devices are exceptions that require special attention and dimensioning. Following this grid-layout principle simplifies drafting and subsequent machine operations in board manufacture and assembly.

Drilled holes must be employed if the stated requirements for punched holes cannot be met, or if the

material is not of a punching grade. Drilling is less detrimental to the laminate surrounding the hole; punching may cause crazing or separation of the laminate layers.

The diameter of punched holes in circuit boards should not be less than $\frac{2}{3}$ the thickness of the base material.

The distance between punched holes or between holes and the edge of the material should not be less than the material thickness.

Punched-hole tolerance should not be less than ± 0.005 inch on the diameters.

Hole sizes should not exceed by more than 0.020 inch the diameter of the wire to be inserted in the hole. With smaller holes, hand insertion of the wire is

difficult. Machine insertion requires the larger allowance. Clinching of the lead is desirable if the clearance is larger.

Tolerances with respect to the true-grid location for terminal area centers and for locating edges of boards or other locating features (datums) should not exceed on the board: 0.014 inch diameter for conductor widths and spacings above 0.031 inch; 0.010 inch diameter for conductor widths and spacings 0.010 to 0.031 inch, inclusive. Tolerances on other dimensions (except conductor widths and spacings) may be larger. Closer tolerances may be needed if machine insertion is required.

Terminal area diameters should be at least (A) 0.020 inch larger than the diameter of the flange or projection of the flange on eyelets or standoff terminals, or the diameter of a plated-through hole, and (B) 0.040 inch larger than the diameter of an unsupported hole. Since the terminal area should be unbroken around the finished hole, the diameter should be further increased over the above minimum to allow for the permitted hole-position tolerance.

Conductor widths should be adequate for the current carried. See Fig. 20. For a given conductor-width and copper-thickness intersection, proceed vertically to the allowable temperature-rise line and then horizontally to the left to determine the permissible current. An additional 15% derating is recommended for board thicknesses of $\frac{1}{32}$ inch or less, or for conductors thicker than 0.004 inch (3 oz). The normal ambient temperature surrounding the board plus the allowable temperature rise should not exceed the maximum safe operating temperature of the laminate. For ordinary work, copper conductor widths of 0.060 inch are convenient; with high-grade technique (extra cost), conductor widths as small as 0.010 inch can be readily produced.

Conductor spacing requirements are governed by the applied voltage, the maximum altitude, the conductor protective coating used, and the power-source size. The guide in Table 20 is suggested.

Preparation of Artwork

In preparing the master artwork for printed circuits, careful workmanship and accuracy are important. When circuits are reproduced by photographic means, much retouching time is saved if care is taken with the original artwork.

Artwork should be prepared on a dimensionally stable material. Tracing paper and bristol board are now outmoded, and specially treated (toothed) polyethylene terephthalate (Mylar, Cronar) base drafting films are used for most printed-circuit layouts. The layout pattern may be produced by one of the following methods:

1. Hand application of opaque, permanent black ink
2. Pressure-sensitive tape
3. Hand-cut stencil made from self-adhesive opaque film
4. Preformed self-adhesive layout patterns

Artwork should be prepared to a scale that is two to five times oversize. Photographic reduction to final negative size should be possible, however, in one step.

Avoid the use of sharp corners when laying out the circuit. See Fig. 21.

The centers of holes to be manually drilled or punched in the circuit board should be indicated by a circle of $\frac{1}{32}$ -inch diameter (final size after reduction). See Fig. 22. This feature is not needed on each board if templates or numerically controlled machine tools are

TABLE 20. CONDUCTOR SPACINGS*

Voltage Between Conductors	Minimum Spacing Between Conductors (Inches)		
	Uncoated Boards		Conformal Coated Boards
	Sea Level to 10 000 Ft	Over 10 000 Ft	All Altitudes
0-30	0.025	0.025	0.010
0-50	0.025	0.025	—
0-150	0.025	—	—
31-50	0.025	0.025	0.015
51-100	0.025	0.060	0.020
101-170	—	0.125	0.030
101-300	—	—	0.030
151-300	0.050	—	0.030
171-250	0.050	0.250	0.030
251-500	—	0.500	—
301-500	0.100	0.500	0.060
Above 500	0.0002	0.001	0.00012
	per volt	per volt	per volt

* From MIL-STD-275C

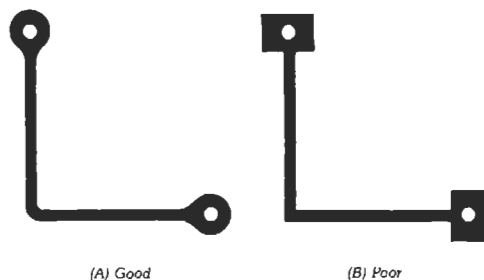


Fig. 21. Design of bends for printed-circuit conductors.

used for hole preparation; however, it is still a convenience for checking drawings, master artwork, and photographic negatives.

When drawing the second side of a printed-circuit board, corresponding centers should be taken directly from the back of the drawing of the first side.

In addition to the illustration of the circuit pattern, the trim line, registration marks, and two scale dimensions at right angles should be shown. Nomenclature, reference designations, operating instructions, and other information may also be added.

Assembly

All components should be inserted on one side of the board if practicable. In the case of boards with the circuit on one side only, the parts should be inserted on the side opposite the circuit. This allows all connections to be soldered simultaneously by dip-soldering.

Dip-soldering consists of applying a flux, usually a rosin-alcohol mixture, to the circuit pattern and then placing the board in contact with molten solder. Slight agitation of the board will ensure good fillets around the wire leads. In good present technique, the circuit board with its components assembled (on one side only) has its conductor pattern passed through the crest of a "wave" of molten solder; all junctions are soldered as



Fig. 22. Indication for hole.

the board progresses through the wave. The flux, board temperature, solder temperature, and immersion time are interrelated and must be adjusted for best results. Long exposure to hot solder is detrimental to the insulating material and to the adhesive that joins the copper foil to the insulation. For hand dipping, a 5-second dip in a 60/40 tin-lead solder bath maintained at a temperature of 450 degrees Fahrenheit will give satisfactory results.

After solder-dipping, the residual flux should be removed by a suitable solvent. Be sure the solvent is compatible with the materials used in the component parts mounted on the board; solvents frequently dissolve cements or plastics and marking inks, or cause severe stress cracking of plastics.

To secure the advantages of machine assembly:

(A) Components should be of similar size and shape, or separate inserting heads will be required for each different shape of item.

(B) Components of the same size and shape must be mounted using the same terminal lead spacing at all points.

(C) Different values of a part, or even different parts of similar shape and sizes (if axial-lead style) may be specially sequenced in a lead-taped package for insertion by one programmed insertion head.

(D) A few oddly sized or shaped components may be economically inserted by hand after the machine insertion work is completed.

Reference

An excellent reference on microelectronic printed-circuit techniques is Scarlett, J. A. *Printed Circuit Boards for Microelectronics*. New York: Van Nostrand Reinhold Co., 1970.

6 Fundamentals of Networks

Inductance of Single-Layer Solenoids	6-3
Approximate Equation	
General Remarks	
Decrease of Solenoid Inductance by Shielding	
Q of Unshielded Solenoid	
Reactance Charts	6-5
Impedance Formulas	6-6
Skin Effect	6-6
Symbols	
Skin Depth	
General Considerations	
Equations for Simple R, L, and C Networks	6-14
Self-Inductance of Circular Ring of Round Wire at Radio	
Frequencies, for Nonmagnetic Materials	
Capacitance	
T- π or Y- Δ Transformation	
Transients—Elementary Cases	6-15
Time Constant	
Capacitor Charge and Discharge	
Two Capacitors	
Inductor Charge and Discharge	
Charge and Discharge of Series R-L-C Circuit	
Series R-L-C Circuit With Sinusoidal Applied Voltage	

Transients—Operational Calculus and Laplace Transforms *6-19*

Example

Circuit Response Related to Unit Impulse

Circuit Response Related to Unit Step

Heaviside Expansion Theorem

Application to Linear Networks

INDUCTANCE OF SINGLE-LAYER SOLENOIDS

The approximate value of the low-frequency inductance of a single-layer solenoid is*

$$L = Fn^2d$$

where,

- L = inductance in microhenrys,
- F = form factor, a function of the ratio d/l (value of F may be read from Fig. 1),
- n = number of turns,
- d = diameter of coil (inches) between centers of conductors,
- l = length of coil (inches) = n times the distance between centers of adjacent turns.

The equation is based on the assumption of a uniform current sheet, but the correction due to the use of spaced round wires is usually negligible for practical purposes. For higher frequencies, skin effect alters the inductance slightly. This effect is not readily calculated, but is often negligibly small. However, it must be borne in mind that the equation gives approximately the true value of inductance. In contrast, the apparent value is affected by the shunting effect of the distributed capacitance of the coil.

Example: Required, a coil of 100 microhenrys inductance, wound on a form 2 inches in diameter by 2 inches winding length. The $d/l = 1.00$, and $F = 0.0173$ in Fig. 1.

$$\begin{aligned} n &= (L/Fd)^{1/2} \\ &= [100/(0.0173 \times 2)]^{1/2} \\ &= 54 \text{ turns} \end{aligned}$$

Reference to Table 1 will assist in choosing a desirable size of wire, allowing for a suitable spacing between turns according to the application of the coil. A slight correction may then be made for the increased diameter (diameter of form, plus two times radius of wire), if this small correction seems justified.

Approximate Equation

For single-layer solenoids of the proportions normally used in radio work, the inductance in microhenrys is given to an accuracy of about 1 percent by the formula

$$L = n^2[r^2/(9r + 10l)]$$

where $r = d/2$ and the other quantities are as defined for the previous inductance formula.

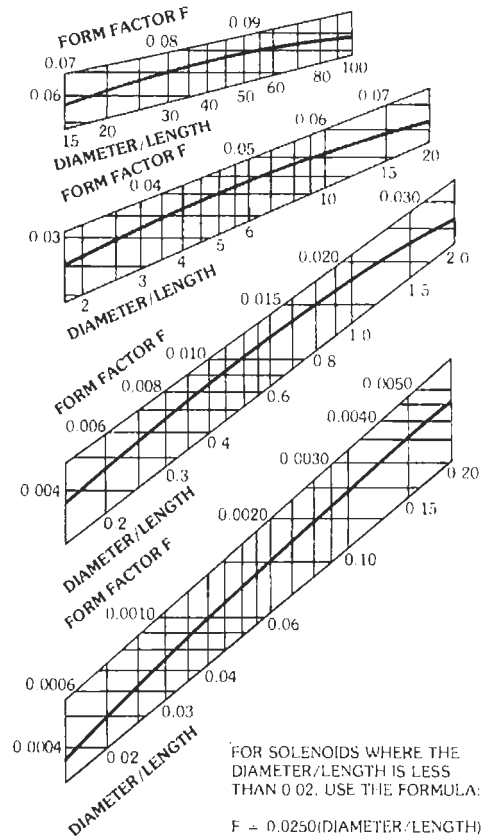


Fig. 1. Chart showing inductance of a single-layer solenoid, form factor = F.

General Remarks

In the use of various charts, tables, and calculators for designing inductors, the following relationships are useful in extending the range of the devices. They apply to coils of any type or design.

(A) If all dimensions are held constant, inductance is proportional to n^2 .

(B) If the proportions of the coil remain unchanged, then for a given number of turns the inductance is proportional to the dimensions of the coil. A coil with all dimensions m times those of a given coil (having the same number of turns) has m times the inductance of the given coil. That is, inductance has the dimensions of length.

Decrease of Solenoid Inductance by Shielding*

When a solenoid is enclosed in a cylindrical shield, the inductance is reduced by a factor given in Fig. 2. This effect has been evaluated by considering the shield

* Equations and Fig. 1 are derived from equations and tables in Bureau of Standards Circular No. C74.

* RCA Application Note No. 48; 12 June 1935.

TABLE 1. MAGNET-WIRE DATA

AWG B & S Gauge									Bare		Enameled	
	Bare Nom Diam (in.)	Enam Nom Diam (in.)	SCC* Diam (in.)	DCC* Diam (in.)	SCE* Diam (in.)	SSC* Diam (in.)	DSC* Diam (in.)	SSE* Diam (in.)	Min Diam (in.)	Max Diam (in.)	Min Diam (in.)	Diam* (in.)
10	0.1019	0.1039	0.1079	0.1129	0.1104	—	—	—	0.1009	0.1029	0.1024	0.1044
11	0.0907	0.0927	0.0957	0.1002	0.0982	—	—	—	0.0898	0.0917	0.0913	0.0932
12	0.0808	0.0827	0.0858	0.0903	0.0882	—	—	—	0.0800	0.0816	0.0814	0.0832
13	0.0720	0.0738	0.0770	0.0815	0.0793	—	—	—	0.0712	0.0727	0.0726	0.0743
14	0.0641	0.0659	0.0691	0.0736	0.0714	—	—	—	0.0634	0.0647	0.0648	0.0664
15	0.0571	0.0588	0.0621	0.0666	0.0643	0.0591	0.0611	0.0613	0.0565	0.0576	0.0578	0.0593
16	0.0508	0.0524	0.0558	0.0603	0.0579	0.0528	0.0548	0.0549	0.0503	0.0513	0.0515	0.0529
17	0.0453	0.0469	0.0503	0.0548	0.0523	0.0473	0.0493	0.0493	0.0448	0.0457	0.0460	0.0473
18	0.0403	0.0418	0.0453	0.0498	0.0472	0.0423	0.0443	0.0442	0.0399	0.0407	0.0410	0.0422
19	0.0359	0.0374	0.0409	0.0454	0.0428	0.0379	0.0399	0.0398	0.0355	0.0363	0.0366	0.0378
20	0.0320	0.0334	0.0370	0.0415	0.0388	0.0340	0.0360	0.0358	0.0316	0.0323	0.0326	0.0338
21	0.0285	0.0299	0.0335	0.0380	0.0353	0.0305	0.0325	0.0323	0.0282	0.0287	0.0292	0.0303
22	0.0253	0.0266	0.0303	0.0343	0.0320	0.0273	0.0293	0.0290	0.0251	0.0256	0.0261	0.0270
23	0.0226	0.0238	0.0276	0.0316	0.0292	0.0246	0.0266	0.0262	0.0223	0.0228	0.0232	0.0242
24	0.0201	0.0213	0.0251	0.0291	0.0266	0.0221	0.0241	0.0236	0.0199	0.0203	0.0208	0.0216
25	0.0179	0.0190	0.0224	0.0264	0.0238	0.0199	0.0219	0.0213	0.0177	0.0181	0.0186	0.0193
26	0.0159	0.0169	0.0204	0.0244	0.0217	0.0179	0.0199	0.0192	0.0158	0.0161	0.0166	0.0172
27	0.0142	0.0152	0.0187	0.0227	0.0200	0.0162	0.0182	0.0175	0.0141	0.0144	0.0149	0.0155
28	0.0126	0.0135	0.0171	0.0211	0.0183	0.0146	0.0166	0.0158	0.0125	0.0128	0.0132	0.0138
29	0.0113	0.0122	0.0158	0.0198	0.0170	0.0133	0.0153	0.0145	0.0112	0.0114	0.0119	0.0125
30	0.0100	0.0108	0.0145	0.0185	0.0156	0.0120	0.0140	0.0131	0.0099	0.0101	0.0105	0.0111
31	0.0089	0.0097	0.0134	0.0174	0.0144	0.0109	0.0129	0.0119	0.0088	0.0090	0.0094	0.0099
32	0.0080	0.0088	0.0125	0.0165	0.0135	0.0100	0.0120	0.0110	0.0079	0.0081	0.0085	0.0090
33	0.0071	0.0078	0.0116	0.0156	0.0125	0.0091	0.0111	0.0100	0.0070	0.0072	0.0075	0.0080
34	0.0063	0.0069	0.0108	0.0148	0.0116	0.0083	0.0103	0.0091	0.0062	0.0064	0.0067	0.0071
35	0.0056	0.0061	0.0101	0.0141	0.0108	0.0076	0.0096	0.0083	0.0055	0.0057	0.0059	0.0063
36	0.0050	0.0055	0.0090	0.0130	0.0097	0.0070	0.0090	0.0077	0.0049	0.0051	0.0053	0.0057
37	0.0045	0.0049	0.0085	0.0125	0.0091	0.0065	0.0085	0.0071	0.0044	0.0046	0.0047	0.0051
38	0.0040	0.0044	0.0080	0.0120	0.0086	0.0060	0.0080	0.0066	0.0039	0.0041	0.0042	0.0046
39	0.0035	0.0038	0.0075	0.0115	0.0080	0.0055	0.0075	0.0060	0.0034	0.0036	0.0036	0.0040
40	0.0031	0.0034	0.0071	0.0111	0.0076	0.0051	0.0071	0.0056	0.0030	0.0032	0.0032	0.0036
41	0.0028	0.0031	—	—	—	—	—	—	0.0027	0.0029	0.0029	0.0032
42	0.0025	0.0028	—	—	—	—	—	—	0.0024	0.0026	0.0026	0.0029
43	0.0022	0.0025	—	—	—	—	—	—	0.0021	0.0023	0.0023	0.0026
44	0.0020	0.0023	—	—	—	—	—	—	0.0019	0.0021	0.0021	0.0024

* Nominal bare diameter plus maximum additions.

For additional data on copper wire, see Chapters 4 and 13.

to be a short-circuited single-turn secondary. The curves in Fig. 2 are reasonably accurate provided the clearance between each end of the coil winding and the corresponding end of the shield is at least equal to the radius of the coil. For square shield cans, take the equivalent shield diameter (for Fig. 2) as being 1.2 times the width of one side of the square.

Example: Let the coil winding length be 1.5 inches and its diameter 0.75 inch, while the shield diameter is

1.25 inches. What is the reduction of inductance due to the shield? The proportions are

$$(\text{winding length})/(\text{winding diameter}) = 2.0$$

$$(\text{winding diameter})/(\text{shield diameter}) = 0.6$$

Referring to Fig. 2, the actual inductance in the shield is 72 percent of the inductance of the coil in free space.

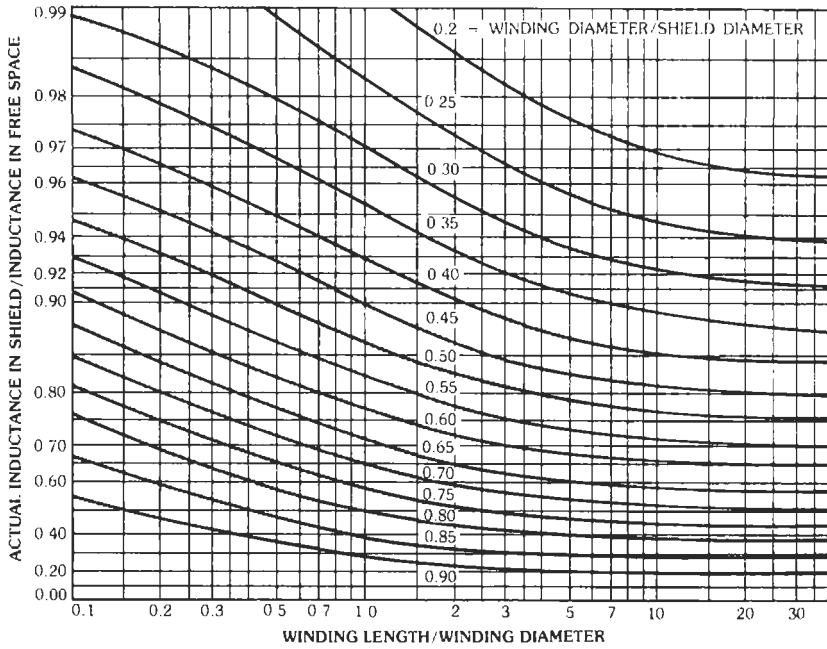


Fig. 2. Inductance decrease when a solenoid is shielded. By permission of RCA, copyright proprietor.

Q of Unshielded Solenoid

Fig. 3 can be used to obtain the unloaded *Q* of an unshielded solenoid.

value may be determined in terms of two others by use of a straightedge laid across the correct chart for the frequency under consideration.

Example: Given a capacitance of 0.001 microfarad, find the reactance at 50 kilohertz and inductance required to resonate. Place a straightedge through these values and read the intersections on the other scales, giving approximately 3200 ohms and 11 millihenrys. See Fig. 5.

REACTANCE CHARTS

Figs. 4, 5, and 6 give the relationships of capacitance, inductance, reactance, and frequency. Any one

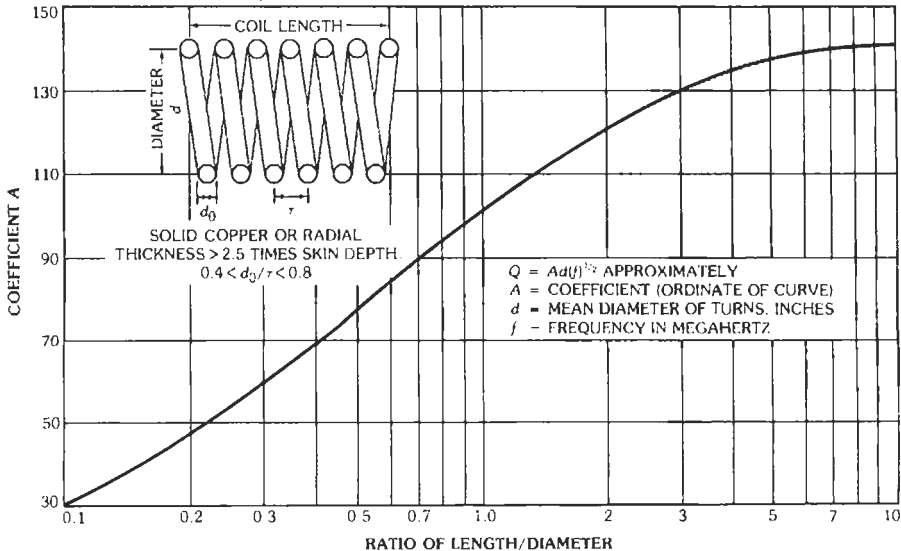


Fig. 3. *Q* of unshielded coil.

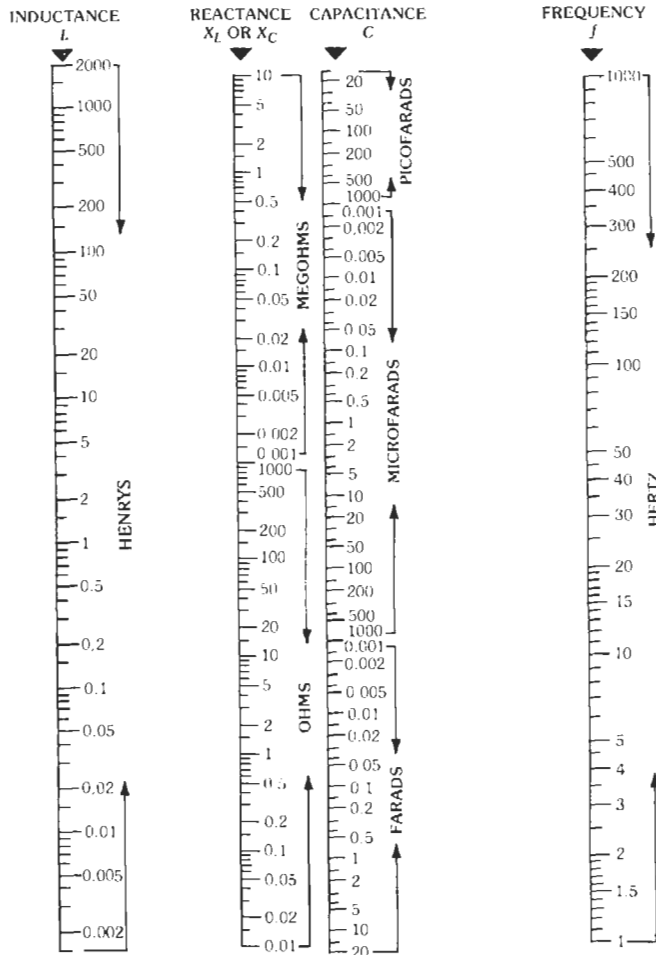


Fig. 4. Chart covering 1 hertz to 1000 hertz.

IMPEDANCE FORMULAS

Impedance and admittance formulas for various combinations of L , C , and R are given in Tables 2 and 3.

SKIN EFFECT

Symbols

- A = correction coefficient
- D = diameter of conductor in inches
- f = frequency in hertz
- R_{ac} = resistance at frequency f
- R_{dc} = direct-current resistance
- R_{sq} = resistance per square
- T = thickness of tubular conductor in inches
- T_1 = depth of penetration of current
- δ = skin depth
- λ = free-space wavelength in meters
- μ_r = relative permeability of conductor material

($\mu_r = 1$ for copper and other nonmagnetic materials)

ρ = resistivity of conductor material at any temperature

ρ_c = resistivity of copper at 20 °C = 1.724 microhms-centimeter

Skin Depth

The skin depth is that distance below the surface of a conductor where the current density has diminished to $1/e$ of its value at the surface. The thickness of the conductor is assumed to be several (perhaps at least three) times the skin depth. Imagine the conductor replaced by a cylindrical shell of the same surface shape but of thickness equal to the skin depth, with uniform current density equal to that which exists at the surface of the actual conductor. Then the total current in the shell and its resistance are equal to the corresponding values in the actual conductor.

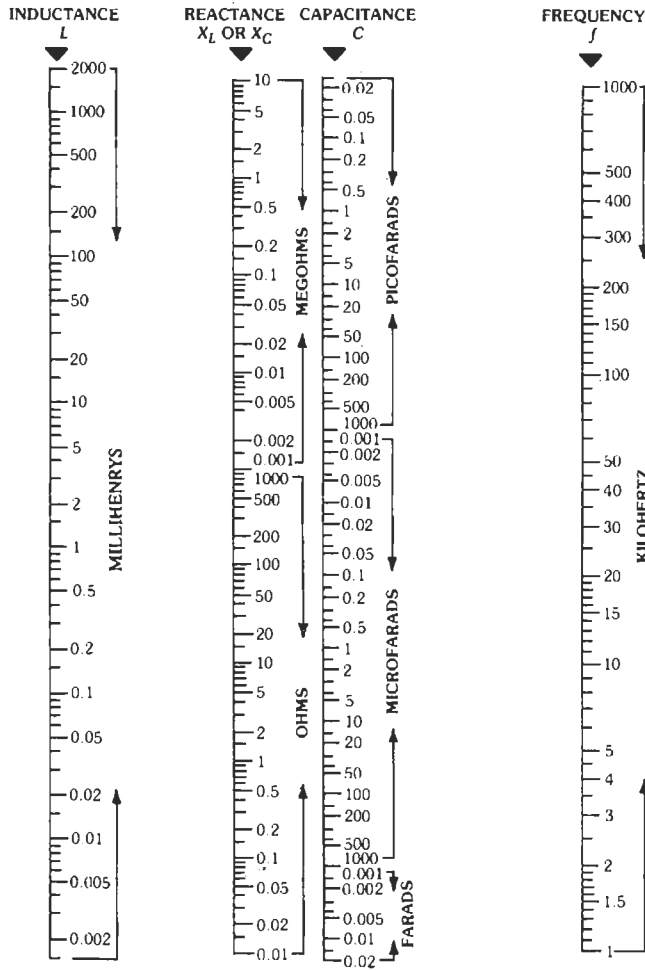


Fig. 5. Chart covering 1 kilohertz to 1000 kilohertz.

The skin depth and the resistance per square (of any size), in meter-kilogram-second (rationalized) units, are

$$\delta = (\lambda / \pi \sigma \mu c)^{1/2}$$

$$R_{sq} = 1 / \delta \sigma$$

where,

- δ = skin depth in meters,
- R_{sq} = resistance per square in ohms,
- c = velocity of light *in vacuo*
= 2.998×10^8 meter/second,
- $\mu = 4\pi \times 10^{-7} \mu_r$ henry/meter,
- $1/\sigma = 1.724 \times 10^{-8} \rho / \rho_c$ ohm-meter.

For numerical computations

$$\begin{aligned} \delta &= (3.82 \times 10^{-4} \lambda^{1/2}) k_1 \\ &= (6.61 / f^{1/2}) k_1 \text{ centimeter} \end{aligned}$$

$$\begin{aligned} \delta &= (1.50 \times 10^{-4} \lambda^{1/2}) k_1 \\ &= (2.60 / f^{1/2}) k_1 \text{ inch} \\ \delta_m &= (2.60 / f_{mc}^{1/2}) k_1 \text{ mil} \\ R_{sq} &= (4.52 \times 10^{-3} / \lambda^{1/2}) k_2 \\ &= (2.61 \times 10^{-7} f^{1/2}) k_2 \text{ ohm} \\ k_1 &= [(1/\mu_r) \rho / \rho_c]^{1/2} \\ k_2 &= (\mu_r \rho / \rho_c)^{1/2} \\ k_1, k_2 &= \text{unity for copper} \end{aligned}$$

Example: What is the resistance/foot of a cylindrical copper conductor of diameter D inches?

$$\begin{aligned} R &= (12 / \pi D) R_{sq} \\ &= (12 / \pi D) \times 2.61 \times 10^{-7} (f^{1/2}) \\ &= 0.996 \times 10^{-6} (f^{1/2}) / D \text{ ohm/foot} \end{aligned}$$

If $D = 1.00$ inch and $f = 100 \times 10^6$ hertz, then $R = 0.996 \times 10^{-6} \times 10^4 \approx 1 \times 10^{-2}$ ohm/foot.

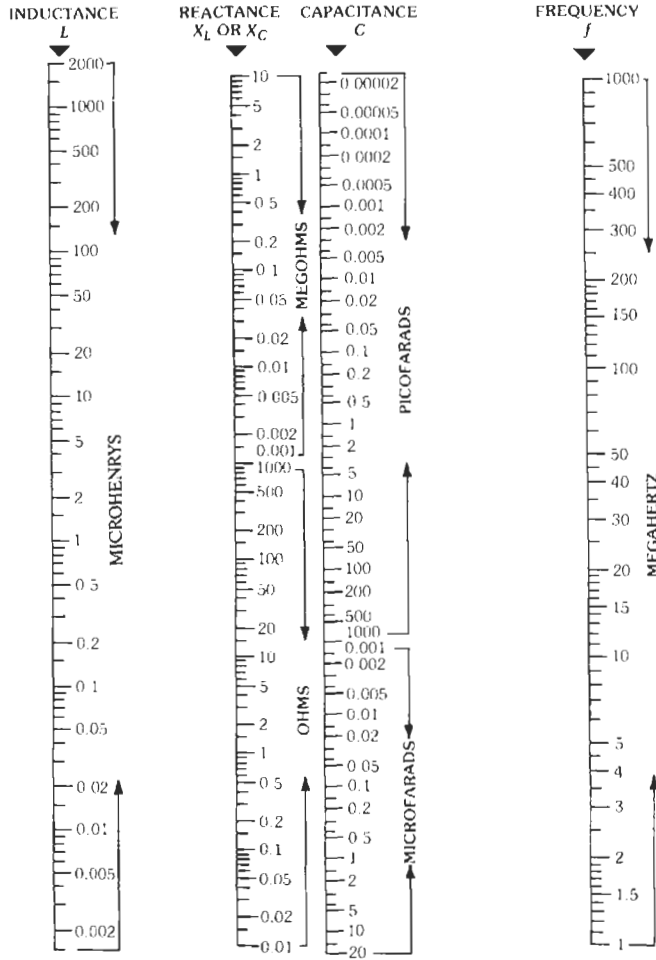


Fig. 6. Chart covering 1 megahertz to 1000 mcgahertz.

General Considerations

Fig. 7 shows the relationship of R_{ac}/R_{dc} versus $D(f^{1/2})$ for copper, or versus $D(f^{1/2})(\mu_r \rho_c / \rho)^{1/2}$ for any conductor material, for an isolated straight solid conductor of circular cross section. Negligible error in the equations for R_{ac} results when the conductor is spaced at least $10D$ from adjacent conductors. When the spacing between axes of parallel conductors carrying the same current is $4D$, the resistance R_{ac} is increased about 3 percent, when the depth of penetration is small. The equations are accurate for concentric lines due to their circular symmetry.

For values of $D(f^{1/2})(\mu_r \rho_c / \rho)^{1/2}$ greater than 40

$$R_{ac}/R_{dc} = 0.0960 D(f^{1/2})(\mu_r \rho_c / \rho)^{1/2} + 0.26 \quad (\text{Eq. 1})$$

The high-frequency resistance in ohms/foot of an isolated straight conductor (either solid or tubular for $T < D/8$ or $T_1 < D/8$) is given in Eq. (2). If the current is

along the inside surface of a tubular conductor, D is the inside diameter.

$$R_{ac} = A[(f^{1/2})/D][\mu_r(\rho/\rho_c)]^{1/2} \times 10^{-6} \quad (\text{Eq. 2})$$

The values of the correction coefficient A for solid conductors and for tubular conductors are given in Table 4.

The value of $T(f^{1/2})(\mu_r \rho_c / \rho)^{1/2}$ that just makes $A = 1$ indicates the penetration of the currents below the surface of the conductor. Thus, approximately

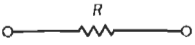
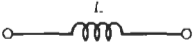
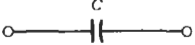
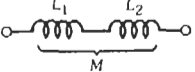

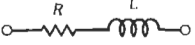


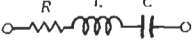
$$T_1 = [3.5/(f^{1/2})](\rho/\mu_r \rho_c)^{1/2} \quad (\text{Eq. 3})$$

where T_1 is in inches.

When $T_1 < D/8$, the value of R_{ac} as given by Eq. (2) (but not the value of R_{ac}/R_{dc} in Table 4, "Tubular Conductors") is correct for any value $T \geq T_1$.

Under the limitation that the radius of curvature of all parts of the cross section is appreciably greater than T_1 ,

TABLE 2. IMPEDANCE AND ADMITTANCE FORMULAS—SERIES AND PARALLEL COMBINATIONS

	Impedance $Z = R + jX$	Magnitude $ Z = (R^2 + X^2)^{1/2}$	Phase Angle $\phi = \tan^{-1}(X/R)$	Admittance $Y = 1/Z$
	R	R	0	$1/R$
	$j\omega L$	ωL	$+\pi/2$	$-j/\omega L$
	$-j/\omega C$	$1/\omega C$	$-\pi/2$	$j\omega C$
	$j\omega(L_1 + L_2 \pm 2M)$	$\omega(L_1 + L_2 \pm 2M)$	$+\pi/2$	$-j/\omega(L_1 + L_2 \pm 2M)$
	$-(j/\omega)(1/C_1 + 1/C_2)$	$(1/\omega)(1/C_1 + 1/C_2)$	$-\pi/2$	$j\omega C_1 C_2 / (C_1 + C_2)$
	$R + j\omega L$	$(R^2 + \omega^2 L^2)^{1/2}$	$\tan^{-1}(\omega L/R)$	$(R - j\omega L) / (R^2 + \omega^2 L^2)$
	$R - j/\omega C$	$(1/\omega C)(1 + \omega^2 C^2 R^2)^{1/2}$	$-\tan^{-1}(1/\omega CR)$	$(R + j/\omega C) / (R^2 + 1/\omega^2 C^2)$
	$j(\omega L - 1/\omega C)$	$(\omega L - 1/\omega C)$	$\pm\pi/2$	$j\omega C / (1 - \omega^2 LC)$
	$R + j(\omega L - 1/\omega C)$	$[R^2 + (\omega L - 1/\omega C)^2]^{1/2}$	$\tan^{-1}[(\omega L - 1/\omega C)/R]$	$\frac{R - j(\omega L - 1/\omega C)}{R^2 + (\omega L - 1/\omega C)^2}$

Continued on next page.

TABLE 2 (CONT). IMPEDANCE AND ADMITTANCE FORMULAS—SERIES AND PARALLEL COMBINATIONS

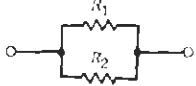
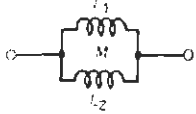
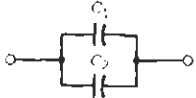
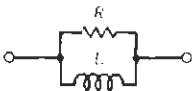
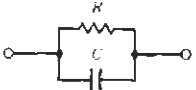
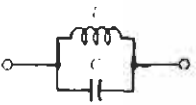
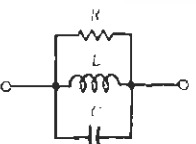
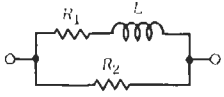
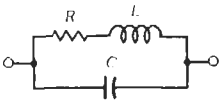
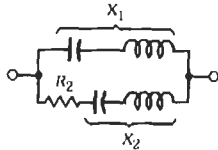
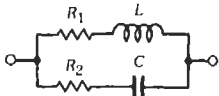
	Impedance $Z = R + jX$	Magnitude $ Z = (R^2 + X^2)^{1/2}$	Phase Angle $\phi = \tan^{-1}(X/R)$	Admittance $Y = 1/Z$
	$R_1 R_2 / (R_1 + R_2)$	$R_1 R_2 / (R_1 + R_2)$	0	$(1/R_1 + 1/R_2)$
	$j\omega \left(\frac{L_1 L_2 - M^2}{L_1 + L_2 \mp 2M} \right)$	$\omega \left(\frac{L_1 L_2 - M^2}{L_1 + L_2 \mp 2M} \right)$	$+\pi/2$	$-(j/\omega) \left(\frac{L_1 + L_2 \mp 2M}{L_1 L_2 - M^2} \right)$
	$-j/\omega(C_1 + C_2)$	$1/\omega(C_1 + C_2)$	$-\pi/2$	$j\omega(C_1 + C_2)$
	$\omega L R \left(\frac{\omega L + jR}{R^2 + \omega^2 L^2} \right)$	$\omega L R / (R^2 + \omega^2 L^2)^{1/2}$	$\tan^{-1}(R/\omega L)$	$1/R - j/\omega L$
	$R(1 - j\omega C R) / (1 + \omega^2 C^2 R^2)$	$R / (1 + \omega^2 C^2 R^2)^{1/2}$	$-\tan^{-1} \omega C R$	$1/R + j\omega C$
	$j\omega L / (1 - \omega^2 LC)$	$\omega L / (1 - \omega^2 LC)$	$\pm \pi/2$	$j(\omega C - 1/\omega L)$
	$\frac{1/R - j(\omega C - 1/\omega L)}{(1/R)^2 + (\omega C - 1/\omega L)^2}$	$[(1/R)^2 + (\omega C - 1/\omega L)^2]^{-1/2}$	$\tan^{-1} R(1/\omega L - \omega C)$	$1/R + j(\omega C - 1/\omega L)$

TABLE 3. IMPEDANCE AND ADMITTANCE FORMULAS—SERIES-PARALLEL COMBINATIONS

	Impedance Z	$R_2 \frac{R_1(R_1 + R_2) + \omega^2 L^2 + j\omega L R_2}{(R_1 + R_2)^2 + \omega^2 L^2}$
	Magnitude $ Z $	$R_2 \left[\frac{R_1^2 + \omega^2 L^2}{(R_1 + R_2)^2 + \omega^2 L^2} \right]^{1/2}$
	Phase Angle ϕ	$\tan^{-1} \frac{\omega L R_2}{R_1(R_1 + R_2) + \omega^2 L^2}$
	Admittance Y	$\frac{R_1(R_1 + R_2) + \omega^2 L^2 - j\omega L R_2}{R_2(R_1^2 + \omega^2 L^2)}$
	Impedance Z	$\frac{R + j\omega[L(1 - \omega^2 LC) - CR^2]}{(1 - \omega^2 LC)^2 + \omega^2 C^2 R^2}$
	Magnitude $ Z $	$\left[\frac{R^2 + \omega^2 L^2}{(1 - \omega^2 LC)^2 + \omega^2 C^2 R^2} \right]^{1/2}$
	Phase Angle ϕ	$\tan^{-1} \left\{ \omega[L(1 - \omega^2 LC) - CR^2]/R \right\}$
	Admittance Y	$\frac{R - j\omega[L(1 - \omega^2 LC) - CR^2]}{R^2 + \omega^2 L^2}$

Continued on next page.

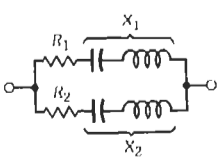
TABLE 3 (CONT). IMPEDANCE AND ADMITTANCE FORMULAS—SERIES-PARALLEL COMBINATIONS

	Impedance Z	$X_1 \frac{X_1 R_2 + j[R_2^2 + X_2(X_1 + X_2)]}{R_2^2 + (X_1 + X_2)^2}$
	Magnitude $ Z $	$X_1 \left[\frac{R_2^2 + X_2^2}{R_2^2 + (X_1 + X_2)^2} \right]^{1/2}$
	Phase Angle ϕ	$\tan^{-1} \frac{R_2^2 + X_2(X_1 + X_2)}{X_1 R_2}$
	Admittance Y	$\frac{R_2 X_1 - j(R_2^2 + X_2^2 + X_1 X_2)}{X_1 (R_2^2 + X_2^2)}$
	Impedance Z	$\frac{R_1 R_2 (R_1 + R_2) + \omega^2 L^2 R_2 + R_1 / \omega^2 C^2}{(R_1 + R_2)^2 + (\omega L - 1/\omega C)^2} + j \frac{\omega L R_2^2 - R_1^2 / \omega C - (L/C)(\omega L - 1/\omega C)}{(R_1 + R_2)^2 + (\omega L - 1/\omega C)^2}$
	Magnitude $ Z $	$\left[\frac{(R_1^2 + \omega^2 L^2)(R_2^2 + 1/\omega^2 C^2)}{(R_1 + R_2)^2 + (\omega L - 1/\omega C)^2} \right]^{1/2}$
	Phase Angle ϕ	$\tan^{-1} \left[\frac{\omega L R_2^2 - R_1^2 / \omega C - (L/C)(\omega L - 1/\omega C)}{R_1 R_2 (R_1 + R_2) + \omega^2 L^2 R_2 + R_1 / \omega^2 C^2} \right]$
	Admittance Y	$\frac{R_1 + \omega^2 C^2 R_1 R_2 (R_1 + R_2) + \omega^4 L^2 C^2 R_2}{(R_1^2 + \omega^2 L^2)(1 + \omega^2 C^2 R_2^2)} + j \omega \left[\frac{C R_1^2 - L + \omega^2 L C (L - C R_2^2)}{(R_1^2 + \omega^2 L^2)(1 + \omega^2 C^2 R_2^2)} \right]$

Note: When $R_1 = R_2 = \sqrt{L/C}$, then $Z = R_1 = R_2$, a pure resistance at any frequency where the given conditions hold.

Continued on next page.

TABLE 3 (CONT). IMPEDANCE AND ADMITTANCE FORMULAS—SERIES-PARALLEL COMBINATIONS

	Impedance Z	$\frac{(R_1 R_2 - X_1 X_2) + j(R_1 X_2 + R_2 X_1)}{(R_1 + R_2) + j(X_1 + X_2)}$
	Magnitude $ Z $	$\left[\frac{(R_1^2 + X_1^2)(R_2^2 + X_2^2)}{(R_1 + R_2)^2 + (X_1 + X_2)^2} \right]^{1/2}$
	Phase Angle ϕ	$\tan^{-1}(X_1/R_1) + \tan^{-1}(X_2/R_2) - \tan^{-1}[(X_1 + X_2)/(R_1 + R_2)]$
	Admittance Y	$1/(R_1 + jX_1) + 1/(R_2 + jX_2)$

$$Z = R + jX$$

$$|Z| = (R^2 + X^2)^{1/2}$$

$$\phi = \tan^{-1}(X/R)$$

$$Y = 1/Z$$

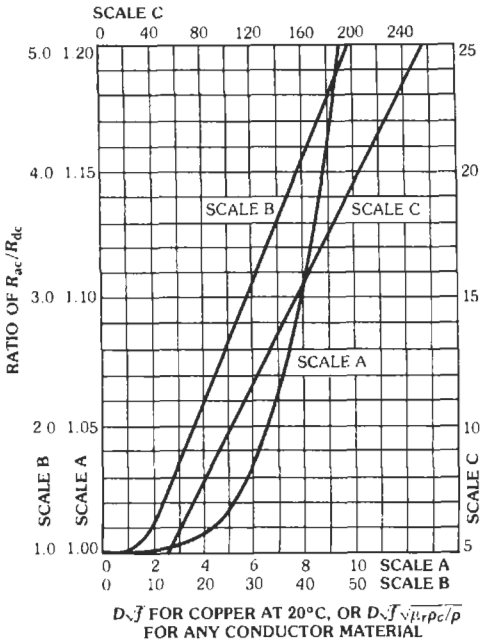


Fig. 7. Resistance ratio for isolated straight solid conductors of circular cross section.

Eqs. (2) and (3) hold for isolated straight conductors of any shape. In this case the term $D = (\text{perimeter of cross section})/\pi$.

Examples:

(A) At 100 megahertz, a copper conductor has a depth of penetration $T_1 = 0.00035$ inch.

(B) A steel shield with 0.005-inch copper plate, which is practically equivalent in R_{ac} to an isolated copper conductor 0.005-inch thick, has a value of $A = 1.23$ at 200 kilohertz. This 23-percent increase in resistance over that of a thick copper sheet is satisfactorily low as regards its effect on the losses of the components within the shield. By comparison, a thick aluminum sheet has a resistance $(\rho/\rho_c)^{1/2} = 1.28$ times that of copper.

EQUATIONS FOR SIMPLE R, L, AND C NETWORKS*

Self-Inductance of Circular Ring of Round Wire at Radio Frequencies, for Nonmagnetic Materials

$$L = (a/100) \times [7.353 \log_{10}(16a/d) - 6.386]$$

* Many equations for computing capacitance, inductance, and mutual inductance will be found in Bureau of Standards Circular No. C74, obtainable from the Superintendent of Documents, Government Printing Office, Washington, D.C. 20402.

TABLE 4. SKIN-EFFECT CORRECTION COEFFICIENT A FOR SOLID AND TUBULAR CONDUCTORS

Solid Conductors		
$D(f^{1/2})[\mu_r(\rho_c/\rho)]^{1/2}$	A	
> 370	1.000	
220	1.005	
160	1.010	
98	1.02	
48	1.05	
26	1.10	
13	1.20	
9.6	1.30	
5.3	2.00	
< 3.0	$R_{ac} \approx R_{dc}$	
$R_{dc} = (10.37/D^2)(\rho/\rho_c) \times 10^{-6}$ ohm/foot		
Tubular Conductors		
$T(f^{1/2})[\mu_r(\rho_c/\rho)]^{1/2}$	A	R_{ac}/R_{dc}
= B where $B > 3.5$	1.00	0.384B
3.5	1.00	1.35
3.15	1.01	1.23
2.85	1.05	1.15
2.60	1.10	1.10
2.29	1.20	1.06
2.08	1.30	1.04
1.77	1.50	1.02
1.31	2.00	1.00
= B where $B < 1.3$	2.60/B	1.00

where,

- L = self-inductance in microhenrys,
- a = mean radius of ring in inches,
- d = diameter of wire in inches,
- $a/d > 2.5$.

Capacitance

For Parallel-Plate Capacitor:

$$C = 0.0885\epsilon_r[(N - 1)A]/t$$

$$= 0.225\epsilon_r[(N - 1)A''/t'']$$

where,

- C = capacitance in picofarads,
- A = area of one side of one plate in square centimeters,
- A'' = area in square inches,
- N = number of plates,
- t = thickness of dielectric in centimeters,
- t'' = thickness in inches,
- ϵ_r = dielectric constant relative to air.

This equation neglects “fringing” at the edges of the plates.

For Coaxial Cylindrical Capacitor (Fig. 8):

$$C = 2\pi\epsilon_r\epsilon_v / [\log_e(b/a)]$$

$$= \{(5 \times 10^6 \epsilon_r) / [c^2 \log_e(b/a)]\}$$

where,

- C = capacitance per unit axial length in farads/meter,
- c = velocity of light in vacuo, meters per second
= 2.998×10^8 ,
- ϵ_r = dielectric constant relative to air,
- ϵ_v = permittivity of free space in farads/meter
= 8.85×10^{-12}

Alternative forms of the equation are:

$$C = 0.2416\epsilon_r / [\log_{10}(b/a)] \text{ picofarads/centimeter}$$

$$= 0.614\epsilon_r / [\log_{10}(b/a)] \text{ picofarads/inch}$$

$$= 7.36\epsilon_r / [\log_{10}(b/a)] \text{ picofarads/foot}$$

When $1.0 < (b/a) < 1.4$, then with accuracy of 1 percent or better, the capacitance in picofarads/foot is:

$$C = 8.50\epsilon_r \frac{(b/a) + 1}{(b/a) - 1}$$

T-π or Y-Δ Transformation

The two networks (Fig. 9) are equivalent, as far as conditions at the terminals are concerned, provided the listed equations are satisfied (either the impedance equations or the admittance equations may be used)

$$Y_1 = 1/Z_1 \quad Y_c = 1/Z_c, \text{ etc.}$$

TRANSIENTS—ELEMENTARY CASES

The complete transient in a linear network is, by the principle of superposition, the sum of the individual transients due to the store of energy in each inductor and capacitor and to each external source of energy connected to the network. To this is added the steady-state condition due to each external source of energy. The transient may be computed as starting from any arbitrary time $t = 0$ when the initial conditions of the energy of the network are known.

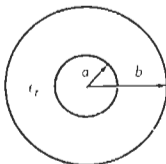


Fig. 8. Coaxial cylindrical capacitor.

Time Constant

The time constant (designated T) of the discharge of a capacitor through a resistor is the time $t_2 - t_1$ required for the voltage or current to decay to $1/\epsilon$ of its value at time t_1 . For the charge of a capacitor, the same definition applies, the voltage “decaying” toward its steady-state value. The time constant of discharge or charge of the current in an inductor through a resistor follows an analogous definition.

Stored energy, in joules (watt-seconds), is:

$$\text{Energy stored in a capacitor} = CE^2/2$$

$$\text{Energy stored in an inductor} = LI^2/2$$

In the equations in this section, T and t are in seconds, R is in ohms, L is in henrys, C is in farads, E is in volts, and I is in amperes. Values of the numerical constants are:

$$\epsilon = 2.718 \quad 1/\epsilon = 0.3679 \quad \log_{10}\epsilon = 0.4343$$

Capacitor Charge and Discharge

Closing of switch (Fig. 10) occurs at time $t = 0$. Initial conditions (at $t = 0$): Battery = E_b ; $e_c = E_0$. Steady state (at $t = \infty$): $i = 0$; $e_c = E_b$. Transient:

$$i = [(E_b - E_0)/R]\exp(-t/RC)$$

$$= I_0 \exp(-t/RC)$$

$$\log_{10}(i/I_0) = -(0.4343/RC)t$$

$$e_c = E_0 + C^{-1} \int_0^t i dt$$

$$= E_0 \exp(-t/RC) + E_b[1 - \exp(-t/RC)]$$

Time constant:

$$T = RC$$

Fig. 11 shows current:

$$i/I_0 = \exp(-t/T)$$

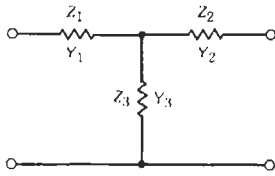
Fig. 11 shows discharge (for $E_b = 0$):

$$e_c/E_0 = \exp(-t/T)$$

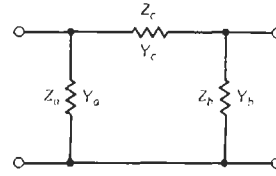
Fig. 12 shows charge (for $E_0 = 0$):

$$e_c/E_b = 1 - \exp(-t/T)$$

These curves are plotted for a wider range in Fig. 13.



(A) T or Y.



(B) π or Δ .

Impedance Equations:

$$Z_c = (Z_1 Z_2 + Z_1 Z_3 + Z_2 Z_3) / Z_3$$

$$Z_a = (Z_1 Z_2 + Z_1 Z_3 + Z_2 Z_3) / Z_2$$

$$Z_b = (Z_1 Z_2 + Z_1 Z_3 + Z_2 Z_3) / Z_1$$

$$Z_1 = Z_a Z_c / (Z_a + Z_b + Z_c)$$

$$Z_2 = Z_b Z_c / (Z_a + Z_b + Z_c)$$

$$Z_3 = Z_a Z_b / (Z_a + Z_b + Z_c)$$

Admittance Equations:

$$Y_c = Y_1 Y_2 / (Y_1 + Y_2 + Y_3)$$

$$Y_a = Y_1 Y_3 / (Y_1 + Y_2 + Y_3)$$

$$Y_b = Y_2 Y_3 / (Y_1 + Y_2 + Y_3)$$

$$Y_1 = (Y_a Y_b + Y_a Y_c + Y_b Y_c) / Y_b$$

$$Y_2 = (Y_a Y_b + Y_a Y_c + Y_b Y_c) / Y_a$$

$$Y_3 = (Y_a Y_b + Y_a Y_c + Y_b Y_c) / Y_c$$

(C) Equations.

Fig. 9. T or Y network and π or Δ network.

Two Capacitors

Closing of switch (Fig. 14) occurs at time $t = 0$.
Initial conditions (at $t = 0$):

$$e_1 = E_1; \quad e_2 = E_2$$

Steady state (at $t = \infty$):

$$e_1 = E_f; \quad e_2 = -E_f; \quad i = 0$$

$$E_f = (E_1 C_1 - E_2 C_2) / (C_1 + C_2)$$

$$C' = C_1 C_2 / (C_1 + C_2)$$

Transient:

$$i = [(E_1 + E_2) / R] \exp(-t / RC')$$

$$e_1 = E_f + (E_1 - E_f) \exp(-t / RC')$$

$$= E_1 - (E_1 + E_2)(C' / C_1)$$

$$\times [1 - \exp(-t / RC')]$$

$$e_2 = -E_f + (E_2 + E_f) \exp(-t / RC')$$

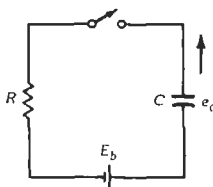


Fig. 10. Circuit for capacitor charge and discharge.

$$= E_2 - (E_1 + E_2)(C' / C_2)$$

$$\times [1 - \exp(-t / RC')]$$

Original energy (joules) = $(C_1 E_1^2 + C_2 E_2^2) / 2$

Final energy (joules) = $(C_1 + C_2) E_f^2 / 2$

Loss of energy (joules) = $\int_0^\infty i^2 R dt$
 $= C' (E_1 + E_2)^2 / 2$

(Loss is independent of the value of R.)

Inductor Charge and Discharge

Initial conditions (at $t = 0$) in Fig. 15:

$$\text{Battery} = E_b; \quad i = I_0$$

Steady state (at $t = \infty$):

$$i = I_f = E_b / R$$

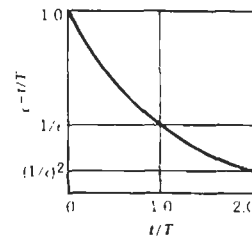


Fig. 11. Capacitor discharge.

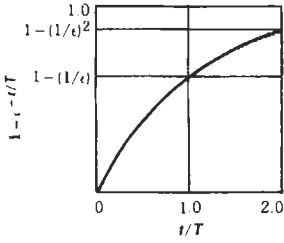


Fig. 12. Capacitor charge.

Transient plus steady state:

$$i = I_f [1 - \exp(-Rt/L)] + I_0 \exp(-Rt/L)$$

$$e_L = -L di/dt$$

$$= -(E_b - RI_0) \exp(-Rt/L)$$

Time constant:

$$T = L/R$$

Fig. 11 shows discharge (for $E_b = 0$):

$$i/I_0 = \exp(-t/T)$$

Fig. 12 shows charge (for $I_0 = 0$)

$$i/I_f = [1 - \exp(-t/T)]$$

These curves are plotted for a wider range in Fig. 13.

Charge and Discharge of Series R-L-C Circuit

Initial conditions (at $t = 0$) in Fig. 16:

$$\text{Battery} = E_b; \quad e_c = E_0; \quad i = I_0$$

Use exponential $\exp(-t/T)$ for charge or discharge of capacitor or discharge of inductor:

$$(\text{current at time } t)/(\text{initial current})$$

Discharge of capacitor:

$$(\text{voltage at time } t)/(\text{initial voltage})$$

Use exponential $1 - \exp(-t/T)$ for charge of capacitor:

$$(\text{voltage at time } t)/(\text{battery or final voltage})$$

Charge of inductor:

$$(\text{current at time } t)/(\text{final current})$$

Fig. 13. Exponential functions $\exp(-t/T)$ and $1 - \exp(-t/T)$ applied to transients in R-C and L-R circuits.

Steady state (at $t = \infty$):

$$i = 0; \quad e_c = E_b$$

Differential equation:

$$E_b - E_0 - C^{-1} \int_0^t i dt - Ri - L(di/dt) = 0$$

when

$$L(d^2i/dt^2) + R(di/dt) + (i/C) = 0$$

Solution of equation:

$$i = \exp(-Rt/2L) \left[\frac{2(E_b - E_0) - RI_0}{R(D^{1/2})} \sinh(Rt/2L) \times (D^{1/2}) + I_0 \cosh(Rt/2L)(D^{1/2}) \right]$$

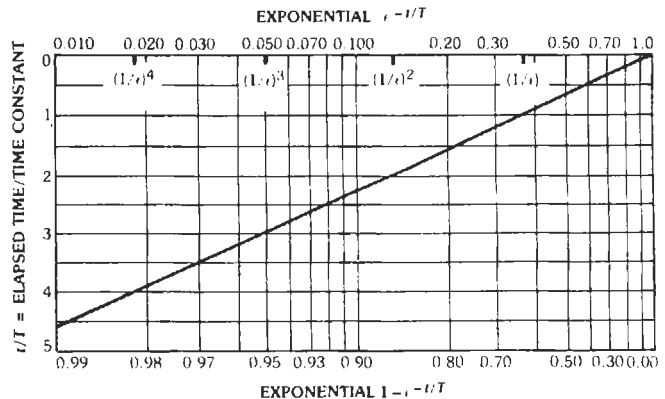
where $D = 1 - (4L/R^2C)$.

Case 1: When L/R^2C is small

$$i = (1 - 2A - 2A^2)^{-1} \left\{ \left[\frac{E_b - E_0}{R} - I_0(A + A^2) \right] \times \exp\left(-\frac{t}{RC}(1 + A + 2A^2)\right) + \left[I_0(1 - A - A^2) - \frac{E_b - E_0}{R} \right] \times \exp\left(-\frac{Rt}{L}(1 - A - A^2)\right) \right\}$$

where $A = L/R^2C$.

For practical purposes, the terms A^2 can be neglected when $A < 0.1$. The terms A may be neglected when $A < 0.01$.



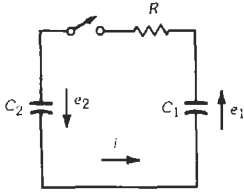


Fig. 14. Circuit for two capacitors.

Case 2: When $4L/R^2C < 1$ for which $D^{1/2}$ is real

$$i = \frac{\exp(-Rt/2L)}{D^{1/2}} \left\{ \left[\frac{E_b - E_0}{R} - \frac{1}{2} I_0 (1 - D^{1/2}) \right] \times \exp\left(\frac{Rt}{2L} D^{1/2}\right) + \left[\frac{1}{2} I_0 (1 + D^{1/2}) - \frac{E_b - E_0}{R} \right] \exp\left(-\frac{Rt}{2L} D^{1/2}\right) \right\}$$

Case 3: When D is a small positive or negative quantity

$$i = \exp(-Rt/2L) \left\{ \frac{2(E_b - E_0)}{R} \left[\frac{Rt}{2L} + \frac{1}{6} \left(\frac{Rt}{2L} \right)^3 D \right] + I_0 \left[1 - \frac{Rt}{2L} + \frac{1}{2} \left(\frac{Rt}{2L} \right)^2 D - \frac{1}{6} \left(\frac{Rt}{2L} \right)^3 D \right] \right\}$$

This equation may be used for values of D up to ± 0.25 , at which values the error in the computed current i is approximately 1 percent of I_0 or of $(E_b - E_0)/R$.

Case 3A: When $4L/R^2C = 1$ for which $D = 0$, the equation reduces to

$$i = \exp(-Rt/2L) \left[\frac{E_b - E_0}{R} \frac{Rt}{L} + I_0 \left(1 - \frac{Rt}{2L} \right) \right]$$

or $i = i_1 + i_2$, plotted in Fig. 17. For practical purposes, this equation may be used when $4L/R^2C = 1 \pm 0.05$ with errors of 1 percent or less.

Case 4: When $4L/R^2C > 1$ for which $D^{1/2}$ is imaginary

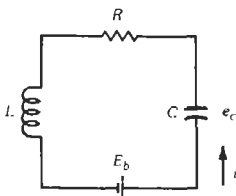


Fig. 16. Series R-L-C circuit.

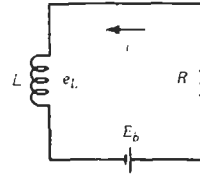


Fig. 15. Circuit for inductor charge and discharge.

$$i = \exp(-Rt/2L) \left[\left(\frac{E_b - E_0}{\omega_0 L} - \frac{RI_0}{2\omega_0 L} \right) \times \sin \omega_0 t + I_0 \cos \omega_0 t \right] = I_m \exp(-Rt/2L) \sin(\omega_0 t + \Psi)$$

where,

$$\omega_0 = [(LC)^{-1} - (R^2/4L^2)]^{1/2}$$

$$I_m = (\omega_0 L)^{-1} \{ [E_b - E_0 - (RI_0)/2]^2 + \omega_0^2 L^2 I_0^2 \}^{1/2}$$

$$\Psi = \tan^{-1} \{ \omega_0 L I_0 / [E_b - E_0 - (RI_0)/2] \}$$

The envelope of the voltage wave across the inductor is

$$\pm \exp(-Rt/2L) [\omega_0 (LC)^{1/2}]^{-1} \times \{ [E_b - E_0 - (RI_0)/2]^2 + \omega_0^2 L^2 I_0^2 \}^{1/2}$$

Example: Relay with transient-suppressing capacitor (Fig. 18). The switch is closed until time $t = 0$, then opened.

Let $L = 0.10$ henry, $R_1 = 100$ ohms, and $E = 10$ volts.

Suppose we choose $C = 10^{-6}$ farad and $R_2 = 100$ ohms.

Then $R = 200$ ohms, $I_0 = 0.10$ ampere, $E_0 = 10$ volts, $\omega_0 = 3 \times 10^3$, and $f_0 = 480$ hertz.

Maximum peak voltage across L (envelope at $t = 0$) is approximately 30 volts. Time constant of decay of envelope is 0.001 second.

Nonoscillating Condition: It is preferable that the circuit be just nonoscillating (Case 3A) and that it present a pure resistance at the switch terminals for any frequency.

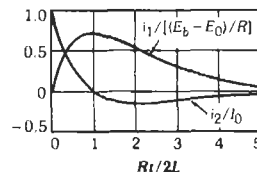


Fig. 17. Transients for $4L/R^2C = 1$.

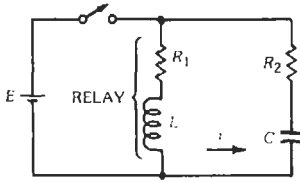


Fig. 18. Equivalent circuit of relay with transient-suppressing capacitor.

$$R_2 = R_1 = R/2 = 100 \text{ ohms}$$

$$4L/R^2C = 1$$

$$C = 10^{-5} \text{ farad} = 10 \text{ microfarads}$$

At the instant of opening the switch, the voltage across the parallel circuit is $E_0 - R_2 I_0 = 0$.

Series R-L-C Circuit With Sinusoidal Applied Voltage

By the principle of superposition, the transient and steady-state conditions are the same for the actual circuit and the equivalent circuit shown in Fig. 19, the closing of the switch occurring at time $t = 0$. In the equivalent circuit, the steady state is due to the source e acting continuously from time $t = -\infty$, while the transient is due to short-circuiting the source $-e$ at time $t = 0$.

Source:

$$e = E \sin(\omega t + \alpha)$$

Steady state:

$$i = (E/Z) \sin(\omega t + \alpha - \phi)$$

where,

$$Z = \{R^2 + [\omega L - (1/\omega C)]^2\}^{1/2}$$

$$\tan \phi = (\omega^2 LC - 1)/\omega CR$$

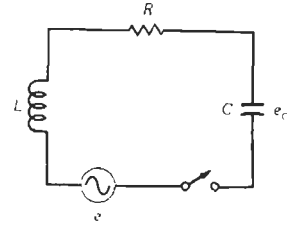
The transient is found by determining current $i = I_0$ and capacitor voltage $e_c = E_0$ at time $t = 0$, due to the source $-e$. These values of I_0 and E_0 are then substituted in the equations of Case 1, 2, 3, or 4, above, according to the values of R , L , and C .

At time $t = 0$, due to the source $-e$:

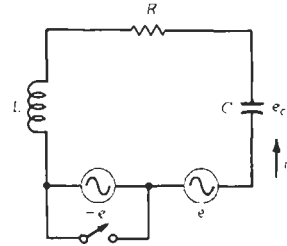
$$i = I_0 = -(E/Z) \sin(\alpha - \phi)$$

$$e_c = E_0 = (E/\omega CZ) \cos(\alpha - \phi)$$

This form of analysis may be used for any periodic applied voltage e . The steady-state current and the



(A) Actual circuit.



(B) Equivalent circuit.

Fig. 19. Series R-L-C circuit.

capacitor voltage for an applied voltage $-e$ are determined, the periodic voltage being resolved into its harmonic components for this purpose, if necessary. Then the instantaneous values $i = I_0$ and $e_c = E_0$ at the time of closing the switch are easily found, from which the transient is determined. It is evident, from this method of analysis, that the waveform of the transient need bear no relationship to that of the applied voltage, depending only on the constants of the circuit and the hypothetical initial conditions I_0 and E_0 .

TRANSIENTS—OPERATIONAL CALCULUS AND LAPLACE TRANSFORMS

Among the various methods of operational calculus used to solve transient problems, one of the most efficient makes use of the Laplace transform.

If we have a function $v = f(t)$, then by definition the Laplace transform is $\mathcal{L} [f(t)] = F(p)$, where

$$F(p) = \int_0^\infty \exp(-pt)f(t) dt \tag{Eq. 4}$$

The inverse transform of $F(p)$ is $f(t)$. Most of the mathematical functions encountered in practical work fall in the class for which Laplace transforms exist. Transforms of several functions are given in Chapter 46.

In the following, an abbreviated symbol such as $\mathcal{L} [i]$ is used instead of $\mathcal{L} [i(t)]$ to indicate the Laplace transform of the function $i(t)$.

The electrical (or other) system for which a solution of the differential equation is required is considered

only in the time domain $t \geq 0$. Any currents or voltages existing at $t = 0$, before the driving force is applied, constitute initial conditions. Driving force is assumed to be 0 when $t < 0$.

Example

Take the circuit of Fig. 20, in which the switch is closed at time $t = 0$. Before the closing of the switch, suppose the capacitor is charged; then at $t = 0$, we have $v = V_0$. It is required to find the voltage v across capacitor C as a function of time.

Since $i = dq/dt = C(dv/dt)$, the differential equation of the circuit in terms of voltage is:

$$e(t) = v + Ri = v + RC(dv/dt) \tag{Eq. 5}$$

where $e(t) = E_b$.

Referring to the table of transforms, the applied voltage is E_b multiplied by unit step, or $E_b S_{-1}(t)$; the transform for this is E_b/p . The transform of v is $\mathcal{L}[v]$. That of $RC(dv/dt)$ is $RC[p \mathcal{L}[v] - v(0)]$, where $v(0) = V_0 =$ value of v at $t = 0$. Then the transform of Eq. (5) is:

$$E_b/p = \mathcal{L}[v] + RC[p \mathcal{L}[v] - V_0]$$

Rearranging and resolving into partial fractions:

$$\begin{aligned} \mathcal{L}[v] &= \frac{E_b}{p(1 + RCp)} + \frac{RCV_0}{1 + RCp} \\ &= E_b[p^{-1} - (p + 1/RC)^{-1}] + \frac{V_0}{p + 1/RC} \tag{Eq. 6} \end{aligned}$$

Now we must determine the equation that would transform into Eq. (6). The inverse transform of $\mathcal{L}[v]$ is v , and those of the terms on the right-hand side are found in the table of transforms. Then, in the time domain $t \geq 0$

$$v = E_b[1 - \exp(-t/RC)] + V_0 \exp(-t/RC) \tag{Eq. 7}$$

This solution is also well known by classical methods. However, the advantages of the Laplace transform method become more and more apparent in reducing the labor of solution as the equations become more involved.

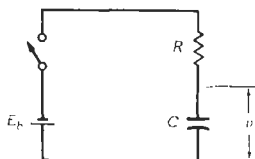


Fig. 20. Series R-C circuit.

Circuit Response Related to Unit Impulse

Unit impulse (see Laplace transforms) has the dimensions of time⁻¹. For example, suppose a capacitor of 1 microfarad is suddenly connected to a battery of 100 volts, with the circuit inductance and resistance negligibly small. Then the current is 10⁻⁴ coulomb multiplied by unit impulse.

The general transformed equation of a circuit or system may be written

$$\mathcal{L}[i] = \psi(p) \mathcal{L}[e] + \psi(p) \tag{Eq. 8}$$

Here $\mathcal{L}[i]$ is the transform of the required current (or other quantity) and $\mathcal{L}[e]$ is the transform of the applied voltage or driving force $e(t)$. The transform of the initial conditions, at $t = 0$, is included in $\psi(p)$.

First considering the case when the system is initially at rest, $\psi(p) = 0$. Writing i_a for the current in this case

$$\mathcal{L}[i_a] = \phi(p) \mathcal{L}[e] \tag{Eq. 9}$$

Now apply unit impulse $S_0(t)$ (multiplied by 1 volt-second), and designate the circuit current in this case by $B(t)$ and its transform by $\mathcal{L}[B]$. The Laplace transform of $S_0(t)$ is 1, so

$$\mathcal{L}[B] = \phi(p) \tag{Eq. 10}$$

Equation (9) becomes, for any driving force

$$\mathcal{L}[i_a] = \mathcal{L}[B] \mathcal{L}[e] \tag{Eq. 11}$$

Applying the convolution function (Laplace transform)

$$\begin{aligned} i_a &= \int_0^t B(t - \lambda) e(\lambda) d\lambda \\ &= \int_0^t B(\lambda) e(t - \lambda) d\lambda \tag{Eq. 12} \end{aligned}$$

To this there must be added the current i_0 due to any initial conditions that exist. From (8)

$$\mathcal{L}[i_0] = \psi(p) \tag{Eq. 13}$$

Then i_0 is the inverse transform of $\psi(p)$.

Circuit Response Related to Unit Step

Unit step is defined and designated $S_{-1}(t) = 0$ for $t < 0$ and equals unity for $t > 0$. It has no dimensions. Its Laplace transform is $1/p$. Let the circuit current be designated $A(t)$ when the applied voltage is $e = S_{-1}(t) \times (1 \text{ volt})$. Then, the current i_a for the case when the

system is initially at rest, and for any applied voltage $e(t)$, is given by any of

$$\begin{aligned}
 i_a &= A(t) e(0) + \int_0^t A(t-\lambda) e'(\lambda) d\lambda \\
 &= A(t) e(0) + \int_0^t A(\lambda) e'(t-\lambda) d\lambda \\
 &= A(0) e(t) + \int_0^t A'(t-\lambda) e(\lambda) d\lambda \\
 &= A(0) e(t) + \int_0^t A'(\lambda) e(t-\lambda) d\lambda \quad (\text{Eq. 14})
 \end{aligned}$$

where A' is the first derivative of A and similarly for e' of e .

As an example, consider the problem of Fig. 20 and Eqs. (5) to (7) above. Suppose $V_0 = 0$, and that the battery is replaced by a linear source

$$e(t) = Et/T_1$$

where T_1 is the duration of the voltage rise in seconds. By Eq. (7), setting $E_b = 1$

$$A(t) = 1 - \exp(-t/RC)$$

Then using the first equation in (14) and noting that $e(0) = 0$, and $e'(t) = E/T_1$ when $0 \leq t \leq T_1$, the solution is

$$v = (Et/T_1) - (ERC/T_1)[1 - \exp(-t/RC)]$$

This result can, of course, be found readily by direct application of the Laplace transform to Eq. (5) with $e(t) = Et/T_1$.

Heaviside Expansion Theorem

When the system is initially at rest, the transformed equation is given by Eq. (9) and may be written

$$\mathcal{L}[i_a] = [M(p)/G(p)] \mathcal{L}[e] \quad (\text{Eq. 15})$$

$M(p)$ and $G(p)$ are rational functions of p . In the following, $M(p)$ must be of lower degree than $G(p)$, as is usually the case. The roots of $G(p) = 0$ are p_r , where $r = 1, 2, \dots, n$, and there must be no repeated roots. The response may be found by application of the Heaviside expansion theorem.

For a force $e = E_{\max} \exp(j\omega t)$ applied at time $t = 0$

$$\begin{aligned}
 \frac{i_a(t)}{E_{\max}} &= \frac{M(j\omega)}{G(j\omega)} \exp(j\omega t) + \sum_{r=1}^n \frac{M(p_r) \exp(p_r t)}{(p_r - j\omega)G'(p_r)} \\
 &\quad (\text{Eq. 16a})
 \end{aligned}$$

$$= \frac{\exp(j\omega t)}{Z(j\omega)} + \sum_{r=1}^n \frac{\exp(p_r t)}{(p_r - j\omega)Z'(p_r)} \quad (\text{Eq. 16b})$$

The first term on the right-hand side of either form of (16) gives the steady-state response, and the second term gives the transient. When $e = E_{\max} \cos\omega t$, take the real part of (16), and similarly for $\sin\omega t$ take the imaginary part. $Z(p)$ is defined in Eq. (19). If the applied force is the unit step, set $\omega = 0$ in Eq. (16).

Application to Linear Networks

The equation for a single mesh is of the form

$$\begin{aligned}
 A_n(d^n i/dt^n) + \dots + A_1(di/dt) \\
 + A_0 i + B \int i dt = e(t) \quad (\text{Eq. 17})
 \end{aligned}$$

System Initially at Rest—Then, Eq. (17) transforms into

$$(A_n p^n + \dots + A_1 p + A_0 + B_p^{-1}) [i] = [e] \quad (\text{Eq. 18})$$

where the expression in parentheses is the operational impedance, equal to the alternating-current impedance when we set $p = j\omega$.

If there are m meshes in the system, we get m simultaneous equations like (17) with m unknowns i_1, i_2, \dots, i_m . The m algebraic equations like (18) are solved for $\mathcal{L}[i_1]$, etc., by means of determinants, yielding an equation of the form of (15) for each unknown, with a term on the right-hand side for each mesh in which there is a driving force. Each such driving force may of course be treated separately and the responses added.

If any two meshes are designated by the letters h and k , the driving force $e(t)$ being in either mesh and the mesh current $i(t)$ in the other, then the fraction $M(p)/G(p)$ in (15) becomes

$$M_{hk}(p)/G(p) = 1/Z_{hk}(p) = Y_{hk}(p) \quad (\text{Eq. 19})$$

where $Y_{hk}(p)$ is the operational transfer admittance between the two meshes. The determinant of the system is $G(p)$, and $M_{hk}(p)$ is the cofactor of the row and column that represent $e(t)$ and $i(t)$.

System Not Initially at Rest—The transient due to the initial conditions is solved separately and added to the above solution. The driving force is set equal to zero in Eq. (17), $e(t) = 0$, and each term is transformed according to

$$\begin{aligned}
 \mathcal{L}[d^n i/dt^n] &= p^n \mathcal{L}[i] - \sum_{r=1}^n p^{n-r} [d^{r-1} i/dt^{r-1}]_{t=0} \\
 &\quad (\text{Eq. 20a})
 \end{aligned}$$

$$\mathcal{L}\left[\int_0^t i dt\right] = p^{-1} \mathcal{L}[i] + p^{-1} \left[\int i dt\right]_{t=0} \quad (\text{Eq. 20b})$$

where the last term in each equation represents the initial conditions. For example, in Eq. (20b) the last term would represent, in an electrical circuit, the quantity of electricity existing on a capacitor at time $t = 0$, the instant when the driving force $e(t)$ begins to act.

Resolution into Partial Fractions—The solution of the operational form of the equations of a system involves rational fractions that must be simplified before the inverse transform is found. Let the fraction be $h(p)/g(p)$ where $h(p)$ is of lower degree than $g(p)$, for example $(3p + 2)/(p^2 + 5p + 8)$. If $h(p)$ is of equal or higher degree than $g(p)$, it can be reduced by division.

The reduced fraction can be expanded into partial fractions. Let the factors of the denominator be $(p - p_r)$ for the n nonrepeated roots p_r of the equation $g(p) = 0$, and $(p - p_a)$ for a root p_a repeated m times.

$$\frac{h(p)}{g(p)} = \sum_{r=1}^n \frac{A_r}{p - p_r} + \sum_{r=1}^m \frac{B_r}{(p - p_a)^{m-r+1}} \quad (\text{Eq. 21a})$$

There is a summation term for each root that is repeated. The constant coefficients A_r and B_r can be evaluated by reforming the fraction with a common denominator. Then the coefficients of each power of p in $h(p)$ and the reformed numerator are equated and the resulting equations solved for the constants. More formally, they may be evaluated by

$$A_r = \frac{h(p_r)}{g'(p_r)} = \left[\frac{h(p)}{g(p)/(p - p_r)} \right]_{p=p_r} \quad (\text{Eq. 21b})$$

$$B_r = [1/(r - 1)!] f^{(r-1)}(p_a) \quad (\text{Eq. 21c})$$

where

$$f(p) = (p - p_a)^m [h(p)/g(p)]$$

and $f^{(r-1)}(p_a)$ indicates that the $(r - 1)$ th derivative of $f(p)$ is to be found, after which we set $p = p_a$.

Fractions of the form $(A_1p + A_2)/(p^2 + \omega^2)$ or, more generally

$$\frac{A_1p + A_2}{p^2 + 2ap + b} = \frac{A(p + a) + B\omega}{(p + a)^2 + \omega^2} \quad (\text{Eq. 22a})$$

where $b > a^2$ and $\omega^2 = b - a^2$ need not be reduced further. From the Laplace transforms the inverse transform of (22a) is

$$\exp(-at)(A \cos \omega t + B \sin \omega t) \quad (\text{Eq. 22b})$$

where

$$A = \frac{h(-a + j\omega)}{g'(-a + j\omega)} + \frac{h(-a - j\omega)}{g'(-a - j\omega)} \quad (\text{Eq. 22c})$$

$$B = j \left[\frac{h(-a + j\omega)}{g'(-a + j\omega)} - \frac{h(-a - j\omega)}{g'(-a - j\omega)} \right] \quad (\text{Eq. 22d})$$

Similarly, the inverse transform of the fraction

$$[A(p + a) + B\alpha]/[(p + a)^2 - \alpha^2]$$

is $\exp(-at)(A \cosh \alpha t + B \sinh \alpha t)$, where A and B are found by (22c) and (22d), except that $j\omega$ is replaced by α and the coefficient j is omitted in the expression for B .

7

Fourier Waveform Analysis

*Revised by
David C. Munson, Jr.*

Fourier Transform of a Function 7-2

Fourier Series 7-2

 Real Form of Fourier Series

 Complex Form of Fourier Series

 Average Power

 Odd and Even Functions

 Odd or Even Harmonics

Pulse-Train Analysis 7-9

Spectral Analysis 7-10

FOURIER TRANSFORM OF A FUNCTION

The Fourier transform, F , of function f is defined by the integral (where x and y are real variables)

$$F(y) = \int_{-\infty}^{\infty} f(x) \exp(-j2\pi xy) dx$$

provided this integral exists.

A sufficient, but not necessary, existence condition is that f be absolutely integrable; that is

$$\int_{-\infty}^{\infty} |f(x)| dx < \infty$$

An important example of a function that has a Fourier transform even though it is not absolutely integrable is $(\sin x)/x$.

In general, F and f may be complex. Letting $f(x) = f_r(x) + jf_i(x)$, where f_r and f_i are real-valued, one has

$$\begin{aligned} F(y) &= \int_{-\infty}^{\infty} [f_r(x) \cos 2\pi yx + f_i(x) \sin 2\pi yx] dx \\ &\quad - j \int_{-\infty}^{\infty} [f_r(x) \sin 2\pi yx - f_i(x) \cos 2\pi yx] dx \end{aligned}$$

Conversely, the function f , whose Fourier transform is a given function F , is given by the integral (inverse Fourier transform)

$$f(x) = \int_{-\infty}^{\infty} F(y) \exp(j2\pi xy) dy$$

where it is assumed that, at points of discontinuity of the integral (if any), the function $f(x)$ is given the value

$$f(x) = [f(x^+) + f(x^-)]/2$$

The functions $f(x^+)$ and $f(x^-)$ are the limits of $f(x+t)$ as t approaches 0 through positive and negative values, respectively.

Letting $F(y) = F_r(y) + jF_i(y)$, one has

$$\begin{aligned} f(x) &= \int_{-\infty}^{\infty} [F_r(y) \cos 2\pi yx - F_i(y) \sin 2\pi yx] dy \\ &\quad + j \int_{-\infty}^{\infty} [F_r(y) \sin 2\pi yx + F_i(y) \cos 2\pi yx] dy \end{aligned}$$

In many engineering applications it is customary to denote the variable y as "frequency"; in most cases x represents time or space.

If the radian frequency $\omega = 2\pi y$ is introduced as a variable, the definitions of the Fourier transform and of its inverse are written as

$$F(\omega/2\pi) = F_1(\omega)$$

$$= \int_{-\infty}^{\infty} f(x) \exp(-j\omega x) dx$$

$$f(x) = (2\pi)^{-1} \int_{-\infty}^{\infty} F_1(\omega) \exp(j\omega x) d\omega$$

Properties of the Fourier transform are listed in Table 1. Table 2 contains normalized graphs of the Fourier transform (or its magnitude) for a number of common pulse shapes.

FOURIER SERIES

Real Form of Fourier Series

A periodic function with period T , defined by its values in the intervals $-T/2$ to $+T/2$ or 0 to T can be written in a Fourier series expansion as

$$\begin{aligned} f(x) &= (1/2)A_0 \\ &\quad + \sum_{n=1}^{\infty} [A_n \cos(n\omega_0 x) + B_n \sin(n\omega_0 x)] \end{aligned}$$

where $\omega_0 = 2\pi/T$, and the coefficients are given by

$$A_n = (2/T) \int_{-T/2}^{T/2} f(x) \cos(n\omega_0 x) dx$$

$$= (2/T) \int_0^T f(x) \cos(n\omega_0 x) dx$$

$$B_n = (2/T) \int_{-T/2}^{T/2} f(x) \sin(n\omega_0 x) dx$$

$$= (2/T) \int_0^T f(x) \sin(n\omega_0 x) dx$$

for $n = 0, 1, 2, \dots$.

An alternative form of the above expansion is

$$f(x) = (1/2)C_0 + \sum_{n=1}^{\infty} C_n \cos(n\omega_0 x + \phi_n)$$

where

$$\begin{aligned} C_0 &= A_0 \\ C_n &= (A_n^2 + B_n^2)^{1/2} \\ \cos \phi_n &= A_n/C_n \\ \sin \phi_n &= -B_n/C_n \end{aligned}$$

Complex Form of Fourier Series

The Fourier series can be written more concisely as

$$f(x) = \sum_{n=-\infty}^{\infty} D_n \exp(jn\omega_0 x)$$

TABLE 1. PROPERTIES OF FOURIER TRANSFORM*

	Function	Fourier Transform
1. Definition	$f(x)$	$F(y) = \int_{-\infty}^{+\infty} f(x) \exp(-2\pi jxy) dx$
2. Inverse transform	$f(x) = \int_{-\infty}^{+\infty} F(y) \exp(2\pi jxy) dy$	$F(y)$
3. Linearity	$af(x)$ $f_1(x) \pm f_2(x)$	$aF(y)$ $F_1(y) \pm F_2(y)$
4. Translation or shifting theorem	$g(x) = f(x - x_0)$, $x_0 = \text{real const.}$ $g(x) = \exp(2\pi jy_0x)f(x)$, $y_0 = \text{real const.}$	$G(y) = \exp(-2\pi jx_0y)F(y)$ $G(y) = F(y - y_0)$
5. Change of scale	$g(x) = f(x/a)$, $a = \text{real const.}$	$G(y) = a F(ay)$
6. Frequency shifting and change of scale	$g(x) = \exp(2\pi jy_0x)f(x/a)$, y_0 and $a = \text{real const.}$	$G(y) = a F[a(y - y_0)]$
7. Interchange of function and transform	$g(x) = F(x)$	$G(y) = f(-y)$
8a. Convolution in x -space (product of Fourier transforms)	$h = f * g = g * f$ i.e., $h(x) = \int_{-\infty}^{+\infty} f(x - \tau)g(\tau) d\tau$ $= \int_{-\infty}^{+\infty} f(\tau)g(x - \tau) d\tau$	$H = F \cdot G$
8b. Convolution in y -space (product of inverse Fourier transforms)	$h = f \cdot g$	$H = F * G$
9. Unit pulse (or Dirac function)	$\delta(x)$ $f(x) = 1$ (for all x)	$F(y) = 1$ (for all y) $\delta(y)$

Continued on next page.

TABLE 1 (CONT). PROPERTIES OF FOURIER TRANSFORM*

	Function	Fourier Transform
10. Periodic train of equal pulses	$A \sum_{n=-\infty}^{+\infty} \delta(x-nT)$	$(A/T) \sum_{n=-\infty}^{+\infty} \delta(y-n/T)$
11a. Derivative in x -space	$g(x)=d^n f/dx^n$	$G(y)=(2\pi j)^n F(y)$, if $G(y)$ exists
11b. Derivative in y -space	$g(x)=(-2\pi jx)^n f(x)$	$G(y)=d^n F/dy^n$
12a. Integral in x -space	$g(x)=\int_{-\infty}^x f(x)dx$	$G(y)=[1/(2\pi j)]F(y)+F(0)\delta(y)$ where $F(0)=\int_{-\infty}^{\infty} f(x)dx$
12b. Integral in y -space	$g(x)=-[1/(2\pi jx)]f(x)$	$G(y)=\int_{-\infty}^y F(y)dy$
13. Symmetry	$g(x)=f(-x)$ f even: $f(x)=f(-x)$ f odd: $f(x)=-f(-x)$	$G(y)=F(-y)$ F even: $F=2 \int_0^{\infty} f(x) \cos(2\pi xy)dx$ F odd: $F=-2j \int_0^{\infty} f(x) \sin(2\pi xy)dx$
14. Complex conjugate	$g(x)=f^*(x)$ Hence, if $f(x)$ is real	$G(y)=F^*(-y)$ $F(-y)=F^*(y)$
15. Area under the curve	$\int_{-\infty}^{+\infty} f(x)dx=F(0)$	$\int_{-\infty}^{+\infty} F(y)dy=f(0)$
16a. Parseval's theorem	$\int_{-\infty}^{+\infty} f^*(x)g(x)dx$	$= \int_{-\infty}^{+\infty} F^*(y)G(y)dy$

Continued on next page.

TABLE 1 (CONT). PROPERTIES OF FOURIER TRANSFORM*

	Function	Fourier Transform
16b. Alternative forms	$\int_{-\infty}^{+\infty} f(x)g(x)dx$	$= \int_{-\infty}^{+\infty} F(-y)G(y)dy$
	$\int_{-\infty}^{+\infty} f(u)G(u)du$	$= \int_{-\infty}^{+\infty} F(u)g(u)du$
16c. "Energy" relation	$\int_{-\infty}^{+\infty} f(x) ^2 dx$	$= \int_{-\infty}^{+\infty} F(y) ^2 dy$
17. Initial value theorem	<p>If $f(x) = 0$ for $x < 0$, is real, and contains no pulses</p> $f(0^+) = 2 \int_{-\infty}^{\infty} F_r(y)dy$ <p>where $F_r(y) = \text{Re}[F(y)]$</p>	
18. Relationships between $F_r(y)$ and $F_i(y)$	<p>a) If $f(x)=0$ for $x < 0$, is real, and contains no pulses</p> $\int_{-\infty}^{\infty} F_r^2(y)dy = \int_{-\infty}^{\infty} F_i^2(y)dy$ <p>where $F_i(y) = \text{Im}[F(y)]$</p> <p>b) The following integral relationships apply (Hilbert transforms)</p> $F_r(y) = 2 \int_{-\infty}^{\infty} [F_i(\tau)/2\pi(y-\tau)]d\tau$ $F_i(y) = -2 \int_{-\infty}^{\infty} [F_r(\tau)/2\pi(y-\tau)]d\tau$ <p>in which the Cauchy principal values of the integrals are taken</p>	

* In the table, functions of x are denoted by lower-case letters and their transforms by the corresponding capital letters.

TABLE 2. TIME AND FREQUENCY FUNCTIONS FOR COMMONLY ENCOUNTERED PULSE SHAPES*

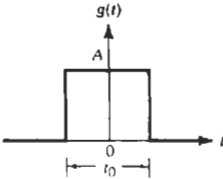
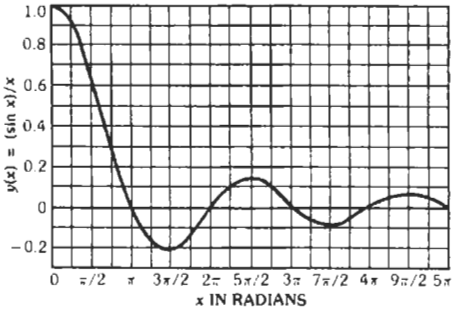
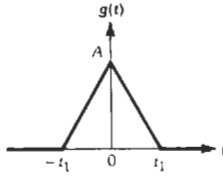
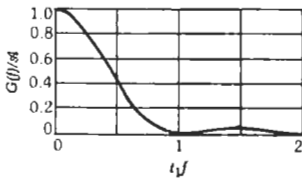
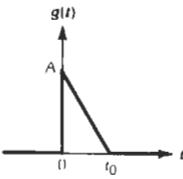
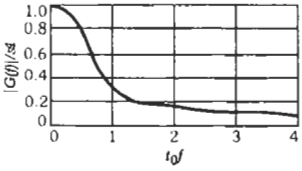
Time Function	Frequency Function
<p>A. Rectangular pulse</p>  <p style="text-align: center;">$g(t) = A$ for $-(1/2)t_0 < t < (1/2)t_0$ $= 0$ otherwise</p> <p>Area $\mathcal{A} = At_0$</p>	 <p style="text-align: center;">$G(f) = \mathcal{A} (\sin \alpha) / \alpha$</p> <p>where $\alpha = \pi_0 f$</p> <p>[See curve $(\sin x)/x$ above.]</p>
<p>B. Isosceles-triangle pulse</p>  <p style="text-align: center;">$g(t) = A[1 - (t/t_1)], 0 \leq t < t_1$ $= A[1 + (t/t_1)], -t < t \leq 0$ $= 0$, otherwise</p> <p>Area $\mathcal{A} = At_1$</p>	 <p style="text-align: center;">$G(f) = \mathcal{A} [(\sin \alpha) / \alpha]^2$</p> <p>where $\alpha = \pi_1 f$</p>
<p>C. Sawtooth pulse</p>  <p style="text-align: center;">$g(t) = A[1 - (t/t_0)], 0 < t < t_0$ $= 0$, otherwise</p> <p>Area $\mathcal{A} = (1/2)At_0$</p>	 <p style="text-align: center;">$G(f) = \mathcal{A} (j/\alpha) \{[(\sin \alpha) / \alpha] \exp(-j\alpha) - 1\}$ $= \mathcal{A} \frac{1 - \exp(-2j\alpha) - 2j\alpha}{2\alpha^2}$</p> <p>where $\alpha = \pi_0 f$</p>

TABLE 2 (CONT). TIME AND FREQUENCY FUNCTIONS FOR COMMONLY ENCOUNTERED PULSE SHAPES*

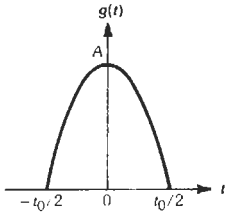
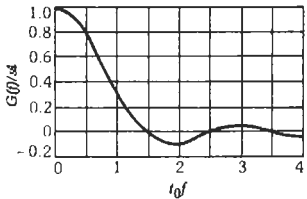
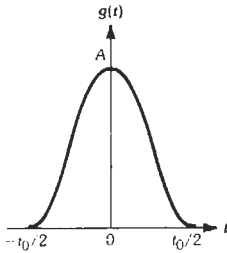
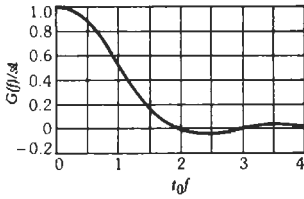
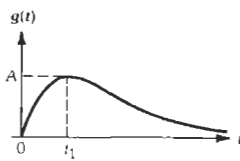
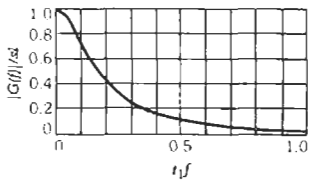
Time Function	Frequency Function
<p>D. Any pulse of polygonal form may be represented as a linear combination of waveforms such as A, B, and C above eventually after some shifts in time. The pulse spectrum is the same linear combination of the corresponding spectra (eventually modified according to property 4, Table 1).</p>	
<p>E. Cosine pulse</p>	
	
$g(t) = A \cos \pi(t/t_0), \quad -(1/2)t_0 < t < (1/2)t_0$ $= 0, \text{ otherwise}$	$G(f) = \mathcal{A} \{ [\cos(\pi/2)\alpha] / (1 - \alpha^2) \}$
<p>Area $\mathcal{A} = (2/\pi)At_0$</p>	<p>where $\alpha = 2t_0f$ For $\alpha = 1$, $G(f) = \mathcal{A}\pi/4$</p>
<p>F. Cosine-squared pulse</p>	
	
$g(t) = A \cos^2 \pi(t/t_0)$ $= (1/2)A [1 + \cos 2\pi(t/t_0)] \quad \left. \vphantom{g(t)} \right\} \quad -(1/2)t_0 < t < (1/2)t_0$ $= 0, \text{ otherwise}$	$G(f) = \mathcal{A} \frac{\sin \pi \alpha}{\pi \alpha (1 - \alpha^2)}$
<p>Area $\mathcal{A} = (1/2)At_0$</p>	<p>where $\alpha = t_0f$ For $\alpha = 1$, $G(f) = (1/2)\mathcal{A}$</p>
<p>G. Gaussian pulse</p>	
<p>Use curve of Fig. 1, Chapter 44 with standard deviation</p>	<p>Use curve of Fig. 1, Chapter 44 with standard deviation</p>
$\sigma = t_1$ $= (2 \ln 2)^{-1/2} t_{6dB}; \quad \Delta t_{6dB} \equiv 2t_{6dB}$ $g(t) = A \exp[-(1/2)(t/t_1)^2]$ $= A \exp[-(\ln 2)(t/t_{6dB})^2]$	$\sigma = f_1 = 1/2\pi t_1$ $= (2 \ln 2)^{1/2} / \pi \Delta t_{6dB}; \quad f_{3dB} = 2^{1/2} \ln 2 / \pi \Delta t_{6dB}$ $G(f) = \mathcal{A} \exp[-(1/2)(f/f_1)^2]$ $= (\ln 2/2\pi)^{1/2} (A/f_{3dB})$ $\times \exp[-(1/2)(\ln 2)(f/f_{3dB})^2]$ $= (1/2)(\pi/\ln 2)^{1/2} A \Delta t_{6dB}$ $\times \exp[-(\pi^2/4 \ln 2)(f \Delta t_{6dB})^2]$
<p>Area $\mathcal{A} = (2\pi)^{1/2} At_1$</p> $= (1/2)(\pi/\ln 2)^{1/2} A \Delta t_{6dB}$	

TABLE 2 (CONT). TIME AND FREQUENCY FUNCTIONS FOR COMMONLY ENCOUNTERED PULSE SHAPES*

Time Function	Frequency Function
<p>H. Critically damped exponential pulse</p>  $g(t) = Ae(t/t_1) \exp(-t/t_1), t > 0$ $= 0, t \leq 0$ <p>$e = 2.71828 \dots$</p> <p>Area $\mathcal{A} = Aet_1$</p>	 $G(f) = \mathcal{A} [1/(1+j\alpha)^2]$ <p>where $\alpha = 2\pi t_1 f$</p>

*For an extensive tabulation of the Fourier transform and its inverse, see Campbell, G. A., and Foster, R. M. *Fourier Integrals for Practical Applications*. New York: D. Van Nostrand Co., Inc., 1948. See also Erdélyi, A., ed. *Tables of Integral Transforms*, Vol. 1, Bateman Manuscript Project. New York: McGraw-Hill Book Co., 1954.

where

$$D_n = T^{-1} \int_0^T f(x) \exp(-jn \omega_0 x) dx$$

and

$$D_0 = (1/2)A_0 = (1/2)C_0$$

$$D_n = (1/2)(A_n - jB_n)$$

$$D_{-n} = (1/2)(A_n + jB_n)$$

$$n = 1, 2, \dots$$

Average Power

The average power of the periodic waveform $f(x)$ is given by

$$T^{-1} \int_0^T |f(x)|^2 dx = \sum_{n=-\infty}^{+\infty} |D_n|^2$$

$$= (1/4)C_0^2 + (1/2) \sum_{n=1}^{\infty} C_n^2$$

$$= (1/4)A_0^2 + (1/2) \sum_{n=1}^{\infty} (A_n^2 + B_n^2)$$

Odd and Even Functions

If $f(x)$ is an odd function, i.e.,

$$f(x) = -f(-x)$$

then all the coefficients of the cosinc terms (A_n) vanish, and the Fourier series consists of sine terms alone.

If $f(x)$ is an even function, i.e.,

$$f(x) = f(-x)$$

then all the coefficients of the sine terms (B_n) vanish, and the Fourier series consists of cosinc terms alone, and a possible constant.

The Fourier expansions of functions in general include both cosinc and sine terms. Every function capable of Fourier expansion consists of the sum of an even and an odd part.

$$f(x) = \underbrace{(1/2)A_0 + \sum_{n=1}^{\infty} A_n \cos n\omega_0 x}_{\text{even}} + \underbrace{\sum_{n=1}^{\infty} B_n \sin n\omega_0 x}_{\text{odd}}$$

To separate a general function $f(x)$ into its odd and even parts, use

$$f(x) \equiv \underbrace{(1/2) [f(x) + f(-x)]}_{\text{even}} + \underbrace{(1/2) [f(x) - f(-x)]}_{\text{odd}}$$

In some cases by suitable selection of the origin, the function may be made either odd or even, thus simplifying the expansion.

Odd or Even Harmonics

An odd or even function may contain odd or even harmonics. A condition that causes a function $f(x)$ of

period T to have only odd harmonics in its Fourier expansion is

$$f(x) = -f(x + T/2)$$

A condition that causes a function $f(x)$ of period T to have only even harmonics in the Fourier expansion is

$$f(x) = f(x + T/2)$$

To separate a general function $f(x)$ into its odd and even harmonics, use

$$f(x) \equiv \underbrace{(1/2)[f(x) + f(x + T/2)]}_{\text{even harmonics}} + \underbrace{(1/2)[f(x) - f(x + T/2)]}_{\text{odd harmonics}}$$

A periodic function may sometimes be changed from odd to even (and vice versa) by a shift of the origin, but

the presence of particular odd or even harmonics is unchanged by such a shift.

PULSE-TRAIN ANALYSIS

If the pulse defined by the function $g(t)$ is repeated every interval T , a periodic waveform

$$y(t) = \sum_{n=-\infty}^{+\infty} g(t-nT)$$

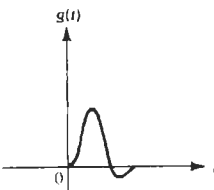
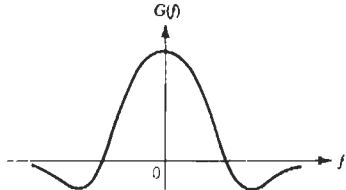
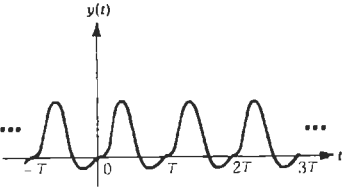
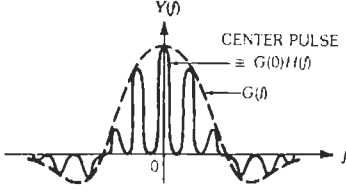
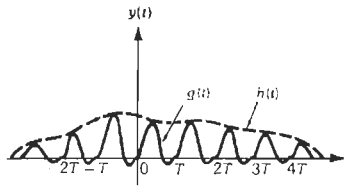
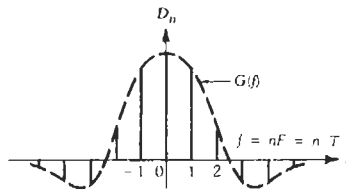
results with period T and repetition frequency $F = 1/T$ (see A and B in Table 3).

This pulse train may be expressed as a convolution product

$$y(t) = \left[\sum_{n=-\infty}^{+\infty} \delta(t - nT) \right] * g(t)$$

and, from application of properties 8 and 10 (Table 1), its Fourier transform is

TABLE 3. THE SPECTRUM FOR PULSE TRAIN*

Waveform	Spectrum
<p>A. Single pulse</p> 	
<p>B. Infinite periodic pulse train</p> 	
<p>C. Limited pulse train</p> 	

* Spectra are in general complex functions. They are represented here by real curves only to simplify the illustration.

$$Y(f) = (1/T) \left[\sum_{n=-\infty}^{+\infty} \delta(f - nF) \right] G(f)$$

The function $y(t)$ is represented by the Fourier series

$$y(t) = \sum_{n=-\infty}^{+\infty} D_n \exp(jn\omega_0 t)$$

where

$$D_n = (1/T)G(nF)$$

The coefficients D_n are obtained by sampling the pulse spectrum at frequencies that are multiples of the repetition frequency.

The amplitude C_n of the n th harmonic in the real representation (see subsection headed "Real Form of Fourier Series" earlier in this chapter) is

$$C_n = 2|D_n| = (2/T) |G(nF)|$$

The constant term of the series

$$D_0 = A_0/2 = C_0/2$$

is the average amplitude

$$A_{av} = \mathcal{A} T = G(0)/T$$

where

$$\mathcal{A} = \int_0^T g(t) dt$$

is the area under one pulse.

If the pulses do not overlap (i.e., if function $g(t)$ is zero outside of some period a to $a + T$), the energy in a pulse is

$$E = \int_a^{a+T} g^2(t) dt = \int_{-\infty}^{+\infty} |G(f)|^2 df$$

The root-mean-square amplitude is

$$A_{rms} = (E/T)^{1/2}$$

The average power of the pulse train is

$$\begin{aligned} E/T = A_{rms}^2 &= \sum_{n=-\infty}^{+\infty} |D_n|^2 \\ &= (1/4) C_0^2 + (1/2) \sum_{n=1}^{+\infty} C_n^2 \end{aligned}$$

A pulse train of finite extent, where all pulses have the same shape and are spaced periodically, may be represented as a product

$$y(t) = h(t) \sum_{n=-\infty}^{+\infty} g(t - nT)$$

The function $h(t)$ defines the envelope of the pulse train. The Fourier transform

$$Y(f) = (1/T) \sum_{n=-\infty}^{+\infty} G(nF)H(f - nF)$$

is given by a weighted sum of shifted Fourier transforms of $h(t)$. If $h(t) = 1$, then $H(f)$ is the δ function and the pulse train is a periodic waveform having a line spectrum as in B of Table 3. However, if $h(t)$ does not have infinite support, then this line spectrum is broadened. If $h(t)$ has support over a large number of periods T , so that $H(f)$ is narrow compared to F , then $Y(f)$ is approximately as shown in C of Table 3 where the envelope is $G(f)$ and each pulse is a scaled, shifted version of $H(f)$.

The Fourier series coefficients for a number of commonly encountered pulse trains are given in Table 4.

SPECTRAL ANALYSIS

If $g(t)$ is band-limited and also nearly time-limited, then approximate samples of its Fourier transform $G(f)$ can be computed from samples of $g(t)$ by using the discrete Fourier transform (DFT). The DFT of the sampled sequence $\{g(nT_s)\}_{n=0}^{N-1}$ is defined by

$$G_m = \sum_{n=0}^{N-1} g(nT_s) \exp[-j(2\pi/N)nm]$$

$$m = 0, \dots, N - 1$$

If the sampling frequency $1/T_s$ is greater than twice the highest frequency contained in $g(t)$, and if $g(t)$ is nearly zero outside the interval $[0, (N - 1)T_s]$, then it can be shown that

$$G_m \approx \begin{cases} (1/T_s) G\left(\frac{m}{NT_s}\right) & 0 \leq m < N/2 \\ (1/T_s) G\left(\frac{m - N}{NT_s}\right) & N/2 \leq m \leq N - 1 \end{cases}$$

where the argument of the Fourier transform G is in hertz. Ordinarily $g(t)$ is not time-limited, and only a finite-length section of $g(t)$ can be sampled and processed. This places a limit on achievable spectral resolution.

For additional discussion, including implementation of the DFT using the fast Fourier transform (FFT) algorithm, see Chapter 28.

TABLE 4. PERIODIC WAVEFORMS AND FOURIER SERIES

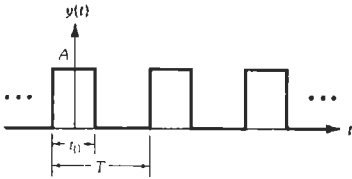
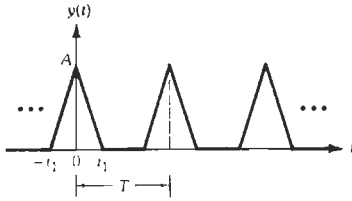
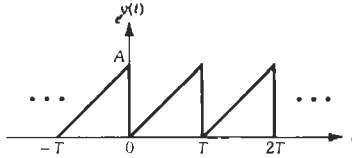
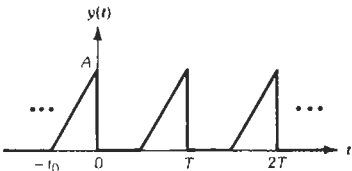
Waveform	Coefficient of Fourier Series
<p>A. Rectangular wave</p> 	$C_n = 2D_n = 2A_{Av} \left \frac{\sin(n\pi t_0/T)}{n\pi t_0/T} \right $ <p>Can be read off curve of $(\sin x)/x$, A in Table 2, by sampling at $n\pi t_0/T$</p> <p>Example: If $T = 2t_0$</p> $y(t) = 2A_{Av} [\frac{1}{2} + (2/\pi) \cos\theta - (2/3\pi) \cos 3\theta + \dots]$ <p>where $\theta = 2\pi t/T$</p>
<p>Derived from rectangular pulse, A in Table 2</p> $A_{Av} = A(t_0/T) \quad A_{rms} = A(t_0/T)^{1/2}$	
<p>B. Isosceles-triangle wave</p> 	$C_n = 2A_{Av} \left(\frac{\sin(n\pi t_1/T)}{n\pi t_1/T} \right)^2$ <p>Example: If $T = 2t_1$</p> $y(t) = 2A_{Av} [\frac{1}{2} + (2/\pi)^2 \cos\theta + (2/3\pi)^2 \cos 3\theta + \dots]$ <p>where $\theta = 2\pi t/T$</p>
<p>Derived from triangular pulse, B in Table 2</p> $A_{Av} = A(t_1/T) \quad A_{rms} = A(2t_1/3T)^{1/2}$	
<p>C. Sawtooth wave</p> 	$C_n = 2A_{Av}(1/n\pi)$ $y(t) = 2A_{Av} [\frac{1}{2} - (1/\pi) \sin\theta - (1/2\pi) \sin 2\theta - \dots]$
<p>Derived from triangular pulse, C in Table 2</p> $A_{Av} = A/2 \quad A_{rms} = A(3^{-1/2})$	
<p>D. Clipped sawtooth wave</p> 	$C_n = 2A_{Av}(1/\alpha^2) [\sin^2\alpha + \alpha(\alpha - \sin 2\alpha)]^{1/2}$ <p>where $\alpha = n\pi t_0/T$</p>
<p>Derived from triangular pulse, C in Table 2</p> $A_{Av} = A(t_0/2T) \quad A_{rms} = A(t_0/3T)^{1/2}$	

TABLE 4 (CONT). PERIODIC WAVEFORMS AND FOURIER SERIES

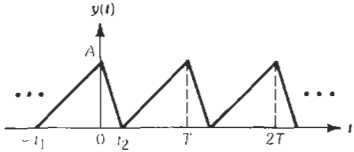
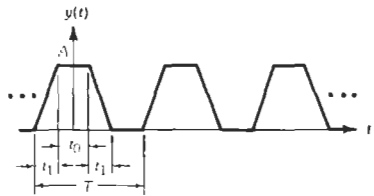
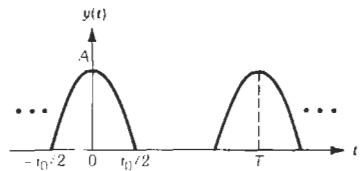
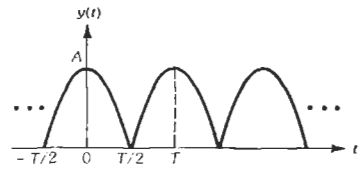
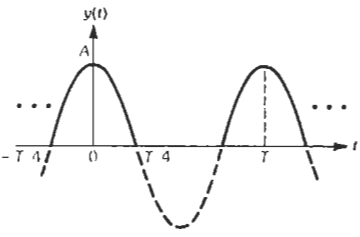
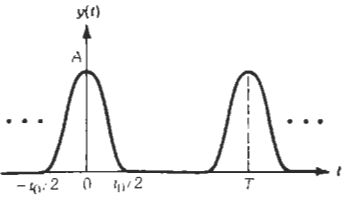
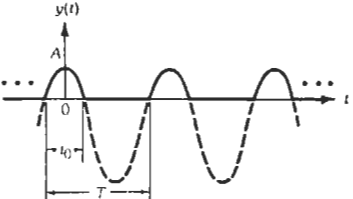
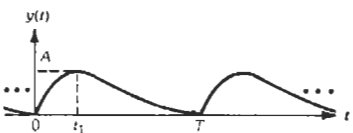
Waveform	Coefficient of Fourier Series
<p>E. Sawtooth wave</p>  <p>Derived from the sum of two triangular pulses, C in Table 2</p> $A_{Av} = A/2 \quad A_{rms} = A(3^{-1/2})$	$C_n = 2A_{Av}(T^2/\pi^2 n^2 t_1 t_2) \sin(n\pi t_1/T) $ <p>where $t_1 + t_2 = T$</p>
<p>F. Symmetrical trapezoidal wave</p>  <p>Derived as in D in Table 2</p> $A_{Av} = A[(t_0 + t_1)/T]$ $A_{rms} = A[(3t_0 + 2t_1)/3T]^{1/2}$	$D_n = A_{Av} \frac{\sin(\pi n_1/T) \sin[\pi n(t_1 + t_0)/T]}{\pi n_1/T \quad \pi n(t_1 + t_0)/T}$ $C_n = 2 D_n $
<p>G. Train of cosine pulses</p>  <p>Derived from cosine pulse, E in Table 2</p> $A_{Av} = (2/\pi)A(t_0/T) \quad A_{rms} = A(t_0/2T)^{1/2}$	$C_n = 2A_{Av} \left \frac{\cos(n\pi t_0/T)}{1 - (2nt_0/T)^2} \right $ <p>For $nt_0T = 1/2$, this becomes $\pi A_{Av}/2$</p>
<p>H. Full-wave-rectified sine wave</p>  <p>Derived from cosine pulse, E in Table 2 (same as G in Table 4 with $t_0 = T$)</p> $A_{Av} = (2/\pi)A \quad A_{rms} = A/(2^{1/2})$	$C_0 = 2A_{Av}$ $C_n = 2A_{Av}(4n^2 - 1)^{-1}, \text{ for } n \neq 0$ $y(t) = 2A_{Av} [1/2 + (1/3) \cos\theta - (1/15) \cos 2\theta + (1/35) \cos 3\theta \dots - (-1)^n (4n^2 - 1)^{-1} \cos n\theta \dots]$ <p>where $\theta = 2\pi t/T$</p>

TABLE 4 (CONT). PERIODIC WAVEFORMS AND FOURIER SERIES

Waveform	Coefficient of Fourier Series
<p>I. Half-wave-rectified sine wave</p> 	$C_0 = 2A_{Av}$ $C_{2n+1} = 0, \text{ except for } C_1 = 2A_{Av}(\pi/4)$ $C_{2n} = 2A_{Av}(4n^2 - 1)^{-1}, \text{ for } n \neq 0$ $y(t) = 2A_{Av}[1/2 + (\pi/4) \cos\theta + (1/3) \cos 2\theta - (1/15) \cos 4\theta + \dots - (-1)^n(4n^2 - 1)^{-1} \cos 2n\theta \dots]$
<p>Derived from cosine pulse, E in Table 2 (same as G in Table 4 with $t_0 = T/2$)</p> $A_{Av} = (1/\pi)A \quad A_{rms} = A/2$	
<p>J. Train of cosine-squared pulses</p> 	$C_n = 2A_{Av} \left \frac{\sin(n\pi t_0/T)}{(n\pi t_0/T)[1 - (nt_0/T)^2]} \right $
<p>Derived from cosine-squared pulse, F in Table 2</p> $A_{Av} = (1/2)A(t_0/T) \quad A_{rms} = (1/2)4(3t_0/2T)^{1/2}$	
<p>K. Fractional sine wave</p> 	$C_n = 2A_{Av} \left \frac{\sin n\alpha \cos\alpha - n \sin\alpha \cos n\alpha}{n(n^2 - 1)(\sin\alpha - \alpha \cos\alpha)} \right $
$A_{Av} = \frac{A \sin\alpha - \alpha \cos\alpha}{\pi (1 - \cos\alpha)}$ $A_{rms} = \frac{A}{(2\pi)^{1/2}} \frac{[2\alpha + \alpha \cos 2\alpha - (3/2) \sin 2\alpha]^{1/2}}{1 - \cos\alpha}$ <p>where $\alpha = \pi t_0/T$</p>	<p>when $n = 1,$</p> $C_1 = A_{Av} \left \frac{\alpha - \cos\alpha \sin\alpha}{\sin\alpha - \alpha \cos\alpha} \right $
<p>L. Critically damped exponential wave</p> 	$C_n = 2A_{Av}[1 + (2\pi n t_1/T)^2]^{-1}$ $= 2A_{Av} \cos^2 \theta_n$ <p>where $\tan \theta_n = 2\pi n t_1/T$</p>
<p>Derived from exponential pulse, H in Table 2 (period $T \gg$ period t_1 to make overlap negligible)</p> $A_{Av} = Ae(t_1/T) \quad A_{rms} = (Ae/2)(t_1/T)^{1/2}$ <p>$e = 2.71828 \dots$</p>	

8

Filters, Simple Bandpass Design

Coefficient of Coupling 8-2

Gain at Resonance 8-2

 Single Circuit

 Pair of Coupled Circuits

Selectivity Far From Resonance 8-3

 For Inductive Coupling

 For Capacitive Coupling

Selectivity of Single- and Double-Tuned Circuits Near
Resonance 8-3

 Selectivity and Phase Shift of Single-Tuned Circuits

Q Determination by 3-Decibel Points

 Selectivity and Phase Shift of Pairs of Coupled Tuned Circuits

Node Input Impedance or Mesh Input Admittance of a Double-Tuned
Circuit 8-10

COEFFICIENT OF COUPLING

Several types of coupled circuits are shown in Table 1 (pages 8-4 and 8-5) and Figs. 1 and 2, together with equations for the coefficient of coupling. Also shown is the dependence of bandwidth on resonance frequency. This dependence is only a rough approximation to show the trend and may be altered radically if L_m , M , or C_m is adjusted in tuning to various frequencies.

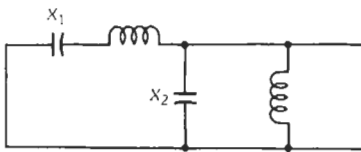
- K_{12} = coefficient of coupling between resonant circuits
- X_{10} = reactance of inductor (or capacitor) of first circuit at f_0
- X_{20} = reactance of similar element of second circuit at f_0
- $(bw)_C$ = bandwidth with capacitive tuning
- $(bw)_L$ = bandwidth with inductive tuning

GAIN AT RESONANCE

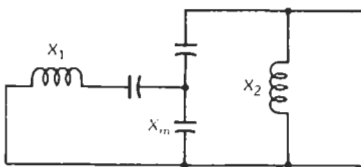
Single Circuit

In Table 1A

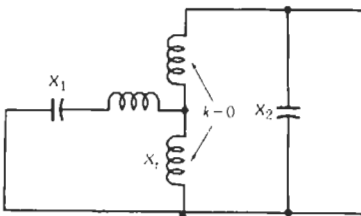
$$E_0/E_g = -g_m |X_{10}| Q$$



$$K_{12} = (X_2/X_1)^{1/2} = (L_2/L_1)^{1/2} = (C_1/C_2)^{1/2}$$



$$K_{12} = X_m / (X_1 X_2)^{1/2}$$



$$K_{12} = X_1 / (X_1 X_2)^{1/2}$$

where,

- E_0 = output volts at resonance frequency f_0 ,
- E_g = input volts to active device,
- g_m = transconductance of active device.

Pair of Coupled Circuits

Pairs of coupled circuits are illustrated in Figs. 3 and 4.

In B through F in Table 1

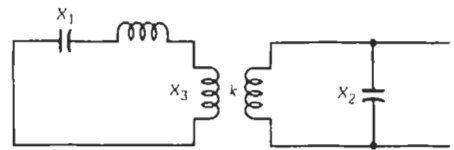
$$E_0/E_g = jg_m(X_{10}X_{20})^{1/2} Q [K_{12}Q / (1 + K_{12}^2Q^2)]$$

This is maximum at critical coupling, where $K_{12}Q = 1$.

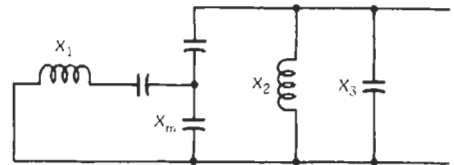
$Q = (Q_1Q_2)^{1/2}$ = geometric-mean Q for the two circuits, as loaded with active device input and output impedances

For circuits with critical coupling and overcoupling, the approximate gain is

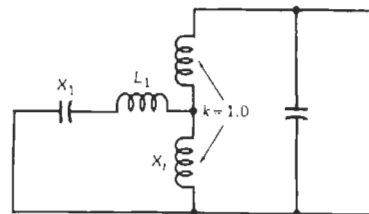
$$|E_0/E_g| \approx 0.1g_m / (C_1C_2)^{1/2} (bw)$$



$$K_{12} = k(X_3/X_1)^{1/2}$$



$$K_{12} = [X_m / (X_1 X_2)^{1/2}] (1 - X_2/X_3)$$



$$K_{12} = (X_1/X_1)^{1/2} = (L_1/L_1)^{1/2}$$

Fig. 1. Additional coefficient-of-coupling configurations (the node resonator is tuned to the desired midfrequency with the mesh resonator open-circuited, or the mesh resonator is tuned to this midfrequency with the node resonator short-circuited).

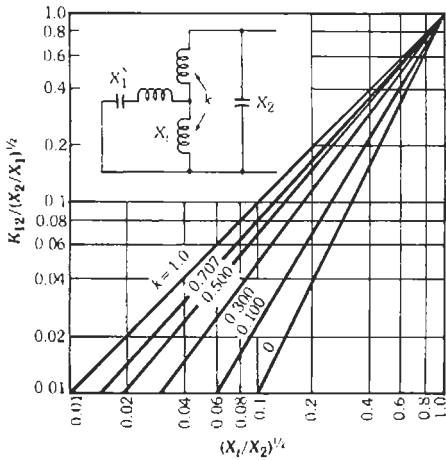


Fig. 2. Coefficient of coupling for configuration shown.

where,

- (bw) is the useful passband in megahertz,
- g_m is in micromhos,
- C is in picofarads.

SELECTIVITY FAR FROM RESONANCE

The selectivity curves of Fig. 5 are based on the presence of only a single type of coupling between the circuits. The curves are useful beyond the peak region treated in the section on selectivity near resonance.

In the equations for selectivity in Table 1

E = output volts at signal frequency f for same value of E_0 as that producing E_0

For Inductive Coupling

$$A = \frac{Q^2}{1 + K_{12}^2 Q^2} \left[\left(\frac{f}{f_0} - \frac{f_0}{f} \right)^2 - K_{12}^2 \left(\frac{f}{f_0} \right)^2 \right]$$

$$\approx \frac{Q^2}{1 + K_{12}^2 Q^2} \left(\frac{f}{f_0} - \frac{f_0}{f} \right)^2$$

For Capacitive Coupling

A similar equation for A applies, except that the neglected term is $-K_{12}^2(f_0/f)^2$. The 180-degree phase shift far from resonance is indicated by the minus sign in the expression for E_0/E .

Example. The use of the curves in Figs. 5, 6, and 7 is indicated by the following example. Given the circuit of Table 1C with input to PB across capacitor C_1 . Let $Q = 50$, $K_{12}Q = 1.50$, and $f_0 = 16.0$ megahertz. Required is the response at $f = 8.0$ megahertz.

Hence, $f/f_0 = 0.50$ and curve C , Fig. 5, gives -75 decibels. Then applying the corrections from Figs. 6 and 7 for Q and $K_{12}Q$, we find

$$\text{Response} = -75 + 12 + 4 = -59 \text{ decibels}$$

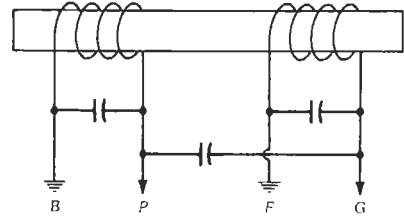


Fig. 3. Connection wherein k_m opposes k_c (k_c may be due to stray capacitance). Peak of attenuation is at $f = f_0 (-k_m/k_c)^{1/2}$. Reversing connections or winding direction of one coil causes k_m to aid k_c .

SELECTIVITY OF SINGLE- AND DOUBLE-TUNED CIRCUITS NEAR RESONANCE

Equations and curves are presented for the selectivity and phase shift of n single-tuned circuits and of m pairs of coupled tuned circuits. The conditions assumed are

- (A) All circuits are tuned to the same frequency, f_0 .
- (B) All circuits have the same Q , or each pair of circuits includes one circuit having Q_1 and the other having Q_2 .
- (C) Otherwise the circuits need not be identical.
- (D) Each successive circuit or pair of circuits is isolated from the preceding and following ones by active devices, with no regeneration around the system.

Certain approximations have been made to simplify the equations. In most actual applications of the types of circuits treated, the error involved is negligible from a practical standpoint. Over the narrow frequency band in question, it is assumed that

- (A) The reactance around each circuit is equal to $2X_0\Delta f/f_0$.
- (B) The resistance of each circuit is constant and equal to X_0/Q .
- (C) The coupling between two circuits of a pair is reactive and constant. (When an untuned link is used to couple the two circuits, this condition frequently is far from satisfied, resulting in a lopsided selectivity curve.)
- (D) The equivalent input voltage, taken as being in series with the tuned circuit (or the first of a

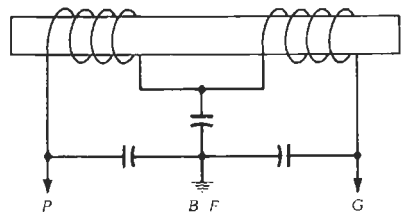
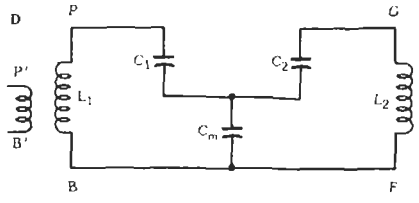


Fig. 4. Connection wherein k_m aids k_c . If mutual-inductance coupling is reversed, k_m will oppose k_c and there will be a transfer minimum at $f = f_0(-k_m/k_c)^{1/2}$.

TABLE I. SEVERAL TYPES OF COUPLED CIRCUITS, SHOWING COEFFICIENT OF COUPLING AND SELECTIVITY EQUATIONS

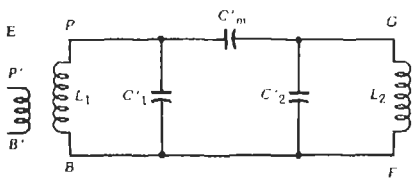
Diagram	Coefficient of Coupling	Approximate Bandwidth Variation With Frequency	Selectivity Far From Resonance	
			Equation*	Curve in Fig. 5
<p>A</p>			Input to PB or to P'B': $E_0/E = jQ[(f/f_0) - (f_0/f)]$	A
<p>B</p>	$K_{12} = L_m / [(L_1 + L_m)(L_2 + L_m)]^{1/2}$ $= \omega_0^2 L_m (C_1 C_2)^{1/2}$ $\approx L_m / (L_1 L_2)^{1/2}$	$(bw)_C \propto f_0$ $(bw)_L \propto f_0^3$	Input to PB: $E_0/E = -A(f/f_0)$ Input to P'B': $E_0/E = -A(f_0/f)$	C D
<p>C</p>	$K_{12} = M / (L_1 L_2)^{1/2}$ $= \omega_0^2 M (C_1 C_2)^{1/2}$ <p>M may be positive or negative</p>	$(bw)_C \propto f_0$ $(bw)_L \propto f_0^3$	Input to PB: $E_0/E = -A(f/f_0)$ Input to P'B': $E_0/E = -A(f_0/f)$	C D



$$K_{12} = -\{C_1 C_2 / [(C_1 + C_m)(C_2 + C_m)]\}^{1/2} \quad (\text{bw})_C \propto 1/f_0 \quad \text{Input to } PB \text{ or to } P'B': \quad D$$

$$\approx -1/\omega_0^2 C_m (L_1 L_2)^{1/2} \quad (\text{bw})_L \propto f_0 \quad E_0/E = -A(f_0/f)$$

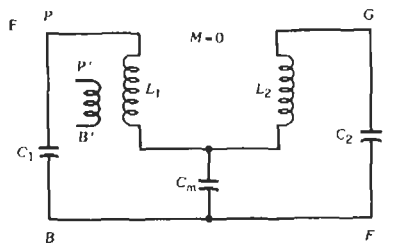
$$\approx -(C_1 C_2)^{1/2} / C_m$$



$$K_{12} = -C_m' / [(C_1' + C_m')(C_2' + C_m')]^{1/2} \quad (\text{bw})_C \propto f^3 \quad \text{Input to } PB \text{ or to } P'B': \quad D$$

$$\approx -\omega_0^2 C_m' (L_1 L_2)^{1/2} \quad (\text{bw})_L \propto f \quad E_0/E = -A(f_0/f)$$

$$\approx -C_m' / (C_1' C_2')^{1/2}$$



$$K_{12} = -\{C_1 C_2 / [(C_1 + C_m)(C_2 + C_m)]\}^{1/2} \quad (\text{bw})_C \propto 1/f_0 \quad \text{Input to } PB: \quad B$$

$$\approx -1/\omega_0^2 C_m (L_1 L_2)^{1/2} \quad (\text{bw})_L \propto f_0 \quad E_0/E = -A(f/f_0)^3$$

$$\approx -(C_1 C_2)^{1/2} / C_m \quad \text{Input to } P'B': \quad C$$

$$E_0/E = -A(f/f_0)$$

* Where $A = [Q^2 / (1 + K_{12}^2 Q^2)] / [(f/f_0) - (f_0/f)]^2$

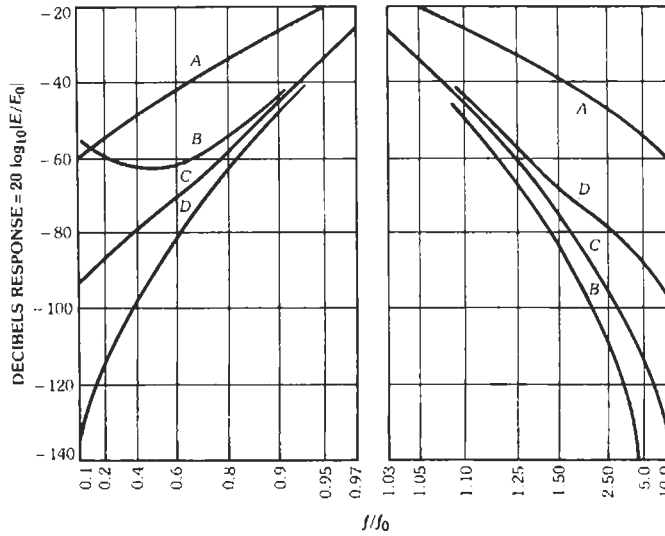


Fig. 5. Selectivity for frequencies far from resonance. $Q = 100$ and $|K_{12}| Q = 1.0$.

pair), is assumed to bear a constant proportionality to the input voltage of the active device or other driving source, at all frequencies in the band.

- (E) Likewise, the output voltage across the circuit (or the final circuit of a pair) is assumed to be proportional only to the current in the circuit.

The following symbols are used in the equations in addition to those defined previously:

- $\Delta f/f_0 = (f - f_0)/f_0 =$ (deviation from resonance frequency)/(resonance frequency)
- (bw) = bandwidth = $2\Delta f$
- X_0 = reactance at f_0 of inductor in tuned circuit
- n = number of single-tuned circuits
- m = number of pairs of coupled circuits
- ϕ = phase shift of signal at f relative to shift at f_0 as signal passes through cascade of circuits

$p = K_{12}^2 Q^2$ or $p = K_{12}^2 Q_1 Q_2$, a parameter determining the form of the selectivity curve of coupled circuits

$$B = p - (1/2)[(Q_1/Q_2) + (Q_2/Q_1)]$$

Selectivity and Phase Shift of Single-Tuned Circuits

The configuration of a single-tuned circuit is shown in Fig. 8.

$$E/E_0 = \{[1 + (2Q\Delta f/f_0)^2]^{-1/2}\}^n$$

$$\Delta f/f_0 = \pm(2Q)^{-1}[(E_0/E)^{2/n} - 1]^{1/2}$$

Decibel response = $20 \log_{10}(E/E_0)$

(dB response of n circuits) = $n \times$ (dB response of single circuit)

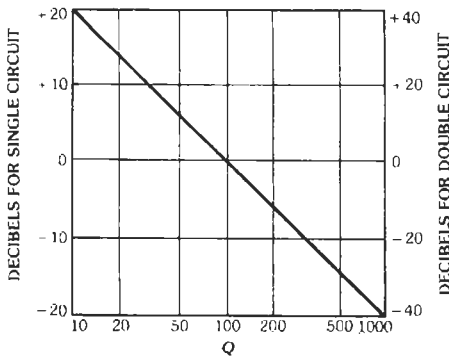


Fig. 6. Correction for $Q \neq 100$.

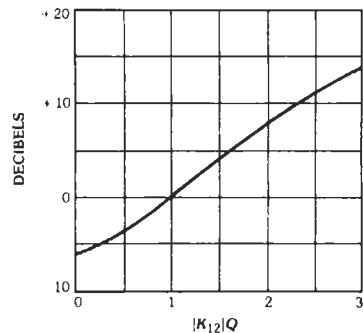


Fig. 7. Correction for $|K_{12}| Q \neq 1.0$.

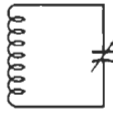


Fig. 8. Single-tuned circuit.

$$\phi = n \tan^{-1} (-2Q\Delta f/f_0)$$

These equations are plotted in Figs. 9 and 10.

Example of the use of Figs. 9 and 10: Suppose there are three single-tuned circuits ($n = 3$). Each circuit has a Q of 200 and is tuned to 1000 kilohertz. The results are shown in Table 2.

Q Determination by 3-Decibel Points

For a single-tuned circuit, when

$$E/E_0 = 0.707 \text{ (3 decibels down)}$$

$$Q = f_0/2\Delta f$$

$$= (\text{resonance frequency})/(\text{bandwidth})_{3\text{dB}}$$

Selectivity and Phase Shift of Pairs of Coupled Tuned Circuits

An example of a pair of coupled tuned circuits is shown in Fig. 11.

CASE 1: When $Q_1 = Q_2 = Q$:

These equations can be used with reasonable accuracy when Q_1 and Q_2 differ by ratios up to 1.5 or even 2 to 1. In such cases, use the value $Q = (Q_1Q_2)^{1/2}$.

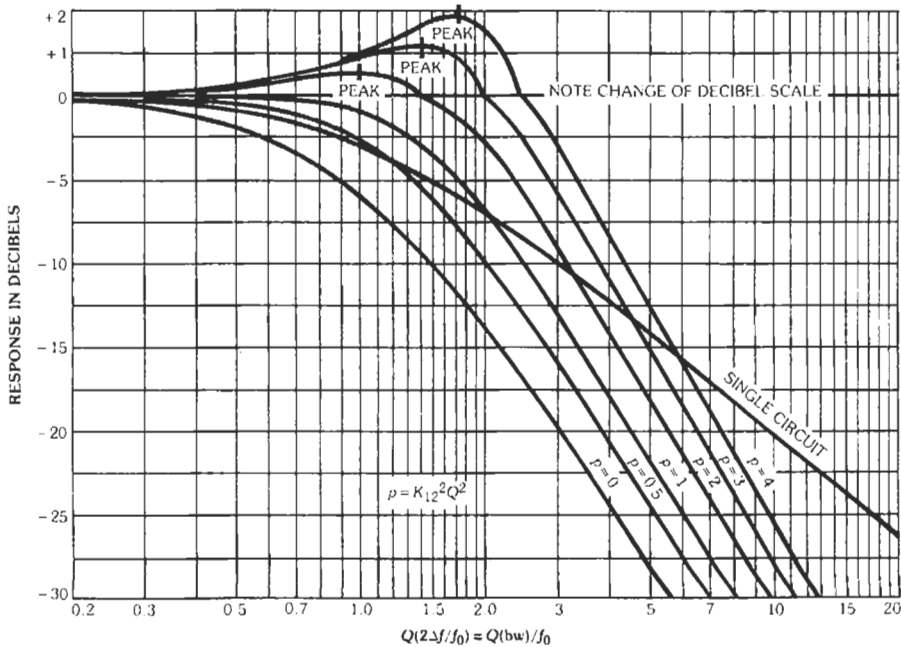
$$E/E_0 = \left[\frac{p + 1}{\{(2Q\Delta f/f_0)^2 - (p - 1)\}^2 + 4p} \right]^m$$

$$\Delta f/f_0 = \pm(2Q)^{-1}$$

$$\times \{(p - 1) \pm [(p + 1)^2(E_0/E)^{2/m} - 4p]^{1/2}\}^{1/2}$$

For very small values of E/E_0 ,

$$E/E_0 = [(p + 1)/(2Q\Delta f/f_0)^2]^m$$



The selectivity curves are symmetrical about the axis $Q\Delta f/f_0 = 0$ for practical purposes. Extrapolation beyond lower limits of chart:

Δ Response for Doubling Δf	Circuit	Useful Limit	
		at (bw)/f₀	Error becomes
- 6 dB	- single - -	0.6	1 to 2 dB
- 12 dB	- pair - -	0.4	3 to 4 dB

Fig. 9. Selectivity curves showing response of a single circuit, $n = 1$, and a pair of coupled circuits, $m = 1$.

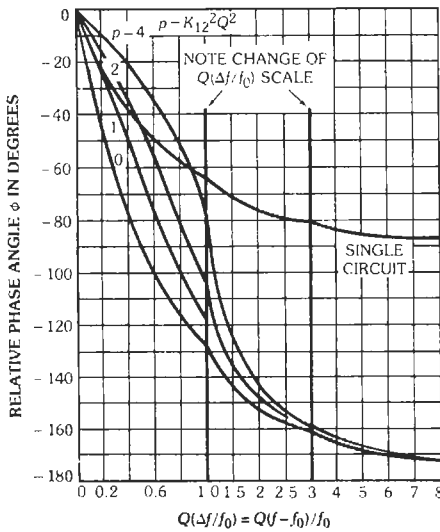


Fig. 10. Phase-shift curves for a single circuit, $n = 1$, and a pair of coupled circuits, $m = 1$. For $f > f_0$, ϕ is negative, while for $f < f_0$, ϕ is positive. The numerical value is identical in either case for the same $|f - f_0|$.

Decibel response = $20 \log_{10}(E/E_0)$
 (dB response of m pairs of circuits) = $m \times$ (dB response of one pair)

$$\phi = m \tan^{-1} \left[\frac{-4Q\Delta f/f_0}{(p+1) - (2Q\Delta f/f_0)^2} \right]$$

As p approaches zero, the selectivity and phase shift approach the values for n single circuits, where $n = 2m$ (gain also approaches zero).

The above equations are plotted in Figs. 9 and 10.

For Overcoupled Circuits ($p > 1$):

Location of peaks:

$$(f_{\text{peak}} - f_0)/f_0 = \pm(2Q)^{-1}(p - 1)^{1/2}$$

Amplitude of peaks:

$$E_{\text{peak}}/E_0 = [(p + 1)/2(p^{1/2})]^m$$

Phase shift at peaks:

$$\phi_{\text{peak}} = m \tan^{-1} [\mp(p - 1)^{1/2}]$$

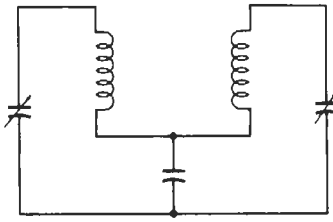


Fig. 11. One of several types of coupling.

Approximate passband (where $E/E_0 = 1$) is

$$(f_{\text{unity}} - f_0)/f_0 = \sqrt{2} [(f_{\text{peak}} - f_0)/f_0] \\ = \pm Q^{-1} [(1/2)(p - 1)]^{1/2}$$

CASE 2: General equation for any Q_1 and Q_2 :

$$E/E_0 = \left[\frac{p + 1}{\{(2Q\Delta f_0/f_0)^2 - B\}^2 + (p + 1)^2 - B^2} \right]^m$$

$$B = p - (1/2)[(Q_1/Q_2) + (Q_2/Q_1)]$$

$$\Delta f/f_0 = \pm(2Q)^{-1} \\ \times \{B \pm [(p + 1)^2(E_0/E)^{2/m} - (p + 1)^2 + B^2]^{1/2}\}^{1/2}$$

$$\phi = m \tan^{-1}$$

$$\times \left(-\frac{2Q \Delta f/f_0 [(Q_1/Q_2)^{1/2} + (Q_2/Q_1)^{1/2}]}{(p + 1) - (2Q \Delta f/f_0)^2} \right)$$

For Overcoupled Circuits:

Location of peaks:

$$(f_{\text{peak}} - f_0)/f_0 = \pm B^{1/2}/2Q \\ = \pm (1/2)[K_{12}^2 - (1/2)(1/Q_1^2 + 1/Q_2^2)]^{1/2}$$

Amplitude of peaks:

$$E_{\text{peak}}/E_0 = \{(p + 1)/[(p + 1)^2 - B^2]^{1/2}\}^m$$

CASE 3: Peaks just converged to a single peak:

$$\text{Here } B = 0 \text{ or } K_{12}^2 = (1/2)(1/Q_1^2 + 1/Q_2^2)$$

$$E/E_0 = \{2/[(2Q'\Delta f/f_0)^4 + 4]^{1/2}\}^m$$

TABLE 2. EXAMPLE OF USE OF FIGS. 9 AND 10

Abscissa $Q(\text{bw})/f_0$	Bandwidth (kilohertz)	Ordinate dB Response for $n=1$	Decibels Response for $n=3$	ϕ^* for $n=1$	ϕ for $n=3$
1.0	5.0	-3.0	-9	$\mp 45^\circ$	$\mp 135^\circ$
3.0	15	-10.0	-30	$\mp 71\frac{1}{2}^\circ$	$\mp 215^\circ$
10.0	50	-20.2	-61	$\mp 84^\circ$	$\mp 252^\circ$

* ϕ is negative for $f > f_0$, and vice versa.

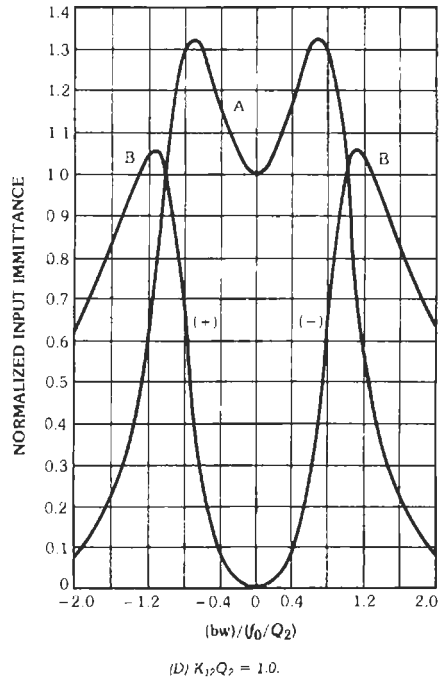
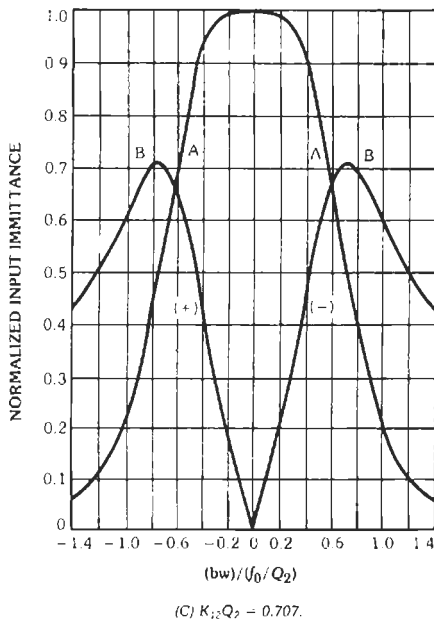
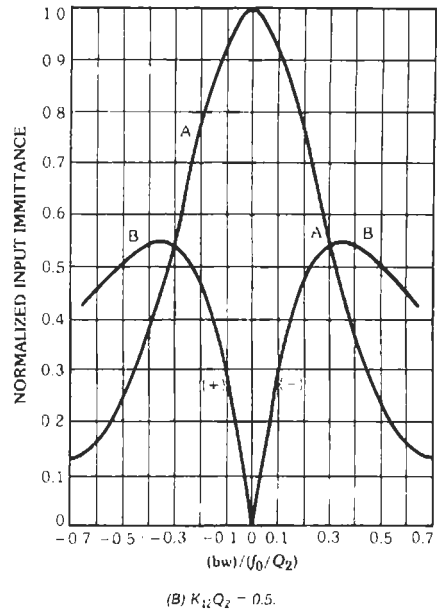
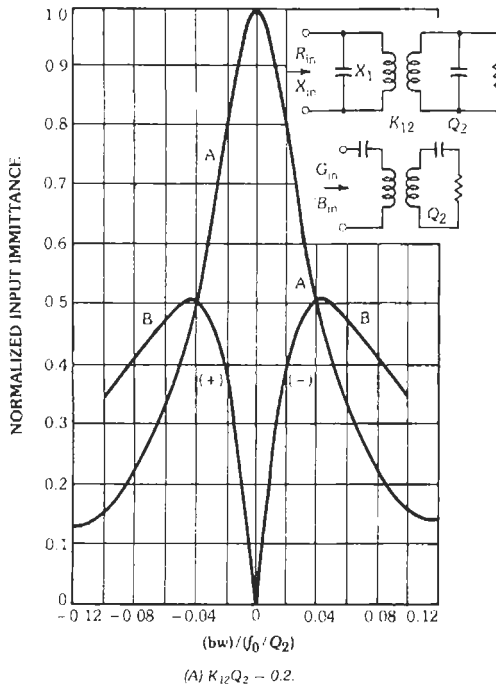


Fig. 12. Normalized input immittance versus normalized frequency of double-tuned circuits. $A = R_{in}/(X_1/K_{12}^2Q_2)$ or $G_{in}/(B_1/K_{12}^2Q_2)$. $B = X_{in}/(X_1/K_{12}^2Q_2)$ or $B_{in}/(B_1/K_{12}^2Q_2)$.

where,

$$Q' = 2Q_1Q_2/(Q_1 + Q_2)$$

$$\Delta f/f_0 = \pm(\sqrt{2}/4)(1/Q_1 + 1/Q_2)[(E_0/E)^{2/m} - 1]^{1/4}$$

$$\phi = m \tan^{-1} \left[-\frac{4Q' \Delta f/f_0}{2 - (2Q' \Delta f/f_0)^2} \right]$$

The curves of Figs. 9 and 10 may be applied to this case, using the value $p = 1$ and substituting Q' for Q .

**NODE INPUT IMPEDANCE OR
MESH INPUT ADMITTANCE OF
A DOUBLE-TUNED CIRCUIT**

Fig. 12 gives the normalized input immittance versus the normalized frequency of double-tuned circuits.

9 Filters, Modern-Network- Theory Design

Ford Shepherd

Introduction	9-2
Normalization	9-2
The Reference Low-Pass	9-3
Cauer-Parameter Low-Pass Filters	9-4
Chebyshev Low-Pass Filters	9-14
Butterworth Low-Pass Filters	9-14
Low-Pass to High-Pass Transformation	9-14
Low-Pass to Bandpass Transformation	9-14
Low-Pass to Band-Stop Transformation	9-19
Numerical Examples	9-19

In contrast to filter design by image-parameter methods, the design of filters by the use of modern network theory is a domain for specialists with digital computers because of the complex calculations required. There are sufficient advantages of filter circuits computed by this method, however, to warrant making some application of them easy and straightforward. This chapter focuses on a very limited subset of the limitless possibilities for low-pass networks and simple transform methods to allow calculation of high-pass, bandpass, and band-stop circuits.

The design information is drawn from experience in the application of modern network theory to the design of electric wave filters. As stated above, only limited design results are supplied, and a concentrated study of the cited references is essential to gain a working knowledge of the synthesis process through which these results were computed. References 1, 4, and 6 provide details of the design theory. Reference 5 provides a concise summary of the theory with graphs and tables to enable an engineer to compute filter circuits with the help of a computer program.

Reference 2 provides a much larger tabulation of Cauer-parameter and Chebyshev filter networks up to degree 9, and Saal has also produced another volume that extends to degree 15. Reference 3 presents many practical ideas (drawn from a 25-year career in Europe and the USA) on designing, testing, and manufacturing filters and mentions two of the computer programs that are available. Many books and articles written on this subject since the work of Cauer and Darlington in the late 1930s are more than worthy of mention here. However, the scope of this chapter and the space available do not necessitate nor permit a detailed discussion. Reference 7 is an example of some of the work done toward practical implementation of these filter networks with standard-value capacitors.

No attempt has been made to present details of the theory and formulation involved, of the approximation of transfer polynomials to performance requirements, or the very useful but less frequently required topics such as zero and infinite terminations, phase and delay performance, and the effects of and compensation for the losses in real coils and capacitors.

INTRODUCTION

Filter networks continue to be of great importance in the design of electrical equipment, especially in communication engineering. Unlike previous methodologies to design spectrum-shaping networks, modern network theory enables the engineer to design filter networks that are based on the actual requirements for signal transmission. While image-parameter design is rather simplistic, only a very limited approximation of the specific requirements can be achieved. Today's methods are not so straightforward and are generally considered beyond the scope of nonspecialists, since the mathematical design process does not directly parallel physical conceptions and the calculations are complex

and extensive. The development of digital computers has led to the capture of much of the knowledge of the specialist and the complex calculation algorithms into programs that allow many engineers to design some of their own networks. Still, the sophisticated requirements necessitate the special expert.

Many of the less complex requirements can be satisfied by designs that can be done with the procedures described here and without access to the computer programs. To accomplish this task, the requirements for all the filter types considered here (low-pass, high-pass, bandpass and band-stop) are transformed to a set of reference low-pass requirements. From the reference low-pass requirements, the network complexity is determined, and a normalized reference low-pass is selected. Then the suitable transformation is applied to this low-pass to arrive at the network to satisfy the initial requirements.

Fig. 1 shows a "reasonably typical" requirement for a filter network. To satisfy this requirement with the minimum network, it is necessary to use one of the available computer programs. However, it is possible to satisfy the modified requirements, as shown in Fig. 2, by the simple computations of this chapter. The requirements depicted by the tolerance plot of Fig. 2 can be transformed to the basic reference low-pass tolerance plot shown by Fig. 3 with $A_p = A_1$ and $A_s = A_6$. The algorithm for this transform is covered in the section headed "The Reference Low-Pass."

NORMALIZATION

To simplify the calculations, frequencies are normalized with respect to some reference frequency, f_{ref} , and the impedances are normalized with respect to a refer-

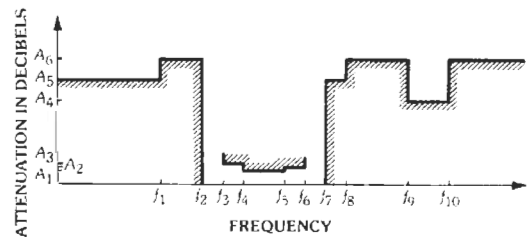


Fig. 1. Requirement for a filter network.

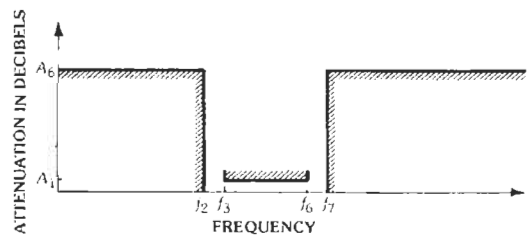


Fig. 2. Modified requirements for filter network.

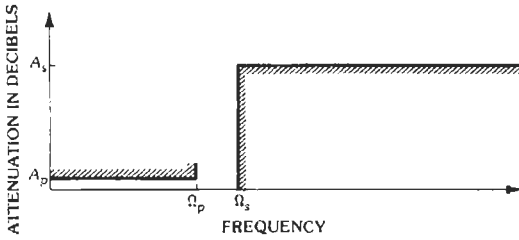


Fig. 3. Basic reference low-pass tolerance plot.

ence resistance, R_{ref} . Hence, normalized frequencies are defined by $\Omega_i = f_i/f_{ref}$, and the normalized resistances are defined by $r_i = R_i/R_{ref}$.

In the case of low-pass filters, it is appropriate to use the upper edge of the passband (or cutoff frequency), f_p , as the reference frequency (hence, $\Omega_p = 1.0$) and use the input resistance, R_1 , as the reference resistance ($r_1 = 1.0$).

Conversely, all normalized data of a circuit can be converted to any frequency range and any impedance level by selection of f_{ref} and R_{ref} . By defining

$$L_{ref} = R_{ref}/2\pi f_{ref}$$

and

$$C_{ref} = 1/2\pi f_{ref}R_{ref}$$

the entire circuit can be computed for the appropriate range of frequency and impedance.

$$R_i = r_i \cdot R_{ref} \text{ (ohms)}$$

$$f_i = \Omega_i \cdot f_{ref} \text{ (hertz)}$$

$$L_i = l_i \cdot L_{ref} \text{ (henrys)}$$

$$C_i = c_i \cdot C_{ref} \text{ (farads)}$$

THE REFERENCE LOW-PASS

Fig. 4 shows a normal tolerance plot of the requirements for a low-pass filter. Also shown are the associated formulas to calculate the quantities for the reference low-pass of Fig. 3. With $f_{ref} = f_p$ ($\Omega_p = 1.0$), three

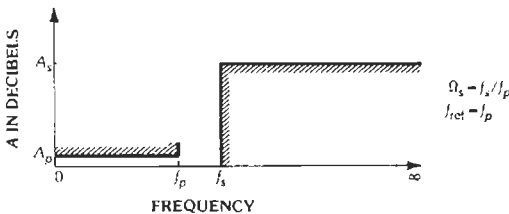


Fig. 4. Requirements for a low-pass filter.

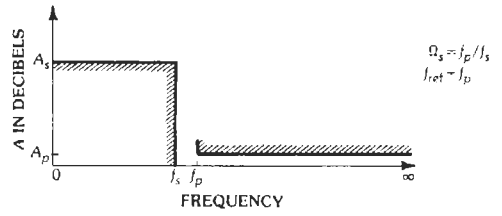


Fig. 5. Requirements for a high-pass filter.

parameters are required: (1) $\Omega_s = f_s/f_{ref}$; (2) A_p (decibels), the maximum acceptable ripple in the passband; and (3) A_s (decibels), the minimum attenuation in the stopband.

Fig. 5 shows a tolerance plot of requirements for a high-pass filter. Also shown are the equations to determine the reference low-pass parameters.

Fig. 6 shows a tolerance plot of requirements for a bandpass filter, together with the equations to resolve the parameters for the reference low-pass. Note that geometric symmetry is required for both stopband and passband limits. Generally, an "over-design" will result for one or more of the frequencies.

Fig. 7 shows the tolerance plot of the requirements for a band-stop network, together with the associated equations to calculate the reference low-pass parameters. Again, symmetry is required as for the bandpass case.

In both the bandpass and band-stop equations, a transformation factor, a , is defined and is needed for transforming the reference low-pass into the desired circuit.

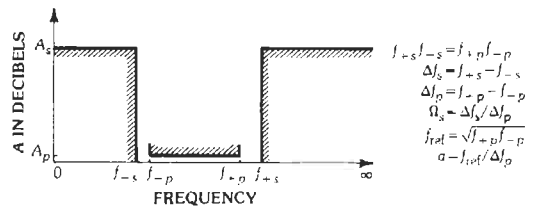


Fig. 6. Requirements for a bandpass filter.

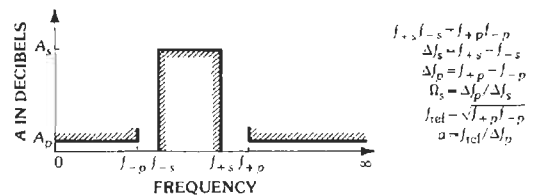


Fig. 7. Requirements for a band-stop network.

CAUER-PARAMETER LOW-PASS FILTERS

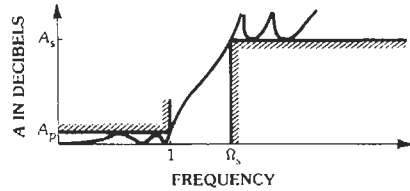
Probably the most important type of low-pass filter is the elliptic-function, or Cauer-parameter, network, which provides equal attenuation maxima in the passband region and equal attenuation minima in the stopband. Fig. 8 shows the attenuation versus frequency performance and the two possible circuit configurations for this type of filter when the degree is odd.

While the attenuation maximum (A_p) in the passband region is one of the parameters necessary to determine the reference low-pass, practical experience is that filter networks are terminated with other transmission networks within a system rather than pure resistances. The transmission quality through these networks connected to filters is adversely affected by excessive variations in the impedance of the filter network. The maximum variation of the input impedance, and consequently the related reflection coefficient (ρ), in the passband of a filter designed with elliptic functions is directly related to the variation in attenuation (A_p) by

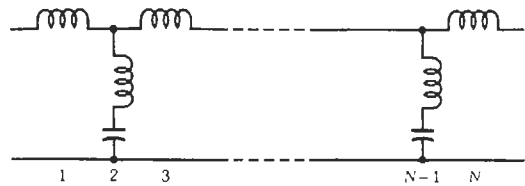
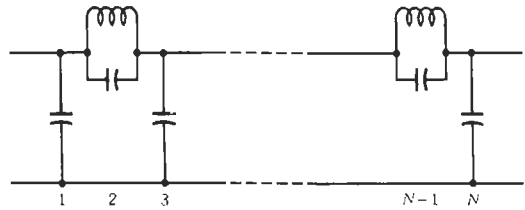
$$A_p = -\ln \sqrt{1 - \rho^2}$$

Normally, the attenuation maximum related to the allowable reflection coefficient is so small that it cannot be measured by practical means and is, in fact, masked by the component losses. Table 1 shows the relationships among A_p , ρ , and return loss in networks with equal passband variations of attenuation.

Tables 2 through 8 include normalized Cauer-parameter low-pass filters of degree 3 through 9 with reflection coefficients of 2, 5, 10, and 25 percent. For each reflection coefficient, ten different cutoff rates, and



(A) Performance.



(B) Configurations.

Fig. 8. Cauer-parameter low-pass filters—odd degree.

TABLE 1. VALUES OF ρ , A_p , AND RETURN LOSS

ρ (%)	A_p (dB)	Return Loss (dB)
0.5	0.000109	46.021
1.0	0.000434	40.000
2.0	0.001738	33.979
3.0	0.003910	30.458
4.0	0.006954	27.959
5.0	0.01087	26.021
6.0	0.01566	24.437
8.0	0.02788	21.938
10	0.04365	20.000
12	0.06299	18.416
15	0.09883	16.478
20	0.17729	13.979
25	0.28029	12.041
30	0.40959	10.458
35	0.56753	9.119
40	0.75721	7.959
45	0.98269	6.936
50	1.24939	6.021
60	1.93820	4.437
70	2.92430	3.098

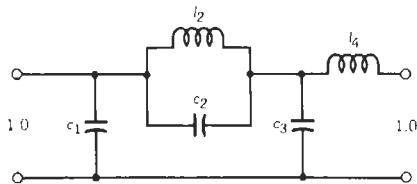
TABLE 2. CAUER-PARAMETER LOW-PASS, DEGREE=3

ρ	θ	Ω_s	A_s (dB)	c_1	c_2	l_2	Ω_2	c_3
5	03	19.1072	74.9	0.6381	0.0021	0.9761	22.059	0.6381
5	05	11.4737	61.6	0.6354	0.0059	0.9711	13.242	0.6354
5	07	8.20550	52.8	0.6314	0.0116	0.9636	9.4661	0.6314
5	09	6.39245	46.3	0.6261	0.0193	0.9536	7.3700	0.6261
5	11	5.24084	41.0	0.6194	0.0291	0.9411	6.0377	0.6194
5	13	4.44541	36.6	0.6113	0.0412	0.9261	5.1166	0.6113
5	16	3.62795	31.1	0.5968	0.0640	0.8991	4.1688	0.5968
5	20	2.92380	25.2	0.5728	0.1043	0.8544	3.3505	0.5728
5	24	2.45859	20.4	0.5438	0.1585	0.8004	2.8079	0.5438
5	29	2.06266	15.4	0.5009	0.2526	0.7205	2.3438	0.5009
10	04	14.3355	73.5	0.8510	0.0033	1.100	16.548	0.8510
10	06	9.56675	62.9	0.8479	0.0075	1.094	11.039	0.8479
10	10	5.75877	49.5	0.8380	0.0211	1.078	6.6370	0.8380
10	12	4.80974	44.8	0.8313	0.0306	1.066	5.5386	0.8313
10	14	4.13356	40.7	0.8233	0.0420	1.052	4.7552	0.8233
10	17	3.42030	35.6	0.8090	0.0630	1.028	3.9277	0.8090
10	21	2.79043	30.0	0.7857	0.0991	0.9885	3.1951	0.7857
10	25	2.36620	25.3	0.7576	0.1458	0.9408	2.7000	0.7576
10	30	2.00000	20.4	0.7160	0.2230	0.8701	2.2701	0.7160
10	36	1.70130	15.4	0.6570	0.3536	0.7701	1.9164	0.6570
25	03	19.1072	89.2	1.344	0.0018	1.140	22.059	1.344
25	06	9.56675	71.1	1.340	0.0072	1.134	11.039	1.340
25	09	6.39245	60.5	1.332	0.0164	1.125	7.3700	1.332
25	13	4.44541	50.9	1.318	0.0345	1.108	5.1166	1.318
25	16	3.62795	45.4	1.305	0.0528	1.090	4.1688	1.305
25	19	3.07155	40.8	1.288	0.0754	1.069	3.5224	1.288
25	23	2.55930	35.7	1.261	0.1128	1.036	2.9256	1.261
25	28	2.13005	30.4	1.221	0.1729	0.9853	2.4230	1.221
25	34	1.78829	25.1	1.163	0.2688	0.9121	2.0198	1.163
25	40	1.55573	20.6	1.094	0.3990	0.8255	1.7423	1.094
ρ	θ	Ω_s	A_s (dB)	l_1	l_2	c_2	Ω_2	l_3

hence ten different stopband minima, are tabulated. In fact, any finer graduation of cutoff sharpness can be computed, and in several publications (Reference 2) the fine graduation is by integer degree of the modular angle, θ ($0^\circ \leq \theta \leq 90^\circ$). For odd degree, $\theta = \arcsin(1/\Omega_s)$.

To determine the degree of the reference low-pass for a Cauer-parameter filter, Fig. 9 (and for sharp cutoff rates, Fig. 10) is used. On the grid of A_p and A_s , an intersection point is found where the two lines meet. From that point, a horizontal line is drawn to the left until the value of Ω_s is reached, and the next curve

TABLE 3. CAUER-PARAMETER LOW-PASS, DEGREE=4



ρ	θ	Ω_s	A_r (dB)	c_1	c_2	l_2	Ω_2	c_3	l_4
2	07	9.60283	75.2	0.4738	0.0091	0.9888	10.546	0.9978	0.4828
2	10	6.73180	62.7	0.4654	0.0187	0.9745	7.3903	0.9929	0.4838
2	14	4.82164	50.9	0.4494	0.0377	0.9476	5.2892	0.9839	0.4857
2	16	4.22622	46.2	0.4393	0.0500	0.9306	4.6338	0.9782	0.4869
2	19	3.56970	40.2	0.4215	0.0726	0.9010	3.9107	0.9685	0.4890
2	22	3.09392	35.0	0.4003	0.1007	0.8661	3.3861	0.9572	0.4914
2	25	2.73382	30.4	0.3756	0.1355	0.8260	2.9885	0.9446	0.4941
2	29	2.37170	25.0	0.3367	0.1953	0.7644	2.5880	0.9257	0.4980
2	33	2.09956	20.3	0.2908	0.2757	0.6938	2.2863	0.9052	0.5021
2	38	1.84278	15.2	0.2223	0.4207	0.5939	2.0006	0.8784	0.5068
5	09	7.47571	74.4	0.6236	0.0125	1.183	8.2082	1.196	0.6360
5	13	5.18854	61.5	0.6112	0.0266	1.162	5.6929	1.188	0.6371
5	18	3.76399	50.1	0.5891	0.0523	1.124	4.1248	1.174	0.6391
5	20	3.39504	46.3	0.5782	0.0654	1.106	3.7182	1.168	0.6401
5	23	2.96326	41.3	0.5595	0.0886	1.074	3.2419	1.156	0.6418
5	27	2.53903	35.6	0.5301	0.1269	1.025	2.7732	1.139	0.6443
5	31	2.22655	30.6	0.4954	0.1754	0.9675	2.4272	1.119	0.6471
5	36	1.93651	25.1	0.4440	0.2552	0.8844	2.1050	1.091	0.6508
5	41	1.72025	20.2	0.3830	0.3650	0.7889	1.8637	1.060	0.6545
5	47	1.52527	15.0	0.2957	0.5609	0.6591	1.6446	1.021	0.6578
10	08	8.40608	84.5	0.7912	0.0088	1.340	9.2310	1.348	0.7999
10	11	6.12348	73.4	0.7838	0.0167	1.327	6.7213	1.344	0.8003
10	15	4.50395	62.5	0.7704	0.0314	1.304	4.9396	1.335	0.8011
10	21	3.23721	50.6	0.7422	0.0634	1.256	3.5441	1.316	0.8027
10	24	2.84368	45.9	0.7245	0.0844	1.226	3.1099	1.305	0.8037
10	28	2.45229	40.3	0.6968	0.1183	1.179	2.6772	1.287	0.8052
10	32	2.16100	35.5	0.6645	0.1604	1.125	2.3545	1.267	0.8068
10	37	1.88831	30.1	0.6171	0.2270	1.047	2.0514	1.239	0.8090
10	42	1.68355	25.3	0.5615	0.3146	0.9568	1.8226	1.206	0.8111
10	48	1.49794	20.1	0.4829	0.4606	0.8337	1.6137	1.164	0.8129
25	10	6.73180	84.9	1.134	0.0123	1.493	7.3903	1.505	1.146
25	14	4.82164	73.1	1.122	0.0243	1.474	5.2892	1.498	1.146
25	20	3.39504	60.6	1.097	0.0505	1.433	3.7182	1.482	1.146
25	26	2.63260	51.2	1.062	0.0877	1.377	2.8767	1.460	1.145
25	30	2.29663	46.0	1.034	0.1197	1.332	2.5049	1.442	1.144
25	35	1.98762	40.3	0.9923	0.1691	1.266	2.1619	1.417	1.143
25	40	1.75892	35.3	0.9437	0.2313	1.189	1.9070	1.387	1.142
25	45	1.58398	30.8	0.8875	0.3102	1.101	1.7108	1.354	1.140
25	51	1.42310	25.8	0.8095	0.4356	0.9820	1.5289	1.309	1.136
25	58	1.28267	20.4	0.7018	0.6497	0.8227	1.3678	1.250	1.129

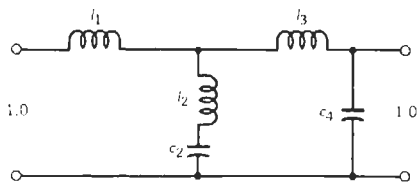


TABLE 4. CAUER-PARAMETER LOW-PASS, DEGREE=5

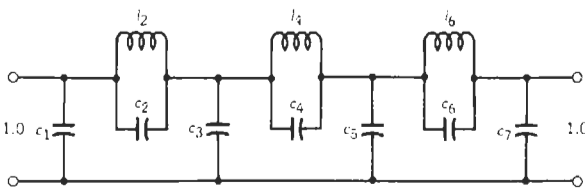
ρ	θ	Ω_3	A_s (dB)	c_1	c_2	l_2	Ω_2	c_3	c_4	l_4	Ω_4	c_5
2	11	5.24084	85.7	0.5764	0.0110	1.158	8.8625	1.349	0.0292	1.128	5.5057	0.5587
2	15	3.86370	72.1	0.5691	0.0206	1.148	6.4997	1.332	0.0557	1.092	4.0559	0.5360
2	19	3.07155	61.7	0.5595	0.0335	1.135	5.1318	1.310	0.0922	1.046	3.2212	0.5059
2	24	2.45859	51.3	0.5441	0.0544	1.114	4.0631	1.277	0.1551	0.9728	2.5743	0.4575
2	27	2.20269	46.0	0.5329	0.0697	1.099	3.6119	1.254	0.2045	0.9212	2.3038	0.4223
2	30	2.00000	41.3	0.5203	0.0874	1.082	3.2508	1.230	0.2652	0.8638	2.0892	0.3821
2	34	1.78829	35.5	0.5010	0.1151	1.056	2.8683	1.197	0.3693	0.7789	1.8646	0.3200
2	38	1.62427	30.3	0.4787	0.1480	1.026	2.5662	1.163	0.5110	0.6851	1.6901	0.2470
2	42	1.49448	25.6	0.4532	0.1875	0.9896	2.3213	1.132	0.7118	0.5836	1.5515	0.1610
2	47	1.36733	20.1	0.4163	0.2492	0.9334	2.0733	1.103	1.111	0.4495	1.4150	0.0304
5	13	4.44541	86.4	0.7556	0.0137	1.295	7.4994	1.558	0.0365	1.259	4.6684	0.7336
5	18	3.23607	72.1	0.7456	0.0266	1.280	5.4168	1.531	0.0716	1.212	3.3946	0.7033
5	23	2.55930	61.2	0.7321	0.0441	1.261	4.2397	1.496	0.1211	1.149	2.6807	0.6628
5	29	2.06266	50.8	0.7112	0.0718	1.231	3.3629	1.444	0.2038	1.056	2.1556	0.6004
5	32	1.88708	46.3	0.6987	0.0887	1.214	3.0476	1.415	0.2574	1.002	1.9695	0.5632
5	36	1.70130	40.8	0.6798	0.1148	1.187	2.7089	1.373	0.3456	0.9214	1.7721	0.5070
5	40	1.55573	35.8	0.6581	0.1456	1.156	2.4377	1.328	0.4590	0.8332	1.6170	0.4427
5	45	1.41421	30.2	0.6265	0.1920	1.110	2.1660	1.271	0.6535	0.7125	1.4654	0.3494
5	49	1.32501	26.0	0.5973	0.2372	1.067	1.9881	1.225	0.8756	0.6091	1.3693	0.2628
5	55	1.22078	20.0	0.5460	0.3242	0.9858	1.7689	1.166	1.419	0.4467	2.2561	0.1073
10	15	3.86370	86.2	0.9590	0.0175	1.352	6.4997	1.760	0.0464	1.310	4.0559	0.9312
10	20	2.92380	73.5	0.9479	0.0315	1.337	4.8753	1.727	0.0843	1.262	3.0653	0.8983
10	27	2.20269	60.1	0.9265	0.0587	1.307	3.6119	1.666	0.1607	1.173	2.3038	0.8363
10	33	1.83608	50.9	0.9024	0.0899	1.273	2.9553	1.601	0.2534	1.076	1.9154	0.7678
10	37	1.66164	45.6	0.8833	0.1155	1.247	2.6356	1.552	0.3337	1.001	1.7299	0.7140
10	41	1.52425	40.7	0.8613	0.1454	1.216	2.3781	1.498	0.4338	0.9195	1.5833	0.6533
10	46	1.39016	35.1	0.8296	0.1902	1.172	2.1187	1.427	0.5976	0.8075	1.4396	0.5670
10	51	1.28676	30.0	0.7924	0.2451	1.119	1.9095	1.352	0.8269	0.6858	1.3279	0.4679
10	56	1.20622	25.1	0.7488	0.3141	1.056	1.7368	1.278	1.168	0.5567	1.2402	0.3534
10	61	1.14335	20.4	0.6973	0.4039	0.9778	1.5913	1.207	1.723	0.4234	1.1708	0.2197
25	13	4.44541	101	1.444	0.0137	1.293	7.4994	2.244	0.0362	1.268	4.6684	1.422
25	18	3.23607	86.3	1.433	0.0266	1.281	5.4168	2.209	0.0703	1.234	3.3946	1.391
25	24	2.45859	73.5	1.416	0.0480	1.261	4.0631	2.151	0.1282	1.177	2.5743	1.341
25	32	1.88708	60.5	1.383	0.0879	1.225	3.0476	2.050	0.2390	1.079	1.9695	1.253
25	40	1.55573	50.1	1.340	0.1430	1.177	2.4377	1.923	0.4007	0.9543	1.6170	1.139
25	44	1.43956	45.5	1.313	0.1775	1.148	2.2154	1.851	0.5083	0.8830	1.4926	1.072
25	49	1.32501	40.2	1.275	0.2288	1.106	1.9881	1.752	0.6786	0.7859	1.3693	0.9789
25	54	1.23607	35.2	1.231	0.2914	1.057	1.8021	1.646	0.9067	0.6808	1.2728	0.8749
25	59	1.16663	30.5	1.179	0.3693	0.9986	1.6467	1.534	1.228	0.5689	1.1966	0.7586
25	64	1.11260	25.9	1.118	0.4695	0.9290	1.5142	1.417	1.714	0.4518	1.1364	0.6280
ρ	θ	Ω_f	A_s (dB)	l_1	l_2	c_2	Ω_2	l_3	l_4	c_4	Ω_4	l_5

The diagram shows a ladder network starting with an inductor l_1 in series. This is followed by a shunt capacitor c_2 , then a series inductor l_2 . Next is a shunt capacitor c_4 , followed by a series inductor l_3 , then a shunt capacitor c_4 , and finally a series inductor l_5 . The input and output ports are both labeled 1.0.

TABLE 5. CAUER-PARAMETER LOW-PASS, DEGREE=6

ρ	θ	Ω_3	A_s (dB)	c_1	c_2	l_2	Ω_2	c_3	c_4	l_4	Ω_4	c_5	l_6
2	18	3.45698	86.1	0.5417	0.0337	1.179	5.0202	1.398	0.0571	1.365	3.5830	1.179	0.5746
2	23	2.72832	73.1	0.5221	0.0562	1.150	3.9348	1.361	0.0959	1.307	2.8252	1.150	0.5762
2	29	2.19194	60.6	0.4916	0.0924	1.106	3.1279	1.307	0.1597	1.219	2.2666	1.105	0.5785
2	35	1.84554	50.2	0.4518	0.1421	1.049	2.5897	1.242	0.2479	1.111	1.9051	1.048	0.5812
2	38	1.71560	45.6	0.4302	0.1707	1.019	2.3975	1.206	0.3043	1.049	1.7692	1.015	0.5828
2	41	1.60614	41.3	0.4050	0.2050	0.9846	2.2258	1.168	0.3714	0.9836	1.6545	0.9795	0.5845
2	45	1.48509	35.8	0.3673	0.2591	0.9341	2.0329	1.115	0.4829	0.8877	1.5274	0.9261	0.5867
2	49	1.38624	30.8	0.3242	0.3252	0.8780	1.8716	1.060	0.6302	0.7833	1.4233	0.8658	0.5890
2	53	1.30476	26.0	0.2749	0.4073	0.8157	1.7349	1.005	0.8333	0.6712	1.3371	0.7976	0.5911
2	58	1.22215	20.3	0.2027	0.5435	0.7271	1.5907	0.9410	1.228	0.5219	1.2491	0.6984	0.5932
5	21	2.97737	85.9	0.6896	0.0408	1.321	4.3068	1.527	0.0706	1.488	3.0843	1.321	0.7294
5	27	2.34339	72.4	0.6643	0.0692	1.283	3.3569	1.477	0.1206	1.411	2.4244	1.283	0.7305
5	33	1.94627	61.4	0.6319	0.1068	1.235	2.7538	1.414	0.1882	1.314	2.0103	1.235	0.7319
5	40	1.64063	50.6	0.5844	0.1651	1.165	2.2802	1.326	0.2969	1.178	1.6907	1.167	0.7337
5	44	1.51304	45.1	0.5520	0.2071	1.118	2.0777	1.270	0.3789	1.089	1.5568	1.121	0.7347
5	48	1.40916	40.0	0.5153	0.2572	1.067	1.9093	1.210	0.4810	0.9922	1.4475	1.070	0.7357
5	52	1.32371	35.1	0.4738	0.3174	1.009	1.7671	1.146	0.6119	0.8874	1.3571	1.013	0.7365
5	56	1.25292	30.5	0.4267	0.3908	0.9453	1.6453	1.081	0.7852	0.7749	1.2820	0.9498	0.7372
5	60	1.19408	26.0	0.3729	0.4825	0.8741	1.5398	1.015	1.026	0.6556	1.2191	0.8794	0.7373
5	65	1.13432	20.5	0.2941	0.6368	0.7722	1.4260	0.9343	1.505	0.4984	1.1546	0.7785	0.7363
10	23	2.72832	87.1	0.8548	0.0458	1.412	3.9348	1.625	0.0793	1.581	2.8252	1.412	0.8994
10	30	2.12408	72.7	0.8231	0.0801	1.365	3.0250	1.558	0.1399	1.483	2.1959	1.366	0.8996
10	37	1.75640	61.1	0.7820	0.1266	1.305	2.4610	1.475	0.2240	1.360	1.8119	1.308	0.8996
10	44	1.51304	51.2	0.7304	0.1883	1.230	2.0777	1.375	0.3401	1.213	1.5568	1.237	0.8993
10	48	1.40916	46.0	0.6957	0.2323	1.181	1.9093	1.312	0.4268	1.118	1.4475	1.190	0.8989
10	53	1.30476	40.0	0.6461	0.2989	1.112	1.7349	1.227	0.5656	0.9890	1.3371	1.125	0.8981
10	57	1.23718	35.4	0.6007	0.3643	1.049	1.6175	1.155	0.7121	0.8773	1.2652	1.068	0.8969
10	61	1.18098	30.9	0.5494	0.4445	0.9794	1.5156	1.080	0.9079	0.7585	1.2050	1.004	0.8951
10	66	1.12397	25.4	0.4747	0.5753	0.8802	1.4053	0.9831	1.273	0.6009	1.1434	0.9150	0.8917
10	70	1.08733	21.0	0.4040	0.7195	0.7880	1.3281	0.9065	1.752	0.4691	1.1032	0.8337	0.8872
25	21	2.97737	100	1.238	0.0365	1.479	4.3068	1.795	0.0597	1.760	3.0843	1.480	1.274
25	27	2.34339	86.7	1.213	0.0613	1.447	3.3569	1.742	0.1011	1.684	2.4244	1.448	1.272
25	35	1.84554	72.5	1.167	0.1073	1.389	2.5897	1.649	0.1775	1.553	1.9051	1.394	1.268
25	43	1.54246	60.7	1.110	0.1681	1.318	2.1248	1.535	0.2854	1.390	1.5877	1.326	1.263
25	51	1.34357	50.6	1.036	0.2515	1.227	1.8005	1.398	0.4401	1.196	1.3782	1.244	1.256
25	55	1.26941	45.9	0.9928	0.3042	1.173	1.6741	1.322	0.5440	1.088	1.2995	1.198	1.251
25	60	1.19408	40.2	0.9313	0.3842	1.098	1.5398	1.219	0.7133	0.9434	1.2191	1.134	1.243
25	64	1.14518	35.8	0.8752	0.4638	1.029	1.4474	1.132	0.8969	0.8195	1.1664	1.078	1.236
25	69	1.09581	30.3	0.7939	0.5924	0.9309	1.3466	1.016	1.233	0.6553	1.1126	1.000	1.224
25	73	1.06444	25.8	0.7171	0.7331	0.8386	1.2753	0.9195	1.665	0.5172	1.0777	0.9305	1.211
ρ	θ	Ω_c	A_s (dB)	l_1	l_2	c_2	Ω_2	l_3	l_4	c_4	Ω_4	l_5	c_6

TABLE 6. CAUER-PARAMETER LOW-PASS, DEGREE=7



ρ	θ	Ω_s	A_s (dB)	c_1	c_2	l_2	Ω_2	c_3	c_4	l_4	Ω_4	c_5	c_6	l_6	Ω_6	c_7
2	26	2.28117	85.2	0.6054	0.0316	1.247	5.0387	1.464	0.1331	1.379	2.3339	1.424	0.1075	1.137	2.8596	0.5358
2	32	1.88708	72.0	0.5919	0.0490	1.228	4.0756	1.405	0.2100	1.281	1.9282	1.348	0.1717	1.063	2.3410	0.4851
2	38	1.62427	60.8	0.5749	0.0713	1.204	3.4121	1.334	0.3130	1.163	1.6571	1.259	0.2597	0.9729	1.9895	0.4215
2	44	1.43956	51.0	0.5537	0.0995	1.174	2.9248	1.252	0.4523	1.029	1.4660	1.161	0.3820	0.8674	1.7371	0.3430
2	48	1.34563	45.0	0.5369	0.1224	1.151	2.6647	1.191	0.5745	0.9295	1.3685	1.092	0.4919	0.7887	1.6056	0.2807
2	51	1.28676	40.7	0.5225	0.1421	1.131	2.4948	1.141	0.6882	0.8504	1.3072	1.040	0.5957	0.7253	1.5213	0.2278
2	55	1.22078	35.3	0.5006	0.1727	1.100	2.2946	1.071	0.8829	0.7387	1.2383	0.9708	0.7758	0.6352	1.4245	0.1472
2	58	1.17918	31.3	0.4818	0.1997	1.073	2.1606	1.016	1.077	0.6505	1.1946	0.9214	0.9567	0.5638	1.3616	0.7749
2	62	1.13257	26.1	0.4527	0.2429	1.030	1.9990	0.9374	1.444	0.5277	1.1455	0.8623	1.299	0.4639	1.2883	-0.319
2	66	1.09464	21.0	0.4178	0.2984	0.9761	1.8530	0.8556	2.045	0.4003	1.1052	0.8182	1.851	0.3598	1.2253	-1.677
5	29	2.06266	86.3	0.7776	0.0364	1.353	4.5080	1.609	0.1603	1.403	2.1090	1.558	0.1227	1.231	2.5729	0.6988
5	36	1.70130	72.4	0.7607	0.0578	1.328	3.6093	1.530	0.2589	1.281	1.7366	1.454	0.2000	1.141	2.0933	0.6379
5	43	1.46628	60.6	0.7387	0.0859	1.296	2.9970	1.435	0.3945	1.136	1.4936	1.332	0.3081	1.031	1.7741	0.5612
5	50	1.30541	50.1	0.7107	0.1226	1.255	2.5494	1.324	0.5859	0.9699	1.3266	1.195	0.4628	0.9015	1.5482	0.4657
5	53	1.25214	45.4	0.6963	0.1417	1.234	2.3913	1.271	0.6935	0.8925	1.2711	1.133	0.5505	0.8396	1.4709	0.4180
5	57	1.19236	40.5	0.6744	0.1712	1.202	2.2039	1.197	0.8734	0.7840	1.2085	1.049	0.6973	0.7512	1.3817	0.3467
5	61	1.14335	35.3	0.6488	0.2068	1.165	2.0378	1.117	1.116	0.6696	1.1569	0.9643	0.8942	0.6561	1.3056	0.2647
5	65	1.10338	30.2	0.6183	0.2505	1.120	1.8883	1.032	1.464	0.5500	1.1145	0.8819	1.173	0.5543	1.2402	0.1690
5	69	1.07114	25.0	0.5813	0.3065	1.063	1.7515	0.9419	2.012	0.4260	1.0800	0.8059	1.599	0.4462	1.1838	0.0540
5	72	1.05146	21.1	0.5478	0.3613	1.010	1.6553	0.8719	2.692	0.3314	1.0587	0.7575	2.103	0.3616	1.1467	-0.0505
10	32	1.88708	86.0	0.9740	0.0435	1.386	4.0756	1.752	0.1983	1.356	1.9282	1.684	0.1459	1.251	2.3410	0.8812
10	39	1.58902	73.1	0.9551	0.0667	1.359	3.3209	1.661	0.3090	1.232	1.6206	1.564	0.2286	1.160	1.9418	0.8165
10	47	1.36733	60.5	0.9270	0.1021	1.319	2.7259	1.536	0.4856	1.064	1.3910	1.402	0.3614	1.034	1.6362	0.7236
10	54	1.23607	50.6	0.8951	0.1431	1.273	2.3422	1.409	0.7085	0.8972	1.2543	1.243	0.5293	0.9020	1.4473	0.6230
10	58	1.17918	45.3	0.8728	0.1725	1.242	2.1606	1.329	0.8827	0.7938	1.1946	1.146	0.6596	0.8177	1.3616	0.5557
10	62	1.13257	40.0	0.8466	0.2076	1.205	1.9990	1.243	1.112	0.6852	1.1455	1.046	0.8291	0.7267	1.2883	0.4795
10	65	1.10338	36.2	0.8238	0.2391	1.173	1.8883	1.174	1.340	0.6006	1.1145	0.9707	0.9946	0.6537	1.2402	0.4153
10	69	1.07114	31.1	0.7879	0.2905	1.122	1.7515	1.079	1.771	0.4840	1.0800	0.8701	1.297	0.5500	1.1838	0.3176
10	73	1.04569	25.8	0.7433	0.3584	1.058	1.6241	0.9763	2.481	0.3639	1.0524	0.7732	1.769	0.4386	1.1352	0.2008
10	77	1.02630	20.4	0.6856	0.4564	0.9706	1.5025	0.8684	3.875	0.2428	1.0310	0.6866	2.618	0.3195	1.0933	0.0540
25	29	2.06266	101	1.456	0.0377	1.306	4.5080	2.204	0.1777	1.265	2.1090	2.137	0.1242	1.216	2.5729	1.375
25	36	1.70130	86.6	1.437	0.0597	1.285	3.6093	2.106	0.2841	1.167	1.7366	2.005	0.1989	1.147	2.0933	1.312
25	44	1.43956	73.2	1.409	0.0932	1.254	2.9248	1.970	0.4504	1.033	1.4660	1.822	0.3158	1.050	1.7371	1.223
25	53	1.25214	60.2	1.367	0.1449	1.207	2.3913	1.785	0.7231	0.8560	1.2711	1.579	0.5055	0.9143	1.4709	1.096
25	60	1.15470	50.9	1.323	0.2000	1.158	2.0776	1.618	1.044	0.7013	1.1689	1.366	0.7233	0.7893	1.3235	0.9741
25	64	1.11260	45.7	1.292	0.2402	1.124	1.9243	1.513	1.303	0.6071	1.1243	1.237	0.8942	0.7093	1.2556	0.8937
25	68	1.07853	40.6	1.256	0.2897	1.084	1.7847	1.401	1.660	0.5090	1.0880	1.102	1.121	0.6225	1.1972	0.8031
25	72	1.05146	35.4	1.210	0.3530	1.034	1.6553	1.280	2.190	0.4074	1.0587	0.9634	1.440	0.5280	1.1467	0.6996
25	76	1.03061	30.0	1.153	0.4387	0.9705	1.5326	1.149	3.075	0.3031	1.0358	0.8226	1.934	0.4248	1.1032	0.5779
25	79	1.01872	25.7	1.098	0.5291	0.9081	1.4426	1.043	4.276	0.2237	1.0225	0.7168	2.543	0.3404	1.0747	0.4689
ρ	θ	Ω_s	A_s (dB)	l_1	l_2	c_2	Ω_2	l_3	l_4	c_4	Ω_4	l_5	l_6	c_6	Ω_6	l_7

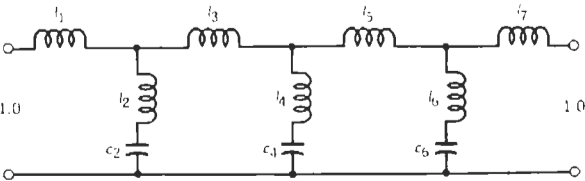
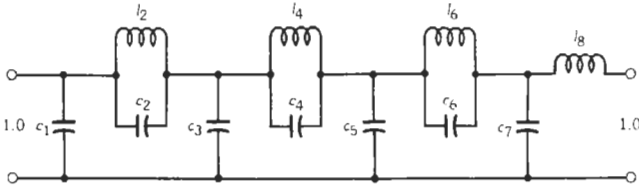


TABLE 7. CAUER-PARAMETER LOW-PASS, DEGREE=8



ρ	θ	Ω_s	A_s (dB)	c_1	c_2	l_2	Ω_2	c_3	c_4	l_4	Ω_4	c_5	c_6	l_6	Ω_6	c_7	l_8
2	33	1.89679	86.6	0.5447	0.0724	1.205	3.3870	1.397	0.1969	1.364	1.9298	1.416	0.1458	1.358	2.2472	1.183	0.6139
2	40	1.60249	72.2	0.5116	0.1113	1.157	2.7865	1.310	0.3075	1.226	1.6281	1.303	0.2270	1.251	1.8769	1.125	0.6155
2	46	1.42779	61.4	0.4761	0.1549	1.107	2.4152	1.223	0.4367	1.091	1.4486	1.193	0.3208	1.141	1.6528	1.065	0.6170
2	53	1.28108	50.0	0.4245	0.2221	1.036	2.0853	1.106	0.6510	0.9127	1.2973	1.050	0.4728	0.9928	1.4595	0.9798	0.6186
2	56	1.23189	45.3	0.3982	0.2582	1.000	1.9680	1.051	0.7755	0.8300	1.2464	0.9845	0.5590	0.9222	1.3928	0.9377	0.6193
2	59	1.18924	40.8	0.3688	0.3001	0.9613	1.8617	0.9938	0.9303	0.7438	1.2022	0.9179	0.6634	0.8474	1.3337	0.8918	0.6198
2	62	1.15228	36.3	0.3357	0.3496	0.9184	1.7648	0.9335	1.129	0.6542	1.1637	0.8505	0.7934	0.7679	1.2812	0.8409	0.6203
2	66	1.11071	30.5	0.2845	0.4316	0.8536	1.6475	0.8485	1.504	0.5299	1.1202	0.7613	1.027	0.6545	1.2199	0.7641	0.6205
2	69	1.08464	26.1	0.2392	0.5107	0.7976	1.5668	0.7815	1.931	0.4337	1.0927	0.6969	1.275	0.5635	1.1796	0.6975	0.6202
2	73	1.05591	20.1	0.1658	0.6551	0.7095	1.4668	0.6880	2.921	0.3034	1.0622	0.6198	1.797	0.4338	1.1325	0.5919	0.6188
5	37	1.71384	86.1	0.6891	0.0834	1.318	3.0164	1.479	0.2376	1.386	1.7423	1.446	0.1711	1.435	2.0178	1.294	0.7686
5	44	1.48002	72.9	0.6538	0.1236	1.266	2.5280	1.381	0.3577	1.239	1.5023	1.322	0.2566	1.317	1.7203	1.234	0.7689
5	51	1.31802	61.1	0.6094	0.1764	1.203	2.1707	1.267	0.5242	1.070	1.3355	1.181	0.3732	1.177	1.5088	1.161	0.7688
5	58	1.20279	50.3	0.5534	0.2474	1.125	1.8960	1.137	0.7680	0.8803	1.2162	1.023	0.5387	1.015	1.3526	1.074	0.7683
5	61	1.16402	45.8	0.5249	0.2856	1.085	1.7962	1.077	0.9115	0.7935	1.1759	0.9512	0.6330	0.9376	1.2980	1.031	0.7676
5	64	1.13047	41.4	0.4928	0.3304	1.042	1.7046	1.013	1.092	0.7034	1.1409	0.8777	0.7484	0.8559	1.2494	0.9843	0.7669
5	68	1.09287	35.5	0.4433	0.4039	0.9755	1.5931	0.9231	1.425	0.5785	1.1014	0.7777	0.9510	0.7393	1.1926	0.9146	0.7650
5	71	1.06945	31.1	0.3995	0.4739	0.9182	1.5159	0.8517	1.791	0.4817	1.0766	0.7025	1.162	0.6450	1.1552	0.8547	0.7635
5	75	1.04394	25.0	0.3287	0.5999	0.8274	1.4194	0.7512	2.600	0.3493	1.0493	0.6046	1.588	0.5096	1.1115	0.7616	0.7599
5	78	1.02877	20.2	0.2615	0.7379	0.7432	1.3503	0.6718	3.764	0.2490	1.0329	0.5369	2.135	0.3994	1.0828	0.6767	0.7554
10	40	1.60249	86.2	0.8520	0.0934	1.380	2.7865	1.542	0.2723	1.386	1.6281	1.451	0.1901	1.493	1.8769	1.355	0.9406
10	48	1.38060	72.0	0.8075	0.1420	1.318	2.3115	1.418	0.4227	1.207	1.4000	1.299	0.2936	1.345	1.5913	1.285	0.9388
10	55	1.24752	60.9	0.7582	0.1987	1.251	2.0058	1.292	0.6099	1.029	1.2626	1.146	0.4196	1.192	1.4141	1.211	0.9363
10	62	1.15228	50.4	0.6959	0.2751	1.167	1.7648	1.149	0.8888	0.8309	1.1637	0.9771	0.6004	1.015	1.2812	1.123	0.9324
10	65	1.12034	46.0	0.6639	0.3167	1.125	1.6757	1.082	1.057	0.7408	1.1303	0.9002	0.7050	0.9309	1.2344	1.080	0.9301
10	69	1.08464	40.1	0.6146	0.3844	1.060	1.5668	0.9874	1.360	0.6159	1.0927	0.7942	0.8860	0.8111	1.1796	1.017	0.9261
10	72	1.06248	35.6	0.5714	0.4483	1.003	1.4911	0.9121	1.686	0.5189	1.0692	0.7127	1.070	0.7143	1.1436	0.9628	0.9225
10	75	1.04394	31.0	0.5206	0.5292	0.9379	1.4194	0.8326	2.167	0.4192	1.0493	0.6303	1.325	0.6107	1.1115	0.9021	0.9178
10	78	1.02877	26.2	0.4591	0.6380	0.8597	1.3503	0.7483	2.953	0.3174	1.0329	0.5482	1.708	0.4992	1.0829	0.8320	0.9113
10	81	1.01681	21.0	0.3813	0.7981	0.7617	1.2825	0.6588	4.480	0.2147	1.0197	0.4691	2.367	0.3778	1.0575	0.7473	0.9023
25	37	1.71384	100	1.255	0.0776	1.417	3.0164	1.756	0.2208	1.492	1.7423	1.545	0.1432	1.715	2.0178	1.398	1.330
25	45	1.45323	85.3	1.211	0.1197	1.369	2.4704	1.631	0.3461	1.329	1.4748	1.405	0.2243	1.569	1.6858	1.343	1.324
25	53	1.28108	72.2	1.155	0.1757	1.309	2.0853	1.484	0.5223	1.138	1.2973	1.239	0.3369	1.393	1.4595	1.277	1.315
25	61	1.16402	60.0	1.082	0.2519	1.230	1.7962	1.313	0.7850	0.9213	1.1759	1.049	0.5002	1.187	1.2980	1.197	1.303
25	67	1.10156	51.2	1.013	0.3297	1.155	1.6200	1.168	1.090	0.7436	1.1105	0.8922	0.6813	1.009	1.2060	1.126	1.291
25	71	1.06945	45.3	0.9576	0.3976	1.094	1.5159	1.062	1.396	0.6181	1.0766	0.7803	0.8529	0.8785	1.1552	1.072	1.281
25	74	1.04974	40.8	0.9084	0.4616	1.041	1.4429	0.9777	1.724	0.5204	1.0556	0.6929	1.028	0.7730	1.1218	1.026	1.271
25	77	1.03346	36.0	0.8507	0.5430	0.9768	1.3732	0.8874	2.211	0.4198	1.0380	0.6024	1.271	0.6597	1.0920	0.9747	1.259
25	80	1.02044	31.0	0.7801	0.6534	0.8985	1.3051	0.7901	3.017	0.3162	1.0238	0.5091	1.641	0.5368	1.0656	0.9158	1.244
25	83	1.01057	25.4	0.6892	0.8209	0.7963	1.2369	0.6835	4.639	0.2102	1.0127	0.4133	2.297	0.4007	1.0423	0.8443	1.223
ρ	θ	Ω_s	A_s (dB)	l_1	l_2	c_2	Ω_2	l_3	l_4	c_4	Ω_4	l_5	l_6	c_6	Ω_6	l_7	c_8

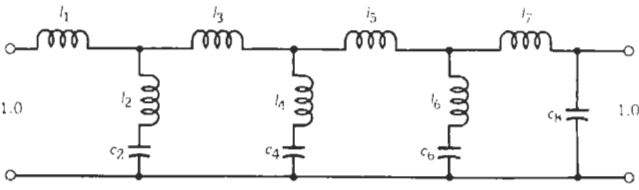


TABLE 8. CAUER-PARAMETER LOW-PASS, DEGREE=9

ρ	θ	Ω_c	A_r (dB)	c_1	c_2	l_2	Ω_2	c_3	c_4	l_4	Ω_4	c_5	c_6	l_6	Ω_6	c_7	c_8	l_8	Ω_8	c_9
2	34	1.78829	101	0.6236	0.0332	1.289	4.8248	1.498	0.1723	1.423	2.0198	1.479	0.2254	1.353	1.8111	1.413	0.1236	1.162	2.6394	0.5415
2	40	1.55573	87.0	0.6123	0.0476	1.273	4.0634	1.438	0.2489	1.324	1.7423	1.368	0.3280	1.231	1.5740	1.325	0.1805	1.096	2.2480	0.4964
2	47	1.36733	72.9	0.5955	0.0689	1.249	3.4103	1.355	0.3674	1.188	1.5135	1.219	0.4907	1.068	1.3815	1.207	0.2707	1.002	1.9200	0.4306
2	54	1.23607	60.2	0.5736	0.0968	1.217	2.9140	1.257	0.5320	1.032	1.3498	1.053	0.7263	0.8855	1.2469	1.075	0.3999	0.8896	1.6766	0.3469
2	60	1.15470	50.0	0.5493	0.1282	1.182	2.5690	1.159	0.7331	0.8806	1.2446	0.8987	1.032	0.7158	1.1632	0.9530	0.5632	0.7753	1.5133	0.2558
2	63	1.12233	45.0	0.5345	0.1476	1.160	2.4161	1.105	0.8673	0.7990	1.2013	0.8192	1.250	0.6269	1.1297	0.8902	0.6752	0.7111	1.4432	0.2011
2	66	1.09464	40.0	0.5172	0.1704	1.135	2.2732	1.046	1.037	0.7131	1.1631	0.7389	1.540	0.5358	1.1009	0.8270	0.8193	0.6415	1.3793	0.1381
2	69	1.07114	35.1	0.4970	0.1978	1.106	2.1381	0.9828	1.259	0.6225	1.1296	0.6592	1.948	0.4430	1.0764	0.7643	1.013	0.5661	1.3208	0.0642
2	72	1.05146	30.1	0.4726	0.2315	1.070	2.0089	0.9141	1.569	0.5266	1.1002	0.5817	2.570	0.3490	1.0558	0.7039	1.287	0.4841	1.2669	-0.0247
2	74	1.04030	26.6	0.4534	0.2592	1.041	1.9249	0.8650	1.857	0.4594	1.0827	0.5326	3.204	0.2864	1.0440	0.6662	1.546	0.4253	1.2332	-0.0958
5	37	1.66164	102	0.7940	0.0370	1.385	4.4144	1.625	0.2023	1.415	1.8692	1.548	0.2644	1.337	1.6820	1.521	0.1360	1.247	2.4282	0.7041
5	44	1.43956	86.7	0.7798	0.0545	1.364	3.6664	1.546	0.3014	1.293	1.6018	1.402	0.3972	1.189	1.4554	1.404	0.2042	1.167	2.0480	0.6501
5	51	1.28676	73.5	0.7614	0.0773	1.336	3.1119	1.453	0.4347	1.151	1.4137	1.236	0.5807	1.021	1.2990	1.270	0.2973	1.070	1.7727	0.5826
5	58	1.17918	61.3	0.7374	0.1072	1.301	2.6775	1.343	0.6207	0.9884	1.2767	1.051	0.8480	0.8349	1.1884	1.120	0.4290	0.9531	1.5640	0.4977
5	64	1.11260	51.3	0.7103	0.1416	1.260	2.3675	1.234	0.8514	0.8322	1.1880	0.8812	1.201	0.6641	1.1196	0.9810	0.5941	0.8334	1.4212	0.4052
5	67	1.08636	46.3	0.6935	0.1631	1.235	2.2274	1.173	1.008	0.7480	1.1514	0.7931	1.457	0.5754	1.0923	0.9082	0.7071	0.7654	1.3592	0.3513
5	70	1.06418	41.4	0.6737	0.1890	1.206	2.0945	1.107	1.210	0.6594	1.1193	0.7037	1.804	0.4849	1.0691	0.8339	0.8531	0.6911	1.3023	0.2883
5	73	1.04567	36.3	0.6498	0.2208	1.171	1.9667	1.036	1.484	0.5658	1.0912	0.6141	2.309	0.3930	1.0497	0.7586	1.051	0.6092	1.2498	0.2142
5	76	1.03061	31.1	0.6202	0.2616	1.127	1.8420	0.9578	1.884	0.4665	1.0668	0.5256	3.114	0.3005	1.0337	0.6835	1.338	0.5181	1.2011	0.1245
5	79	1.01872	25.6	0.5821	0.3172	1.069	1.7175	0.8713	2.538	0.3601	1.0459	0.4412	4.597	0.2087	1.0209	0.6109	1.801	0.4159	1.1555	0.0102
10	40	1.55573	101	0.9894	0.0430	1.410	4.0634	1.754	0.2443	1.348	1.7423	1.618	0.3195	1.264	1.5740	1.626	0.1567	1.262	2.2489	0.8867
10	48	1.34563	85.0	0.9710	0.0651	1.383	3.3315	1.650	0.3753	1.205	1.4868	1.429	0.4957	1.092	1.3593	1.473	0.2419	1.168	1.8808	0.8199
10	55	1.22078	72.5	0.9498	0.0909	1.353	2.8522	1.543	0.5334	1.059	1.3303	1.241	0.7152	0.9223	1.2312	1.319	0.3452	1.068	1.6470	0.7465
10	62	1.13257	60.7	0.9218	0.1253	1.313	2.4659	1.417	0.7579	0.8936	1.2152	1.034	1.043	0.7376	1.1403	1.147	0.4917	0.9466	1.4658	0.6547
10	68	1.07853	50.7	0.8894	0.1657	1.267	2.1824	1.292	1.046	0.7351	1.1403	0.8439	1.494	0.5694	1.0841	0.9870	0.6783	0.8214	1.3397	0.5551
10	71	1.05762	45.8	0.8689	0.1920	1.238	2.0514	1.221	1.251	0.6496	1.1095	0.7454	1.837	0.4825	1.0622	0.9023	0.8091	0.7492	1.2844	0.4953
10	74	1.04030	40.6	0.8440	0.2244	1.203	1.9249	1.145	1.525	0.5592	1.0827	0.6454	2.328	0.3941	1.0440	0.8140	0.9838	0.6683	1.2332	0.4267
10	77	1.02630	35.3	0.8130	0.2661	1.159	1.8006	1.060	1.922	0.4634	1.0595	0.5447	3.093	0.3052	1.0291	0.7254	1.228	0.5793	1.1856	0.3424
10	79	1.01872	31.6	0.7873	0.3019	1.123	1.7175	0.9987	2.309	0.3958	1.0459	0.4781	3.903	0.2458	1.0209	0.6645	1.462	0.5124	1.1555	0.2761
10	82	1.00983	25.6	0.7369	0.3771	1.050	1.5893	0.8954	3.290	0.2874	1.0284	0.3807	6.205	0.1576	1.0113	0.5730	2.029	0.3983	1.1124	0.1503
25	44	1.43956	101	1.452	0.0572	1.301	3.6664	2.091	0.3414	1.141	1.6018	1.849	0.4468	1.057	1.4554	1.910	0.2068	1.153	2.0480	1.318
25	52	1.26902	85.9	1.428	0.0847	1.274	3.0437	1.958	0.5122	1.008	1.3914	1.610	0.6778	0.8994	1.2808	1.713	0.3101	1.066	1.7395	1.238
25	60	1.15470	72.2	1.397	0.1223	1.238	2.5690	1.800	0.7565	0.8534	1.2446	1.338	1.022	0.7235	1.1632	1.487	0.4562	0.9572	1.5133	1.137
25	67	1.08636	60.6	1.358	0.1686	1.196	2.2274	1.635	1.079	0.6991	1.1514	1.077	1.505	0.5570	1.0923	1.264	0.6444	0.8400	1.3592	1.024
25	73	1.04569	50.6	1.312	0.2258	1.145	1.9667	1.470	1.523	0.5514	1.0912	0.8373	2.230	0.4069	1.0497	1.052	0.8937	0.7163	1.2498	0.9005
25	76	1.03061	45.3	1.281	0.2653	1.111	1.8420	1.375	1.865	0.4712	1.0668	0.7127	2.836	0.3299	1.0337	0.9388	1.078	0.6431	1.2011	0.8245
25	78	1.02234	41.7	1.256	0.2982	1.084	1.7592	1.306	2.176	0.4149	1.0525	0.6282	3.423	0.2781	1.0248	0.8596	1.240	0.5887	1.1704	0.7652
25	81	1.01247	35.8	1.209	0.3637	1.031	1.6328	1.190	2.878	0.3251	1.0339	0.4995	4.859	0.2001	1.0142	0.7350	1.590	0.4956	1.1265	0.6615
25	83	1.00751	31.5	1.166	0.4258	0.9844	1.5447	1.102	3.665	0.2606	1.0233	0.4128	6.637	0.1481	1.0087	0.6467	1.959	0.4230	1.0985	0.5748
25	85	1.00382	26.7	1.108	0.5174	0.9194	1.4499	0.9994	5.088	0.1911	1.0143	0.3257	10.24	0.0968	1.0045	0.5537	2.578	0.3381	1.0711	0.4640
ρ	θ	Ω_c	A_r (dB)	l_1	l_2	c_2	Ω_2	l_3	l_4	c_4	Ω_4	l_5	l_6	c_6	Ω_6	l_7	l_8	c_8	Ω_8	l_9

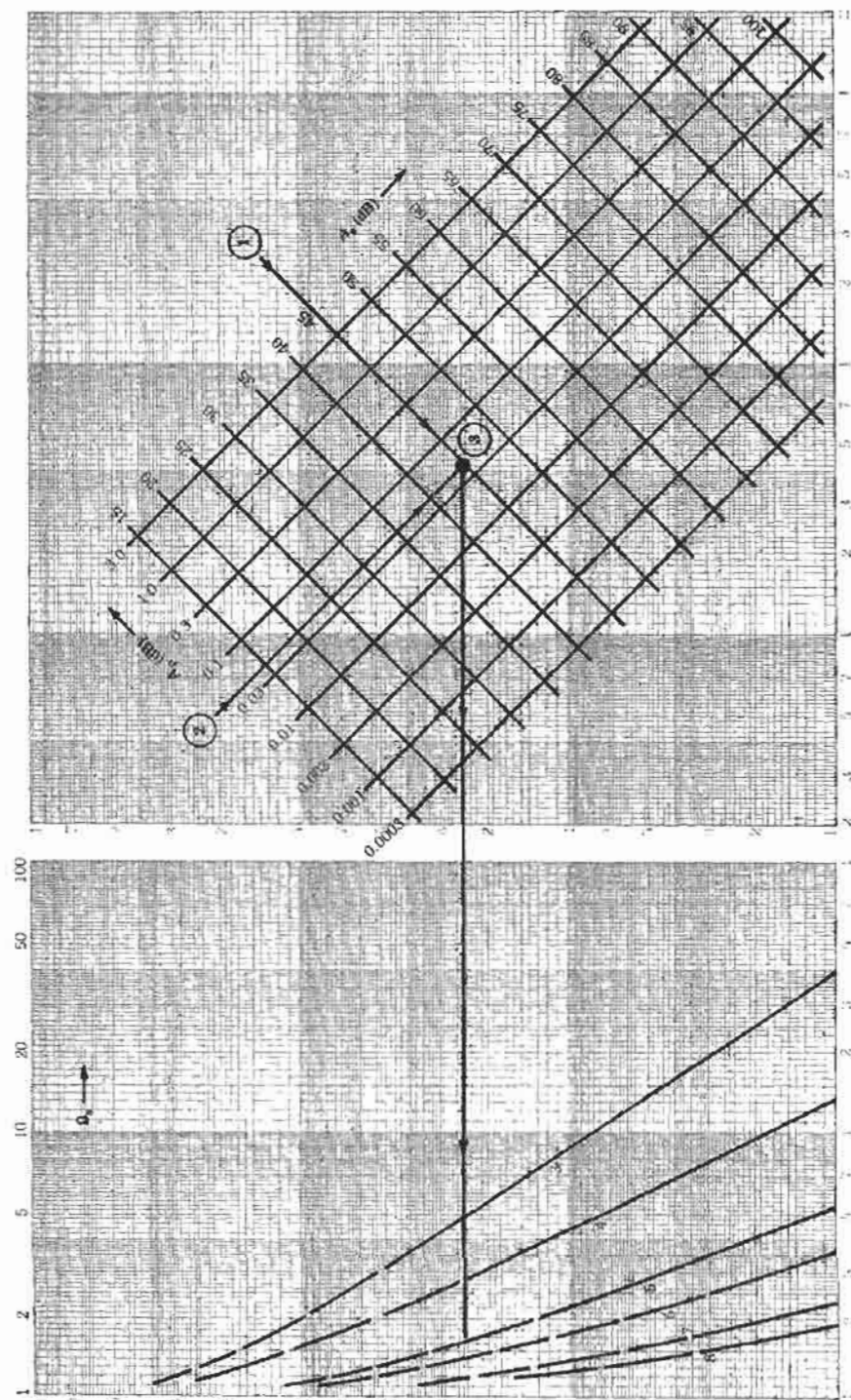


Fig. 9. Selecting the degree for a Cauer-parameter low-pass. (Example 1 is drawn in to illustrate use of the chart.)

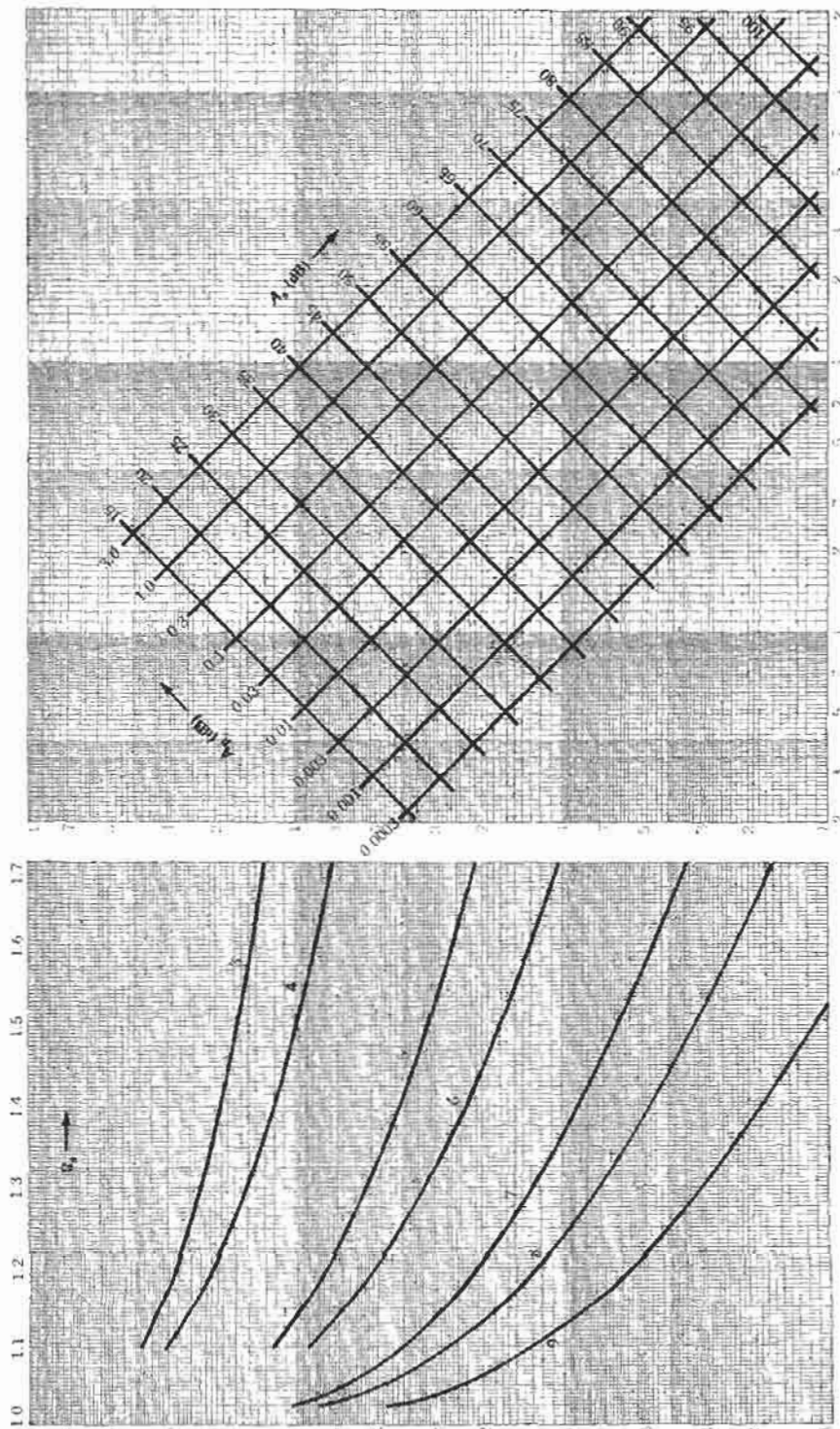


Fig. 10. Selecting the degree for a Cauer-parameter low-pass—expanded.

beyond that intersection specifies the necessary degree. Example 1 in the section headed "Numerical Examples" demonstrates this application.

Chebyshev Low-Pass Filters

If the passband requirement is for a variation of attenuation between zero and a maximum and the stop requirement is only to increase monotonically to infinity, a Chebyshev (sometimes spelled "Tchebyscheff") filter can be used. Fig. 11 shows this performance characteristic and the circuit configurations for both odd- and even-degree filters.

The discussion of impedance variation in the preceding section also applies for this network type.

Table 9 contains data for a number of normalized Chebyshev low-pass filters for degree 2 through 9. For each degree, one network is tabulated for reflection coefficients of 2, 5, 10, and 25 percent. Since stopband attenuation increases monotonically, there is only one circuit for a given A_p .

Figs. 12 and 13 are used in the same way as described for Cauer-parameter filters to ascertain the order of low-pass required.

Butterworth Low-Pass Filters

Whenever a filter performance is desired with monotonically increasing attenuation from frequency equal to zero to infinity, Butterworth filters may be applicable. Fig. 14 shows the performance characteristic and circuit configurations for Butterworth low-pass filters of odd and even degree.

Table 10 contains data for a number of normalized Butterworth low-pass filters of degree 2 through 9. For each degree, circuit elements are tabulated for attenuation maxima at f_{ref} of 0.011, 0.044, 0.28, and 3.0 dB.

Figs. 15 and 16 are used in the same way as described for Cauer-parameter filters to determine the degree of the low-pass to satisfy the requirements.

Low-Pass to High-Pass Transformation

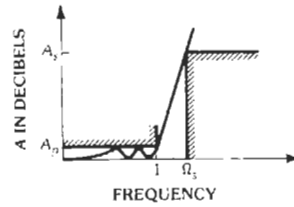
The frequency performance of a low-pass filter can be transformed into high-pass performance by the frequency transformation

$$\Omega_{HP} = 1/\Omega_{LP}$$

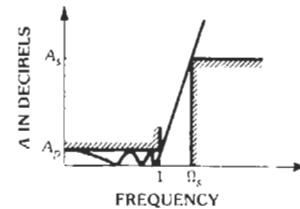
The transformation of normalized circuit-element values is shown in Fig. 17.

Low-Pass to Bandpass Transformation

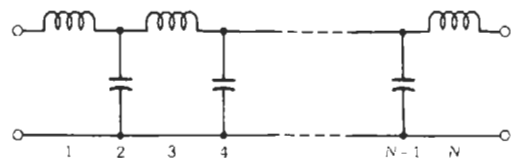
By the transformation



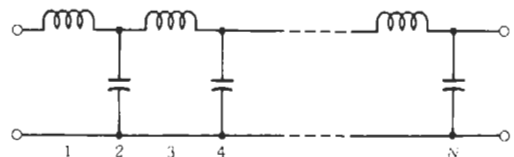
(A) Odd-degree performance.



(B) Even-degree performance.



(C) Odd-degree configuration.



(D) Even-degree configuration.

Fig. 11. Chebyshev low-pass filters.

$$\Omega = a(\eta - 1/\eta)$$

the frequency characteristic of a low-pass filter is converted into that of a bandpass filter. Two normalized bandpass frequencies, $\eta_p = f_{+p}/f_{ref}$ and $\eta_{-p} = 1/\eta_p = f_{-p}/f_{ref}$, correspond to the normalized low-pass Ω_p . Similarly, the geometrically symmetrical stopband limits, $\eta_s = f_s/f_{ref}$ and $\eta_{-s} = f_{-s}/f_{ref}$ are defined as shown in Fig. 6, and:

$$\Omega_p = a(\eta_{-p} + \eta_{+p}) = 1$$

$$\Omega_s = a(\eta_{-s} + \eta_s)$$

The relation between stopband limits and passband limits is

$$\Omega_s = (\eta_s - \eta_{-s})/(\eta_p - \eta_{-p})$$

TABLE 9. CHEBYSHEV LOW-PASS, DEGREE 2-9

N	ρ	c_1	l_2	c_3	l_4	c_5	l_6	c_7	l_8	c_9	r_l
2	2	0.2858	0.2746								0.9608
	5	0.4588	0.4151								0.9048
	10	0.6667	0.5455								0.8182
	25	1.155	0.6928								0.6000
3	2	0.4520	0.7838	0.4520							1.000
	5	0.6396	0.9789	0.6396							1.000
	10	0.8534	1.104	0.8534							1.000
	25	1.345	1.141	1.345							1.000
4	2	0.5379	1.041	1.084	0.5168						0.9608
	5	0.7231	1.207	1.334	0.6543						0.9048
	10	0.9332	1.292	1.580	0.7636						0.8182
	25	1.420	1.260	2.100	0.8519						.6000
5	2	0.5849	1.169	1.369	1.169	0.5849					1.000
	5	0.7664	1.310	1.588	1.310	0.7664					1.000
	10	0.9732	1.372	1.803	1.372	0.9732					1.000
	25	1.456	1.307	2.283	1.307	1.456					1.000
6	2	0.6126	1.240	1.505	1.446	1.290	0.5885				0.9608
	5	0.7912	1.365	1.700	1.538	1.508	0.7159				0.9048
	10	0.9958	1.413	1.895	1.550	1.727	0.8147				0.8182
	25	1.476	1.330	2.353	1.412	2.216	0.8855				0.6000
7	2	0.6301	1.282	1.579	1.575	1.579	1.282	0.6301			1.000
	5	0.8068	1.400	1.757	1.634	1.757	1.400	0.8068			1.000
	10	1.010	1.437	1.941	1.622	1.941	1.437	1.010			1.000
	25	1.488	1.343	2.388	1.451	2.388	1.343	1.488			1.000
8	2	0.6419	1.309	1.623	1.643	1.710	1.560	1.363	0.6167		0.9608
	5	0.8171	1.417	1.791	1.684	1.861	1.621	1.566	0.7392		0.9048
	10	1.019	1.452	1.968	1.657	2.025	1.610	1.774	0.8336		0.8182
	25	1.496	1.351	2.407	1.470	2.450	1.444	2.252	0.8977		0.6000
9	2	0.6501	1.328	1.652	1.684	1.779	1.684	1.652	1.328	0.6501	1.000
	5	0.8242	1.431	1.813	1.712	1.913	1.712	1.813	1.431	0.8242	1.000
	10	1.025	1.462	1.985	1.677	2.066	1.677	1.985	1.462	1.025	1.000
	25	1.502	1.357	2.420	1.481	2.480	1.481	2.420	1.357	1.502	1.000
N	ρ	l_1	c_2	l_3	c_4	l_5	c_6	l_7	c_8	l_9	g_l

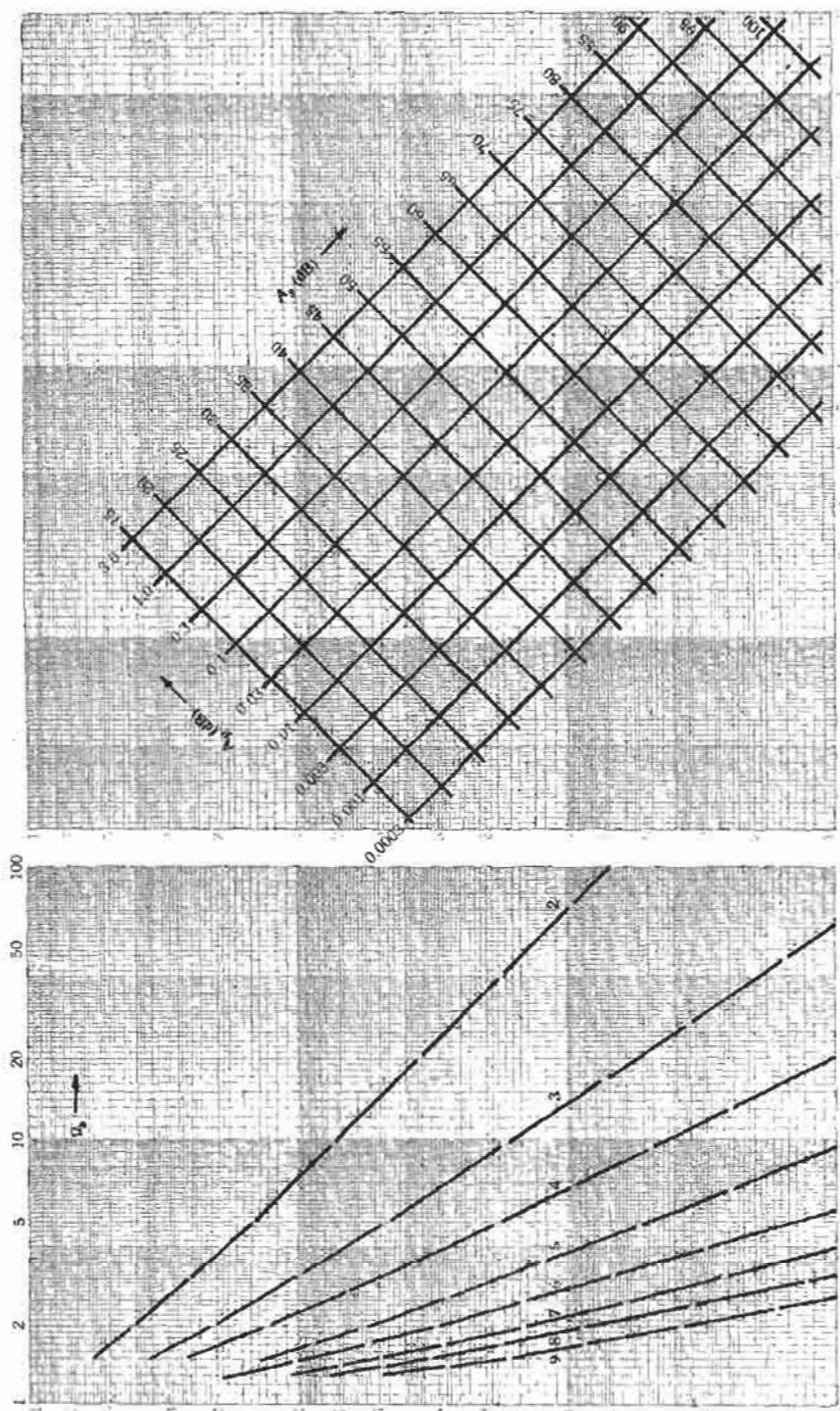


Fig. 12. Selecting the degree for a Chebyshev low-pass.

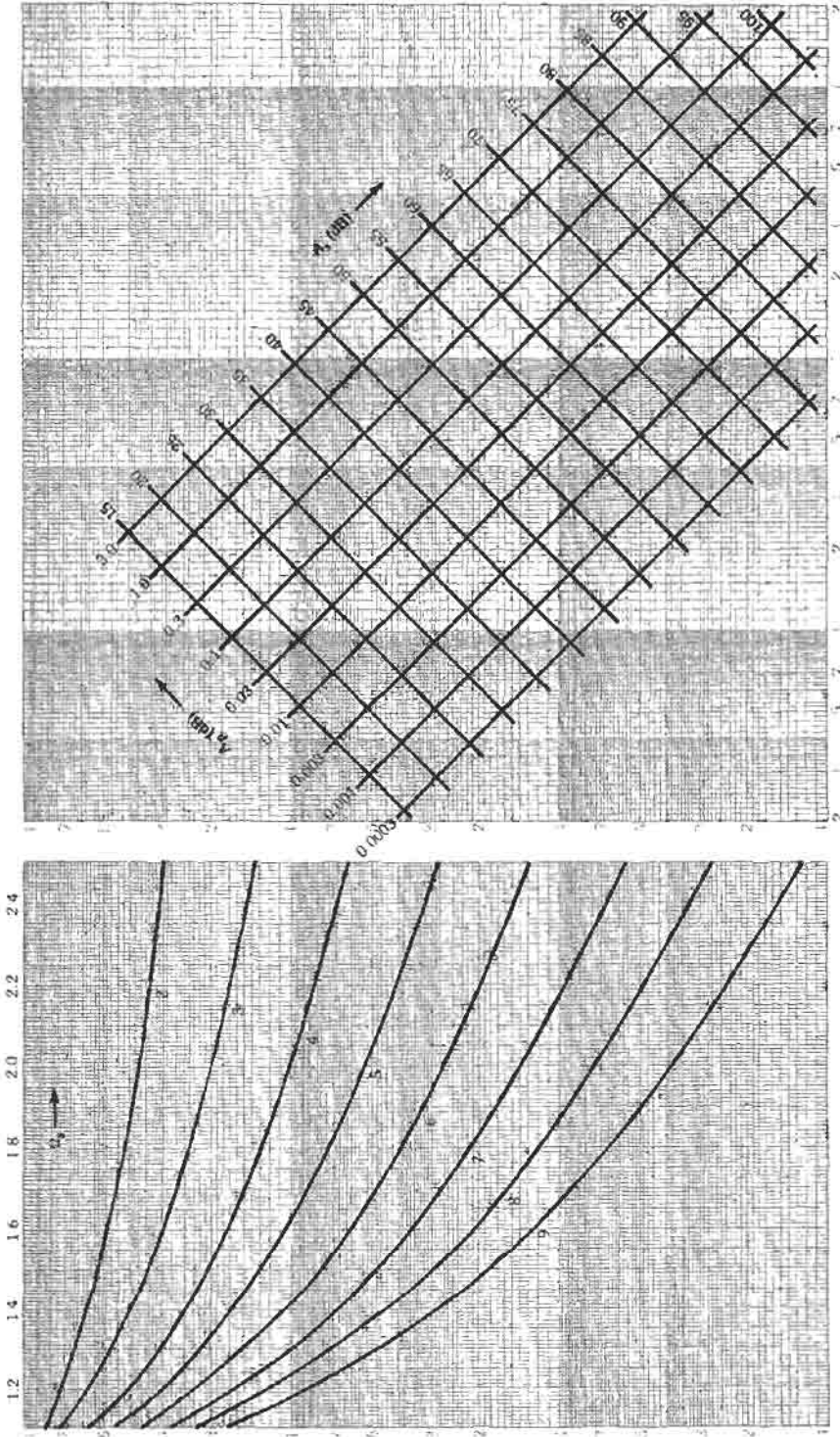
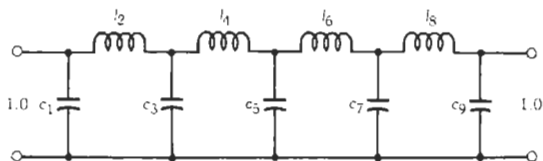
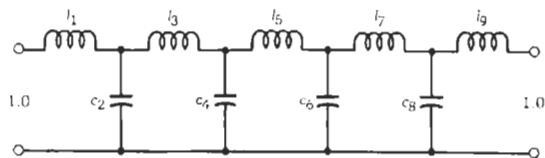


Fig. 13. Selecting the degree for a Chebyshev low-pass—expanded.

TABLE 10. BUTTERWORTH LOW-PASS, DEGREE 2-9



N	A_p (dB)	c_1	l_2	c_3	l_4	c_5	l_6	c_7	l_8	c_9
2	0.01087	0.3164	0.3164							
	0.04365	0.4483	0.4483							
	0.28029	0.7186	0.7186							
	3.0000	1.413	1.413							
3	0.01087	0.3686	0.7371	0.3686						
	0.04365	0.4649	0.9299	0.4649						
	0.28029	0.6368	1.274	0.6368						
	3.0000	1.000	2.000	1.000						
4	0.01087	0.3620	0.8740	0.8740	0.3620					
	0.04365	0.4309	1.040	1.040	0.4309					
	0.28029	0.5456	1.317	1.317	0.5456					
	3.0000	0.7649	1.847	1.847	0.7649					
5	0.01087	0.3396	0.8890	1.099	0.8890	0.3396				
	0.04365	0.3903	1.022	1.263	1.022	0.3903				
	0.28029	0.4714	1.234	1.526	1.234	0.4714				
	3.0000	0.6177	1.617	2.000	1.617	0.6177				
6	0.01087	0.3143	0.8586	1.173	1.173	0.8586	0.3143			
	0.04365	0.3530	0.9643	1.317	1.317	0.9643	0.3530			
	0.28029	0.4131	1.129	1.542	1.542	1.129	0.4131			
	3.0000	0.5174	1.414	1.931	1.931	1.414	0.5174			
7	0.01087	0.2901	0.8130	1.175	1.304	1.175	0.8130	0.2901		
	0.04365	0.3205	0.8981	1.298	1.440	1.298	0.8981	0.3205		
	0.28029	0.3668	1.028	1.485	1.648	1.485	1.028	0.3668		
	3.0000	0.4449	1.247	1.801	2.000	1.801	1.247	0.4449		
8	0.01087	0.2684	0.7642	1.144	1.349	1.349	1.144	0.7642	0.2684	
	0.04365	0.2928	0.8338	1.248	1.472	1.472	1.248	0.8338	0.2928	
	0.28029	0.3294	0.9381	1.404	1.656	1.656	1.404	0.9381	0.3294	
	3.0000	0.3901	1.111	1.662	1.961	1.961	1.662	1.111	0.3901	
9	0.01087	0.2490	0.7170	1.098	1.347	1.434	1.347	1.098	0.7170	0.2490
	0.04365	0.2690	0.7747	1.187	1.456	1.549	1.456	1.187	0.7747	0.2690
	0.28029	0.2988	0.8603	1.318	1.617	1.721	1.617	1.318	0.8603	0.2988
	3.0000	0.3472	1.000	1.532	1.879	2.000	1.879	1.532	1.000	0.3472
	A_p (dB)	l_1	c_2	l_3	c_4	l_5	c_6	l_7	c_8	l_9



The transformation of normalized circuit-element values is shown in Fig. 18.

LOW-PASS TO BAND-STOP TRANSFORMATION

The transformation of a low-pass filter into a band-stop filter is similar to the bandpass transform.

$$\Omega = 1/a(\eta - 1/\eta)$$

The resulting transformation of circuit-element values is shown in Fig. 19.

NUMERICAL EXAMPLES

Example 1

A low-pass filter with input impedance of 600 ohms and output impedance less than or equal to 600 ohms is needed to pass frequencies up to 3.4 kHz with less than 0.05 dB attenuation and to attenuate frequencies at 8.0 kHz and above by at least 45 dB (Fig. 20). Design both a Cauer-parameter and a Chebyshev low-pass to meet the requirement.

The value of Ω_s for the reference low-pass is calculated

$$\Omega_s = 8000/3400 = 2.353$$

Cauer-Parameter Design—To determine the Cauer-parameter low-pass, use Fig. 9 and proceed as follows.

Find the 45-dB (A_s) line on the grid (1), and move to the crossing with the $A_p = 0.05$ line (2). From this intersection (3), move to the left to intersect with the curves of degree versus cutoff rate. Note that for degree 3, Ω_s is $\cong 5.0$; for degree 4, $\Omega_s \cong 2.9$. Not until the curve for degree 5 is reached is an $\Omega_s < 2.353$ attained ($\Omega_s \cong 1.65$).

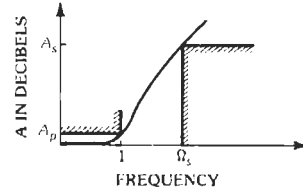
So a fifth-order low-pass with $A_s = 45$ dB, $A_p = 0.05$ dB and $\Omega_s = 1.65$ can be designed to meet the specified requirements. The tables included here are limited, however, so they must be checked to see what is available. From Table 4 (for filters of degree 5), $\theta = 37^\circ$ with 10% reflection (0.044 dB) and $\Omega_s = 1.66164$ meets the requirements. The normalized low-pass is shown in Fig. 21.

Since only an $\Omega_s \leq 2.353$ is required, some over-design can be accomplished. Generally, the losses erode the passband performance more rapidly than the stopband performance. As a result, the reference frequency (f_{ref}) can be chosen such that $f_s = 8000$ Hz.

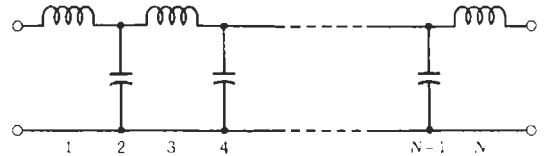
$$f_{ref} = 4.6245 \text{ kHz}$$

and since $R_{ref} = 600$ ohms,

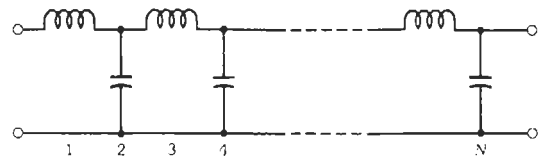
$$L_{ref} = 20.6492 \text{ mH}$$



(A) Performance.



(B) Odd-degree configuration.



(C) Even-degree configuration.

Fig. 14. Butterworth low-pass.

$$C_{ref} = 57.3588 \text{ nF}$$

The actual low-pass circuit (Fig. 21) then is:

$$C_1 = 50.665 \text{ nF}$$

$$C_2 = 6.625 \text{ nF} \quad L_2 = 25.75 \text{ mH}$$

$$C_3 = 89.021 \text{ nF} \quad L_4 = 20.67 \text{ mH}$$

$$C_4 = 19.141 \text{ nF} \quad f_2 = 12.188 \text{ kHz}$$

$$C_5 = 40.954 \text{ nF} \quad f_4 = 8.000 \text{ kHz}$$

Chebyshev Design—To determine the Chebyshev low-pass, use Fig. 12 and proceed as for the Cauer-parameter design. Here $N = 5$ results in $\Omega_s \cong 2.7$, and $N = 6$ provides $\Omega_s \cong 2.1$. Thus a sixth-order filter is required. From Table 9 for sixth degree and 10% reflection coefficient the normalized low-pass is as shown in Fig. 22.

With

$$f_{ref} = 3.8 \text{ kHz}$$

$$L_{ref} = 25.1297 \text{ mH}$$

$$C_{ref} = 69.8048 \text{ nF}$$

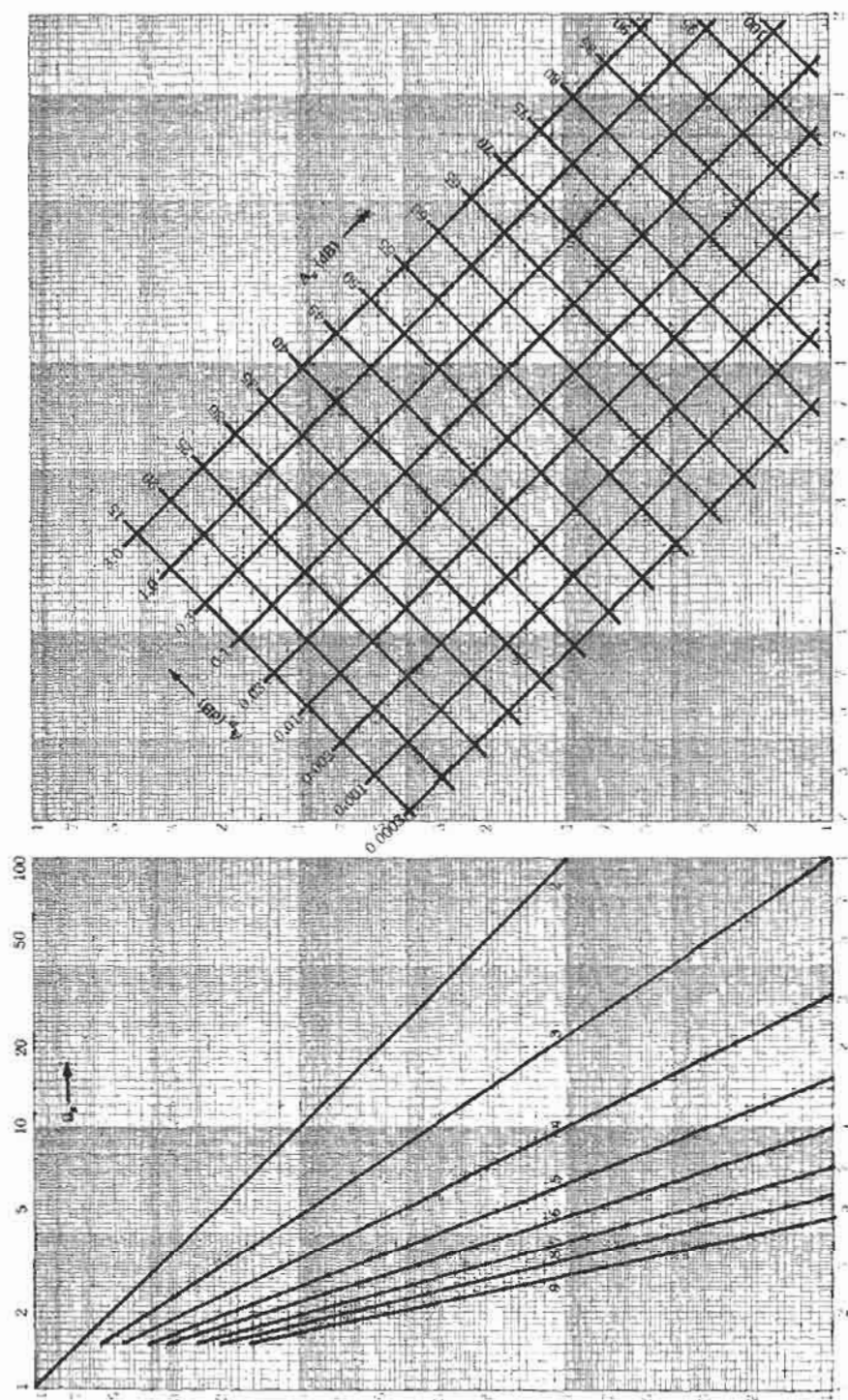


Fig. 15. Selecting the degree for a Butterworth low-pass.

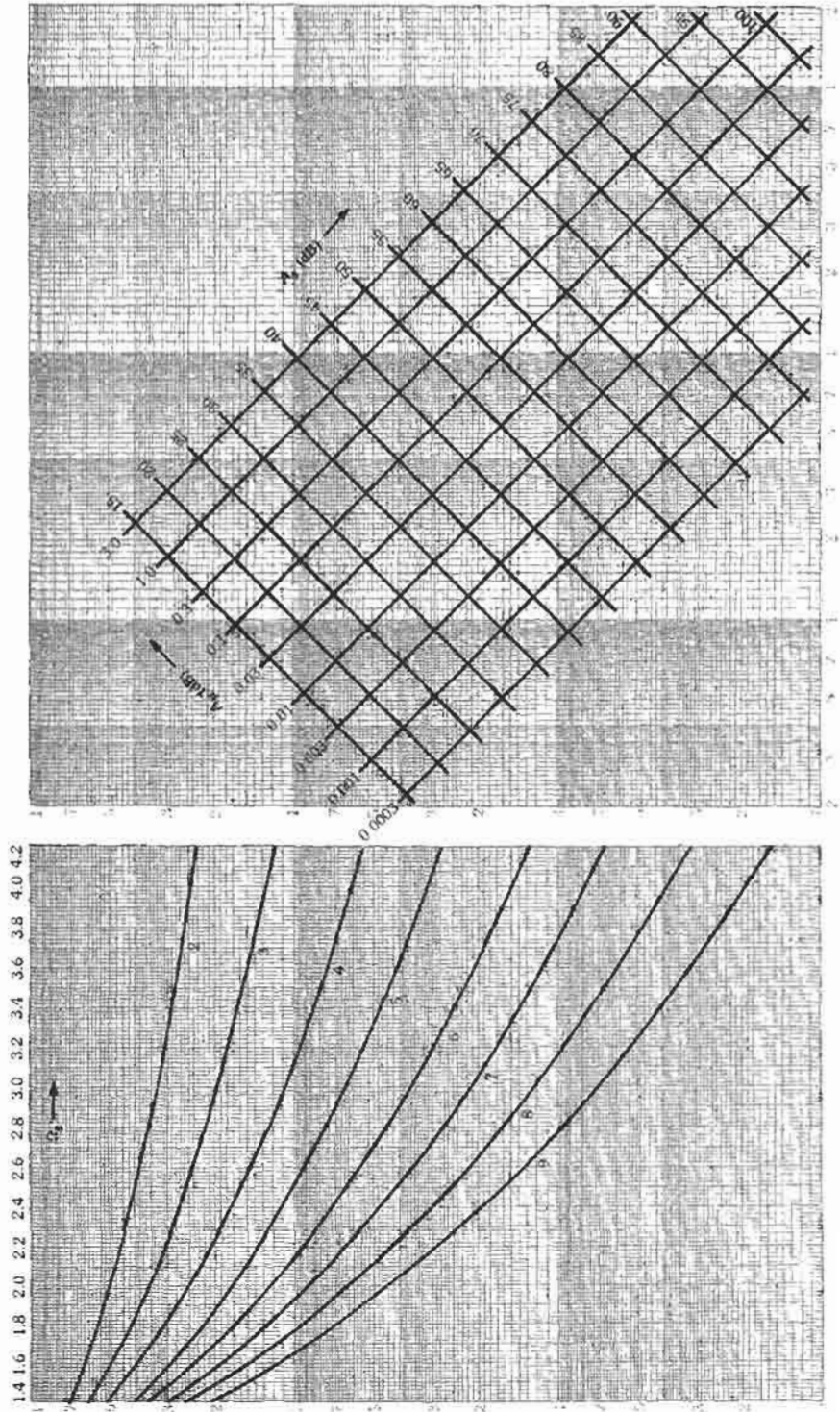


Fig. 16. Selecting the degree for a Butterworth low-pass —expanded.

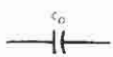
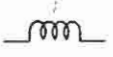

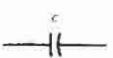
LOW-PASS	HIGH-PASS	TRANSFORMATION
		$l = 1/c_o$
		$c = 1/l_o$

Fig. 17. Low-pass to high-pass transform.

the actual circuit elements are:

$$C_1 = 69.512 \text{ nF}$$

$$L_2 = 35.508 \text{ mH}$$

$$C_3 = 132.28 \text{ nF}$$

$$L_4 = 38.951 \text{ mH}$$

$$C_5 = 120.55 \text{ nF}$$

$$L_6 = 20.473 \text{ mH}$$

$$R_7 = 490.9 \text{ ohms}$$

than 0.4-dB passband variation. (See Fig. 23.) From Fig. 5, $\Omega_s = 64/16 = 4.0$.

To determine the reference low-pass, use Fig. 15, which shows that for degree 5 an Ω_s of 3.2 is attainable. In Table 10, however, the choice is more limited, and with $N = 5$ and $A_p = 0.28$, an Ω_s of 3.3 is found (return to Fig. 15 with $A_p = 0.28$ dB).

The normalized reference low-pass is shown in Fig. 24. From the transform of Fig. 17, the normalized high-pass shown in Fig. 25 is obtained. With $f_{ref} = 64.000$ kHz and $R_{ref} = 600$ ohms,

$$L_{ref} = 1.4921 \text{ mH}$$

$$C_{ref} = 4.1447 \text{ nF}$$

and the actual circuit becomes

$$L_2 = 1.209 \text{ mH}$$

$$C_1 = 8.791 \text{ nF}$$

$$L_4 = 1.209 \text{ mH}$$

$$C_3 = 2.716 \text{ nF}$$

$$C_5 = 8.791 \text{ nF}$$

Example 2

A high-pass filter is needed to pass frequencies of 64 kHz and greater and to attenuate frequencies below 16 kHz by 40.0 dB with input and output impedances of 600 ohms. A Butterworth filter is preferred with less

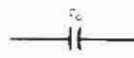
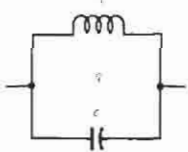
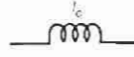
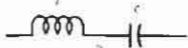
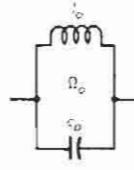
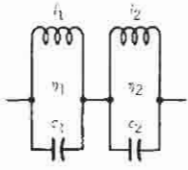
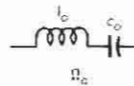
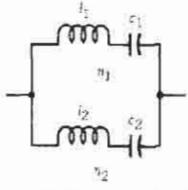
LOW-PASS	BANDPASS	TRANSFORMATION
		$\eta = 1.0 \quad c = 1/l \pi a c_o$
		$\eta = 1.0 \quad l = 1/c \pi a l_o$
		$c_1 = 1/l_2 - a c_o (1 + \eta_2^2)$ $c_1 = 1/l_1 - a c_o (1 + \eta_1^2)$
		$\eta_1 \eta_2 = \sqrt{(\Omega_o/2a)^2 + 1} = \Omega_o/2a$ $l_1 = 1/c_2 - a l_o (1 + \eta_2^2)$ $l_2 = 1/c_1 - a l_o (1 + \eta_1^2)$

Fig. 18. Low-pass to bandpass transform.

LOW-PASS	BAND-STOP	TRANSFORMATION
		$\eta = 1.0 \quad c = 1/l = c_0/a$
		$\eta = 1.0 \quad l = 1/c = l_0/a$
		$c_1 = 1/l_2 = (a/l_0)(1 + \eta_2^2)$ $c_2 = 1/l_1 = (a/l_0)(1 + \eta_1^2)$ $\eta_1 \eta_2 = \sqrt{(1/2a\Omega_0)^2 + 1} = 1/2a\Omega_0$
		$l_1 = 1/c_2 = (a/c_0)(1 + \eta_2^2)$ $l_2 = 1/c_1 = (a/c_0)(1 + \eta_1^2)$

Fig. 19. Low-pass to band-stop transform.

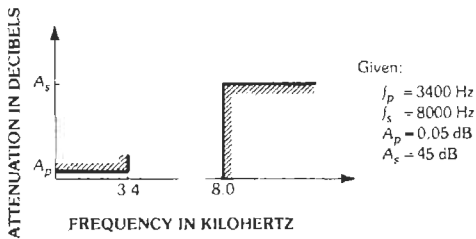


Fig. 20. Requirement for Example 1.

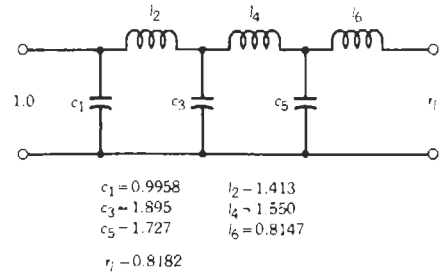


Fig. 22. Normalized Chebyshev low-pass for Example 1.

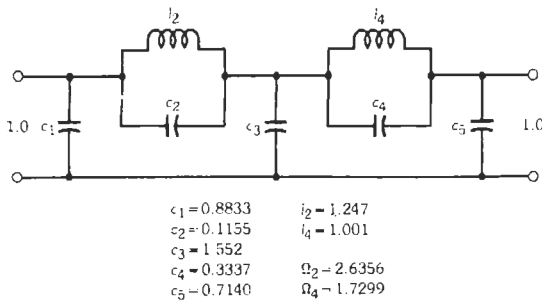


Fig. 21. Normalized Caucr-parameter low-pass for Example 1.

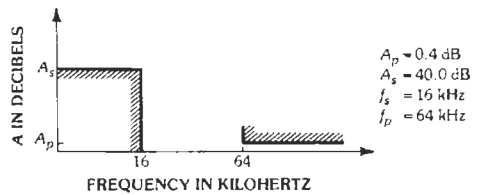


Fig. 23. Requirements for Example 2.

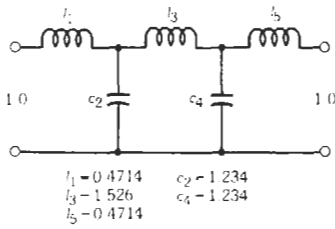


Fig. 24. Normalized reference low-pass for Example 2.

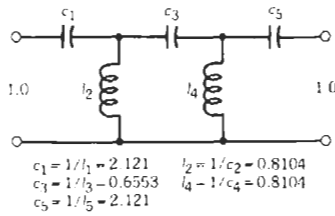


Fig. 25. Normalized high-pass for Example 2.

Example 3

A bandpass as shown by the tolerance plot of Fig. 26 with input and output impedances of 300 ohms is needed with stopband attenuation of 58 dB and a reflection coefficient of 5%. From Fig. 6 we might compute a reference low-pass $\Omega_s = (104 - 72)/(96 - 80) = 2.0$. However, geometric symmetry is required for both passband and stopband, so more calculation is necessary. From Fig. 6, it can be found that:

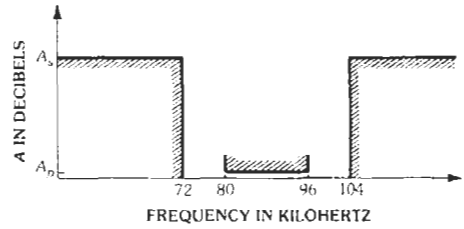
$$f_{ref}^2 = f_{+s}f_{-s} = f_{+p}f_{-p}$$

There are, then, two choices for the reference low-pass:

(A) Calculate a new f_{-p}

$$f_{-p} = f_{-s}f_{+s}/f_{+p} = 78.0 \text{ kHz}$$

from which



$A_p = 0.01087 \text{ dB}$ $f_{-p} = 80 \text{ kHz}$
 $A_s = 58 \text{ dB}$ $f_{+p} = 96 \text{ kHz}$
 $f_{-s} = 72 \text{ kHz}$
 $f_{+s} = 104 \text{ kHz}$

Fig. 26. Tolerance plot for Example 3.

$$\Omega_s = (104 - 72)/(96 - 78) = 1.778$$

(B) Calculate a new f_{-s}

$$f_{-s} = f_{+p}f_{-p}/f_{+s} = 73.846 \text{ kHz}$$

from which

$$\Omega_s = (104 - 73.846)/(96 - 80) = 1.885$$

Fig. 9 shows that a seventh-order filter is required to achieve an Ω_s as sharp as 1.885, and, in fact, $\Omega_s \cong 1.475$ can be achieved. From Table 6 with 5% reflection coefficient, $\theta = 43^\circ$ yields $\Omega_s = 1.46628$ and $A_s = 60.6 \text{ dB}$. (For filters of this complexity, only Cauer-parameter networks are generally considered.) The normalized reference low-pass is shown in Fig. 27.

Since this low-pass has a sharper cutoff rate than was required by either choice A or B, the bandpass transformation can be made with some overdesign throughout the tolerance plot. For instance, choose

$$f_{+p} = 97 \text{ kHz}$$

$$\text{with } f_{ref} = \sqrt{104 \times 72} = 86.5332 \text{ kHz}$$

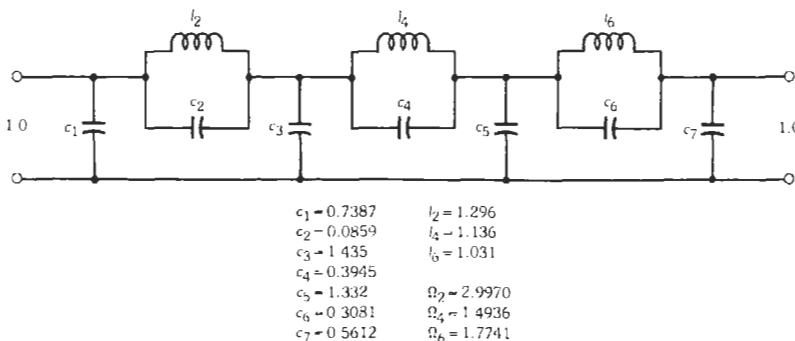


Fig. 27. Normalized reference low-pass for Example 3.

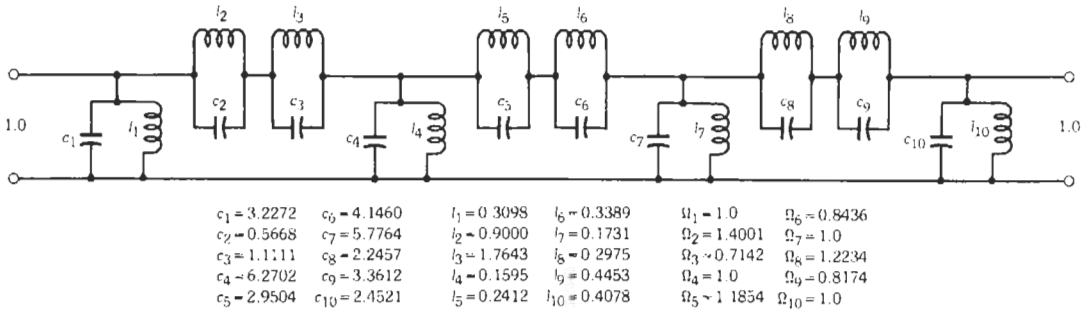


Fig. 28. Normalized bandpass circuit for Example 3.

$$f_{-p} = f_{ref}^2 / f_{+p} = 77.196 \text{ kHz}$$

$$\Delta f_p = 19.804 \text{ kHz}$$

$$\Delta f_s = \Omega_s \cdot \Delta f_p = 29.032 \text{ kHz}$$

Solving $f_{+s}f_{-s} = f_{ref}^2$ with $f_{+s} - f_{-s} = 29.032$

$$f_{+s} = 102.259 \text{ kHz}$$

$$f_{-s} = 73.226 \text{ kHz}$$

and

$$a = 4.369$$

With the transformations of Fig. 18, the normalized bandpass circuit becomes as shown in Fig. 28. (In Fig. 28, $\Omega_i = 1/\sqrt{c_i l_i}$.) Since $R_{ref} = 300$ ohms and $f_{ref} = 86.5332$ kHz,

$$L_{ref} = 551.77 \mu\text{H}$$

$$C_{ref} = 6.131 \text{ nF}$$

Then the actual circuit elements are:

$$C_1 = 19.788 \text{ nF} \quad L_1 = 170.95 \mu\text{H} \quad f_1 = 86.533 \text{ kHz}$$

$$C_2 = 3.475 \text{ nF} \quad L_2 = 496.59 \mu\text{H} \quad f_2 = 121.157 \text{ kHz}$$

$$C_3 = 6.812 \text{ nF} \quad L_3 = 973.48 \mu\text{H} \quad f_3 = 61.804 \text{ kHz}$$

$$C_4 = 38.441 \text{ nF} \quad L_4 = 88.00 \mu\text{H} \quad f_4 = 86.533 \text{ kHz}$$

$$C_5 = 18.088 \text{ nF} \quad L_5 = 133.09 \mu\text{H} \quad f_5 = 102.578 \text{ kHz}$$

$$C_6 = 25.418 \text{ nF} \quad L_6 = 187.01 \mu\text{H} \quad f_6 = 72.998 \text{ kHz}$$

$$C_7 = 35.414 \text{ nF} \quad L_7 = 95.52 \mu\text{H} \quad f_7 = 86.533 \text{ kHz}$$

$$C_8 = 13.768 \text{ nF} \quad L_8 = 164.16 \mu\text{H} \quad f_8 = 105.866 \text{ kHz}$$

$$C_9 = 20.607 \text{ nF} \quad L_9 = 245.70 \mu\text{H} \quad f_9 = 70.731 \text{ kHz}$$

$$C_{10} = 15.034 \text{ nF} \quad L_{10} = 225.02 \mu\text{H} \quad f_{10} = 86.533 \text{ kHz}$$

Example 4

Design a band-stop filter to operate between resistances of 75 ohms with 5% reflection coefficient and $A_s = 40$ dB, per Fig. 29.

Two options exist for computing the reference low-pass (see Fig. 7) with geometric symmetry:

(A) Calculate a new f_{-p} :

$$f_{-p} = (200 \times 248) / 272 = 182.35 \text{ kHz}$$

from which

$$\Omega_s = (272 - 182.35) / (248 - 200) = 1.867$$

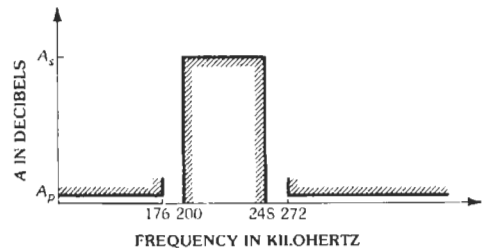
(B) Calculate a new f_{-s} :

$$f_{-s} = (272 - 176) / 248 = 193.03 \text{ kHz}$$

from which

$$\Omega_s = (272 - 176) / (248 - 193.03) = 1.7465$$

Fig. 9 indicates that a fifth-order filter is sufficient to meet the requirements. In fact, $\Omega_s \cong 1.70$ is achievable. From Table 4, for a reflection coefficient of 5%,



- $A_s = 40 \text{ dB}$
- $A_p = 0.01087 \text{ dB}$
- $f_s = 200 \text{ kHz}$
- $f_{-s} = 248 \text{ kHz}$
- $f_{-p} = 176 \text{ kHz}$
- $f_{+p} = 272 \text{ kHz}$

Fig. 29. Tolerance plot for Example 4.

$\theta = 36^\circ$ produces $\Omega_s = 1.7013$ and $A_s = 40.8$ dB. The normalized reference low-pass is shown in Fig. 30. Selecting $f_{ref} = \sqrt{200 \times 248} = 222.711$ kHz and $f_s = 249.000$ kHz gives:

$$f_{-s} = f_{ref}^2 / f_s = 199.197 \text{ and } \Delta f_s = 49.803 \text{ kHz}$$

$$\Delta f_p = \Delta f_s \cdot \Omega_s = 84.730 \text{ kHz}$$

and solving $f_{-p} f_{+p} = f_{ref}^2$ with $\Delta f_p = f_{+p} - f_{-p}$ gives:

$$f_{-p} = 184.339 \text{ kHz}$$

$$f_{+p} = 269.069 \text{ kHz}$$

$$a = 2.6285$$

After application of the transformation of Fig. 19, the normalized bandstop circuit is as shown in Fig. 31.

With $R_{ref} = 75$ ohms and $f_{ref} = 222.711$ kHz

$$L_{ref} = 53.597 \mu\text{H}$$

$$C_{ref} = 9.5284 \text{ nF}$$

and the actual circuit elements are:

$$L_1 = 207.26 \mu\text{H} \quad C_1 = 2.464 \text{ nF} \quad f_1 = 222.711 \text{ kHz}$$

$$L_2 = 11.254 \mu\text{H} \quad C_2 = 39.437 \text{ nF} \quad f_2 = 238.898 \text{ kHz}$$

$$L_3 = 12.950 \mu\text{H} \quad C_3 = 45.378 \text{ nF} \quad f_3 = 207.620 \text{ kHz}$$

$$L_4 = 102.58 \mu\text{H} \quad C_4 = 4.978 \text{ nF} \quad f_4 = 222.711 \text{ kHz}$$

$$L_5 = 8.391 \mu\text{H} \quad C_5 = 49.121 \text{ nF} \quad f_5 = 247.897 \text{ kHz}$$

$$L_6 = 10.397 \mu\text{H} \quad C_6 = 60.859 \text{ nF} \quad f_6 = 200.083 \text{ kHz}$$

$$L_7 = 277.85 \mu\text{H} \quad C_7 = 1.838 \text{ nF} \quad f_7 = 222.711 \text{ kHz}$$

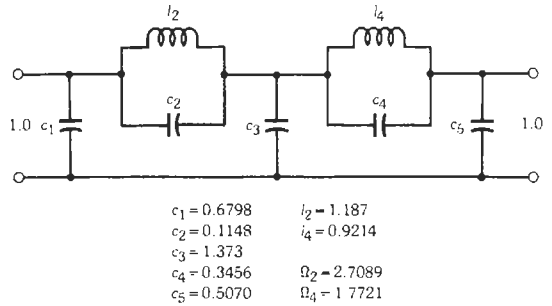


Fig. 30. Normalized reference low-pass for Example 4.

REFERENCES

1. Saal, R., and Ulbrich, E. "On the Design of Filters by Synthesis." *IRE Transactions on Circuit Theory*, December, 1958, pp. 284-327.
2. Saal, R. *The Design of Filters Using the Catalog of Normalized Lowpass Filters*. Telefunken G. M. B. H., Fachbereich Anlagen Weitverkehr und Kabeltechnik, 1963.
3. Christian, E. *LC-Filters Design, Testing and Manufacturing*. New York: John Wiley & Sons, Inc., 1983.
4. Cauer, W. *Synthesis of Linear Communication Networks*. New York: McGraw-Hill Book Co., 1958.
5. Christian, E., and Eisenmann, E. *Filter Design Tables and Graphs*. New York: John Wiley & Sons, Inc., 1966.
6. Darlington, S. "Synthesis of Reactance 4-Poles Which Produce Prescribed Insertion Loss Characteristics." *Journal of Math and Physics*, 1939, pp. 257-353.
7. Wetherhold, E. "Additional Modern Filters and Selected Filter Bibliography." *1984 Radio Amateur's Handbook*, 61st ed. Newington, Conn.: American Radio Relay League; pp. 2-40 to 2-46.

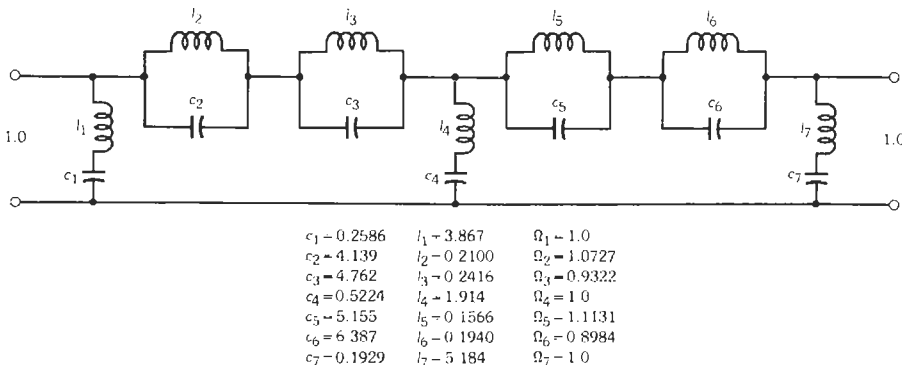


Fig. 31. Normalized band-stop circuit for Example 4.

10 Active Filter Design

Rolf Schaumann

Introduction	10-3
Circuit Elements	10-4
Resistors	
Capacitors	
Operational Amplifiers	
Operational Transconductance Amplifiers	
Fundamentals and Techniques of Active-Filter Design	10-7
The Transfer Function	
General Realization Methods	
Sensitivity	
Frequently Used Building Blocks	10-10
Summers	
Integrators	
The General Impedance Converter (GIC)	
The Gyrator	
The FDNR Element	
First-Order Transfer Functions	
Second-Order Transfer Functions—The Biquads	10-15
Single-Amplifier Filters	
Two-Amplifier Filters	
Three-Amplifier Filters	
Transconductance- C (g_m - C) Filters	

High-Order Transfer Functions *10-23*

- Cascade Realization
- Multiple-Feedback Topologies
- Ladder Simulation

Switched-Capacitor Filters *10-34*

- Sampled-Data Filter Systems
- The Operation of Ideal SC Filters
- SC Integrators
- SC z -Domain Biquadratic Transfer Functions
- Parasitic-Insensitive SC Biquads
- Low-Pass Notch Example
- Other Topologies and Comments

The information in this chapter enables the engineer to design a wide variety of practical active filters for operation in the audio-frequency range, and higher if fast operational amplifiers (op amps) are available. If operational transconductance amplifiers (OTAs) are used, filters even into the range of hundreds of megahertz can be designed. The equations presented permit the user to complete the design and arrive at a fairly comprehensive evaluation of the performance to be expected from the filter, without requiring complicated mathematics. Out of the countless different filters proposed in the technical literature, only those few circuits that have been proven to be practical, state-of-the-art designs are discussed in this chapter. Given the limited space available in a reference volume such as this, sufficient information can be provided only for the design of filters of relatively simple specifications. If system requirements are very stringent, the reader should consult the many excellent books* or papers referred to in the text.

INTRODUCTION

The technology of hybrid and monolithic integrated circuits has profoundly influenced the design and implementation of filters in the audio-frequency range and beyond. Frequency limitations are always set by the bandwidth of the available op amps. Integration has allowed the realization of filters that are small in size, inexpensive, and mass-producible. During the past 25 years, active *RC* networks, typically comprising resistors, capacitors, and operational amplifiers, have been the primary means of hybrid-integrated audio-filter implementation. Active *RC* filters have eliminated the need for the bulky, expensive inductors required in passive implementations, and tuning is simplified and usually involves the adjustment of only resistors. Furthermore, tuning can be automated in manufacture, using commercial laser trimming systems. In addition, active *RC* filters have provided opportunities for standardization and modularity that significantly simplify design and fabrication.

Switched-capacitor (SC) networks† have allowed audio-frequency active filters to be realized with the metal-oxide-semiconductor (MOS) large-scale-integration (LSI) technologies associated with digital networks. Switched-capacitor filters typically contain capacitors, MOSFET switches, and operational amplifiers. The switches are operated by clock signals that are digitally derived from a stable frequency source such as a crystal-controlled oscillator. The characteristics of the filter are then determined by capacitor ratios and the clock frequency, both inherently precise and stable parameters. Hence, SC filters rarely require trimming. The most important attribute of SC filters is that their implementation in silicon is compatible with

digital-circuit integration, enabling digital and analog circuitry to coexist on the same LSI chip.

For applications at much higher frequencies, such as in the read/write channels of magnetic disk recording systems or high-frequency communication systems, active *RC* filters based on operational amplifiers prove insufficient because of the op amps' bandwidth limitations. In these cases, when signals in the tens to hundreds of megahertz, or even gigahertz, must be processed, the designer uses operational transconductance amplifiers (OTAs) as active devices. OTAs are voltage-to-current converters described by their transconductance parameter g_m . It is quite possible to design OTAs with much higher bandwidth than op amps (up to several hundred megahertz and even gigahertz) so that analog active filters in the radio-frequency (RF) range become possible. Because discrete OTAs are not readily available, this latter technology is used mainly for integrated filters, where the use of only OTAs and capacitors enables the designer to implement high-frequency analog g_m -*C* filters compatible with digital CMOS technology. The filter performance depends unavoidably on OTA parameters, which must somehow be tuned. Techniques for handling these problems are becoming available on ICs.‡

At gigahertz frequencies, also passive *LC* filters can be implemented in integrated form because inductors of very small values (typically, nH) can be realized on an IC as single-layer or multiple-layer metal spirals. These inductors are very lossy, with quality factors of the order of only 10, so that active loss compensation techniques are usually employed; also, the spirals consume a relatively large silicon area, but the filters can be shown to have a lower noise and superior dynamic range than those using g_m -*C* techniques. For an overview of this design approach, the reader may consult Reference 28.

Integrated analog filters, especially for high-frequency applications, must be designed in fully differential, balanced form. Normally these filters will share an IC chip with digital circuitry where ground and power-supply lines (ac ground) are noisy, due, for example, to digital switching transients. Referring a signal to ac ground will, therefore, likely result in a severely restricted dynamic range (low signal-to-noise ratio). This problem is greatly reduced by referring two differential signal lines with equal and opposite signal voltages $\pm v_s/2$ to each other: $v_o = v_s/2 - (-v_s/2) = v_s$. Noise voltages on power-supply and ground lines appear as common-mode signals and are rejected in the differential circuit. Because OTA-based filters are intended mainly for IC implementation and for high frequencies, the corresponding examples below will be shown in differential form. Conversions from single-ended to differential form are straightforward.§

* References 1–10.

† Reference 10.

‡ References 6, 7, and 11–14.

§ Reference 6.

Although the design of high-frequency integrated OTA-C filters is, in principle, quite simple, in practice numerous constraints must be considered and these are beyond this introductory discussion. The presentation in this book should make the reader aware of the possibilities, but details are left to References 6, 7, and 10-14.

Active filters are used extensively in instrumentation and communication systems. To design these circuits, whether in active RC (based on op amps or OTAs), g_m -C, or SC techniques, one must first understand what an active filter is and how its performance requirements are specified. An electric filter is a network that transforms an input signal in some specified way into a desired output signal. Although many applications exist where filter requirements are set in terms of time-domain specifications, the majority of filters are designed to satisfy certain frequency-domain criteria. Thus, as shown in Fig. 1, a filter is a two-port network with input voltage V_1 and output voltage V_2 ; the circuit response is described by a transfer function $H(s)$ defined by

$$H(s) = V_2/V_1 = N(s)/D(s) \quad (\text{Eq. 1})$$

where, in steady-state,

$$\begin{aligned} s &= j\omega \text{ is the frequency parameter,} \\ \omega &= 2\pi f \text{ is the radian frequency (rad/s),} \\ f &\text{ is the frequency in Hertz (Hz).} \end{aligned}$$

As indicated, $H(s)$ is a ratio of two polynomials $N(s)$ and $D(s)$. The roots of $N(s)$ are the transmission zeros of the filter, i.e., points of infinite attenuation; the roots of $D(s)$ are its poles. The transfer function is a complex quantity that may be expressed as

$$H(j\omega) = |H(j\omega)| \exp[\phi(\omega)] \quad (\text{Eq. 2})$$

where $|H(j\omega)|$ is the magnitude and $\phi(\omega)$ is the phase. Thus, to specify a transfer function completely, both magnitude and phase must be given at a sufficient number of frequency points. In many cases, the magnitude response is the dominant specification with the phase response either loosely specified or unspecified. In this case, a minimum-phase filter (having all zeros in the right half-plane) is designed to meet the magnitude specification, and whatever phase the design provides is accepted. When both magnitude and phase are specified, a widely accepted design procedure is first to design a minimum-phase filter to meet the magnitude response, as previously mentioned. Then, an all-pass filter is designed, which, when cascaded with the minimum-phase filter, meets the desired phase specification. This all-pass network is often referred to as a phase or delay equalizer. Filtering implies that certain

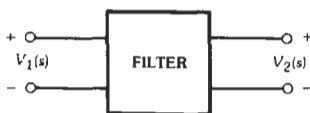


Fig. 1. A general filter representation.

frequency components of the input signal, those in the passband or passbands, are transmitted or passed to the output, whereas those in the stopband(s) are not transmitted. The most frequently used method of identifying the location of passbands and stopbands on the frequency axis is by specifying, versus frequency, the magnitude characteristic via the loss curve in decibels (dB), defined as

$$\alpha(\omega) = -20 \log |H(j\omega)| = -20 \log |V_2/V_1| \quad (\text{Eq. 3})$$

In the stopbands, where $|V_2| \ll |V_1|$, $|H(j\omega)|$ is small and the loss α is large, for example $|H(j\omega)| < 0.01$ or $\alpha > 40$ dB. In the passbands, $|V_2| \approx |V_1|$ or even $|V_2| \geq |V_1|$, so that $|H(j\omega)| \approx 1$ ($\alpha = 0$ dB) or $|H(j\omega)| > 1$ ($\alpha < 0$ dB), i.e., the circuit provides gain, something an active filter can do, whereas a passive filter always provides a loss.

If the phase response is of prime importance, then $\phi(\omega)$ is specified directly in degrees or radians; alternatively, and perhaps more frequently, one prescribes the delay $T(\omega)$ in seconds, defined as

$$T(\omega) = -d\phi(\omega)/d\omega \quad (\text{Eq. 4})$$

For best, distortion-free transmission, the delay should be constant, $T(\omega) = T_0$, i.e., the phase should be linear, $\phi(\omega) = -\omega T_0$, over the frequency range of interest. This is especially important when filtering pulse signals, such as in read/write channels of magnetic disk drives.

Some additional criteria of practical interest in active-filter design are sensitivity, dynamic range, noise, power dissipation, number and range of components, method of fabrication, and cost. All of these specifications place limitations and constraints on an acceptable design. In more cases than one would like, the specifications conflict so that engineering tradeoffs have to be made to resolve the conflict.

In the following sections, the components used for active filters, some important design criteria, and several state-of-the-art practical active filters will be discussed.

CIRCUIT ELEMENTS

Active filters are constructed from resistors, capacitors, and, usually, operational amplifiers for low-frequency applications or operational transconductance amplifiers for applications at high frequencies. A few comments will be made about these components, especially about the op amp because of its serious effect on filter performance. A suitably simple model for the OTA will also be presented. Spiral inductors of a few nanohenries (nH) in size are occasionally employed at the highest (GHz) frequencies in ICs.

Resistors*

Resistors used in active-filter design are carbon composition, metal or carbon film, thin or thick film,

* References 3 and 9.

wirewound, and diffused. The selection depends on cost, on the technology used to implement the filter, and on filter requirements. Carbon composition resistors are the least expensive, but they have large tolerances and temperature coefficients. Further, tracking is not very good so that composition resistors should be used only for noncritical applications. Metal-film and wirewound resistors, although more expensive, are better than composition types in all respects and are the most frequently used resistors in active-filter design today. Wirewound resistors have somewhat larger parasitics (L and C) than metal-film resistors and should not be used for high-frequency applications.

Capacitors*

Of the numerous different types of capacitors available, those commonly used in active filters are ceramic disc, Mylar, polystyrene, Teflon, and thin-film capacitors. Again, the selection depends on factors such as cost, available range, tolerances, temperature coefficients, and dissipation factor (loss). Ceramic and Mylar capacitors are the least expensive types and have the highest loss; they should be used only for noncritical applications. Teflon, thin-film, and especially polystyrene (or for small values, mica) capacitors are more expensive but have much lower dissipation factors and are therefore better suited for critical filter designs. Of course, in integrated implementations, on-chip metal-oxide-semiconductor (MOS) capacitors are being used.

In setting filter parameters, apart from tolerances and temperature coefficient, the dissipation factor (DF) or quality factor (Q_c) is of some importance. If loss is modeled by means of a resistor R_c in parallel with the capacitor C , Q_c and DF are defined as

$$Q_c = 1/DF = \omega CR_c \quad (\text{Eq. 5})$$

where ω is some critical frequency of interest, usually chosen at the passband edge. It should also be remembered that Q_c is a strong function of temperature. Typical values of Q_c range from less than 100 (ceramic) up to several thousand (polystyrene).

Operational Amplifiers

Although, as was mentioned, IC filters preferably use OTAs, the active element used in the vast majority of all discrete active RC filters and integrated SC filters is the operational amplifier. The "op amp" is an integrated circuit with five or more terminals, three of which are used for handling the signal (Fig. 2): the inverting (V^-) and noninverting (V^+) input terminals, at which the input voltages are applied, and the output terminal (V_o). The remaining terminals are for power supply and, in some models, for offset compensation

and frequency compensation. The function of the op amp is described by

$$V_o = A(s)(V^+ - V^-) = A(s)V_i \quad (\text{Eq. 6})$$

Thus, the output voltage is obtained by multiplying the differential input voltage, V_i , by $A(s)$, the op amp's open-loop gain.

In the ideal op amp, A is assumed to be infinite (i.e., $A \rightarrow \infty$) and, furthermore, all input impedances are infinite and the output impedance is zero. Thus, the ideal op amp is an ideal voltage-controlled voltage source with infinite gain; the input currents into the inverting and noninverting input terminals are zero. Further, since $A \rightarrow \infty$ and V_o must remain finite in practice, it follows that the input voltage $V_i = V^+ - V^-$ is zero; the operation is defined such that $V_o = 0$ when $V_i = 0$. Note that $V_i = 0$ implies that $V^- = 0$ when $V^+ = 0$. In that case, the V^- terminal is referred to *virtual ground*—the voltage is zero, but also the current is zero because of the op amp's infinite input impedance.

To get a feel for their operation, op-amp circuits often are first analyzed or designed under the assumption of ideal amplifiers ($A \rightarrow \infty$). It should, however, be strongly emphasized at this point that, except for uncritical applications at very low frequencies ($f < 1$ kHz), the operation of an active filter will rarely be satisfactory in practice if its design is based on ideal op amps. The main reason is that the op-amp gain $A(s)$ is a strong function of frequency. Specifically, for most inexpensive op amps used in active-filter design, the gain decreases from a high value of about 100 dB below 10 Hz throughout the useful frequency range by 20 dB/decade to 0 dB at about 1 MHz. This frequency response is required for stability reasons and is achieved by use of internal or external compensation capacitors.

Thus, the most widely used and quite realistic op-amp model is

$$A(s) = A_0 \sigma / (s + \sigma) = \omega_t / (s + \sigma) \quad (\text{Eq. 7})$$

where

$s = j\omega$ (as defined previously),

σ is the open-loop -3 -dB frequency (usually $\leq 2\pi \times 10$ Hz),

A_0 is the dc gain (usually > 100 dB),

$\omega_t = A_0 \sigma$ is the gain-bandwidth product ($\approx 2\pi \times 1$ MHz).

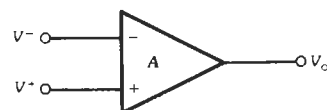


Fig. 2. Op-amp symbol.

* References 3 and 9.

In most practical applications, $|\sigma| \gg \sigma$, so that, instead of Eq. 7, a commonly used model is

$$A(s) \approx \omega_c/s \tag{Eq. 8}$$

Analyzing an active filter with op amps represented by Eqs. 7 or 8 rather than by a simple constant-gain model increases the degree of the describing network function (Eq. 1) by one for each amplifier used. Thus, the filter acquires parasitic poles and zeros in addition to suffering shifts of the nominal pole-zero locations. Stated differently, the two polynomials $N(s)$ and $D(s)$ in Eq. 1 change from their prescribed ideal form. It is for this reason that filters designed with the assumption of ideal amplifiers ($\omega_c \rightarrow \infty$) do not normally have a satisfactory frequency response but show a potentially large deviation from their nominal performance because of finite ω_c .

Occasionally, a filter will behave unpredictably, possibly even oscillate, in spite of a design based on the op-amp model of Eqs. 7 or 8. The reason may often be found in the fact that, whereas the representation of Eq. 7 describes the op-amp gain very well, the phase shift is larger than indicated by Eq. 7. An adequate remedy is to multiply Eq. 7 by an "excess phase factor," $\exp(-j\omega/\omega_2)$, where ω_2 is a normalizing frequency of value $\omega_2 \approx 3\omega_c$, and to use this augmented model in the analysis. Because for all practical purposes ω_2 is much larger than the operating frequency ω , we can set $\exp(-s/\omega_2) \approx 1 - s/\omega_2$ to keep the algebra manageable; i.e., a very accurate op-amp model for use in highly selective filters at fairly high frequencies is

$$A(s) \approx (\omega_c/s)(1 - s/\omega_2) \tag{Eq. 9}$$

In most applications, however, the use of Eq. 8 leads to entirely satisfactory results.

Further op-amp characteristics of concern to filter designers are slew rate and dc offset voltage. Slew rate (SR), given in volts per microsecond, refers to the maximum rate of change of a signal voltage that the amplifier can maintain at its output. Violating slew-rate limitations results in gross signal and/or transfer-function distortion and should be avoided. Thus, if $v_o(t) = V_0 \sin \omega t$ is the amplifier output voltage, one needs to observe $V_0 < SR/\omega$. For example, if $SR = 0.7 \text{ V}/\mu\text{s}$ (a typical value for inexpensive op amps) and the signal frequency is 15 kHz, the signal amplitude must satisfy $V_0 \leq 7.3 \text{ V}$ to avoid slew-rate limitations.

There are two main reasons* for offset voltage. One is the need for dc input bias currents into the input stage of bipolar op amps. (There is no input bias current in MOS op amps.) The second arises from imbalances in the input stage. To be able to provide the input bias, dc paths must exist from the inverting and noninverting input terminals to ground. To minimize offset, the resistances seen from these two terminals back into

the network ought to be equal. Then the voltage drops caused by the two bias currents at the op-amp inputs are equal and the direct differential input voltage is zero, resulting in zero contribution to offset. In practice, things are not quite so simple because imbalances in the op-amp input stage cause the bias currents not to be exactly equal, resulting in a finite differential direct input voltage. This voltage is multiplied by the dc closed-loop gain (see below), leading to an output offset voltage that, however, frequently can be reduced to zero by means of a potentiometer connected to the offset-adjust terminals of the op amp.

All amplifier parameters, especially the important terms A_0 and ω_c , are strong functions of bias voltages and temperature, and in general are not well determined or even predictable from unit to unit. Therefore, one strives to minimize filter dependence on these parameters and always uses op amps in closed-loop feedback configurations (Fig. 3) where the gain dependence on device parameters is reduced. Using straightforward analysis, and remembering that op-amp input currents are zero and that $V^- = -V_2/A$, yields for Fig. 3A the inverting gain

$$V_2/V_1 = -(Z_2/Z_1) / [1 + (1 + Z_2/Z_1)/A] \tag{Eq. 10}$$

and for Fig. 3B the noninverting gain

$$V_2/V_1 = +(1 + Z_2/Z_1) / [1 + (1 + Z_2/Z_1)/A] \tag{Eq. 11}$$

In both circuits, R is chosen to equal $1/[Y_1(0) + Y_2(0)]$ in line with the above discussion about offset minimization so that both op-amp inputs see the same dc impedance back into the circuit. Note, though, that resistor R does not affect the signal gain because no signal cur-

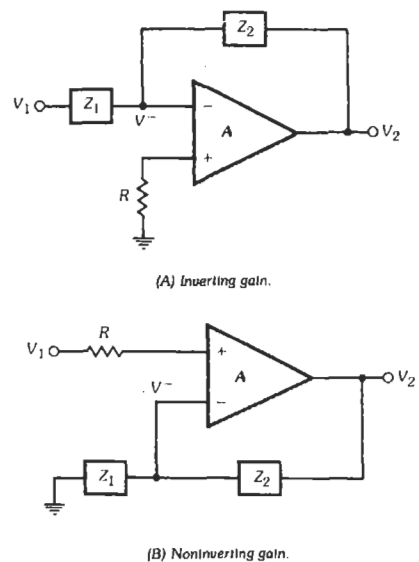


Fig. 3. Closed-loop feedback configurations.

* References 3, 7, and 9.

rent flows through it. Thus, for simplicity, in the remaining discussion in this chapter, R will be neglected; i.e., $R = 0$ is assumed.

For ideal amplifiers, or, in practice, in the frequency range where $|A(j\omega)| \gg |1 + Z_2(j\omega)/Z_1(j\omega)|$, Eqs. 10 and 11 reduce to

$$V_2/V_1 = -Z_2(s)/Z_1(s) \quad (\text{Eq. 12})$$

and

$$V_2/V_1 = 1 + Z_2(s)/Z_1(s) \quad (\text{Eq. 13})$$

respectively. As desired, the closed-loop gain functions are then independent of amplifier parameters and are determined only by presumably accurately adjustable and stable external impedances. The quantities $Z_1(s)$ and $Z_2(s)$ are finally chosen to yield the desired frequency dependence of the gain. For example, setting $Z_1 = R_1$ and $Z_2 = KR_1$ results, of course, in the well-known amplifiers of inverting gain $-K$ (Fig. 3A) and noninverting gain $1 + K$ (Fig. 3B). More will be said later about these important building blocks.

Operational Transconductance Amplifiers

As mentioned previously, an operational transconductance amplifier (OTA) is a voltage-to-current converter described by the conversion parameter g_m , the transconductance. Its circuit symbol is shown in Fig. 4A, and its function is given by

$$I_o = g_m(V^+ - V^-) \quad (\text{Eq. 14})$$

Naturally, at very high frequencies the gain, g_m , of a transconductance decreases because of parasitic poles. Just as in op amps, this decrease can be modeled by

$$g_m(s) = g_{m0}/(s + \sigma) \quad (\text{Eq. 15})$$

but in an OTA, the pole σ is at far higher frequencies (several tens to a few hundred megahertz) so that the effect may be neglected for most applications. Equation 14 indicates that an OTA generates an output current I_o proportional to the differential input voltage $V^+ - V^-$. Ideally, the input and output impedances and the bandwidth of an OTA are infinite. In practice, for well-designed circuits, the input resistor is greater than $10^8 \Omega$ and the output resistor is greater than $10^5 \Omega$, but there

are parasitic input and output capacitors of some small fraction of a picofarad (pF) in size. These capacitors must be accounted for in accurate designs at very high frequencies. The OTA's bandwidth is several hundred megahertz or, depending on technology and design, reaches the gigahertz range. From these few numbers it can be seen that, in contrast to op amps, OTAs can be treated as almost ideal components for most active-filter design tasks. Only at the highest frequencies must OTA nonidealities, in particular the parasitic capacitors and the parasitic pole σ , be taken into account. Fig. 4B shows an OTA with differential output terminals, and Fig. 4C shows an OTA with multiple differential inputs (here two) and differential outputs realizing

$$I_o^+ = I_o^- = g_m[(V_1^+ - V_1^-) + (V_2^+ - V_2^-)] \quad (\text{Eq. 16})$$

FUNDAMENTALS AND TECHNIQUES OF ACTIVE-FILTER DESIGN

The Transfer Function

In Eq. 1 the transfer characteristic of a filter was introduced as a ratio of polynomials, i.e.,

$$H(s) = \frac{N(s)}{D(s)} = \frac{n_m s^m + n_{m-1} s^{m-1} + \dots + n_1 s + n_0}{s^r + d_{r-1} s^{r-1} + \dots + d_1 s + d_0} \quad (\text{Eq. 17})$$

where

- $m \leq r$,
- r is the order of $H(s)$,
- n_i and d_i are real coefficients, with $d_i > 0$.

It is customary to scale the frequency parameter by some convenient normalizing frequency ω_n ; i.e., s in Eq. 17 is a normalized frequency

$$s = j\omega/\omega_n \quad (\text{Eq. 18})$$

Similarly, because a voltage transfer function is independent of the impedance level, all components (R and C) in the filter are scaled by a suitable normalizing resistor R_n so that the circuit is composed of dimensionless normalized resistors R/R_n and capacitors $\omega_n C R_n$. The advantages of this step are that all parameters in Eq. 17

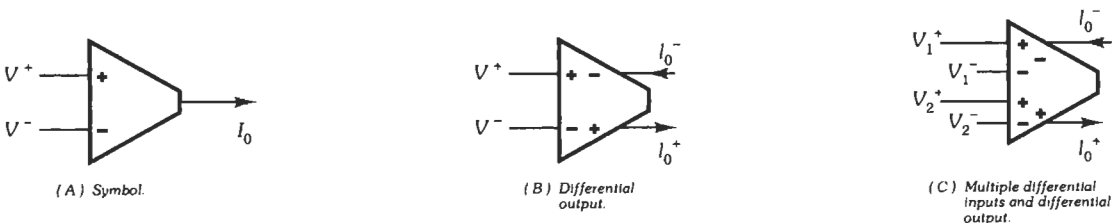


Fig. 4. Operational transconductance amplifiers.

are dimensionless and, more importantly, that the element values of the filter are dimensionless quantities, scaled for easier numerical computation.

In the following, it shall be assumed that $H(s)$ is given and our concern will be how to design an active filter to realize this prescribed transfer function. Space limitations do not permit us to discuss how to obtain $H(s)$ so that it approximates some desired magnitude or phase frequency response. This topic is treated in great detail in many excellent textbooks.*

General Realization Methods

If a high-order function such as Eq. 17 with $r \geq 4$ is prescribed, the engineer has to decide how to find a network structure and element values such that the measured voltage ratio V_2/V_1 is indeed as prescribed in $H(s)$. Numerous different techniques have been developed for this purpose, but they fall essentially into two different groups. One method attempts to break Eq. 17 into simpler functions of lower order, usually first-order

$$T_1(s) = (as + b)/(s + c) \quad c > 0 \quad (\text{Eq. 19})$$

and second-order sections

$$T_2(s) = \frac{as^2 + bs + c}{s^2 + s\omega_0/Q + \omega_0^2} \quad Q > 0, \omega_0 > 0 \quad (\text{Eq. 20})$$

which are then interconnected in a configuration suitable to implement Eq. 17. To keep with standard nomenclature, the denominator coefficients in Eq. 20 have been expressed in terms of the pole frequency ω_0 and the pole quality factor Q .

The second method draws on the vast experience in the many references and the numerous available tables for passive LC ladder filters, and their known excellent performance. In this case, either the structure or the equations describing the LC ladder are simulated via active RC filters; the inconvenient inductors are thereby avoided, but the positive properties of LC ladders are retained.

The details of high-order-function design are discussed in a later section of this chapter. Various popular and practically proven techniques are presented so that well-performing filters of reasonable complexity can be designed. In preparation for this step, realization methods for a number of elementary building blocks, such as first- and second-order sections, summers, integrators, and simulated inductors, are first discussed in the next two sections. Besides being useful in their own right, these components form the essential building blocks of high-order filters.

To decide which of the high-order design methods or which of the numerous available second-order filter sections might be best for a given application, it is necessary to establish practically useful performance criteria that permit the designer to make an informed

selection. In addition to obvious points, such as ease of design effort, number of active and passive components needed, spread and size of element values, and required power consumption, all of which reflect themselves in the final price of the filter, the generally accepted most important criterion for a good filter is that of low sensitivity.

Sensitivity

The response or performance parameters of networks depend in general on some or all of the components in the circuit. For example, the gain of the closed-loop amplifier in Fig. 3A, Eq. 10, depends on the impedances Z_1 and Z_2 and on the op-amp gain $A(s)$. In practice, the components cannot be expected to have or even maintain their ideal, calculated, nominal values. Rather, owing to factors such as fabrication tolerances, aging, and temperature drifts, the element values will vary, and, consequently, so will the circuit response.

Clearly, it is of great interest to the filter designer to get an estimate of the expected magnitude of the response deviation caused by element variations. Preferably, this information should be available before the circuit is ever fabricated so that judgments can be made about the likelihood of the initial and continued operation of the circuit within specifications. The theory of sensitivity addresses this question.

Assume that a circuit response depends on the N elements or parameters k_i , $i = 1, \dots, N$, i.e., $H(j\omega) = H(j\omega, k_1, k_2, \dots, k_N)$. When the elements change from k_i to $k_i + \Delta k_i$ with Δk_i a small change, the total variation of H is obtained from

$$\Delta H \approx \frac{\partial H}{\partial k_1} \Delta k_1 + \frac{\partial H}{\partial k_2} \Delta k_2 + \dots + \frac{\partial H}{\partial k_N} \Delta k_N \quad (\text{Eq. 21})$$

where higher order terms are neglected (Δk_i is assumed small!) and $\partial H / \partial k_i$ is the partial derivative of H with respect to k_i , evaluated at the nominal point. Upon normalizing, Eq. 21 can be rewritten as

$$\Delta H/H \approx \sum_{i=1}^N S_k^H (\Delta k_i/k_i) \quad (\text{Eq. 22})$$

where

$$S_k^H = \frac{k_i}{H} \frac{\partial H}{\partial k_i} = \frac{\partial H/H}{\partial k_i/k_i} = \frac{d \ln H}{d \ln k_i} \quad (\text{Eq. 23})$$

is the classical definition of the *single-parameter small-change sensitivity*; i.e., S_k^H is a measure that indicates the percentage deviation dH/H caused by the percentage variation dk_i/k_i of the element k_i . Thus, with S_k^H and the component tolerance known, for a change in a single parameter the expected response deviation is, from Eqs. 22 and 23,

$$\Delta H/H \approx S_k^H (\Delta k_i/k_i) \quad (\text{Eq. 24})$$

* References 1-8.

Equation 22, of course, gives an indication of the total change in H due to variations in all N parameters k_i . As an example, consider the op-amp stage with positive gain G in Fig. 3B, assuming for now $Z_i(s) = R_i$. From Eq. 11,

$$G = V_2/V_1 = (1 + R_2/R_1) / [1 + (1 + R_2/R_1)/A]$$

Using Eq. 23, it is easy to show that

$$S_{R_1}^G = -S_{R_2}^G = \frac{R_2/R_1}{(1 + R_2/R_1)[1 + (1 + R_2/R_1)/A]}$$

and

$$S_A^G = \frac{1 + R_2/R_1}{1 + (1 + R_2/R_1)/A} \cdot \frac{1}{A}$$

Thus, if the desired value of G is, say, 50 and if in the frequency range of interest $|A(j\omega)| \gg 50$ can be assumed, then $R_2/R_1 = 49$ and

$$S_{R_2}^G = -S_{R_1}^G \approx 49/50 \approx 1 \quad (\text{Eq. 25a})$$

$$S_A^G \approx 50/A \ll 1 \quad (\text{Eq. 25b})$$

Using Eqs. 22, 24, and 25,

$$\frac{\Delta G}{G} \approx -\frac{\Delta R_1}{R_1} + \frac{\Delta R_2}{R_2} + \frac{50}{A} \frac{\Delta A}{A} \quad (\text{Eq. 26})$$

These numbers indicate that, as long as $|A(j\omega)| \gg G$, changes in A have a negligible effect on ΔG , but that the percentage deviation of gain equals that of R_2 for constant R_1 and equals the negative of that of R_1 for constant R_2 . Note also that if R_1 and R_2 track, i.e., if they increase or decrease by the same amount, say, due to temperature changes in a thin-film hybrid circuit, then $\Delta G/G \approx 0$, because the ratio R_2/R_1 stays constant.

Several further observations can be made at this point. First note that since $A(j\omega)$ is a function of frequency, as for example in Eq. 8, G and all sensitivities also are functions of frequency and the sensitivity to $A(s)$ will increase with increasing ω because $|A(j\omega)| \approx \omega_i/\omega$ decreases. Second, note that the frequency limitation imposed on the validity of Eq. 26, the requirement $|A| \gg G = 50$, implies $\omega \ll \omega_i/50$ (e.g., $f \ll 20$ kHz for a 741-type op amp). This reemphasizes the point made earlier that active-network designs based on assumptions of ideal op amps can be expected to give reasonable results only at very low frequencies. Finally, remember that the op-amp gain is a very unreliable parameter: A_0 and ω_i often vary by more than 50% from unit to unit. Thus, this simple example indicates already that the circuit response, here the gain G , is very sensitive (Eq. 26) to A variations unless strong feedback is applied, i.e., $G \ll |A|$.

A number of interesting relationships that simplify sensitivity calculations can be derived.* Two of the more useful ones will be discussed. First, if a network

response function depends on a parameter k_1 that in turn depends on a second variable k_2 ,

$$H = H(k_1)$$

where

$$k_1 = k_1(k_2)$$

then it is easy to show from Eq. 23 that

$$S_{k_2}^H = S_{k_1}^H \cdot S_{k_2}^{k_1} \quad (\text{Eq. 27})$$

For example, if k_1 is a resistor R and k_2 the temperature T (i.e., $R = R(T)$), Eq. 27 permits an easy evaluation of S_T^H .

The second relationship makes use of Eq. 2; if Eq. 2 is inserted in Eq. 23, it should be clear that S_k^H is a complex number equal to

$$S_k^{H(j\omega)} = S_k^{H(j\omega)} + j \frac{\partial \phi(\omega)}{\partial k/k} \quad (\text{Eq. 28})$$

Thus, the real part of the sensitivity of the transfer function $H(j\omega)$ to a parameter k equals the sensitivity of the magnitude (gain) of $H(j\omega)$, and its imaginary part specifies the absolute variation of phase resulting from the relative error in k . Thus, gain and phase deviations can easily be estimated.

Equation 24 describes the change to be expected in $H(s)$ when only one circuit parameter k_i varies. To minimize this change, $\Delta k_i/k_i$ has to be minimized (implying components of tight tolerances and therefore high cost) and/or a design has to be found that minimizes S_k^H . Thus, the importance of low-sensitivity filter circuits should be clear.

Although useful, the insight provided by single-parameter sensitivities is somewhat limited in many cases because the effect of the remaining element variations is not taken into account. Therefore, for a more realistic picture, the deterministic multiparameter sensitivity measure, Eq. 22, ought to be consulted and a performance criterion, such as

$$\sqrt{\sum_{i=1}^N |S_{k_i}^H|^2} \rightarrow \text{Min} \quad (\text{Eq. 29})$$

should be evaluated for a given selection of circuits to arrive at a best choice. Since $\Delta k_i/k_i$ and $S_{k_i}^H$ in general can be of either sign, some cancellations can be expected when $\Delta H/H$ is calculated from Eq. 22. Thus, Eq. 29 gives a somewhat pessimistic worst-case picture of circuit performance. Further, this measure does not take into account that in many modern filter technologies, such as in hybrid or monolithic integrated circuits, the element variations are statistically related and frequently highly correlated. For example, all resistors may track and increase and all capacitors track and decrease with changes in temperature or during fabrication. Treatment of this case is beyond the scope of this chapter; the interested reader is referred to Reference 6.

* Reference 6.

Although the above sensitivity discussion has been concentrated solely on deviations of the transfer function $H(s)$ caused by component changes, it should be understood that sensitivity calculations apply to any network parameter whose value depends on variable circuit components. For example, the center frequency and the selectivity of a bandpass function, or the cutoff frequency of a low-pass filter, depend in general on several resistors and capacitors and possibly on op-amp or OTA gain. Any shift in these parameters resulting from component variations can be estimated by calculating the corresponding sensitivities. Similarly, the precise location of transmission zeros or filter poles depends, of course, on the correct element values. Any variations or tolerances in the latter cause shifts in these critical frequencies and therefore transfer-function errors that can readily be evaluated by sensitivity calculations. All second- and higher-order filters discussed in the remaining four sections of this chapter have been evaluated extensively as to their sensitivity performance and are recognized as the best designs available. If a different active RC filter out of the numerous topologies presented in the literature is to be used, the sensitivities of the contemplated circuit should be carefully investigated to make sure that the design will work satisfactorily in practice. It is noted again that simple postdesign tuning to eliminate fabrication tolerances will in general not suffice because circuit components cannot be expected to retain their values under environmental stresses, such as aging or temperature fluctuations. Thus, low sensitivity is a necessary requirement for any circuit. When evaluating a filter, the reader ought to keep in mind that the sensitivity results presented are valid for *small changes* in element values because of the linearization involved in the analysis: Note from Eq. 21 that the function $H(j\omega, k)$ has been replaced by its slope in the nominal point $k = k_0$ to estimate changes caused by varying k from k_0 to $k_0 + \Delta k$. In particular, this means that a circuit with the desirable property $S_k^H = 0$ may not be entirely independent of the parameter k ; it only says that at the nominal value $k = k_0$ the slope dH/dk is zero, usually implying a quadratic dependence on k . Large variations Δk can still cause unacceptable performance errors!

It is also emphasized that sensitivity expressions are usually functions of frequency so that the results ought to be evaluated in the frequency range of interest (normally in the passband or at the passband edges) when different possible designs are compared.

Finally, a last point is worth noting: Sensitivity is only an intermediate result and by itself can still create a misleading picture of circuit performance. In the ultimate analysis, what is important is the deviation or variability $\Delta H/H$, i.e., by Eqs. 21 or 24, the sensitivity multiplied by the expected component tolerances. Hence, relatively large sensitivities are acceptable when the relevant "component" (e.g., a resistor ratio in a thin-film hybrid realization) can be expected to be very accurate and stable. However, when a component

varies strongly (e.g., op-amp gain-bandwidth product), very low sensitivities must be insisted upon.

FREQUENTLY USED BUILDING BLOCKS

Active filters are generally constructed by interconnecting a number of well-understood elementary building blocks. Understanding the performance of those blocks enables the designer to assemble a filter that works better and more reliably. Thus, in this section, the most important functional blocks are being introduced: the summer, the integrator, the general impedance converter (GIC), the gyrator (which is a circuit that is used to simulate an inductor via a capacitor), the frequency-dependent negative resistor (FDNR), and a circuit for realizing a first-order voltage transfer function. A separate section is devoted to the realization of second-order biquadratic transfer functions because of their special importance.

Summers

As the name implies, a summer is used to add different signals in a filter to form a desired sum signal. The circuit is essentially an extension of Fig. 3A in which various input signals V_i , $i = 1, \dots, n$, are fed through weighting resistors to the virtual ground node V^- of the op amp (Fig. 5). Simple analysis yields

$$V_o = - \left(\sum_{i=1}^n a_i V_i \right) / \left[1 + \frac{1}{A} \left(1 + \sum_{i=1}^n a_i \right) \right] \Big|_{A \rightarrow \infty} = - \sum_{i=1}^n a_i V_i \quad (\text{Eq. 30})$$

Thus, for ideal amplifiers ($A \rightarrow \infty$), the output voltage is simply the negative sum of the input voltages, weighted by resistor ratios a_i . If the op-amp gain is finite and frequency-dependent, e.g., via Eq. 8, the summing coefficients exhibit frequency-dependent errors (specifically low-pass behavior) as is evident from Eq. 30. Whether or not this effect causes unacceptable deviations in a given filter depends on several factors and must be carefully investigated from case to case; no general rule can be given.

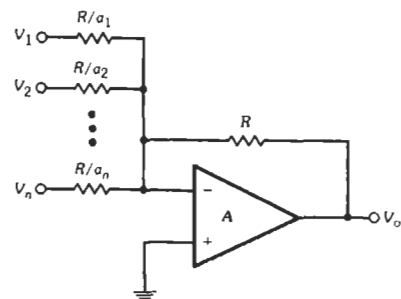


Fig. 5. Resistive op-amp summer.

In an OTA design, voltages are converted at each stage into currents as indicated by Eq. 14. *Current summing* is performed *automatically* at any circuit node via Kirchhoff's current law and so does not cost any elements. The designer is advised, therefore, to construct the circuit such that required summing functions are performed with currents rather than voltages. However, if voltages must be summed, it can be done as shown in differential form in Fig. 6. The circuits realize

$$g_{m1}(V_1^+ - V_1^-) + g_{m2}(V_2^+ - V_2^-) - g_{m3}(V_o^+ - V_o^-) = 0 \tag{Eq. 31}$$

or

$$(V_o^+ - V_o^-) = (g_{m1}/g_{m3})(V_1^+ - V_1^-) + (g_{m2}/g_{m3})(V_2^+ - V_2^-) \tag{Eq. 32}$$

Obviously, if the transconductances are identical, as is assumed in Figs. 6B and 6C, the scaling factors are unity. In Fig. 6B, the two OTAs g_{m1} and g_{m3} are merged into a two-input device; in Fig. 6C, all three devices are merged into one. This merging technique can be used whenever OTAs share a common output node and will generally result in significant savings of electronic components and power. Observe also that the transconductance g_{m3} with output fed back to the input realizes a grounded resistor of value $R = 1/g_{m3}$.

Integrators

One of the most basic functions in analog signal processing is that of integration; i.e., an output signal is obtained by integrating an input signal over time. In the frequency domain, this results in the transfer function

$$V_2/V_1 = \pm 1/(s\tau) \tag{Eq. 33}$$

where τ is the integrator time constant. Equation 33 can be implemented by use of Fig. 3A with $Z_1 = R$ and $Z_2 = 1/(sC)$; i.e., assuming ideal op amps, the simple inverting integrator realizes

$$V_2/V_1 = -1/(sCR) \tag{Eq. 34}$$

so that $\tau = RC$. Constructing filters by use of these very elementary, so-called Miller integrators results invariably in large errors, even at frequencies as low as 1 kHz, because of the overly ideal assumptions made about the op-amp gain. The problem can be appreciated by substituting the more realistic model of Eq. 8 into Eq. 10 with the result, for $\tau \gg 1/\omega$

$$V_2/V_1 = -1/(s\tau + s^2\tau/\omega_c) \tag{Eq. 35}$$

Thus, compared with the ideal case of Eq. 34, the integrator transfer function is no longer purely imaginary (it has acquired a parasitic pole), and its phase is not $\pi/2$ as an ideal integrator would demand. In analogy to the quality (Q) factor of capacitors and inductors, the ratio of imaginary to real part of the denominator of Eq. 35 is defined as the integrator Q -factor

$$Q_1 = -\omega_c/\omega \approx -|A(j\omega)| \tag{Eq. 36}$$

That is, the Q -factor of the Miller integrator is the negative of the op-amp gain at the frequency of interest.

It is apparent from Eq. 35 that the op amp introduces an excess phase lag into the integrator transfer function when compared to the ideal. Since filter behavior is usually particularly sensitive to such phase errors, procedures for their elimination have been devised; they can be classified as passive and active compensation methods, respectively. Fig. 7A shows a Miller integrator with passive phase compensation. Adding a compensation resistor of value $R_c = 1/(\omega_c C)$ into the feedback loop approximately cancels the phase lag introduced by the op amp at some nominal operating condition. The difficulty is that ω_c is not well determined and is variable (e.g., with temperature), and so compensation is difficult to maintain. A better, though more expensive, approach is active compensation (Fig. 7B), which attempts to cancel the errors by use of a second op amp. For well-matched amplifiers, routine analysis shows that the integrator of Fig. 7B has the large Q -factor of value $Q \approx -(\omega/\omega_c)^3 = -|A(j\omega)|^3$, a significant improvement. Two further integrators are shown in Figs. 7C and 7D. They are noninverting (plus sign in Eq. 33) and are very useful in their own right and especially because in filter circuits normally inverting and noninverting integrators are paired. The circuit

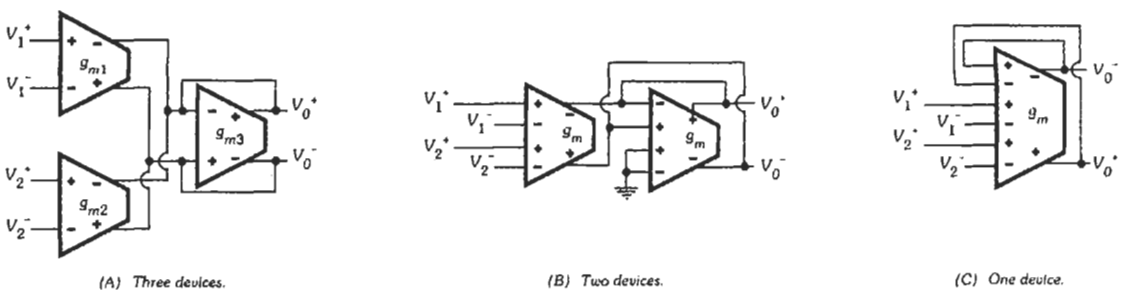
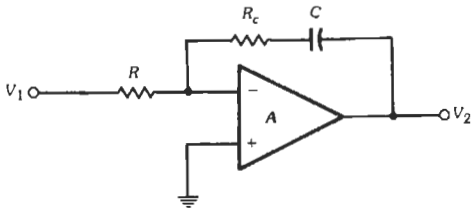
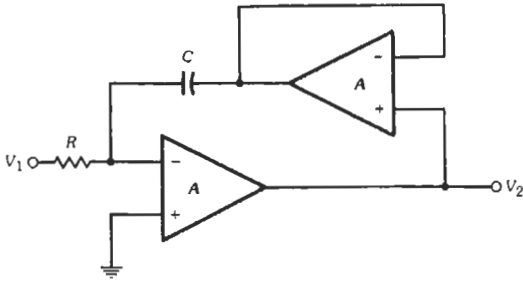


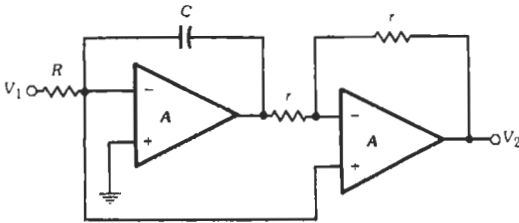
Fig. 6. Voltage summing with OTAs.



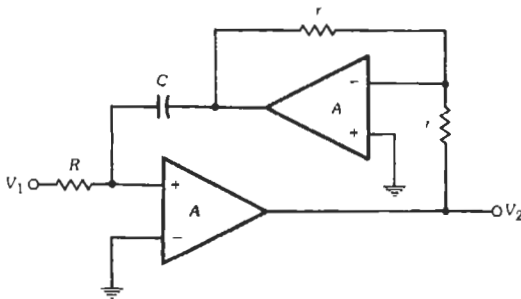
(A) Passive compensation.



(B) Active compensation.



(C) Noninverting, Q negative.



(D) Noninverting, Q positive.

Fig. 7. Different phase-compensated integrator circuits.

in Fig. 7C again has a Q -factor of value $Q \approx -|A(j\omega)|$, whereas Fig. 7D can be shown to have $Q \approx +|A(j\omega)|$. It is noted here especially that the Q -factor of Fig. 7D is positive; i.e., the integrator has a leading phase error. This fact is significant because active-filter structures often consist of two-integrator loops where the total

phase error can be virtually eliminated by pairing the circuit in Fig. 7D together with the regular Miller integrator in Fig. 7A (with $R_c = 0$). The result is much improved filter performance in spite of the relatively large phase error in each individual integrator.

In principle, building an integrator with OTAs requires only two elements: the OTA and a capacitor. The differential circuit in Fig. 8A shows how a voltage is integrated. The OTA converts the input voltage to be integrated into a current, which is then sent through a capacitor to produce a voltage that is the integral of the input:

$$V_o/V_i = g_m/(sC) \quad (\text{Eq. 37})$$

The positions of capacitor and OTA can also be interchanged (Fig. 8B) so that an input current first charges the capacitor whose voltage is sensed by the OTA and converted into an output current. The resulting current transfer function is identical to that in Eq. 37, $I_o/I_i = g_m/(sC)$. Filters based on this type of integration are called *current-mode filters*, which have a few desirable properties but behave essentially the same as voltage mode circuits. In the following, only voltage-mode circuits will be discussed.

As a practical matter, observe that the OTA output resistor $r_o = 1/g_o$ in Fig. 8A is in parallel with the capacitor C , so that the realized integration is lossy as described by

$$V_o/V_i = g_m/(sC + g_o) \quad (\text{Eq. 38})$$

As mentioned earlier, for well-designed OTAs g_o can often be neglected, but the engineer should be aware of this effect and of the need of very small values of g_o (large values of r_o). Fig. 8C illustrates two concepts: (1) the integration can be made lossy on purpose as required for many filter blocks by including a resistor $1/g_{m2}$, and (2) the capacitors can be realized by grounded elements, which are often advantageous for reduced parasitics in IC implementation; the circuit realizes

$$V_o/V_i = g_{m1}/(sC + g_{m2}) \quad (\text{Eq. 39})$$

Because the two OTAs in Fig. 8C share common output terminals, they can be merged into one for potential savings in components and power. The process was illustrated in Fig. 6.

The General Impedance Converter (GIC)

One of the motivations for *active* filters is the elimination of inductors. Thus, finding methods for inductance simulation is an important topic in active-filter design. The best op-amp-based circuit for inductance simulation is the general impedance converter (GIC) shown in Fig. 9. Figure 9A shows the circuit and Fig. 9B shows a symbol that will be

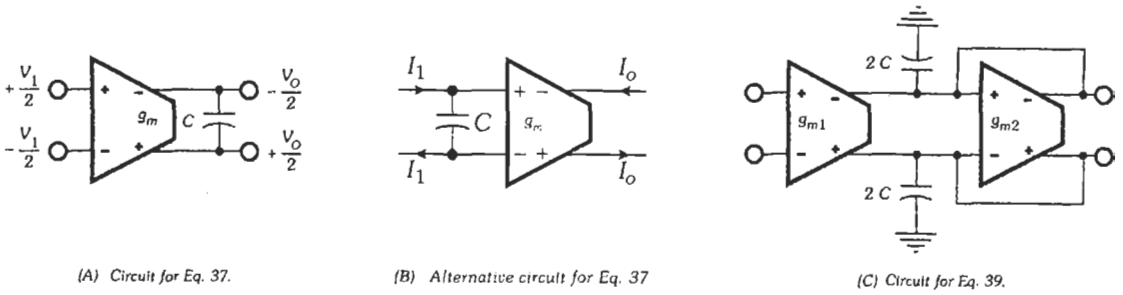


Fig. 8. OTA-C integrators.

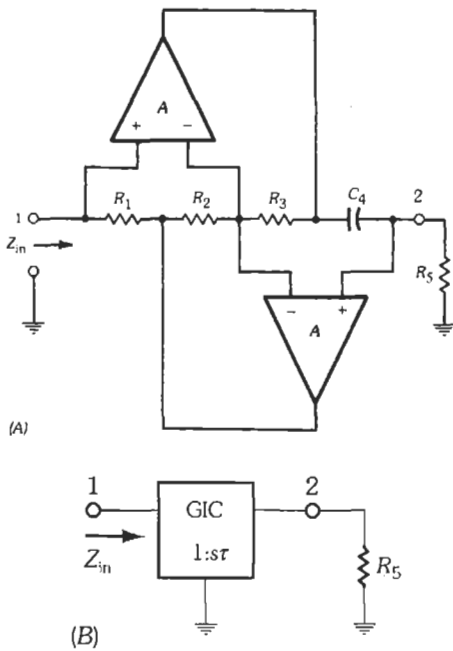


Fig. 9. GIC circuit to simulate a grounded inductor: (A) op-amp circuit, (B) GIC circuit symbol.

used later in the discussion of ladder networks. The circuit realizes between terminals 1 and ground the input impedance

$$Z_{in}(s) = s(C_4 R_1 R_3 / R_2) R_5 = s \tau R_5 \quad (\text{Eq. 40})$$

The definition of the time constant τ is apparent. Evidently, the circuit simulates a grounded inductor of value $L_0 = (C_4 R_1 R_3 / R_2) R_5 = \tau R_5$. To minimize the dependence of L_0 on op-amp performance (finite ω_c), the circuit is carefully optimized by setting $R_2 = R_3$ and $C_4 = 1/(\omega_c R_5)$. Here ω_c is some critical frequency value chosen in the passband or at the passband edge of the filter that uses the inductor. For these conditions, $L_0 = R_1 / \omega_c$. R_1 is chosen to set the inductor value, and

the inductor error and the inductor quality factor can be shown to be

$$\Delta L_0 / L_0|_{\omega=\omega_c} \approx 4/|A(j\omega_c)| \quad (\text{Eq. 41})$$

$$Q = |A(j\omega_c)| / (2\epsilon_r) \quad (\text{Eq. 42})$$

where $\epsilon_r = |R_2 - R_3| / R_2$ represents the fractional mismatch between R_2 and R_3 . Note that $Q = \infty$ can be obtained by tuning the circuit such that $R_2 = R_3$, i.e., $\epsilon_r = 0$.

Fig. 9 simulates a grounded inductor. Although several procedures exist for simulating floating inductors with operational amplifiers, none has ever been found satisfactory in practice. Thus, alternative procedures had to be invented for those (e.g., low-pass) filters whose passive equivalent requires floating inductors. One such method was proposed by Gorski-Popiel.* It uses the fact that a GIC converts a load resistor, R_5 in Fig. 9, into an equivalent inductor of value τR_5 (Eq. 40) at its input terminals. As will be explained below, connecting two GICs back-to-back provides us with a floating inductor. Alternatively, impedance transformations may be used, which lead to the so-called frequency-dependent negative resistor (FDNR). This method will also be discussed below.

The Gyration

A gyration† is a two-port device whose input impedance is inversely proportional to the load impedance, i.e.,

$$Z_{in}(s) = r^2 / Z_{Load}(s) \quad (\text{Eq. 43})$$

where r is the gyration resistance, a parameter set by the implementation of the gyration circuit. Evidently, if the load is a capacitor, $Z_{Load} = 1/(sC)$, the input impedance becomes that of an inductor, $Z_{in}(s) = sr^2C$, i.e., $L = r^2C$, and the circuit is readily realized in IC form. A gyration is, in principle, a combination of an inverting and a noninverting voltage-controlled current source; it is quite difficult to build a good gyration via op amps, but it is almost trivial to achieve such a component with transconductors. The circuit in Fig. 10A shows a

* See References 1, 2, or 7.

† References 2, 6, or 7.

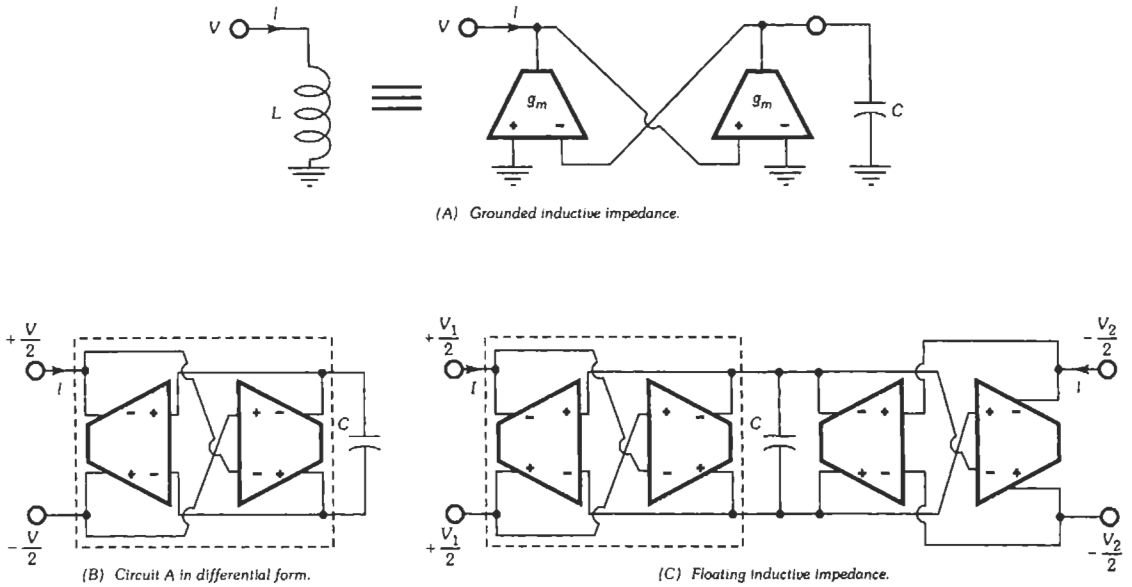


Fig. 10. Grounded and floating inductor simulation with OTA-based gyrators.

gyrator loaded by a capacitor C to realize the grounded inductive impedance

$$V/I = sL = sC/g_m^2 \quad (\text{Eq. 44})$$

For this circuit, the gyration resistance equals $1/g_m$. A floating inductor, as required for example in low-pass filters, can be built with two cascaded gyrators. The dashed box in Fig. 10B shows the circuit of Fig. 10A in differential form, and as an illustration, Fig. 10C contains an implementation of a floating inductor $Z_L = sL = (V_1 - V_2)/I$ in fully differential form. If all g_m values are the same, the inductor value realized is always given by Eq. 44.

The FDNR Element

The use of and motivation for FDNR elements in designing active filters based on passive LC prototypes will be discussed in a later section. Here, only the building block is introduced. The input impedance between terminal n and ground in the circuit of Fig. 11 is

$$Z_{in}(s) = [R_3 / (C_1 R_2 R_4 C_5)] / s^2 = 1 / (D_0 s^2) \quad (\text{Eq. 45})$$

For $s = j\omega$, Z_{in} is a frequency-dependent negative resistor of value $-1/(\omega^2 D_0)$ with $D_0 = C^2 R_2 R_4 / R_3$ where $C_1 = C_5 = C$ was assumed. Again, to minimize the dependence of D_0 on op-amp parameters (finite ω_c), the circuit is carefully optimized by choosing $R_2 = R_3$ and $R_4 = 1/(\omega_c C)$, where ω_c is a critical frequency value chosen in the passband of the filter or at the band edge. The fractional deviation of D_0 is then

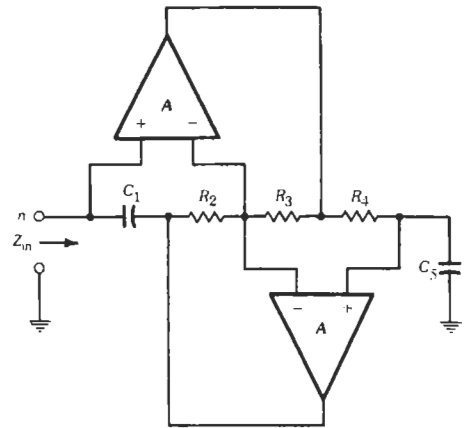


Fig. 11. FDNR element.

$$\Delta D_0 / D_0 = 4 / |A(j\omega_c)| \quad (\text{Eq. 46})$$

and its quality factor, defined as the ratio of real to imaginary part of $Z_{in}(j\omega)$, is

$$Q \approx |A(j\omega_c)| / (2\epsilon_r) \quad (\text{Eq. 47})$$

where ϵ_r is the fractional mismatch between R_2 and R_3 .

First-Order Transfer Functions

It was stated earlier that one technique for realizing a high-order transfer function involves breaking the func-

tion into first- and/or second-order subfunctions. These are then realized by suitable networks that are interconnected in an appropriate fashion. The realization of first-order voltage transfer functions is relatively simple and will be treated in the following. Second-order functions are discussed in the next section.

Fig. 12 shows different realizations of a function of the form of Eq. 19, repeated here as Eq. 48:

$$T(s) = V_2/V_1 = (as + b)/(s + c) \quad c > 0 \quad (\text{Eq. 48})$$

Fig. 12A is an appropriate realization if $0 \leq a \leq 1$ and $0 \leq b \leq c$. Note that in this case $T(s)$ is realized by a passive RC network. The amplifier is a special case of Fig. 3B for $Z_2 = 0$ and $Z_1 = \infty$; it is a buffer amplifier with unity gain and very high input and zero output impedance. The buffer has been included for the case in which the RC network may have a significant load; otherwise it can be omitted. Fig. 12B, a special case of Fig. 3A, realizes $-T(s)$; the minus sign, a 180° phase shift, is immaterial in many applications. Fig. 12C realizes $T(s)$ for the case in which $a \geq 1$ and $0 \leq c \leq b$. Finally, Fig. 12D takes care of the case when $a = 1$ and $b < 0$, i.e., $T(s) = (s - |b|)/(s + c)$, a regime important for all-pass functions (see below).

It is noted here that the circuits in Figs. 12B, 12C, and 12D realize Eq. 48 under the assumption of ideal amplifiers. As discussed earlier, if the op amp is described by Eq. 8, discrepancies from the ideal frequency response must be expected. Fortunately, however, the majority of all cases require only the realization of Eq. 48 with either $a = 0$ or $b = 0$. Fig. 12A can then always be used, possibly after factoring out a multiplying (gain) constant that can be taken care of in the remaining parts of the circuit. The unity-gain buffer amplifier requires for perfect operation only $|A(j\omega)| \gg 1$, as can readily be seen from Eq. 11.

With OTAs, a general and versatile first-order function can be implemented as shown in Fig. 13. The fully differential circuit realizes

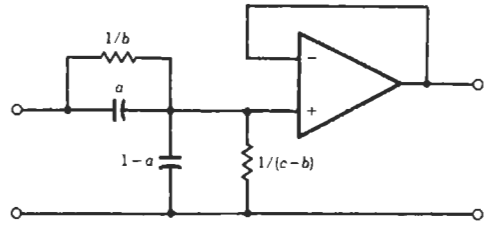
$$V_o/V_i = -(sC_1 + g_{m1})/(sC_2 + g_{m2}) \quad (\text{Eq. 49})$$

Note that depending on the polarity of the transconductances and their connections, the values of g_{m1} and g_{m2} may be positive or negative so that zeros anywhere in the s -plane and both inverting and noninverting transfer functions can be realized.

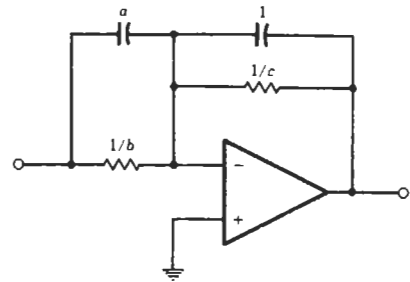
SECOND-ORDER TRANSFER FUNCTIONS—THE BIQUADS

Probably the most important basic building block used in the design of active filters is the biquad, a second-order circuit that realizes the transfer function of Eq. 20, which is repeated here as Eq. 50:

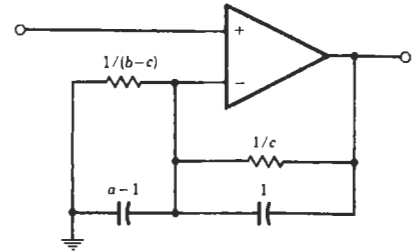
$$T(s) = \frac{N(s)}{D_2(s)} = \frac{as^2 + bs + c}{s^2 + s\omega_0/Q + \omega_0^2} \quad (\text{Eq. 50})$$



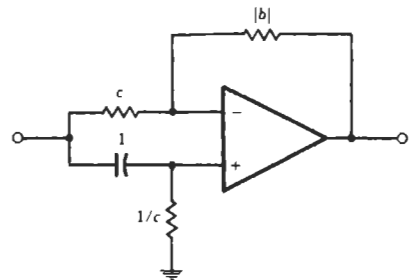
(A) For $0 \leq a \leq 1$ and $0 \leq b \leq c$.



(B) Realizes $-T(s)$.



(C) For $a \geq 1$ and $0 \leq c \leq b$.



(D) For $a = 1$ and $b < 0$.

Fig. 12. Different realizations of $T(s)$ in Eq. 48. The element values are normalized resistors and capacitors.

where $\omega_0 > 0$ and $Q > 0$. It is standard practice to identify the coefficients of the second-order denominator polynomial $D_2(s)$ via the pole frequency ω_0 and the

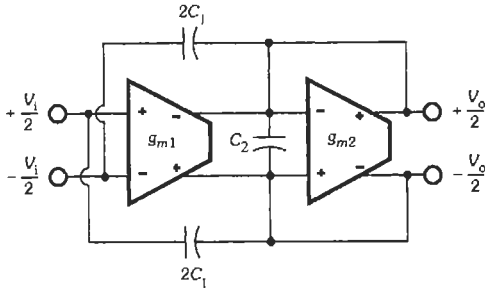
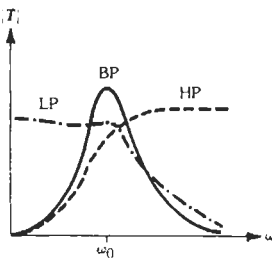
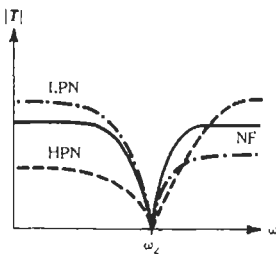


Fig. 13. First-order OTA-C filter section.

pole quality factor Q . The coefficients a , b , and c of the numerator polynomial $N(s)$ specify the type of transfer function realized (Fig. 14): If $N(s) = c$, i.e., $a = b = 0$, $T(s)$ is a low-pass (LP) function; if $N(s) = as^2$, i.e., $b = c = 0$, $T(s)$ is a high-pass (HP) function; if $N(s) = bs$, i.e., $a = c = 0$, $T(s)$ reduces to a bandpass (BP) function. A finite transmission zero at $\omega = \omega_z$ is obtained by setting $b = 0$; i.e., $N(s) = as^2 + c = a(s^2 + \omega_z^2)$ where $\omega_z^2 = c/a$. Such a circuit is called a “notch” filter (NF) and is distinguished further between a high-pass notch (HPN) if $\omega_z \leq \omega_0$ and a low-pass notch (LPN) if $\omega_z \geq \omega_0$. Finally, if $N(s)$ can be brought into the form $N(s) = a(s^2 + sb/a + c/a) = a(s^2 - s\omega_0/Q + \omega_0^2)$, i.e., $a \neq 0$, $c > 0$, $b < 0$, then $|N(j\omega)| = |D_2(j\omega)|$, so that $|T(j\omega)| = 1$.



(A) Pass characteristics.



(B) Notch characteristics.

Fig. 14. Typical plots of the basic second-order transfer functions.

Such an “all-pass” (AP) filter passes signals of all frequencies without attenuation but imposes a frequency-dependent phase shift of value

$$\phi = -2 \arctan \left[\left(\omega \omega_0 / Q \right) / \left(\omega_0^2 - \omega^2 \right) \right] \quad (\text{Eq. 51})$$

To realize the function of Eq. 50 or any of the mentioned special cases, the technical literature contains literally hundreds of different circuits. Fortunately, the choice is simplified for the engineer because closer examination reveals that the vast majority of all circuits have some drawback or another, usually excessive sensitivities, so that only two or three structures are used for practical active filters. These will be discussed in the following.

Single-Amplifier Filters

Fig. 15 contains three state-of-the-art single-amplifier active filters that can be used to realize the function in Eq. 50 with any of its special cases. The circuits are related to each other by complementary and by $RC:CR$ transformations. The transfer-function pole sensitivities, i.e., the ω_0 and Q sensitivities, can be shown to be the same for all three circuits. Specifically, the sensitivities of ω_0 to the passive elements are at their theoretical minima, $S_{k_i}^{\omega_0} = -0.5$, where k_i stands for any of the resistors and capacitors, and the passive sensitivities of Q satisfy

$$|S_k^Q| \approx Q/q - 0.5 \quad (\text{Eq. 52})$$

where q is a free design parameter, defined as a resistor ratio in Fig. 15. Further, it can be shown that ω_0 and Q deviations caused by the finite op-amp gain-bandwidth product ω_1 are

$$\Delta \omega_0 / \omega_0 \approx -q \omega_0 / \omega_1 \quad (\text{Eq. 53})$$

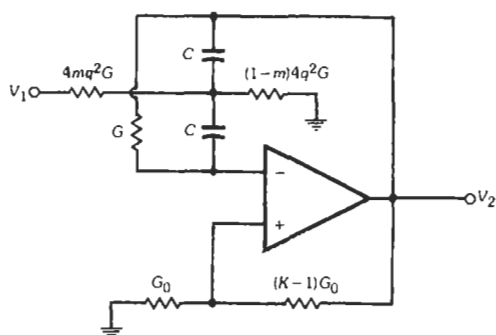
$$\Delta Q / Q \approx q \left(\omega_0 / \omega_1 \right) \left(1 - 2QK \omega_0 / \omega_1 \right) \quad (\text{Eq. 54})$$

where K is defined in Fig. 15. The value of free parameter q is arrived at through a compromise dictated by the technology chosen to implement the filter: It is apparent from Eqs. 52, 53, and 54 that increasing q will reduce the passive Q sensitivities but will increase the ω_0 and Q deviations caused by the active element via finite ω_1 . Taking into consideration that in the design, Eq. 60b below, the value of K depends on the choice of q , we can show that for sensitivity and realizability reasons (see Eq. 59) the value of q must be restricted by

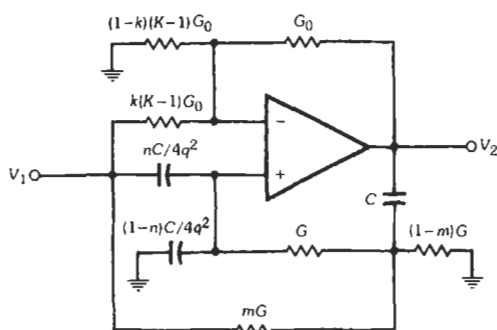
$$\sqrt{6}/2 \leq q \leq Q \quad (\text{Eq. 55})$$

In practice, q in the range of 3 to 5 is a reasonable tradeoff between active and passive sensitivities.

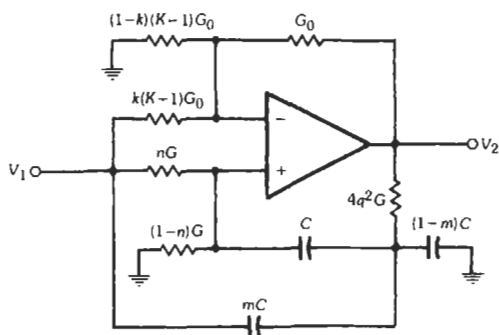
The initial error caused by finite ω_1 can be eliminated by predistortion, i.e., by applying the negative of the deviations, Eqs. 53 and 54, to the nominal values



(A) Bandpass biquad.



(B) Low-pass biquad.



(C) High-pass biquad.

Fig. 15. Single-amplifier biquads.

of ω_0 and Q and by basing the design on the pre-distorted values

$$\omega_{0p} = \omega_0 + \Delta\omega_0 = \omega_0(1 + q\omega_0/\omega_i) \quad (\text{Eq. 56})$$

$$Q_p = Q + \Delta Q = Q[1 - q(\omega_0/\omega_i)(1 - 2QK\omega_0/\omega_i)] \quad (\text{Eq. 57})$$

The transfer functions of the three circuits in Fig. 15A, B, and C can be shown to be

$$T_i(s) = N_i(s)/D(s) \quad i = A, B, C$$

where the denominator equals

$$D(s) = s^2 + s\omega_0/Q + \omega_0^2 \quad (\text{Eq. 58})$$

with $\omega_0 = 2qG/C$ and $Q = q/[1 - 2q^2(K - 1)]$, and the numerators are

$$N_A(s) = -s4Kmq^2G/C \quad (\text{Eq. 59a})$$

$$N_B(s) = K \left\{ s^2(n - p) + s \frac{2G}{C} [n - p(2q^2 + 1)] + \omega_0^2(m - p) \right\} \quad (\text{Eq. 59b})$$

$$N_C(s) = K \left\{ s^2(m - p) + s \frac{2G}{C} [n - p(2q^2 + 1)] + \omega_0^2(n - p) \right\} \quad (\text{Eq. 59c})$$

where the parameter $p = k(K - 1)/K$. To realize a transfer function of the form of Eq. 50 with given pole quality factor Q and pole frequency ω_0 , by Eq. 58 the elements C and $G = 1/R$ and the resistor ratio K in all three circuits in Fig. 15 are determined from

$$RC = 2q/\omega_{0p} \quad (\text{Eq. 60a})$$

$$K = 1 + (1 - q/Q_p)/(2q^2) \quad (\text{Eq. 60b})$$

where q should satisfy Eq. 55 and ω_{0p} and Q_p are replaced by ω_0 and Q if pre-distortion is not used. The resistor $R_0 = 1/G_0$ is arbitrary, because only the ratio $K - 1$ enters the transfer function. Further, either C or $R = 1/G$ can be chosen for convenience because only the RC product is determined by Eq. 60. As a rule, C is chosen because of the more limited number of available standard values.

The values of the circuit parameters m , n , and k (all less than or equal to unity) determine the type of transfer function. They are defined in Fig. 15 and as indicated in Eq. 59 determine the coefficients a , b , and c in Eq. 50. Equation 59 shows that the two circuits in Figs. 15B and 15C can realize a complete biquadratic transfer function as in Eq. 50; however, the transmission zero (ω_z) sensitivities and the notch depth of the circuit in Fig. 15C increase with increasing ratio ω_z/ω_0 . The opposite is true for the circuit in Fig. 15B. Thus, Fig. 15B gives the preferred circuit for a low-pass-notch, $\omega_z \geq \omega_0$ (see Fig. 14), and low-pass applications, whereas Fig. 15C shows the preferred circuit for high-pass-notch ($\omega_z \leq \omega_0$) and high-pass applications. Further, an all-pass filter can be designed using either of the biquads in Figs. 15B and 15C, and a bandpass filter is best built using the circuit in Fig. 15A.

With the help of Eqs. 58 through 60, arbitrary second-order functions can be designed. Specific design

information for the six basic biquads (those in Fig. 14 plus the all-pass filter) is given in Table 1 and Eqs. 50, 51, and 55 through 60. The design of a general biquadratic transfer function can be derived with Eqs. 59. The reader is also referred to References 7 and 15 on which the information in this subsection is based.

A couple of examples will illustrate the use of the information presented. First, assume the bandpass function

$$T(s) = -2s_n / (s_n^2 + 0.1s_n + 1) \quad (\text{Eq. 61})$$

has to be realized, where the frequency parameter s_n is normalized with respect to $\omega_0 = 2\pi \times 3 \text{ kHz}$, i.e., $s_n = s/\omega_0$. Equation 61 is a special case of Eq. 50 with $a = c = 0$, $b = -2\omega_0$; $\omega_0 = 2\pi \times 3 \text{ kHz}$, and $Q = 10$. Note that at $\omega = \omega_0$, i.e., $|s_n| = 1$, $T(s_n) = -20$; the midband gain of the filter is inverting with a magnitude of 20, or 26 dB. From Eq. 55, q should satisfy $1.2 \leq q \leq 10$; a good choice, as stated earlier, is $q = 4$. Assuming the circuit will be built using a 741-type op amp with $\omega_1 = 2\pi \times 900 \text{ kHz}$, the expected sensitivities and deviations are, from Eqs. 52 through 54, $|S_k^Q| \leq 2$, $\Delta\omega_0/\omega_0 \approx -0.013$, and $\Delta Q/Q = 0.013(1 - 0.07K)$, where, from Eq. 60b with $Q_p = Q$, $K = 1.01875$. Thus $\Delta Q \approx 0.12$, giving a 1.2% error. To reduce this error, predistortion may be used, i.e., from Eqs. 56 and 57 the circuit is designed for $Q = Q_p = 9.88$, resulting in $K - 1 = 0.0186$. Further, from Eq. 56, $\omega_{0p} = 2\pi \times 3.04 \text{ kHz}$ so that, from Eq. 60a, $RC = 418.83 \mu\text{s}$. Choosing $C = 10 \text{ nF}$ gives $R = 41.9 \text{ k}\Omega$. Finally, the bandpass design is completed by calculating, from Eq. 59a, $m = |b| / (2Kq\omega_{0p}) = 1/(Kq) = 0.2454$ (note that in this example $b < 0$). Choosing arbitrarily $R_0 = 1 \text{ k}\Omega$, gives the resulting bandpass filter shown in Fig. 16.

As the next example, design a low-pass notch filter realizing, with the maximum possible high-frequency gain H ,

$$T(s) = H(s_n^2 + 2.3) / (s_n^2 + 0.12s_n + 1) \quad (\text{Eq. 62})$$

where s_n is normalized with respect to $\omega_0 = 2\pi \times 2.4 \text{ kHz}$. Comparing Eq. 62 with Eq. 50 results in $\omega_0 = 2\pi$

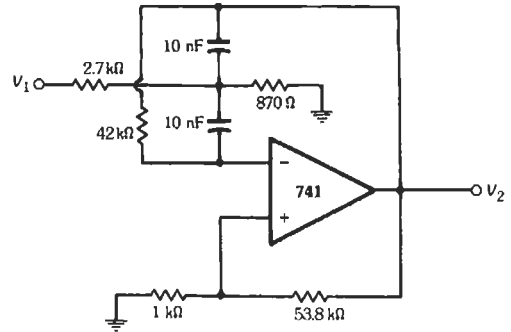


Fig. 16. Bandpass circuit of Fig. 15A realizing the example of Eq. 61.

$\times 2.4 \text{ kHz}$, $Q = 1/0.12 = 8.333$, $a = H$, $b = 0$, and $c = 2.3\omega_0^2 H$, i.e., $\omega_z = 1.5166 \omega_0$. Choosing $q = 3.5$ and assuming $\omega_1 = 2\pi \times 900 \text{ kHz}$ gives $|S_k^Q| \leq 1.881$, $\Delta\omega_0/\omega_0 \approx -0.009$, $\Delta Q/Q \approx 0.009$, and $K \approx 1.0237$, where Eqs. 52 through 54 and 60b were used. The deviations are quite small, so predistortion appears unnecessary; thus, selecting $C = 10 \text{ nF}$, from Eq. 60a, we get $R = 2q / (\omega_0 C) = 46.24 \text{ k}\Omega$. The function will be realized via the circuit in Fig. 15B; thus from Table 1: $k = 1.724H$, $m = 2.289H$, and $n = 1.0167H$. Since k , m , and n must be less than or equal to 1, $H_{\text{max}} = 0.4369$; then $k = 0.7532$, $m = 1$, and $n = 0.4442$. Finally, R_0 is chosen equal to $1 \text{ k}\Omega$, and the final low-pass notch circuit is as shown in Fig. 17.

Two-Amplifier Filters

The literature contains several multiamplifier second-order filters whose main advantage over single-amplifier circuits is that the sensitivities both to the passive components and to the op amps' finite gain-bandwidth product, ω_1 , are reduced. Based on a sensitivity comparison, only a few of these circuits have been shown to result in practical designs. Among

TABLE 1. DESIGN DATA FOR CIRCUITS IN FIG. 15

Filter Type	$N_z(s)$	k	m	n	Best Circuit
Low-Pass	c	0	$c/(K\omega_0^2)$	0	Fig. 15B
High-Pass	as^2	0	a/K	0	Fig. 15C
Bandpass	$-bs$	0	$b/(2Kq\omega_0)$	0	Fig. 15A
Low-Pass Notch	$a(s^2 + \omega_z^2)$ $\omega_z \geq \omega_0$	$a/(1 - q/Q)$	$p[1 + 2q^2(\omega_z/\omega_0)^2]$	$k 1 - q/(KQ) $	Fig. 15B
High-Pass Notch	$a(s^2 + \omega_z^2)$ $\omega_z \leq \omega_0$	$a(\omega_z/\omega_0)^2/(1 - q/Q)$	$p[1 + 2q^2(\omega_0/\omega_z)^2]$	$k 1 - q/(KQ) $	Fig. 15C
All-Pass { $a = \frac{Q-q}{Q+q}$ is fixed }	$a \left(s^2 - s \frac{\omega_0}{Q} + \omega_0^2 \right)$	1	$1 - \frac{2}{K(1+Q/q)}$	$1 - \frac{2}{K(1+Q/q)}$	Either Fig. 15B or Fig. 15C

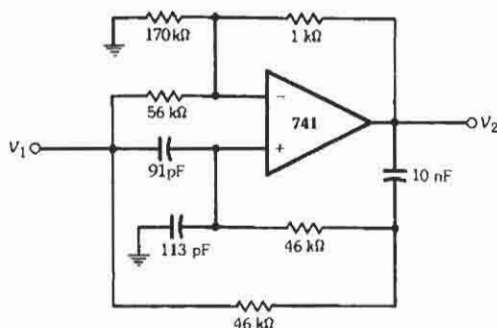


Fig. 17. Realization of low-pass notch function of Eq. 62 using circuit of Fig. 15B.

these, the biquad based on inductance simulation using an alternative version of the general impedance converter (GIC) of Fig. 9 has emerged as one of the best biquads owing to the reduced dependence of filter parameters on ω_1 .

The circuit capable of realizing a general biquadratic transfer function (except low-pass) is shown in Fig. 18. Simple analysis yields, for $\omega_1 \rightarrow \infty$,

$$T(s) = \frac{V_2}{V_1} = \frac{N(s)}{D(s)} = \frac{a_2 s^2 + a_1 s \omega_0 + a_0 \omega_0^2}{s^2 + s \omega_0 (1/Q + R/R_c) + \omega_0^2} \quad (\text{Eq. 63})$$

with

$$\omega_0 = 1/(RC) \quad (\text{Eq. 64a})$$

and

$$a_0 = b \quad (\text{Eq. 64b})$$

$$a_1 = [aH - b(H-1)](1/Q + R/R_c) - aHR/R_c \quad (\text{Eq. 64c})$$

$$a_2 = cH - b(H-1) \quad (\text{Eq. 64d})$$

The resistor R_c is a compensation resistor used only for the case when $H \neq 2$ (see Eq. 68). $R_c = \infty$ for $H = 2$.

Different types of transfer functions are realized by an appropriate choice of the parameters a , b , and c that are defined in Fig. 18. Table 2 gives details. Of course, a , b , and c must be less than or equal to 1, and $H > 1$. The circuit of Fig. 18 cannot realize a good low-pass function because the finite value of ω_1 degrades the performance. Thus, a low-pass filter is built with a slightly modified circuit as shown in Fig. 19; it realizes, with $\omega_0 = 1/(CR)$,

$$T(s) = V_2/V_1 = 2\omega_0^2 / (s^2 + s\omega_0/Q + \omega_0^2) \quad (\text{Eq. 65})$$

A few observations are in order about the performance of GIC-based filters: As was the case with the GIC in Fig. 9, the circuit is optimally insensitive to ω_1 if all GIC-internal resistors are equal. In Fig. 18, this clearly requires $H = 2$, which according to Table 2 implies that the midband gain of the bandpass filter and the high-frequency gain of the high-pass filter are fixed at $H = 2$. However, in many applications, especially in high-order filters based on second-order building blocks (see next section), the gain must be adjustable to maximize dynamic range.

If the required gain is less than 2, the lead-in element (QR for the bandpass filter, C for the high-pass filter), can simply be split into a voltage divider of an appropriate ratio. This step does not upset the desirable value $H = 2$. However, $H > 2$ can be achieved only by setting the two resistors $(H-1)R$ in Fig. 18 to the appropriate value (or using an additional amplifier to raise the signal). The effect of $H > 2$ can be shown to result in a very significant enhancement of Q in addition to a decrease of ω_0 with increasing values of ω_0/ω_1 . The ω_0 error

$$\Delta \omega_0 / \omega_0 = (\omega_s - \omega_0) / \omega_0 = -0.5 [H^2 / (H-1)] \omega_0 / \omega_1 \quad (\text{Eq. 66})$$

TABLE 2. DESIGN DATA FOR THE CIRCUIT IN FIG. 18

Filter Type	$N_2(s)$	a	b	c	Gain at $\omega =$	Comments
Low-Pass*	a_0	$b/2$	b	$b/2$	b , 0	$H = 2, b = a_0$
Bandpass	$a_1 s$	1	0	0	H , ω_0	Set $H = a_1 Q / \omega_0, R_c$ from Eq. 68
High-Pass	$a_2 s^2$	0	0	1	H , ∞	Set $H = a_2, R_c$ from Eq. 68
All-Pass	$s^2 - s\omega_0/Q + \omega_0^2$	0	1	1	1, All	Set $H = 2, R_c = \infty$
Notch	$a_2(s^2 + \omega_z^2)$ $\omega_z^2 = \frac{b\omega_0^2}{2c-b}$	$b/2$	b	b	b , 0	Set $H = 2, R_c = \infty$
					b , ∞	Low-pass notch: $b > c$
					b , 0, ∞	High-pass notch: $b < c$ Notch: $b = c$

* Not realizable well with this circuit (see Fig. 19).

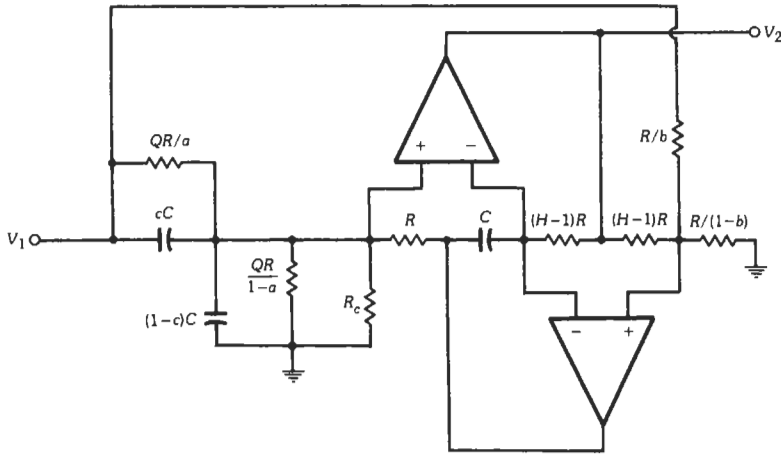


Fig. 18. A GIC-based biquad.

(where ω_a is the actually realized pole frequency) has its minimum at $-2\omega_0/\omega_1$ for $H = 2$. It can be corrected only by predistortion. The Q deviation, expressed as

$$Q_a = Q / \left\{ 1 + \left[1 - 2Q(2-H)/H \right] \Delta\omega_0/\omega_0 + 2(3-4/H)Q(\Delta\omega_0/\omega_0)^2 \right\} \quad (\text{Eq. 67})$$

(where Q_a is the actually realized pole quality factor) also is minimized for $H = 2$ and can be shown to be reduced approximately to zero, i.e., $Q_a \approx Q$, by connecting a compensation resistor

$$R_c \approx R(1+m)^2 / [m(H-2-m)] \quad (\text{Eq. 68})$$

in shunt with the GIC input, as indicated in Fig. 18. In Eq. 68,

$$m = [H/(H-1)]\omega_0/\omega_1 \quad (\text{Eq. 68a})$$

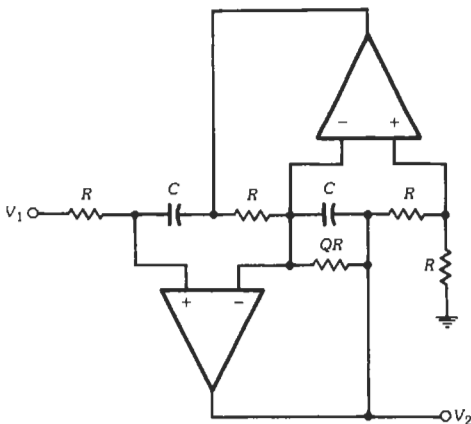


Fig. 19. GIC low-pass filter.

Equation 68 is valid for pole frequencies satisfying (approximately)

$$\omega_0 \leq 0.3\omega_1/H \quad (\text{Eq. 69})$$

For larger values of ω_0 , Eq. 68 becomes increasingly inaccurate, but Q compensation can still be achieved by functional tuning of R_c . A note of caution must be made: The value of R_c is seen to depend on ω_1 , which is process-, temperature-, and bias-dependent, is very inaccurate, and changes from op amp to op amp. Compensating this parameter with a passive resistor can, at best, be only approximate.

The function of R_c is only to eliminate ω_1 -caused Q enhancement for $H > 2$; otherwise, it has no effect on the realized transfer function apart from the slight dependence of a_1 on R_c shown by Eq. 64c. Since by Eq. 69a m is small, as is seen from Eq. 68, R_c is eliminated, i.e., $R_c \rightarrow \infty$ for $H = 2$.

In the all-pass and notch circuits, nothing significant is gained by choosing $H \neq 2$, and no usable compensation exists for $H \neq 2$ in the low-pass filter of Fig. 19. Thus, all GIC-internal resistors are set equal in these circuits.

Finally, it should be mentioned that Q is quite sensitive to capacitor losses. Labeling the capacitor loss resistor R_L and using Eqs. 5 and 64a, we obtain $Q_c = \omega_0 CR = R_L/R$. The actual quality factor, Q_a , can be shown to be

$$Q_a = Q / (1 + 2Q/Q_c) \quad (\text{Eq. 70})$$

where by Q we labeled the ideal quality factor, independent of such effects as finite ω_1 and $H \neq 2$. Thus, high-quality capacitors should be used to build these filters. Further details about GIC-filter performance can be found in References 2, 6, and 7, on which most of this discussion is based.

As a design example, assume that the bandpass function

$$T(s) = 0.8s_n / (s_n^2 + 0.1s_n + 1) \quad (\text{Eq. 71})$$

with $s_n = s/\omega_0$ and $\omega_0 = 2\pi \times 6 \text{ kHz}$ has to be realized. Assume further that op amps with $\omega_t = 2\pi \times 900 \text{ kHz}$ are available. From Eq. 71, pole frequency, quality factor, and midband gain are, respectively, $f_0 = 6 \text{ kHz}$, $Q = 10$, and $H = 8$. From Table 2 for a bandpass filter: $a = 1$ and $b = c = 0$. From Eq. 69 it follows that $f_0 = 6 \text{ kHz} < 33.7 \text{ kHz}$ so that Eq. 68 applies. Thus, $R_c = 22.24R$ will compensate for ω_t -caused Q enhancement that by Eq. 67 would be 79%; i.e., without R_c , $Q_a = 1.79Q = 17.9$. Remember though that this compensation depends on the specific value of ω_t assumed in the design. Further, from Eq. 66, $\Delta f_0 \approx -183 \text{ Hz}$, a -3.0% error. To eliminate this nominal deviation, the filter is designed with predistortion to realize $f_0 = 6.189 \text{ kHz}$ so that a -3% error gives $f_a = 6 \text{ kHz}$. Thus, from Eq. 64a with $C = 10 \text{ nF}$, $R = 2572 \ \Omega$; therefore, $QR = 25.7 \text{ k}\Omega$ and $(H - 1)R = 18.0 \text{ k}\Omega$. The circuit is shown in Fig. 20.

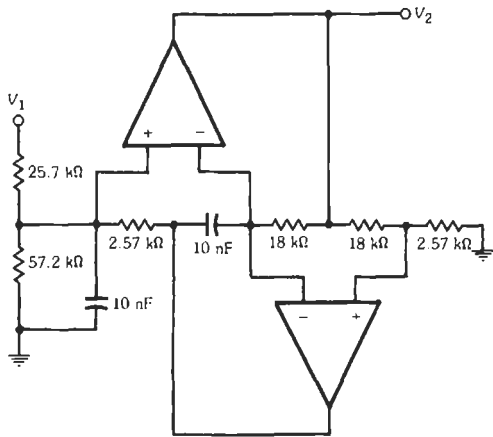


Fig. 20. GIC filter realizing Eq. 71.

Three-Amplifier Filters*

As demonstrated in the previous discussion, using additional amplifiers may lead to design and performance advantages for active filters. Thus, let us briefly present a three-amplifier section, the so-called Åckerberg–Mossberg (AM) biquad depicted in the top part of Fig. 21. The circuit consists of an inverting integrator (amplifier A_1 , C , R , QR , R/k) and a noninverting integrator (A_3 , C , and the inverter, A_2 and R_a) connected in a feedback loop. At the outputs V_L and V_B , respectively, the circuit realizes the inverting low-pass and bandpass transfer functions

$$T_L(s) = -k\omega_0^2/D(s), \quad T_B(s) = -k(s\omega_0/Q)/D(s) \quad (\text{Eq. 72})$$

where

$$D(s) = s^2 + s\omega_0/Q + \omega_0^2, \quad \text{with } \omega_0 = 1/(RC) \quad (\text{Eq. 73})$$

R , C , k , and Q are identified in the figure. If a fourth op amp is available, e.g., in a “quad” (four op amps on an IC chip), a general biquadratic transfer function can readily be obtained by adding the three outputs (V_L , $-V_L$, V_B) available in the AM circuit to its input. This is illustrated in Fig. 21 by the additional summer (built with A_4 and five resistors). The final circuit realizes the function

$$T(s) = \frac{N_2(s)}{D(s)} = \frac{\alpha s^2 + s[\alpha - \beta kQ]\omega_0/Q + [\alpha - (\gamma - \delta)k]\omega_0^2}{s^2 + s\omega_0/Q + \omega_0^2} \quad (\text{Eq. 74})$$

Evidently, any kind of second-order transfer function, as called for in Eq. 50, can be implemented with this circuit by choosing of the values (G_0 , αG_0 , βG_0 , γG_0 , δG_0) of the summing resistors appropriately. Table 3 gives the design details. Of course, low-pass and bandpass functions can be obtained directly from the basic AM circuit (Eq. 72) without any need for a summer.

* References 2, 6, and 7.

TABLE 3. DESIGN DATA FOR CIRCUIT IN FIG. 21*

Filter Type	$N_2(s)$	β	γ	δ
High-Pass	as^2	a	aQ	0
Low-Pass Notch	$a(s^2 + \omega_z^2)$ $\omega_z \geq \omega_0$	a	0	$aQ[1 - (\omega_z/\omega_0)^2]$
Notch	$a(s^2 + \omega_z^2)$ $\omega_z = \omega_0$	a	0	0
High-Pass Notch	$a(s^2 + \omega_z^2)$ $\omega_z \leq \omega_0$	a	$aQ[(\omega_z/\omega_0)^2 - 1]$	0
All-Pass	$a(s^2 + s\omega_0/Q + \omega_0^2)$	$2a$	0	0

* In all circuits $\alpha = a$, $k = 1/Q$, and $\omega_0 = 1/(RC)$.

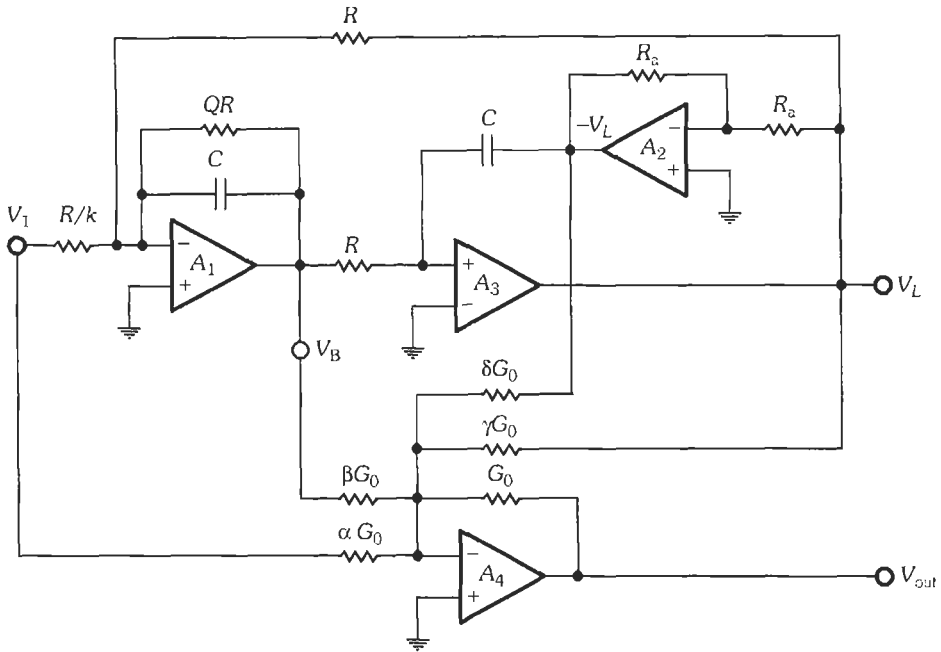


Fig. 21. The Åckerberg-Mossberg biquad with summer.

Transconductance-C (g_m -C) Filters

For applications at high frequencies, especially in integrated form, it is usually necessary to resort to g_m -C implementations of biquads because of the restrictions imposed by op-amp realizations. For these situations, a suitable biquad is shown in Fig. 22. (For case of understanding, a single-ended version is shown; see Reference 6 for a simple and systematic method of

converting a single-ended circuit into fully differential form.) The circuit realizes

$$V_o = \frac{(bV_3)s^2 + [bV_3 - aQ^2V_2](\omega_0/Q)s + (g_{m1}/g_m)\omega_0^2V_1}{s^2 + s\omega_0/Q + \omega_0^2} \tag{Eq. 75}$$

where the capacitors were selected as

$$C_1 = (g_m/\omega_0)Q, \quad C_2 = (g_m/\omega_0)/Q \tag{Eq. 76}$$

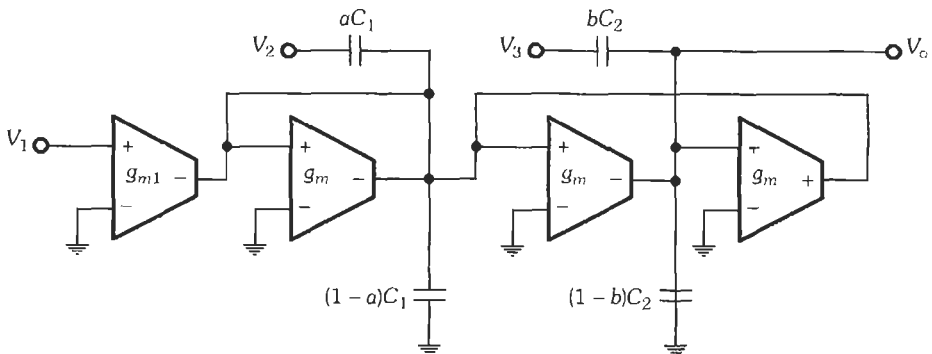


Fig. 22. Second-order g_m -C biquad.

and, therefore,

$$\omega_0 = g_m / \sqrt{C_1 C_2} \quad \text{and} \quad Q = \sqrt{C_1 / C_2} \quad (\text{Eq. 77})$$

Depending on whether the input V_{in} is applied to the terminal V_1 , V_2 , and/or V_3 , the circuit in Fig. 22 can realize any type of second-order function. In particular, a fully biquadratic function can be obtained by setting $V_1 = V_2 = V_3 = V_{in}$.

As an example, a practical differential low-pass circuit is shown in Fig. 23; it is obtained from the circuit in Fig. 22 setting $a = b = 0$ and converting the result to fully differential form. This circuit realizes

$$V_o / V_{in} = (g_{m1} / g_m) \omega_0^2 / (s^2 + s \omega_0 / Q + \omega_0^2) \quad (\text{Eq. 78})$$

Very similar low-pass sections have been used to implement fully integrated seventh-order constant-delay low-pass filters in the frequency range of 10–80 MHz for the read/write channels of magnetic disk storage systems.* Observe that all transconductors with common output terminals can be merged into one transconductor, a step that may save considerable power and electronic components on the integrated circuit.

HIGH-ORDER TRANSFER FUNCTIONS

As indicated earlier, high-order transfer functions are realized by one of three methods: (1) by cascading first- and/or second-order sections, (2) by connecting second-order sections in a suitable feedback topology, and (3) by simulating the elements of a lossless LC ladder or their operation. The primary advantage of Methods 2 and 3 is that they lead to designs with lower sensitivity to component variations than the cascade approach. The price paid for this benefit is a more complicated design and more difficult tuning, but a

* References 16 and 23.

major benefit of Method 3 is the possibility of using numerous design tables† and other filter information.

In the following, the cascade method, “follow-the-leader feedback” (FLF) designs, and the ladder simulation approach are discussed in sufficient detail to permit construction of reliable high-order filters. Derivations, proofs, and further details about these and several other approaches are contained in many excellent textbooks and papers.‡

Cascade Realization

The easiest and most popular method for realizing a high-order function $H(s)$ is via a cascade connection of second-order filter sections as shown in Fig. 24. For many applications cascade circuits are entirely adequate. Assuming no interaction between the sections T_i , $i = 1, \dots, l$, i.e., each filter block $T_i(s)$ has a low output and/or high input impedance, the overall voltage gain of the circuit in Fig. 24 is

$$H(s) = \frac{V_o}{V_i} = \frac{V_1}{V_i} \frac{V_2}{V_1} \dots \frac{V_{l-1}}{V_{l-2}} \frac{V_o}{V_{l-1}} = T_1 T_2 \dots T_l \quad (\text{Eq. 79})$$

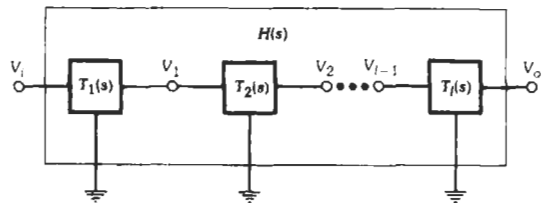


Fig. 24. Cascade connection of blocks $T_i(s)$.

† References 17 and 18.

‡ References 1–3, 6, and 7.

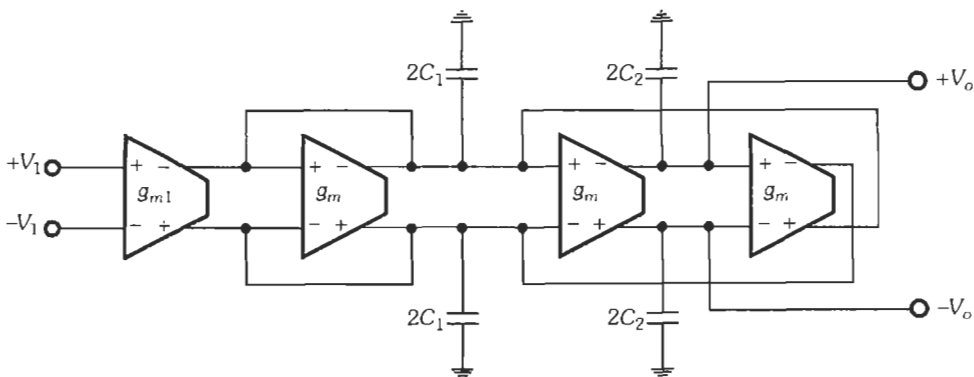


Fig. 23. Fully differential g_m - C second-order low-pass filter.

To this end, Eq. 17 is factored into the product of l second-order transfer functions $T_i(s) = N_i(s)/D_i(s)$ of the form of Eq. 50, such that

$$N(s) = \prod_{i=1}^l N_i(s), \quad D(s) = \prod_{i=1}^l D_i(s) \quad (\text{Eq. 80})$$

This step is accomplished simply by finding the (in general) complex roots of $N(s)$ and $D(s)$ and keeping conjugate complex terms together so that all coefficients in $T_i(s)$ are real.

The notation in Eqs. 17, 79, and 80 has tacitly assumed that $H(s)$ is of even order, i.e., that r is even. Of course, if r is odd, one of the transfer functions $T_i(s)$ in Eq. 79 must be of first order; that is, it must be of the form of Eq. 48, whose realizations are given in Figs. 12 and 13. Such a first-order function can always be cascaded with the remainder of the network. If in Eq. 17 $m < r$, some of the terms $N_i(s)$ in Eq. 80 are frequency independent constants so that the corresponding biquads are low-pass functions. Thus, it should be clear that the discussion can concentrate on the synthesis of even-order functions only; i.e., the question is how to realize Eq. 17 as expressed in Eq. 79 with Eq. 80. This, in turn, implies that practical methods have to be found to realize biquadratic transfer functions of the form of Eq. 50, a topic discussed in the previous section of this chapter.

As a simple example for the steps discussed, assume the function

$$H(s) = \frac{K(s^2 + 4\pi^2 \times 3.3551)}{s^4 + 2.7555s^3 + 3.7964s^2 + 3.099s + 1.4246} \quad (\text{Eq. 81})$$

a low-pass function with the frequency response sketched in Fig. 25, has to be realized. In Eq. 81, s is normalized with respect to the low-pass cutoff frequency $f_{3\text{ dB}} = 4.55$ kHz; i.e., the filter passband is in $0 \leq f \leq f_{3\text{ dB}}$, and $s = j\omega/(2\pi \times f_{3\text{ dB}})$. Realizing Eq. 81 as a cascade of two second-order sections requires factoring the denominator and writing $H(s)$ as a product of two functions, T_1 and T_2 . Simple algebra (root finding) results in

$$H(s) = \frac{K_1}{s^2 + 2.0905s + 1.3544} \cdot \frac{K_2(s^2 + 132.45)}{s^2 + 0.6650s + 1.0518} = T_1(s)T_2(s) \quad (\text{Eq. 82})$$

where $K_1K_2 = K$. The low-pass and low-pass-notch functions $T_1(s)$ and $T_2(s)$ can then be realized as described in the previous section, e.g., by the circuit in Fig. 15B.

In Eq. 82, with equal justification, the numerator factor $s^2 + 132.45$ could have been assigned to $T_1(s)$, the order of T_1 and T_2 could have been interchanged, or any combination of the above could have been per-

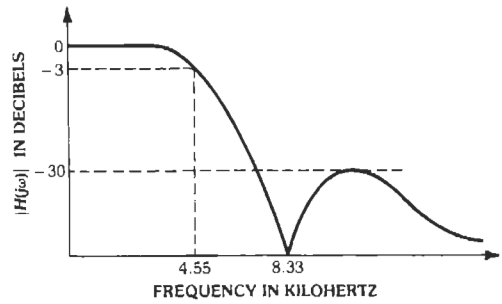


Fig. 25. Frequency-response plot, Eq. 81.

formed. Thus, since only the product of the functions $T_i(s)$ is prescribed, it is clear that considerable freedom exists, especially in functions of high order. Pole-zero pairing, assignments of the gain constants K_i , and ordering of the functions T_i in the cascade are not unique. The final choice of these factors determines such important practical characteristics as sensitivity and, most importantly, dynamic range.

The dynamic range of a filter is a number, usually given in decibels (e.g., 75 dB) that specifies the range of signal voltages the circuit can process without being corrupted by electrical noise at the low end and without causing nonlinear distortion at the high end because of clipping or slew-rate limiting in the amplifier. Clearly, dynamic range ought to be maximized as far as possible and a judicious choice of pole-zero pairing, gain-factor assignment, and section ordering can go a long way toward this goal. Arriving at the optimal choice is fairly complicated and so a computer algorithm is usually needed. If a suitable computer program is not available, a good suboptimal choice can frequently be obtained by

1. assigning the poles with the highest Q factor to the closest transmission zero,
2. choosing a low-pass or bandpass section as the first block and a high-pass or bandpass section as the last block in the cascade (this choice helps to keep high-frequency signals out of the filter, and prevents low-frequency noise from reaching the output),
3. ordering the remaining sections in the order of increasing values of Q , i.e., $Q_i < Q_{i+1}$, so that the voltages at the intermediate section outputs in the passband are as "flat" (that is, independent of frequency) as possible, and, finally,
4. assigning the gain constants K_i such that the voltage maxima at all section outputs are equal.

For critical filter specifications, requiring a transfer function of order 6 to 8 or higher, cascade realizations have been found to be too sensitive to element variations in some cases. That is, the filter response cannot be tuned correctly or maintained within specifications because of fabrication tolerances or later component

changes, such as those caused by aging or temperature drifts. Such sensitivity problems are particularly severe in filters with high- Q (narrow bandwidth or steep roll-off from passband to stopbands). For these cases, frequently, multiple-feedback topologies give useful solutions. For very critical and stringent filter specifications, active simulations of passive LC ladder filters show the best sensitivity performance and are preferred as alternative realizations.

Multiple-Feedback Topologies

Fig. 26 shows one of the many multiple-feedback topologies that have been proposed in the literature. It is recognized as essentially a cascade connection of l (second-order) sections $T(s)$, embedded into a resistive feedback network that ties all section outputs to an input summer. Because in this structure all l second-order sections are identical, it has been labeled the primary-resonator-block (PRB) topology.

The transfer function of the circuit in Fig. 26 can be shown, for an ideal summing amplifier, to be equal to

$$H(s) = \frac{V_i}{V_{in}} = - \frac{F_0 T^l(s)}{1 + \sum_{i=2}^l F_i T^i(s)} \quad (\text{Eq. 83})$$

where

$$F_i = R_j/R_{in}, \quad i = 0, 2, \dots, l.$$

Such multiple-feedback structures are most convenient for the realization of even-order symmetrical all-pole bandpass transfer functions, $H_{BP}(s)$ (Butterworth and Chebyshev filters; see Chapter 9). Such filters are most easily designed by transforming $H_{BP}(s)$, into a prototype low-pass transfer function, $H_{LP}(p)$, via the frequency transformation $p = \hat{Q}(s^2 + 1)/s$. Here p is the frequency parameter of the low-pass filter, normalized to the corner frequency of the low-pass filter's passband and $s = j\omega/\omega_0$ is the normalized frequency parameter of the bandpass function. ω_0 is the center

frequency of the bandpass and $\hat{Q} = \omega_0/\Delta\omega$, where $\Delta\omega$ is the desired bandwidth. The function

$$H_{LP}(p) = - \frac{K_0}{p^l + d_{l-1}p^{l-1} + \dots + d_1p + d_0} \quad (\text{Eq. 84})$$

then corresponds to the to-be-designed bandpass function

$$H_{BP}(s) = - \frac{Ks^{r/2}}{s^r + D_{r-1}s + \dots + D_1s + 1} \quad (r = 2l = \text{even}) \quad (\text{Eq. 85})$$

In Eq. 84, K_0 is the prescribed midband gain of the bandpass filter. In Eqs. 85 and 84, $r = 2l$ is the degree of desired bandpass and l is the degree of prototype low-pass filter. For this situation, the blocks T in the low-pass equivalent of the circuit in Fig. 26 are

$$T(p) = Mk/(p + k) \quad (\text{Eq. 86})$$

To design the circuit, one substitutes Eq. 86 into Eq. 83, brings the resulting equation into the form of Eq. 84, and compares coefficients with the prescribed function. If the definition

$$f_i = F_i M^i \quad (i = 2, \dots, l) \quad (\text{Eq. 87})$$

is used, this process yields the following solution:

$$k = d_{l-1}/l \quad (\text{Eq. 88a})$$

$$f_2 = d_{l-2}/k^2 - l(l-1)/2! \quad (\text{Eq. 88b})$$

$$f_i = \frac{d_{l-i}}{k^i} - \frac{1}{(l-i)!} \left[\frac{l!}{i!} - \sum_{j=2}^{i-1} f_j \frac{(l-j)!}{(i-j)!} \right] \quad (i = 3, \dots, l) \quad (\text{Eq. 88c})$$

Applying now the frequency transformation $p = \hat{Q}(s^2 + 1)/s$ to Eq. 86 results in $T(s)$, the blocks in Fig. 26, being second-order bandpass functions

$$T(s) = (M/Q)s/(s + s/Q + 1) \quad (\text{Eq. 89})$$

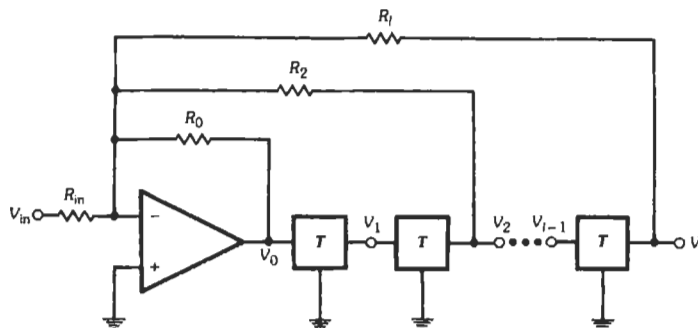


Fig. 26. The FLP topology for the special case of all identical sections $T(s)$.

where $Q = \hat{Q}/d_{i-1}$ is the quality factor of $T(s)$, i.e., ω_0 divided by the 3-dB bandwidth of T . Note that all blocks $T(s)$ are tuned to the same ω_0 because $s = j\omega/\omega_0$. The gain constant M of $T(s)$ should be selected for best dynamic range, such that all filter-internal signal maxima are equal. In general, this requires a computer routine and results in different M factors for each section; a good suboptimal choice, however, which at the same time retains all identical sections, is simply

$$M = \sqrt{1 + (Q/\hat{Q})^2} \quad (\text{Eq. 90})$$

With M and $f_i, i = 2, \dots, l$, known from Eqs. 88 and 90, the actual feedback factors $F_i = R_i/R_{in}$ can then be determined from Eq. 87. Finally, the PRB bandpass design is completed by setting

$$F_0 = R_0/R_{in} = (K_0/M^l)(Q/\hat{Q})^l \quad (\text{Eq. 91})$$

so that the prescribed gain constant K_0 is realized correctly.

An example will illustrate the process: To be designed is a sixth-order bandpass filter with a Butterworth magnitude characteristic, center frequency $f_0 = 4.8$ kHz, center-frequency gain $K_0 = 5$ (14 dB), and 3-dB bandwidth $\Delta f = 600$ Hz. For this case, the prototype third-order Butterworth low-pass transfer function, corresponding to Eq. 84, is

$$H_{LP}(p) = -5/(p^3 + 2p^2 + 2p + 1) \quad (\text{Eq. 92})$$

and the low-pass-to-bandpass transformation, with $\hat{Q} = 4.8 \text{ kHz}/0.6 \text{ kHz} = 8$, is

$$p = 8(s^2 + 1)/s \quad (\text{Eq. 93})$$

where $s = j\omega/(2\pi \times 4.8 \text{ kHz})$. Then, from Eqs. 87 and 92, $k = 2/3, f_2 = 1.5$, and $f_3 = 0.875$. From Eqs. 89 and 90, the second-order sections have a pole quality factor $Q = 8 \times 3/2 = 12$ and a gain constant $M = [1 + (3/2)^2]^{1/2} = 1.8028$; that is, the functions to be realized are

$$T(s) = 0.1502s/(s^2 + s/12 + 1) \quad (\text{Eq. 94})$$

Finally, the feedback resistor ratios F_i are, from Eq. 87, $F_2 = f_2/M^2 = 0.4615$ and $F_3 = f_3/M^3 = 0.1493$ and from Eq. 91, $F_0 = 2.880$.

The single-amplifier bandpass filter of Fig. 15A will be used to realize the function of Eq. 94. The design is completed as follows: From Eqs. 59, 56, and 57 with a choice of $q = 5$, one obtains $K = 1.01149$ and the pre-distorted values, if $f_i = 1$ MHz, $f_{0p} = 4.915$ kHz, and $Q_p = 11.75$. If $C = 12$ nF is chosen, Eq. 58 yields $R = 27.0$ k Ω . Finally, from Table 1, with $|b| = 0.1502\omega_{0p}$, $m = 0.0149 = 1/67.3$. The complete circuit realizing the prescribed bandpass filter is shown in Fig. 27. Note that an inverter is needed in the outer feedback loop because the bandpass sections have an inverting gain and the total loop gain must be negative to assure stability.

For other multiple-loop feedback topologies or for multiple-feedback structures for realizing bandpass fil-

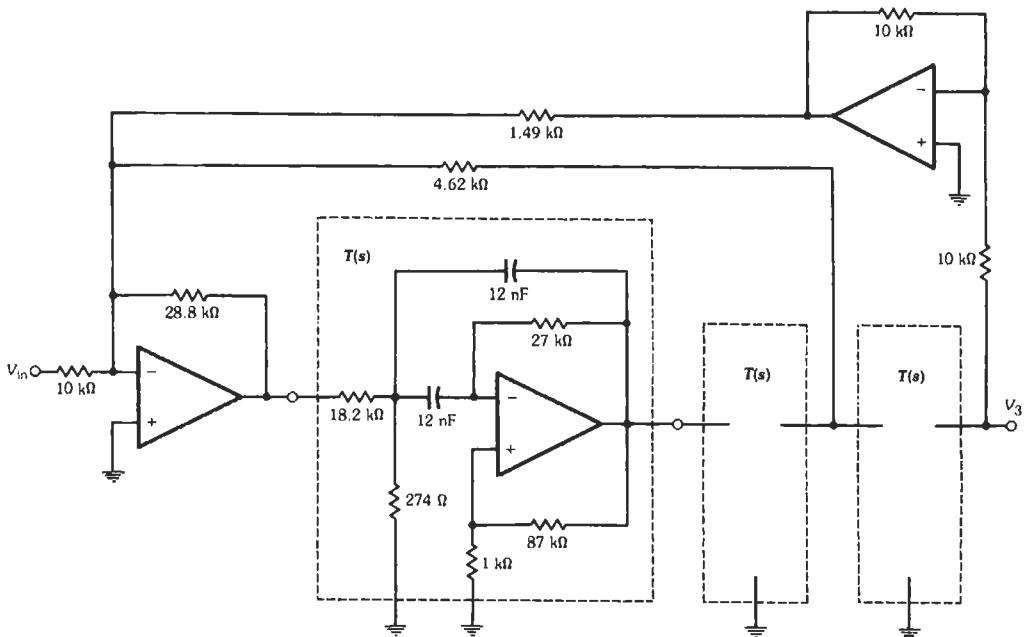


Fig. 27. Realization of the FLF (PRB) example.

ters with finite transmission zeros, the reader is referred to the literature.*

Ladder Simulation

Doubly terminated passive LC ladder filters are known to be optimally insensitive to element variations. Thus, in recent years, many methods have been proposed that simulate the operation, the topology, and/or the elements of LC ladders by active networks. The goal is to eliminate inductors and, at the same time, to retain the ladder's superior sensitivity performance. A few simple but powerful procedures are described in the following; for more details and for the design of very demanding filter circuits, the reader is referred to the literature.†

Operational Simulation—The process of simulating the operation of an LC ladder is best explained by an example from which the general design procedure should become clear.

Consider the simple LC low-pass filter in Fig. 28. Analysis yields the equations

$$V_o = I_3 / (sC_4 + G_L) \tag{Eq. 95a}$$

$$I_3 = (V_C - V_o) / sL_3 \tag{Eq. 95b}$$

$$V_C = (I_1 - I_3) / sC_2 \tag{Eq. 95c}$$

$$I_1 = (V_i - V_C) / (sL_1 + R_S) \tag{Eq. 95d}$$

where the currents and voltages are interpreted as signals to be represented in a signal-flow graph. Since all signals ought to be voltages, Eqs. 95 are then normalized by an arbitrary resistor R, and the notation $R_S/R = r_S$; $R_1/R = r_1$; $I_m R = V_m$; $L_m/R = I_m$, $m = 1, 3$; $C_k R = c_k$, $k = 2, 4$, is introduced. For example,

$$V_o = I_3 R / (sC_4 R + G_L R)$$

becomes

$$V_o = V_3 / (sc_4 + g_L)$$

Using this notation, and using judicious minus signs so that all capacitors become inverting integrators and inductors noninverting integrators, and allowing only summing, we can rewrite Eqs. 95 as follows:

$$V_o = -V_3 [-1 / (sc_4 + g_L)] \tag{Eq. 96a}$$

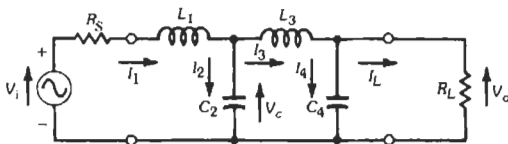


Fig. 28. LC low-pass ladder.

* References 1 and 6.

† References 2, 6, and 7.

$$-V_3 = [(-V_C) + V_o] / (sI_3) \tag{Eq. 96b}$$

$$-V_C = [V_1 + (-V_3)] / (-sC_2) \tag{Eq. 96c}$$

$$V_1 = [V_i + (-V_C)] / (sI_3 + r_S) \tag{Eq. 96d}$$

These equations are represented in the flow graph of Fig. 29, which demonstrates that the operation of the ladder in Fig. 28 is simulated by use of summers and inverting and noninverting integrators. These two functions are easy to realize with active circuits, as was shown earlier in this chapter. Note that all internal integrators are ideal, i.e., lossless, reflecting that the LC prototype ladder is lossless, but that the two end sections are realized via lossy integrators, thereby taking care of the resistive source and load terminations. Note further that each loop contains one inverting and one noninverting integrator (see Fig. 7), thereby assuring inverting loop gains for stability reasons. Stability will always be assured if all inductors are simulated by noninverting and all capacitors by inverting integrators. Lossy integrators are implemented by connecting a resistor in parallel with the integrator capacitor, and the summing operation indicated in Fig. 29 is realized by combining the summer of Fig. 5 and an integrator of Fig. 7 into one circuit. These operations are implemented in Fig. 30 and described by

$$V_o = -\frac{V_1}{sCR_1 + R_1/R_3} - \frac{V_2}{sCR_2 + R_2/R_3} \tag{Eq. 97}$$

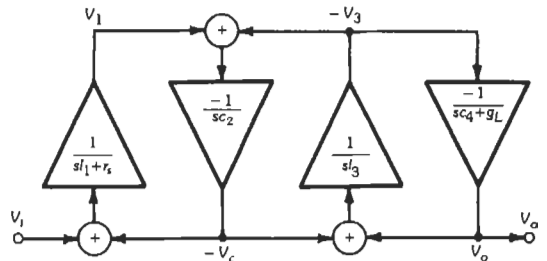


Fig. 29. Signal-flow diagram representing the LC filter shown in Fig. 28.

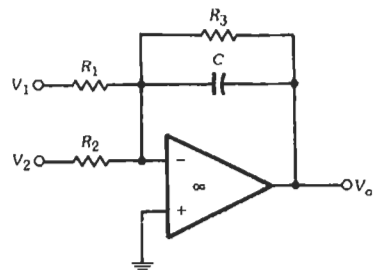


Fig. 30. A summing lossy integrator.

Clearly, if $R_3 = \infty$ the summing integrator is lossless. Of course, to obtain a noninverting summing integrator operation, one of the circuits of Fig. 7C or 7D is used instead of a simple Miller integrator.

The actual circuit implementing the LC ladder of Fig. 28 is shown in Fig. 31; the similarity of this diagram to the flow diagram in Fig. 29 should be apparent. To minimize phase errors caused by the finite

gain-bandwidth product ω_t of the op amps, each loop consists of a Miller integrator with $Q_1 < 0$ (see Eq. 36) and a phase-lead integrator with $Q_1 > 0$ (Fig. 7D). This is the preferred choice of integrator pairing. For convenience, all capacitors are selected to be equal. To determine resistors R_2 and R_9 in Fig. 31, the time constants of the corresponding integrators are set equal to the time constants of the input and output branches of the LC prototype in Fig. 28. Thus

$$CR_2 = L_1/R_S \quad (\text{Eq. 98a})$$

$$C_4R_L = CR_9 \quad (\text{Eq. 98b})$$

or, for some suitable choice of C ,

$$R_2 = L_1/(CR_S) \quad (\text{Eq. 99a})$$

$$R_9 = (C_4/C)R_L \quad (\text{Eq. 99b})$$

The "internal" resistors R_i , $i = 3, \dots, 8$, for Fig. 31 are obtained by equating the loop gains in the active circuit, i.e., the product of any two adjacent integrator time constants, to the product of the two corresponding reactances in the LC prototype. Thus

$$L_1C_2 = CR_3 \times CR_4 \quad (\text{Eq. 100a})$$

$$C_2L_3 = CR_5 \times CR_6 \quad (\text{Eq. 100b})$$

$$L_3C_4 = CR_7 \times CR_8 \quad (\text{Eq. 100c})$$

So, with the reactances of the LC ladder known and some suitable choice of C , as for Eqs. 99, one has, for example,

$$R_3 = R_4 = \sqrt{L_1C_2}/C \quad (\text{Eq. 101a})$$

$$R_5 = R_6 = \sqrt{L_3C_2}/C \quad (\text{Eq. 101b})$$

$$R_7 = R_8 = \sqrt{L_3C_4}/C \quad (\text{Eq. 101c})$$

Finally, resistor R_1 determines the desired dc gain of the filter. In the active network, this can be chosen freely, whereas in the passive prototype it is fixed at $R_1/(R_S + R_L)$. Let the dc gain be called K ; then it can be shown that

$$R_1 = (R_S/R)(R_D/K) \quad (\text{Eq. 102a})$$

if the LC low-pass ladder starts with a series inductor (as in Fig. 28) and

$$R_1 = (R/R_S)(R_D/K) \quad (\text{Eq. 102b})$$

if the LC ladder starts with a shunt capacitor. In Eqs. 102, R_D is the damping resistor of the first integrator (R_2 in Fig. 31), and R is the value of the normalizing resistor that was chosen in connection with Eqs. 96.

In a similar fashion, any type of LC filter can be simulated as an active RC filter by going through a signal-flow graph. For a detailed discussion see References 2, 6, and 7. If the desired filter function is an all-pole bandpass function obtained via the low-pass-to bandpass transformation from a prototype low-pass

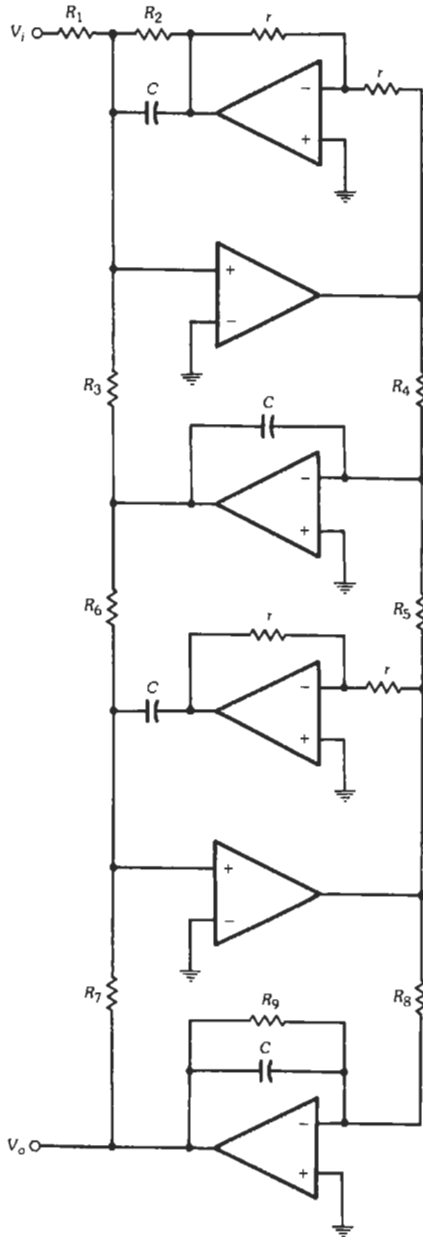


Fig. 31. Active simulation of the LC ladder of Fig. 28.

function, as illustrated in connection with the transformation from Eq. 84 to Eq. 85, the design is particularly simple. One only needs to simulate the low-pass filter by a signal-flow graph, as shown in Fig. 29 for a fourth-order case, and then use the low-pass-to-bandpass transformation to transform each integrator of the form $T_i(p) = a/(bp + c)$ into a second-order bandpass section

$$T_i(s) = \left[sa' / (b\hat{Q}) \right] / \left[s + sc / (b\hat{Q}) + 1 \right] \quad (\text{Eq. 103})$$

similar to the transformation from Eq. 86 to Eq. 89 above. For ideal integrators, the coefficient $c = 0$, implying a bandpass of infinite Q^* for all internal sections of the simulated ladder (compare Fig. 29). After a suitable realization for the second-order sections is obtained as discussed earlier, these are then interconnected in a "leapfrog" topology as shown in Fig. 32, where the resistors of the feedback network are adjusted for the correct loop gains. The "leapfrog" method owes its name to the resulting circuit structure. For details of the design of these very well behaved active ladder simulations, References 1, 3, 6, or 7 should be consulted.

Ladders Using Inductor Simulation—It was pointed out earlier that high-quality grounded inductors can be simulated by the use of general impedance converters implemented with two op amps. Fig. 9 shows the op amp circuit realizing, for $R_2 = R_3$ and $R_5 = 1/(\omega_c C_4)$, an inductor of value $L = R_1/\omega_c$, where ω_c is a frequency in or at the edge of the filter passband (Eq. 40). It stands to reason, therefore, that an LC ladder can be simulated simply by replacing all inductors by the circuit in Fig. 9, provided the inductors are grounded. An illustration of this rather straightforward procedure for grounded inductors is given in Fig. 33. Fig. 33A shows an LC ladder, which may have been obtained from tables and/or filter design handbooks. The filter has only grounded inductors and a finite transmission zero at the frequency $\omega_z = \sqrt{L_2 C_2}$ where C_2 and L_2 resonate.

If floating inductors need to be simulated, as is true for all low-pass and most bandpass filters, a different approach must be used. The strategy, suggested by Gorski-Popiel† is obtained by a different interpretation of the general impedance converter in Fig. 9: If R_5 is

considered the load of the GIC (Fig. 9B), then the GIC with the conversion time constant τ converts a load resistor R_5 into the inductor $L = \tau R_5$ (Eq. 40). Generalizing this idea, we can prove the following fact:

If GICs with the conversion constant τ are inserted between the outside world and all terminals (except ground) of a purely resistive network containing the resistors R_i , (the resistive circuit forms the load of the GICs) then the resulting circuit will appear as an inductive network of the same topology containing the inductors $L_i = \tau R_i$. Conversely, an inductor network with components L_i is replaced by a resistive network with components L_i/τ embedded between GICs with time constant τ .

Fig. 34 shows this method, of which Fig. 9 is the simplest special case. Evidently, in Fig. 34A a floating inductor is simulated and in Fig. 34B an inductive T .

The design proceeds as follows: In the first step the inductive subnetwork(s) are cut from the filter and the cuts replaced by properly connected GICs with conversion constant τ . If then all inductors L_i are replaced by resistors $R_i = L_i/\tau$, the resulting filter behaves like the original LC circuit, bearing in mind, of course, the frequency limitations (not much beyond upper audio range) imposed by the op amps used in the design. Fig. 35 shows Gorski-Popiel's method on a simple low-pass filter.

In the active realization of ladders, often a network transformation can be used to eliminate all floating inductors that would use many op amps for their simulation. A procedure that works particularly well for low-pass filters makes use of frequency-dependent negative resistors (FDNRs).

Ladders Using FDNR Circuits—If the normalized impedances of a circuit are scaled (multiplied) by an impedance scaling factor $1/s$, the result is as follows: A resistor of value $R_i \Omega$ becomes a capacitor of value $(1/R_i) F$, an inductor of $L_i H$ becomes a resistor of $L_i \Omega$, and a capacitor of $C_i F$ is converted into a frequency-dependent negative resistor of value $D_i = C_i$. This was introduced in the section "The FDNR Element." Thus, an impedance $Z = R + sL + 1/(sC)$ is transformed into

$$Z_i = \frac{Z}{s} = \frac{R}{s} + L + \frac{1}{s^2 C} \Rightarrow \frac{1}{sC_i} + R_i + \frac{1}{s^2 D_i} \quad (\text{Eq. 104})$$

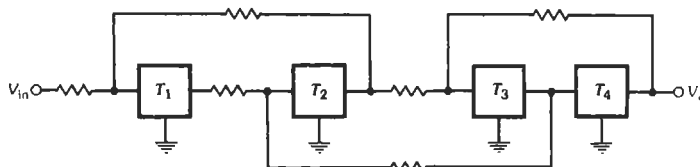


Fig. 32. Leapfrog bandpass filter.

* In practice, a very large but finite Q is sufficient.

† See References 2, 6, or 7.

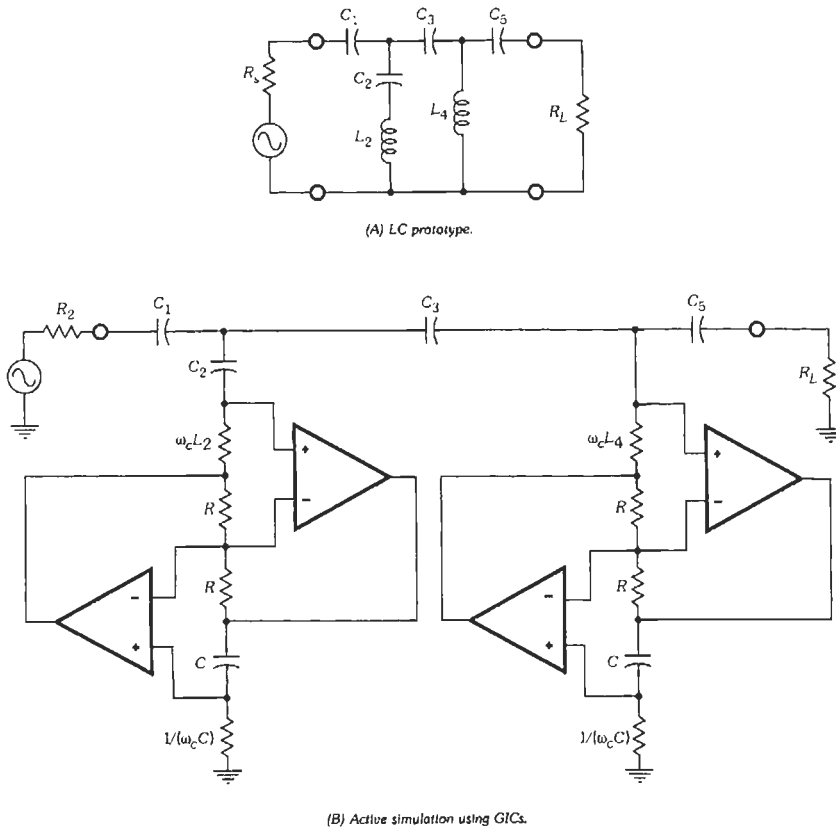


Fig. 33. Simulation of an LC ladder.

The voltage transfer function, being a dimensionless quantity, is of course unaffected by impedance scaling so that the transformed circuit has the same response as the original filter. Hence, according to Eq. 104, inductors are eliminated, and FDNR elements are introduced that can readily be realized as discussed earlier in connection with Fig. 11 (Eq. 45). One limitation is that a good floating FDNR is very difficult to realize. Thus, this approach is most useful for passive LC prototype filters that have no floating capacitors. In this sense, the FDNR and inductor simulation methods are complementary.

A brief example will illustrate the procedure: Fig. 36A shows an LC low-pass filter with finite transmission zeros at $\omega_1 = (L_2C_2)^{-1/2}$ and $\omega_4 = (L_4C_4)^{-1/2}$. Impedance scaling by a factor $1/s$ transforms the circuit into that of Fig. 36B, which shows also the commonly used symbol for FDNR elements. For further design details concerning problems with biasing, buffering input and output, and the selection of element values, the reader is referred to the literature.*

Ladders Using OTAs—It was mentioned before that for active filters at higher frequencies—say above a few hundred kilohertz—op-amp circuits are gener-

ally not useful because of the limited bandwidth of op amps. In contrast, the bandwidth of transconductors extends, depending on the technology, to several hundred megahertz or even gigahertz.† In addition, as will become apparent shortly, the design of ladder simulations with OTAs is fairly simple and systematic and lends itself particularly well to high-frequency integrated filters. As discussed earlier in this section, ladders can be implemented by simulating the inductors in the structure. This can be done for grounded as well as floating inductors via g_m -based gyrators (Eq. 44, Fig. 10). Alternatively, the operation of each branch or mesh of the ladder can be simulated by the signal-flow graph method, using g_m -based integrators (Eq. 37, Fig. 8). Both methods will be illustrated in the following. For further details, the reader is referred to References 6, 7, 11–14, or 20.

Inductor Simulation: The process is illustrated best with a simple example. Consider the third-order elliptic (see Chapter 9) low-pass filter in Fig. 37A. To

* References 2, 6, or 7.

† References 11 and 19.

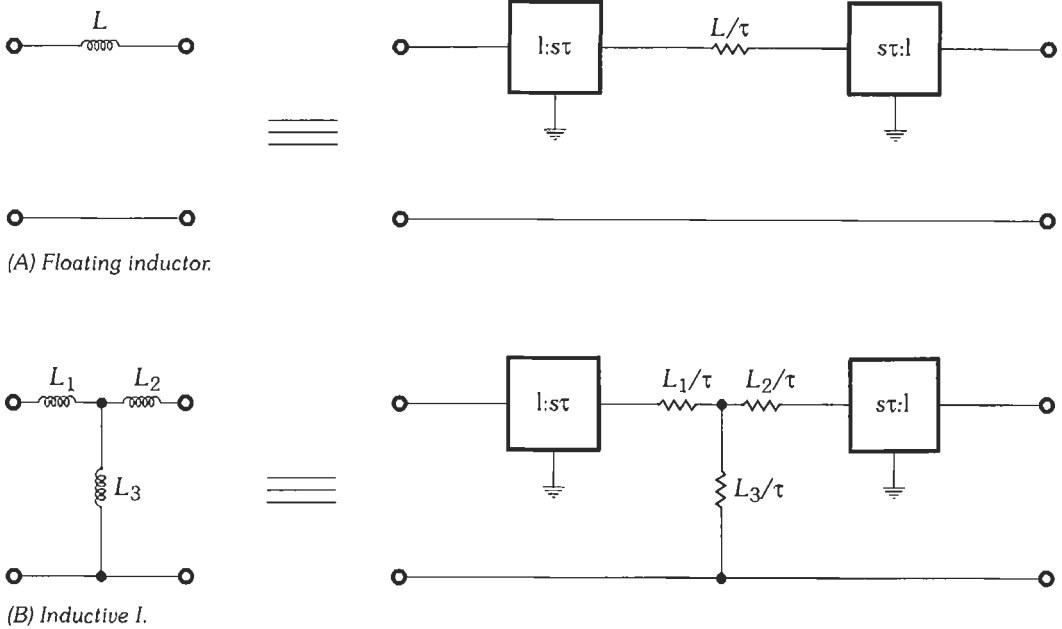


Fig. 34. Simulation of inductor networks via GICs.

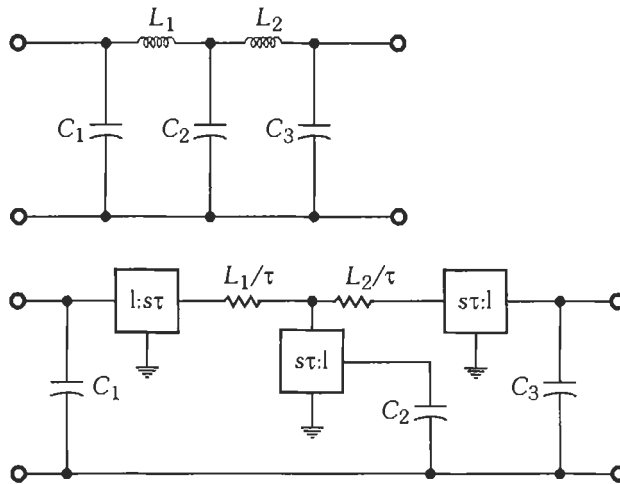


Fig. 35. Simulation of an LC ladder filter.

implement this circuit as an active fully differential integrated filter, the floating inductor must be simulated by the circuit in Fig. 10C. The result is shown in Fig. 37B. OTA 1 converts the input voltage into a current (source transformation); OTAs 2 and 7 simulate the source and load resistors; and OTAs 3 through 6 with the capacitor $C_1/2$ form the floating inductor. Also shown in Fig. 37C is a version that merges

OTAs for savings in power and components as was discussed in connection with Figs. 6B and 6C. OTAs 2 and 3 are merged into OTA *b*, OTAs 4 and 5 into OTA *c*, and OTAs 6 and 7 into OTA *d*. To be able to use identical OTAs, OTA *a* has its inputs connected as shown, thereby doubling the transconductance value to $2g_m$. This, in turn, eliminates the inherent 6-dB loss of the LC prototype in Fig. 37A. (Note that at dc,

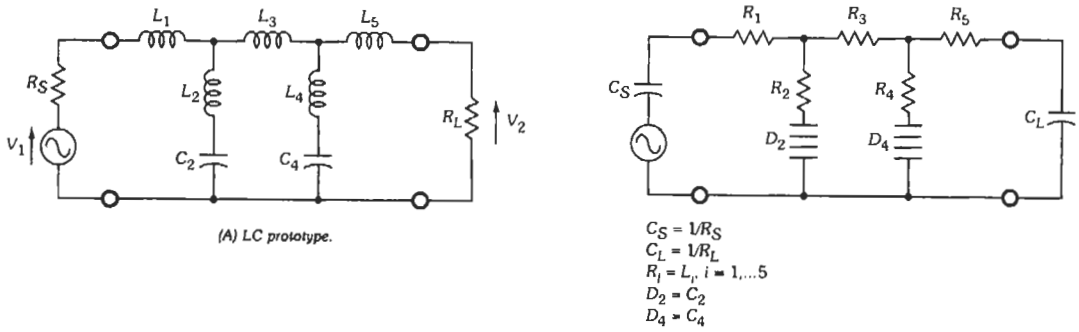
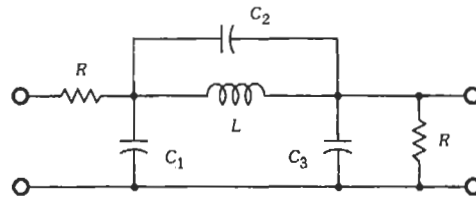
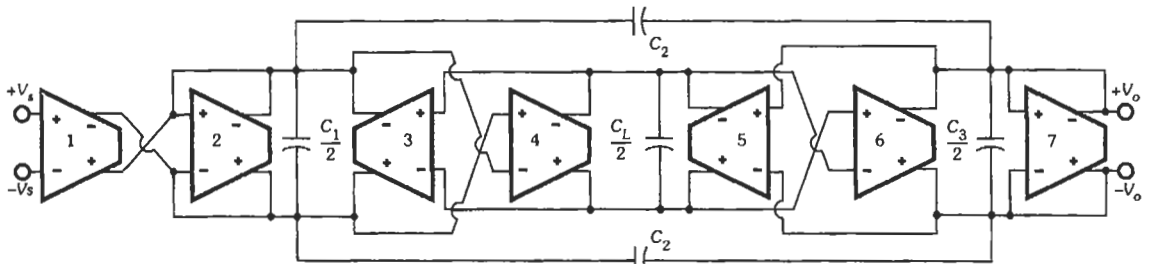


Fig. 36. FDNR ladder implementation.

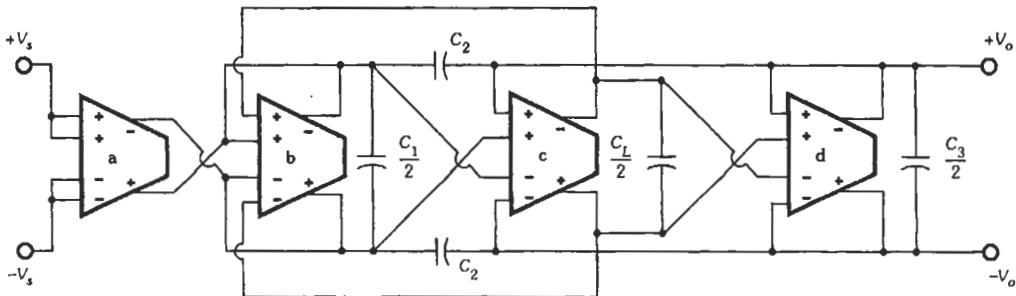
(B) Transformed circuit.



(A) Third-order elliptic LC low-pass filter.



(B) Differential OTA-C simulation.



(C) Differential OTA-C simulation with merged OTAs.

Fig. 37. Low-pass filter using inductor simulation.

the gain of the LC filter is 0.5 because of the resistive voltage divider.) Extensions of this technique to ladders of higher order, such as the one in Fig. 38A, are obvious.

Signal-Flow Graph (SFG) Simulation: In deriving the active circuit to perform an operational or SFG simulation of the passive ladder, the node and mesh equations are written such that each capacitor and

inductor is recognized as performing a signal integration, $1/(sC_i)$ or $1/(sL_i)$. As was done before, a source transformation is first performed on the ladder input so that the source resistor R is in parallel with C_1 , and both the first and the last branch form lossy integrators. Referring to Fig. 38A and denoting $1/R$ by G , one obtains the equations for the five state variables, analogous to Eqs. 95 above:

$$V_1 = \frac{1}{sC_1 + G} [(V_i G - I_{L2}) - sC_2(V_1 - V_3)] \quad (\text{Eq. 105a})$$

$$I_{L2} = [1/(sL_2)](V_1 - V_3) \quad (\text{Eq. 105b})$$

$$V_3 = \frac{1}{sC_3} [(I_{L2} - I_{L4}) + sC_2(V_1 - V_3) - sC_4(V_3 - V_5)] \quad (\text{Eq. 105c})$$

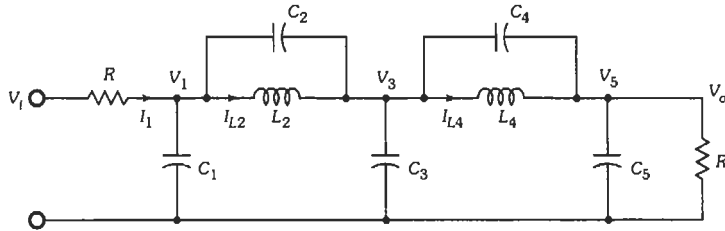
$$I_{L4} = [1/(sL_4)](V_3 - V_5) \quad (\text{Eq. 105d})$$

$$V_5 = \frac{1}{sC_5 + G} [I_{L4} + sC_4(V_3 - V_5)] \quad (\text{Eq. 105e})$$

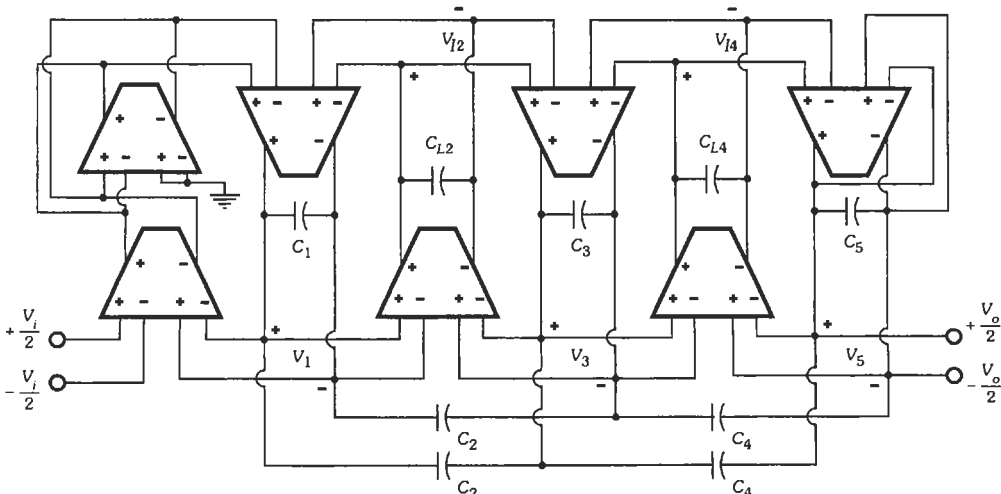
The equations are then scaled by a dimensionless factor to assure that all signal maxima are of equal magnitude. As discussed in detail in, for example, References 2, 6, 21, and 22, this step results in a realization with maximum dynamic range. Finally, as explained earlier, a suitable normalization, here illustrated for Eq. 105b,

$$V_{i2} = \frac{I_{L2}}{g_m} = \frac{1}{sL_2 g_m^2} g_m (V_1 - V_3) \quad (\text{Eq. 106})$$

leaves the currents as voltage signal variables ($I/g_m = V$) and converts the inductors to capacitors ($Lg_m^2 = C_L$); all integrators can then be implemented via transconductors as indicated in Fig. 8. Note that the floating capacitors C_2 and C_4 are not realized via integrators because the currents through these two elements (see Eqs. 105a, 105c, and 105e) can be obtained by connecting C_2 and C_4 directly to the nodes V_1, V_3 , and V_3, V_5 , respectively. At the cost of floating capacitors and the resulting parasitic substrate capacitances, this process saves several transconductors, silicon area, and



(A) LC prototype.



(B) Fully differential g_m -C ladder simulation.

Fig. 38. Fifth-order elliptic low-pass filter using signal-flow graph simulation.

power and reduces noise. Making use of multiple-input transconductors to achieve a fully differential circuit, the five equations (Eqs. 105) are implemented step by step, and interconnected as indicated in Fig. 38B. The normalized signal voltages, representing the state variables V_{12} , V_{14} , V_1 , V_3 , and V_5 , are marked in Fig. 38B for easy reference.

Observe that the ladder simulation technique leads to a very efficient and systematic active realization of the passive prototype. The fifth-order active ladder requires five two-input summing transconductors, plus two for simulating the source and load resistors, and five capacitors, plus two for each floating C in the passive prototype. All transconductors are identical. Again, it should be apparent how the method is extended to filters of higher order.

SWITCHED-CAPACITOR FILTERS

Metal-oxide-semiconductor (MOS) integrated circuit (IC) technology is used widely in industry because of its superior logic density and lower power consumption compared to that achievable with other IC technologies. With very large scale integration (VLSI), hundreds of thousands to millions of MOS transistors can be placed on a single chip. The ubiquitous presence of CMOS chips in appliances, entertainment electronics, and personal computers indicates the economic and social impact of MOS VLSI on modern society.

Switched-capacitor (SC) techniques, in addition to the g_m - C circuits discussed earlier, provide analog signal processing capability that is needed in MOS technology for "mixed-mode" analog/digital signal processing requirements. The marriage of analog SC functions and high-density digital logic on the same piece of silicon extends to mixed-mode systems the same cost and space savings associated with memories and microcomputers. One advantage of properly designed SC filters over g_m - C implementations is that the realized characteristics usually require no trimming and are inherently stable over process and environmental variations. However, because of their sampled-data nature, at the time of this writing, SC filters are commercially useful only at relatively moderate frequencies (from the audio range up to possibly a few hundred kilohertz), whereas continuous-time (c-t) g_m - C circuits can be made to operate commercially at several tens and even hundreds of megahertz. Consequently, these two integrated analog filter design techniques rarely compete in the same applications.

SC filters* consisting of MOS capacitors, switches, and op amps realize infinite impulse response (IIR) analog sampled-data filters, similar topologically to the active RC filters described in previous sections. Narrow and flat passbands can be realized efficiently.

Unlike their c-t active RC and g_m - C counterparts, SC filters are sampled-data systems. This complicates their use and design, but SC filters do indeed take full advantage of the inherent precision achieved by MOS technology. As shown later, the transfer functions are completely determined by precise crystal-controlled clocks and capacitor ratios. Capacitor ratios can be held to an accuracy of about 0.3 percent or less and, with appropriate parasitics insensitive circuit techniques,† capacitors as small as 0.2 pF can be used. Furthermore, MOS capacitors are nearly ideal, with very low dissipation factors and good temperature stability. These properties can be achieved with either NMOS or CMOS processing, but CMOS, with its added flexibility for realizing high-gain, low-noise op amps and low power dissipation, is recognized as the technology of choice.

Sampled-Data Filter Systems

Consider a sampled-data filter system that is suitable for use in a c-t analog (i.e., analog input/analog output) environment. This represents, from a hardware point of view, the most demanding environment for any sampled-data filter, and specifically for a SC filter. The system given in Fig. 39 shows the analog signal being passed through a c-t antialiasing low-pass filter; an input sample-and-hold circuit, (S/H)_i, which samples the band-limited analog input at intervals of $1/f_{s1}$; the SC filter; an output sample-and-hold circuit, (S/H)_o, which resamples the output of the SC filter at intervals $1/f_{s2}$; and a final c-t low-pass reconstruction filter, which smoothes the sharp transitions in the sampled-data output waveform. The SC filter is shown, in general, to be controlled by clocks of multiple frequencies (f_{s2} through $f_{s,n-1}$). Because the capacitor ratios (hence the silicon area) required to realize a given transfer function scale in proportion to the ratio of the clock frequency to the pole and zero frequencies (see Eq. 115), silicon area can be minimized by employing multiple clocks.‡ However, for simplicity and to minimize aliasing of out-of-band signals (including system noise and power-supply feedthrough), SC filter systems are often operated with a single clock ($f_{s1} = f_{s2} = \dots = f_{sn} = f_c$). f_c is typically in the range of 100 to 256 kHz, much higher than the highest signal frequency. Sampling rates (input and output) in this range serve to reduce the complexity of the c-t antialiasing and reconstruction filters. It should be noted that in many SC filters, particularly low-pass filters, the sample-and-hold operations shown in Fig. 39 are inherent to the SC filter. In these cases, the SC filter includes the three blocks enclosed by the dashed box in Fig. 39. Typical frequency responses for the various blocks (shown for a low-pass SC filter) in Fig. 39 are depicted below each system block, assuming the typical case

* References 1, 6, 7, 10, 24, and 25.

† References 1, 6, 7, 10, and 25.

‡ Reference 1.

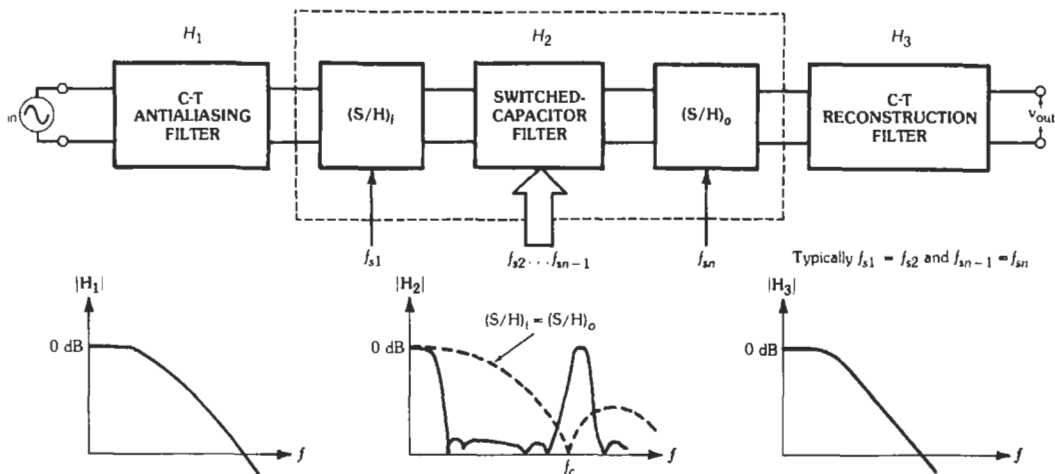


Fig. 39. Sampled-data filter system for continuous-time analog input and smooth analog output.

that all clock frequencies are the same $f_{s_i} = f_s, i = 1, \dots, n$.

It is noted that when an input or output is interfaced with digital or sampled-data circuitry, such as D/A and A/D converters, some of this circuitry is no longer needed. For example, when the output is to be interfaced to another sampled-data function, the reconstruction filter is not required, and a sample-and-hold circuit is typically included at the input of this function. Although the need for c-t filtering is reduced, interfacing to digital and sampled-data circuits requires careful synchronization between the clocks that control the SC filter and those that control the external sampling operations. This is accomplished by passing synchronization pulses between the SC network and the external samplers. One reason for synchronization is to ensure that the SC network output is sampled after all transients have settled and the output is steady.

Another consequence of the sampled-data character of SC networks is the mathematical convenience of the z transform in analyzing or specifying SC filters. The z transform, where $z = \exp(s\tau)$ and τ is the clock period, is covered in Chapter 28 (or refer to References 1, 6-8, or 10) as a mathematical tool for digital filters and discrete-time systems. The sampling instants, the instants at which switches open, are the times at which capacitor charges are updated. These discrete instants of time are the most important times in the operation of the filter. In fact, one can completely describe the behavior of the filter by considering operation at only these discrete instants. This aspect of SC filter behavior is analogous to the operation of digital filters; hence, the mathematical analysis follows in a similar manner.

It has been mentioned that the input and output of an analog sampled-data system are analog signals. Note that the antialiasing filter serves to band-limit the input spectrum to $f_c < f_s/2$ so that the signal can be recon-

structed without error. Here, f_c refers to the highest component allowed in the input spectrum, and $f_s = 1/\tau$ is the sampling frequency. The antialiasing filter will also serve to band-limit high-frequency input noise, which would otherwise be aliased back into the baseband. Inherent in all analog sampled-data systems is the means to provide some sort of analog reconstruction. With $\tau = 1/f_s$, the simplest form of reconstruction is the zero-order hold or sample-and-hold (S/H) circuit. The impulse response of the zero-order hold is described by

$$h_0(t) = \begin{cases} 1/\tau & \text{for } 0 \leq t \leq \tau \\ 0 & \text{elsewhere} \end{cases} \quad (\text{Eq. 107})$$

The sample-and-hold impulse response is sketched in Fig. 40, and an S/H reconstructed signal $x_r(t)$ is shown in Fig. 41. The transfer function for this sample-and-hold is

$$H_0(s) = (1 - e^{-s\tau})/(s\tau) \quad (\text{Eq. 108})$$

and its frequency response is

$$H_0(j\omega) = \frac{1 - e^{j\omega\tau}}{j\omega\tau} = e^{-j\omega\tau/2} \frac{\sin \omega\tau/2}{\omega\tau/2} \quad (\text{Eq. 109})$$

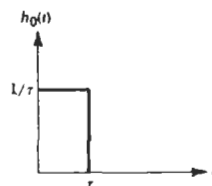


Fig. 40. Impulse response for sample-and-hold.

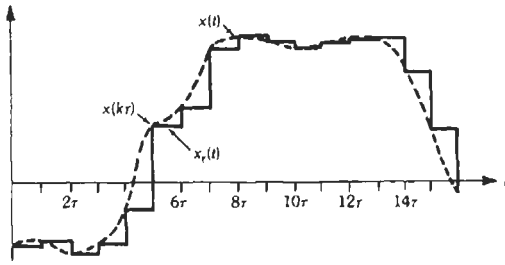


Fig. 41. Reconstruction with sample-and-hold.

It can be seen that $H_0(j\omega)$ has the $(\sin x)/x$ magnitude response shown in Fig. 42A and the phase characteristic shown in Fig. 42B.

Examine the effect on the spectrum of a sampled signal $X^{\#}(j\omega)$, as predicted by z -transform analysis, when the sampled signal is reconstructed with a sample-and-hold. The spectrum of the reconstructed signal $X_r(j\omega)$ is simply the product of the sampled signal spectrum $X^{\#}(j\omega)$ and the sample-and-hold spectrum $H_0(j\omega)$ as shown in Fig. 43. Note that the high-frequency content of the sampled signal outside the baseband has been substantially reduced and the baseband spectrum is slightly altered near the band edge. When $f_s/f_c \gg 1$, as is usually the case in SC filters, this baseband distortion is negligible and

$$X_r(j\omega) = X^{\#}(j\omega) \text{ for } \omega < \omega_s/2 = \pi f_s \quad (\text{Eq. 110})$$

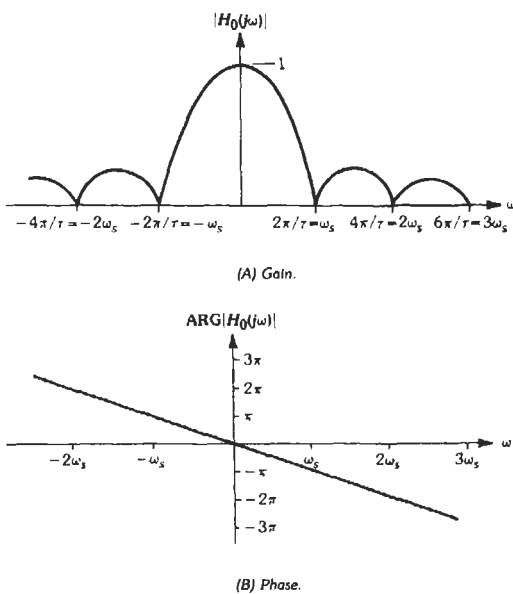


Fig. 42. Sample-and-hold responses.

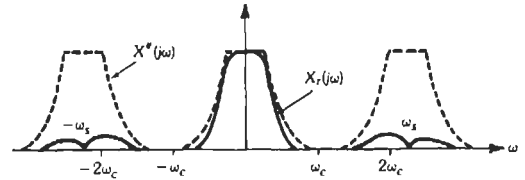


Fig. 43. Spectrum of reconstructed sampled signal.

For lower sampling rates, the sampled-data filter that operates on the signal X is designed with a peak to compensate for the S/H-related band-edge droop. To further smooth $X_r(j\omega)$ (i.e., further attenuate the high-frequency components of X_r) the analog sampled-data system can be followed by a c-t low-pass filter as shown in Fig. 39.

Now that the concept of analog sampled-data systems has been discussed in general terms and the necessary hardware to support the sampled system in a non-band-limited analog environment has been identified, attention will now be focused on the SC-filter portion of this system.

The Operation of Ideal SC Filters

Consider the operation of an ideal SC filter, made up of ideal capacitors, ideal switches, and ideal op amps (infinite gain and infinite bandwidth). It is noted that MOS op amps that settle to within 0.1 percent of final value in less than $2 \mu\text{s}$ and achieve dc gains of greater than 60 dB have been designed.* One subtle difference between SC filters and their active RC counterparts is the effect of op-amp dynamics on filter behavior. Although related, op-amp settling time, rather than op-amp frequency response, is the important op-amp characteristic for SC filters. In fact, it is op-amp settling time that establishes the upper limit for the sampling rate. It has been demonstrated that settling errors can be minimized by carefully orchestrating the timing among the various samples within the filter.† Hence, for sampling rates f_s of no more than 256 kHz and careful timing, a practical MOS op amp usually can be modeled as a voltage-controlled voltage source with gain equal to the dc gain of the op amp.

Consider for the moment the simple SC circuit in Fig. 44A. It consists of a single capacitor with bottom plate grounded and top plate connected between two switches controlled by the clocks ϕ^a and ϕ^b . As shown in Fig. 44B, the clocks ϕ^a and ϕ^b are two-phase, non-overlapping clocks of frequency $f_s = 1/(2T) = 1/\tau$. The clock's duty cycles must be less than 50 percent for the clock pulses to be nonoverlapping. However, to allow maximum time for op-amp settling, duty cycles of

* Reference 24.

† Reference 25.

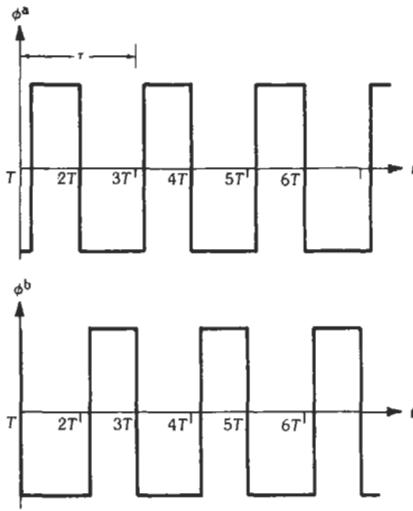
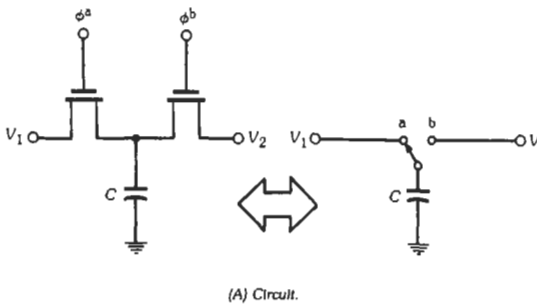


Fig. 44. Simple SC circuit.

greater than 35 percent are typically used. Most SC filters are insensitive to duty cycle, but nonoverlap is essential for the circuit to operate correctly. The fact that op-amp settling typically establishes the upper limit for f_s has already been mentioned. It was also mentioned in the previous section that high f_s results in large capacitor ratios and low f_s results in complex antialiasing/reconstruction filters. Other factors that impact the selection of a sampling rate are noise and capacitor leakage. Evidently the choice of sampling rate involves a compromise among several conflicting factors and no general rule can be stated. One can find in the literature* SC filters with sampling frequencies that range from many megahertz to as low as 8 kHz.

The circuit in Fig. 44A, although simple, provides a rather interesting and important function. Consider the following simple analysis of this circuit and assume initially that the capacitor C is uncharged. When switch ϕ^a closes and switch ϕ^b opens, at $t = 0$, capaci-

tor C charges to the voltage V_1 . When $t = T = \pi/2$, switch ϕ^a opens and switch ϕ^b closes, and capacitor C discharges to the voltage V_2 . The amount of charge that flows into (or out of) V_2 during the time τ is given by

$$\Delta Q = C(V_2 - V_1) \quad (\text{Eq. 111})$$

If the process repeats at $t = \tau, 2\tau, \dots$, etc., the equivalent average current I is

$$I = \frac{\Delta Q}{\tau} = \frac{C}{\tau}(V_2 - V_1) = Cf_s(V_2 - V_1) \quad (\text{Eq. 112})$$

Equation 112 is recognized to be the instantaneous voltage-current relation for a series resistor of value

$$R = 1/(Cf_s) \quad (\text{Eq. 113})$$

Equations 111–113 establish the well-known switched-capacitor-resistor equivalence. An oversimplified application of Eq. 113 is to derive SC filters from active RC prototypes by replacing each resistor with an equivalent switched capacitor (with $C = 1/(Rf_s)$). However, a careful analysis of Fig. 44A reveals that Eq. 113 is an approximation in that the charge transfer described in Eq. 111 involves a delay of T seconds. The approximation improves as T becomes small (f_s becomes large). Hence, the SC filter synthesis method of replacing resistors in an active RC prototype with switched capacitors is not a recommended procedure. In fact, if one were to replace the resistor in the feedback loop of a high- Q active RC biquad with the SC circuit in Fig. 44A, the added delay would no doubt result in instability. This result can be proven rigorously.† Fortunately, all these pitfalls can be avoided by working directly in the z -domain. Such an analysis renders an exact characterization of the ideal SC filter at the sampling instants.

By loading the SC circuit in Fig. 44A with a grounded capacitor, a simple one-pole low-pass filter can be constructed. This circuit is shown in Fig. 45. If we resort again for the purpose of demonstration to the resistor equivalence in Eq. 113, the pole frequency of the equivalent RC filter becomes

$$\omega_p = 1/(R_1C_2) \quad (\text{Eq. 114})$$

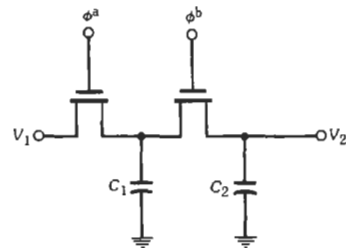


Fig. 45. One-pole low-pass SC filter.

* References 1 and 24.

† Reference 1.

Substituting the SC equivalent for resistor R_1 in Eq. 114 yields the pole frequency for the SC filter, i.e.,

$$\omega_p \approx f_s(C_1/C_2) \text{ for } f_s \gg 2\pi f_p \quad (\text{Eq. 115})$$

Equation 115 provides much of the motivation for SC filters. From Eq. 115, the inherent accuracy of the SC implementation is seen immediately: ω_p no longer depends on the *product* of the components R and C , but on the clock, f_s , and the *ratio* of two C s. It was mentioned previously that MOS capacitor ratios can be held to tight tolerances and f_s is derived from a very precise crystal-controlled master clock. In addition, note that the capacitor ratio C_2/C_1 equals f_s/ω_p . That is, for constant ω_p , C_2/C_1 scales in direct proportion to f_s . As a consequence, the silicon area required to realize the SC filter in Fig. 44 (for a fixed ω_p and minimum capacitance C_1) scales with f_s . For IC realizations, minimum capacitors on the order of 1 pF are typically used.

SC Integrators

Most practical SC filters are realized with one op amp per pole (and zero), and they use SC integrators as basic building blocks. Moreover, SC filters implemented in this manner can be made insensitive to the unavoidable parasitic capacitances that occur in an IC realization and can, therefore, be implemented with minimum-size capacitors.

Let us first consider the inverting SC integrator in Fig. 46. Also shown are the parasitic capacitances associated with each node of the circuit. Capacitors (C_1 and C_2 in Fig. 46) are typically implemented as a thin SiO_2 dielectric sandwiched between two polysilicon plates. The interconnecting lines are typically polysilicon and metal that lie on top of a thick layer of SiO_2 (referred to as field oxide). The switches are usually either polysilicon-gate CMOS or enhancement NMOS devices. In the process of constructing these features in an integrated circuit, parasitic capacitors to the reverse-biased substrate (= signal ground) are unavoidably created. The parasitic capacitances are not controllable, do not track the intended capacitors, and sometimes are nonlinear. Hence, it is important to significantly reduce or eliminate the effect of parasitic

capacitances on the transfer characteristic of the SC filter.

Returning to the integrator in Fig. 46, consider the effect of each of the parasitic capacitances, C_{p1} through C_{p4} . First, the parasitics C_{p1} and C_{p4} shunt low-impedance voltage sources (the output of an op amp and the input voltage possibly derived from another op amp as indicated in the figure). C_{p3} shunts the "virtual ground" of a high-gain op amp and is never charged. These parasitics have no effect on circuit performance. Parasitic C_{p2} , in contrast, parallels C_1 and adds directly to its value. In the absence of parasitics, the z-domain transfer function for the inverting SC integrator in Fig. 45 can be shown to be*

$$H(z) = -(C_1/C_2)z^{-1/2}/(1-z^{-1}) \quad (\text{Eq. 116})$$

If C_{p2} is included, the transfer function becomes

$$H(z) = -\frac{C_1}{C_2} \left(1 + \frac{C_{p2}}{C_1} \right) \frac{z^{-1/2}}{1-z^{-1}} \quad (\text{Eq. 117})$$

From Eq. 117, the parasitic capacitance is seen to cause a gain error factor of $1 + C_{p2}/C_1$.

For the intended use of a minimum-size C_1 this error is large and an alternative circuit configuration is necessary. Consider the SC inverting integrator shown in Fig. 47, which also shows the two important parasitic capacitances. All other parasitic capacitances (shunting either a voltage source or the op-amp virtual ground) have been shown to have no effect. To determine the effect of C_{p1} and C_{p2} , examine the circuit as switches ϕ^a and ϕ^b turn on and off. With ϕ^a on and ϕ^b off, parasitics C_{p1} and C_{p2} shunt a voltage source and the op-amp virtual ground, respectively. When ϕ^a turns off and ϕ^b turns on, both C_{p1} and C_{p2} are shorted harmlessly to ground and do not contribute any charges to the circuit operation. Hence, this SC inverting integrator operates without error. The realization in Fig. 47 is said to be *parasitic insensitive*. The transfer function for this integrator is†

$$V_{out}/V_{in} = -(C_1/C_2)z^{-1/2}/(1-z^{-1}) \quad (\text{Eq. 118})$$

* References 1, 6-8, and 10.

† References 1, 6, 7, and 10.

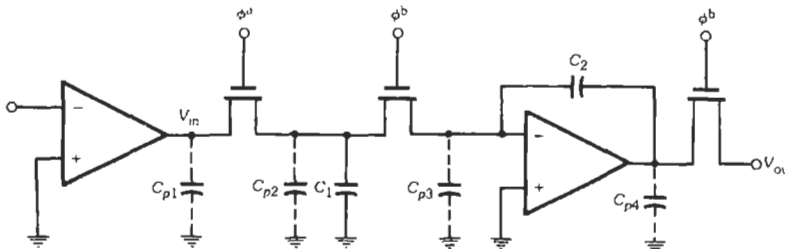


Fig. 46. Parasitic-sensitive SC inverting integrator with associated parasitic capacitors.

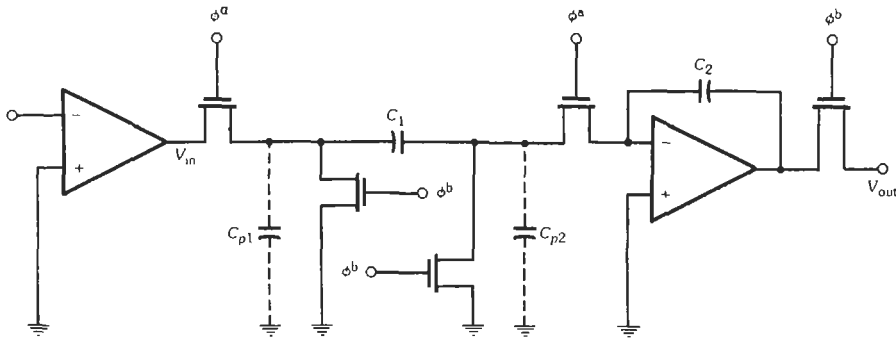


Fig. 47. Parasitic-insensitive SC inverting integrator with associated parasitic capacitances.

Note the presence of the $z^{-1/2}$ (one-half clock period delay) term in Eqs. 116–118. In the digital and SC filter literature, integrators of the form $Kz^{-1/2}/(1 - z^{-1})$ are known as lossless discrete integrators (LDI).

Noninverting integrators can be implemented very efficiently with SC techniques. As shown in Fig. 48, an SC noninverting integrator can be realized with a single op amp. As seen in Figs. 7C and 7D, in the active RC world noninverting integrators require two op amps. With the switching arrangement shown for switched capacitor C_1 , the additional inversion is elegantly performed by C_1 with a simple change of switch phasing. It can be verified that this circuit, like that in Fig. 47, is parasitic insensitive. The transfer function for the noninverting integrator in Fig. 48 is

$$H(z) = (C_1/C_2) z^{-1/2}/(1 - z^{-1}) \quad (\text{Eq. 119})$$

With the inverting and noninverting integrators in Figs. 47 and 48, parasitic-insensitive SC filters of all transfer-function types can be realized. Parasitic insensitivity is not achieved at the expense of design generality. In fact, there is no reason to implement SC filters in any other way.

SC z-Domain Biquadratic Transfer Functions

To design an SC filter, one must first derive a z -domain transfer function from the filter requirements. This can be done by first obtaining the s -domain transfer function according to classic methods. The desired z -domain transfer function can then be obtained from the s -domain function by using the well-known bilinear transform

$$s = (2/\tau)(1 - z^{-1})/(1 + z^{-1}) \quad (\text{Eq. 120})$$

where τ represents one clock period. The bilinear transform is treated in Chapter 28 of this handbook and will not be discussed further here. Alternatively, z -domain transfer functions can be obtained directly from specified filter requirements by using any one of several computational techniques (see Chapter 28).

It was noted in previous sections that active RC biquads play a dominant role in the realization of high-order active RC filters. Although SC filters provide the opportunity to implement filters of much higher order on a single chip, the same substantial advantages are

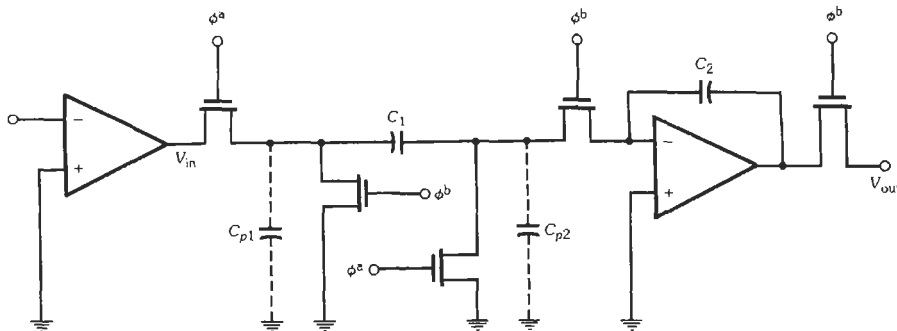


Fig. 48. Parasitic-insensitive SC noninverting integrator with associated parasitic capacitances.

obtained from biquads as in the *RC* case. Using biquads lends regularity to the design process and can reduce considerably the chip layout effort.

As active *RC* biquads realize biquadratic transfer functions in the *s*-domain, SC biquads realize biquadratic *z*-domain transfer functions. Thus, the transfer functions to be realized are of the form

$$H(z) = \frac{N(z)}{D(z)} = \frac{\gamma + \epsilon z^{-1} + \delta z^{-2}}{1 + \alpha z^{-1} + \beta z^{-2}} \quad (\text{Eq. 121})$$

The well-known special cases of Eq. 121, namely low-pass (LP), high-pass (HP), bandpass (BP), low-pass notch (LPN), high-pass notch (HPN), and all-pass (AP), can be derived by applying the bilinear transform to the corresponding second-order *s*-domain transfer functions.

One important property of the bilinear transfer function to be considered at this point is that the *z*-plane zeros at the one-half sampling frequency (i.e., $z = -1$) are mapped by Eq. 120 into *s*-plane zeros at infinity. Such zeros appear in low-pass (LP) and bandpass (BP) functions. Although use of the bilinear transform provides desirable steep cutoff to BP and LP filters in the vicinity of the half-sampling frequency (at $s \rightarrow \infty$), it may not result in the most economical SC realization. The additional attenuation provided at the half-sampling frequency is of little importance and diminishes in importance for relatively small pole, ω_p , and zero, ω_z , frequencies of the biquad (as $\omega_p \tau$ and $\omega_z \tau$ become small for the typical high sampling frequencies). When this relation is taken into account, several alternative LP and BP biquadratic transfer functions can be derived by replacing the zeros factors $1 + z^{-1}$ (i.e., for zeros at $z = -1$) with either 2 or $2z^{-1}$.

In Table 4, the numerators $N(z)$ are listed for the various elementary *z*-domain second-order forms. The LP and BP forms are referred to in this table as LP ij and BP ij , where i and j may have the value 0, 1, or 2. The suffix i denotes the number of $1 + z^{-1}$ factors and j the number of z^{-1} factors.

The recommended design procedure is first to obtain $D(z)$ and $N(z)$ via the bilinear transform. This will ensure the proper placement of the poles and zeros (other than those at $z = -1$). One or more of the zeros at $z = -1$ can then be replaced by either 2 or $2z^{-1}$. To place the gain level accurately, the gain constant K may then require slight alteration.

Parasitic-Insensitive SC Biquads

A general SC biquad that realizes Eq. 121 and its special cases in Table 4 was introduced by Fleischer and Laker, and by others.* This parasitic-insensitive biquad is shown in Fig. 49, in which capital letters A through L , of course, identify capacitors. This circuit is widely used throughout the industry. Because of space limitations, only the salient properties of this circuit will be highlighted. The interested reader is referred to the references for detailed derivations and demonstrations of individual features.

It simplifies the general analysis of this biquad to assume the input, V_{in} , is sampled and held for the full clock period. Although not necessary in every special-case implementation of Fig. 49, the condition is readily arranged. Furthermore, if the input to a cascade of such SC biquads is presented with a full-clock-period sampled-and-held signal, the switch timing in the biquads will propagate this condition through the entire filter. This condition is provided inherently by LP02 and BP01 biquads for any properly band-limited input signal.

Since, depending on the application, the desired output may be either at V or V' , both the corresponding transfer functions are of interest. They are derived as

$$H = V/V_{in} = -N/D \quad \text{and} \quad H' = V'/V_{in} = N'/D \quad (\text{Eq. 122a})$$

* References 1 and 26.

TABLE 4. ELEMENTARY SECOND-ORDER TRANSFER FUNCTIONS

Generic Form	Numerator $N(z)$
LP 20 (bilinear transform)	$K(1 + z^{-1})^2$
LP 11	$Kz^{-1}(1 + z^{-1})$
LP 10	$K(1 + z^{-1})$
LP 02	Kz^{-2}
LP 01	Kz^{-1}
LP 00	K
BP 10 (bilinear transform)	$K(1 - z^{-1})(1 + z^{-1})$
BP 01	$Kz^{-1}(1 - z^{-1})$
BP 00	$K(1 - z^{-1})$
HP	$K(1 - z^{-1})^2$
LPN	$K(1 + \epsilon z^{-1} + z^{-2}), \epsilon > \alpha/\sqrt{\beta}, \beta > 0$
HPN	$K(1 + \epsilon z^{-1} + z^{-2}), \epsilon < \alpha/\sqrt{\beta}, \beta > 0$
AP	$K(\beta + \alpha z^{-1} + z^{-2})$
General	$\gamma + \epsilon z^{-1} + \delta z^{-2}$

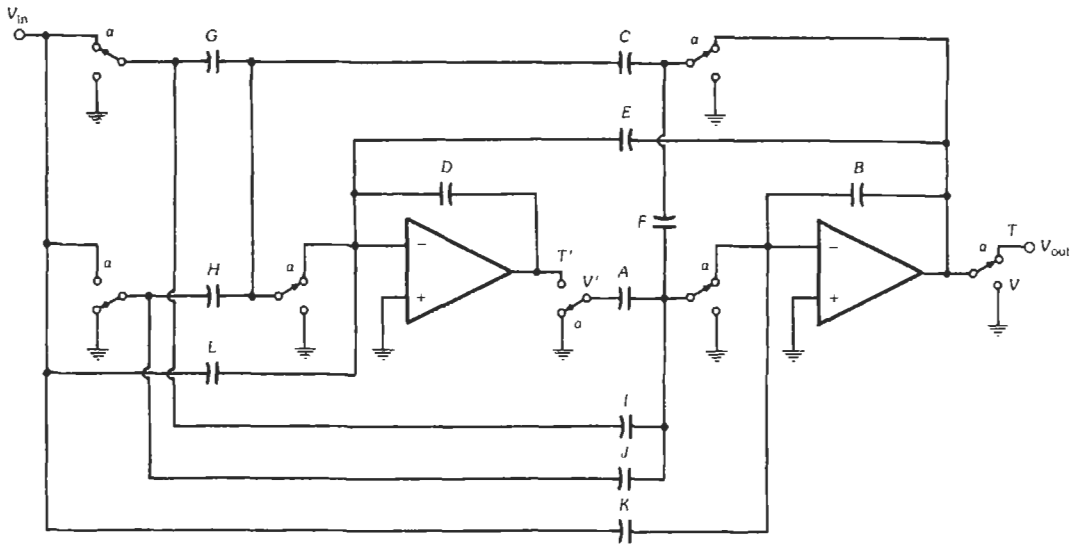


Fig. 49. General active SC biquad topology. (For clarity, the two transistors clocked at ϕ^a and ϕ^b are shown as a simple switch that is in position "a" during ϕ^s .)

where

$$D(z) = D(F + B) + (AC + AE - DF - 2BD)z^{-1} + (DB - AE)z^{-2} \tag{Eq. 122b}$$

is, of course, the same for both transfer functions, and

$$N(z) = D(I + K) + (AG + AL - DI - DJ - 2DK)z^{-1} + (DJ + DK - AH - AL)z^{-2} \tag{Eq. 122c}$$

$$N'(z) = (IC + IE + KC + KE - GF - GB - LF - LB) + (FH + BH + BG + FL + 2BL - JC - JE - IE - KC - 2KE)z^{-1} + (EJ + EK - BH - BL)z^{-2} \tag{Eq. 122d}$$

From Fig. 49 and Eqs. 122, it will be recognized that the transfer-function poles are determined by the feedback loop, consisting of capacitors $A, B, C, D, E,$ and F . Further, the transfer-function zeros are determined by the two feed-forward paths, consisting, respectively, of capacitors G, H, L and I, J, K . As expected, both H and H' share the same denominator; however, the numerators are quite different.

The properties of this general biquad circuit are much more easily discerned if certain simplifications are introduced. Thus, it may be readily observed that the general circuit of Fig. 49 is not minimal, with redundancy occurring in both the feedback and feed-forward paths. In the feedback paths, capacitor E and switched capacitor F provide two means for damping the transfer function poles. Thus, in practice either E or F is used, but not both. Accordingly, it is useful to define an "E-circuit" in which $E \neq 0$ and $F = 0$ and an "F-circuit" in which $E = 0$ and $F \neq 0$.

Also, as noted before, the two zero-forming feed-forward paths consist of six elements. At most, four of these, that is, two elements for each path, are required to realize arbitrary zero locations. Consequently, during the initial design of a biquad, it is convenient to assign $K = I = 0$. This degree of freedom is readily restored to the design by using the element equivalencies shown in Fig. 50. That is, after an initial design is completed, these equivalencies are employed to modify the circuit until an acceptable design is obtained. The z -domain validity of the equivalencies relies on terminals 1 and 2 being connected to a voltage source (independent voltage source or op-amp output) and virtual ground, respectively.*

Because the transfer function of switched-capacitor filters depends only on capacitor ratios, one capacitor

* References 1 and 26.

in each stage may be arbitrarily chosen. It is convenient to exercise these degrees of freedom by setting $B = D = 1$. In addition, it may be shown that one can initially set $A = 1$. The effect of this choice is to relinquish temporarily control of the gain constant associated with the transfer function to the "secondary" output. Once the initial design is completed, the stages may be properly rescaled to restore full generality to the biquad.

In view of the above, it is convenient to set $K = L = 0$ and $A = B = D = 1$ in Eqs. 122 and further to consider the E -circuit and the F -circuit separately. This results in the following equations:

$$H_F(z) = -\frac{I + (G - I - J)z^{-1} + (J - H)z^{-2}}{1 + (C + E - 2)z^{-1} + (1 - E)z^{-2}} \quad (\text{Eq. 123a})$$

$$H'_E(z) = \frac{(IC + IE - G) + (H + G - JC - JE - IE)z^{-1} + (EJ - H)z^{-2}}{1 + (C + E - 2)z^{-1} + (1 - E)z^{-2}} \quad (\text{Eq. 123b})$$

$$H_F(z) = -\frac{\hat{I} + (\hat{G} - \hat{I} - \hat{J})z^{-1} + (\hat{J} - \hat{H})z^{-2}}{(\hat{F} + 1) + (\hat{C} - \hat{F} - 2)z^{-1} + z^{-2}} \quad (\text{Eq. 124a})$$

$$H'_E(z) = -\frac{(\hat{G}\hat{F} + \hat{G} - \hat{I}\hat{C}) + (\hat{J}\hat{C} - \hat{F}\hat{H} - \hat{H} - \hat{G})z^{-1} + \hat{H}z^{-2}}{(\hat{F} + 1) + (\hat{C} - \hat{F} - 2)z^{-1} + z^{-2}} \quad (\text{Eq. 124b})$$

The "hats" are placed on the F -circuit elements to distinguish them from the E -circuit elements.

The synthesis equations for the biquad can be readily derived from Eqs. 121, 123, and 124. To avoid repetition, design equations will be given only for the more frequently used H_E and H_F functions. A complete set of design equations can be found in Reference 26.

For pole placement, the synthesis equations for E , C and \hat{F} , \hat{C} in terms of the z -domain transfer-function coefficients α and β can be stated as follows:

For the "E-circuit"

$$E = 1 - \beta \quad (\text{Eq. 125a})$$

$$C = 1 + \alpha + \beta \quad (\text{Eq. 125b})$$

For the "F-circuit"

$$\hat{F} = (1 - \beta) / \beta \quad (\text{Eq. 126a})$$

$$\hat{C} = (1 + \alpha + \beta) / \beta \quad (\text{Eq. 126b})$$

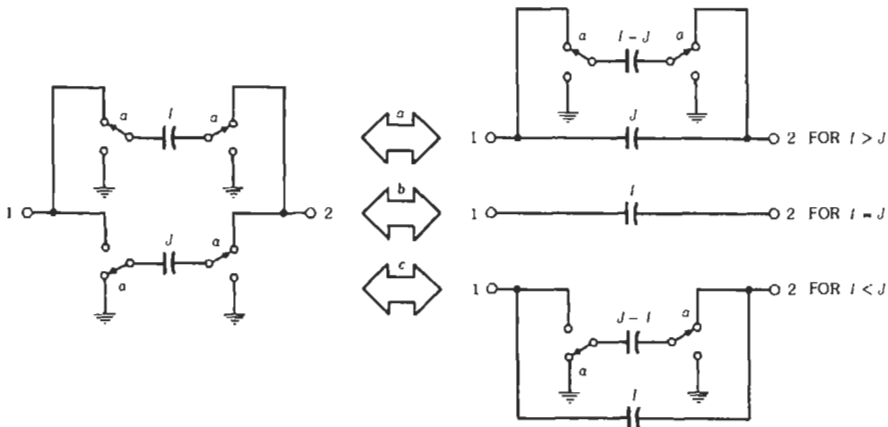


Fig. 50. SC element equivalencies.

The synthesis equations for placing the zeros can be derived in terms of the z -domain transfer-function coefficients γ , ϵ , and δ by comparing Eqs. 123 and 124 with the numerator forms in Table 4. Since in Eqs. 123 and 124 those numerator coefficients are determined by four capacitors, the solution is not unique. Hence, in Table 5 a complete set of design equations is given

TABLE 5. ZERO PLACEMENT FORMULAS FOR H_E AND H_F

Filter Type	Design Equations	Simple Solution
LP 20	$I = K $ $G - I - J = 2 K $ $J - H = K $	$I = J = K $ $G = 4 K , H = 0$
LP 11	$I = 0$ $G - I - J = \pm K $ $J - H = \pm K $	$I = 0, J = K $ $G = 2 K , H = 0$
LP 10	$I = K $ $G - I - J = K $ $J - H = 0$	$I = K , J = 0$ $G = 2 K , H = 0$
LP 02	$I = 0$ $G - I - J = 0$ $J - H = \pm K $	$I = J = 0$ $G = 0, H = K $
LP 01	$I = 0$ $G - I - J = \pm K $ $J - H = 0$	$I = J = 0$ $G = K , H = 0$
LP 00	$I = K $ $G - I - J = 0$ $J - H = 0$	$I = K , J = 0$ $G = K , H = 0$
BP 10	$I = K $ $G - I - J = 0$ $J - H = - K $	$I = K , J = 0$ $G = H = K $
BP 01	$I = 0$ $G - I - J = \pm K $ $J - H = \pm K $	$I = 0, J = K $ $G = H = 0$
BP 00	$I = K $ $G - I - J = - K $ $J - H = 0$	$I = K , J = 0$ $G = H = 0$
HP	$I = K $ $G - I - J = -2 K $ $J - H = K $	$I = J = K $ $G = H = 0$
HPN and LPN	$I = K $ $G - I - J = K \epsilon$ $J - H = K $	$I = J = K $ $G = K \{2 + \epsilon\}, H = 0$
AP ($\beta > 0$)	$I = K \beta$ $G - I - J = K \alpha$ $J - H = K $	$I = K \beta, J = K $ $G = K (1 + \beta + \alpha) = K \alpha$ $H = 0$
General ($\gamma > 0$)	$I = \gamma$ $G - I - J = \epsilon$ $J - H = \delta$	$I = \gamma$ $J = \delta + x$ $G = \gamma + \delta + \epsilon + x$ $H = x \geq 0$

Note: $\hat{G} = G(1 + \hat{F})$, $\hat{H} = H(1 + \hat{F})$, $\hat{I} = I(1 + \hat{F})$, and $\hat{J} = J(1 + \hat{F})$.

for each case. For each of the cases, a "simple" solution is also offered.

The synthesis equations given in the previous paragraphs result in unscaled capacitor values. To complete the synthesis in practice, some scaling is required. The first order of business is to adjust the voltage level at the "secondary" output. If this voltage is too high, amplifier overloads may result; if it is too low, unnecessary noise penalties may be taken.

Although the voltage levels and necessary scaling factors may be obtained by using analysis techniques,*

the simplest procedure is to simulate the unscaled circuit on an analysis program. This method serves also as a confirmation of the correctness of the design.

To adjust the voltage level V' , i.e., the flat gain of H' , without affecting H , only capacitors A and D need to be scaled. More precisely, if it is desired to modify the gain constant associated with V' according to

$$H' \rightarrow \mu H' \quad (\text{Eq. 127})$$

* Reference 1.

then it is only necessary to scale A and D as

$$(A, D) \rightarrow (1/\mu)A, (1/\mu)D \quad (\text{Eq. 128})$$

The gain constant associated with H remains invariant under this scaling. The correctness of this procedure follows directly from signal-flow graph concepts.

In a similar fashion, it can be shown that if the flat gain associated with V is to be modified, i.e.,

$$H \rightarrow \nu H \quad (\text{Eq. 129})$$

the following capacitors must be scaled:

$$(B, C, E, F) \rightarrow (1/\nu)B, (1/\nu)C, (1/\nu)E, (1/\nu)F \quad (\text{Eq. 130})$$

Once satisfactory gain levels have been obtained at both outputs, it is convenient to scale the admittances associated with each stage so that the minimum capacitor value in the circuit becomes unity. This makes it easier to observe the maximum capacitor ratios required to realize a given circuit and also serves to “standardize” different designs so that the total capacitance required can be readily observed. The two groups of capacitors that are to be scaled together are listed below:

Group 1: (C, D, E, G, H, L)

Group 2: (A, B, F, I, J, K)

Note that capacitors in each group are distinguished by the fact that they are all incident on the same input node of one of the operational amplifiers.

This completes the design process for synthesizing practical SC-biquad networks. In the next section, a detailed example is given to demonstrate each step of the design.

Low-Pass Notch Example

The transfer function to be realized will be based on the s -domain low-pass notch function

$$H(s) = \frac{0.891975s^2 + \omega_p^2}{s^2 + s\omega_p/30 + \omega_p^2} \quad (\text{Eq. 131})$$

where

$$\omega_p = 2\pi f_p \text{ with } f_p = 1700 \text{ Hz.}$$

This transfer function provides a notch frequency of $f_n = 1800$ Hz, a peak corresponding to a quality factor $Q_p = 30$, and 0-dB dc gain. The assigned sampling frequency is 128 kHz; i.e., $\tau = 7.8125 \mu\text{s}$.

The z -domain transfer function is conveniently obtained via the bilinear transformation shown in Eq. 120. Because the band-edge frequency of 1700 Hz is much less than the sampling rate, it is not necessary to prewarp (see Chapter 28) $H(s)$. Applying the bilinear transformation to Eq. 131 yields, after some algebraic manipulations,

$$H(z) = 0.89093 \frac{1 - 1.99220z^{-1} + z^{-2}}{1 - 1.99029z^{-1} + 0.99723z^{-2}} \quad (\text{Eq. 132})$$

Note that in obtaining this transfer function a high degree of numerical precision is required. However, this does not result in high sensitivities, since the capacitor ratios define only the departures from -2 and $+1$ in the above terms.

Only the H_E and H_F realizations of the above circuit will be given here, because they are more economical in the number of capacitors required. The synthesis itself is straightforward. Capacitors C, E or \hat{C}, \hat{E} are determined from Eqs. 125 or 126, respectively, and capacitors G, H, I, J or $\hat{G}, \hat{H}, \hat{I}, \hat{J}$ are obtained from the “simple” solution entry in Table 5. Finally, of course, A, B, D or $\hat{A}, \hat{B}, \hat{D}$ are set equal to unity. The resulting unscaled capacitor values are given in the appropriate columns of Table 6. Note that in this table the “hats” are omitted from the F -circuit capacitors for notational convenience. Also note that since $I = J$ these two switched capacitors are replaced by the unswitched capacitor K ($K = I = J$) in accordance with Fig. 50.

As the next step, it is appropriate to simulate the unscaled E - and F -circuits to verify the design. The results of the simulation will confirm the correctness of H_E and H_F , in particular, that the maximum gain in both these realizations is approximately 10.56 dB. However, the maximum gains for H'_E and H'_F are very low. To make the first stage no more susceptible to overloads than the second stage, it is recommended that these gains also be increased to a maximum of 10.56 dB. Specifically, it was found that

$$H'_{E,\text{max}} \approx -11.05 \text{ dB}, H'_{F,\text{max}} \approx -10.96 \text{ dB}$$

Therefore, in accordance with Eq. 127,

$$\mu = 12.0365, \hat{\mu} = 11.9124$$

Using these factors to rescale A, D and \hat{A}, \hat{D} , respectively, as given in Eq. 128, yields the “dynamic range adjusted” capacitor values also shown in Table 6. Finally, the capacitors associated with each operational-amplifier stage are separately rescaled so that the minimum capacitor value becomes 1 pF. Table 6 shows the “final” values.

In comparing the “final” realizations, we note that the F -circuit requires roughly 12 times the total capacitance of the E -circuit, in spite of the fact that the initial values were almost identical. Thus, alternative designs must be carried to completion before they can be meaningfully compared. It should be noted that other practical examples exist in which the F -circuit designs are dramatically more efficient than the corresponding E -circuit designs. The sensitivities for both designs are found to be equivalent.

Other Topologies and Comments

SC filters, as any other type of analog active filter, consist of active inverting and noninverting integrators (see Figs. 46 and 47). It stands to reason, therefore,

TABLE 6. LOW-PASS NOTCH REALIZATION

Capacitor (pF)	E-Circuit			F-Circuit		
	Initial	Dynamic Range Adjusted	Final	Initial	Dynamic Range Adjusted	Final
A	1.0000	0.08308	1.0000	1.0000	0.08395	30.1895
B	1.0000	1.0000	12.0365	1.0000	1.0000	359.629
C	0.00694	0.00694	2.5035	0.00696	0.00696	1.0000
D	1.0000	0.08308	29.9613	1.0000	0.08395	12.0591
E	0.00277	0.00277	1.0000	—	—	—
F	—	—	—	0.00278	0.00278	1.0000
G	0.00694	0.00694	2.5035	0.00696	0.00696	1.0000
H	—	—	—	—	—	—
I	—	—	—	—	—	—
J	—	—	—	—	—	—
$K(I = J)$	0.89093	0.89093	10.7238	0.89340	0.89340	321.293
ΣC (pF)	—	—	59.7	—	—	726.1

that other filter structures, such as ladders, discussed in the earlier parts of this chapter can also be implemented in the SC technique. Because of space limitations and the fact that SC filters are built only in fully integrated form, requiring a considerable design effort and expertise, the reader is referred to the literature* for a discussion of design details. Suffice it here to make the reader aware that analog active SC filters of very high quality and of high and low order have found commercial acceptance for a wide variety of signal-processing and control applications at moderate frequencies. Attention must be paid only to the fact that SC filters are sampled-data circuits, which calls for some care in the system design. Finally, we emphasize again that integrated SC filters, just as g_m -C circuits, should be built in fully differential form for reasons outlined earlier: better dynamic range, power-supply rejection, and clock feedthrough.†

REFERENCES

- Ghausi, M. S., and Laker, K. R. *Modern Filter Design: Active RC and Switched Capacitor*. Englewood Cliffs, NJ: Prentice Hall, 1981.
- Sedra, A. S., and Brackett, P. O. *Filter Theory and Design: Active and Passive*. Portland, OR: Matrix Publishers, Inc., 1978.
- Lindquist, C. S. *Active Network Design with Signal Filtering Applications*. Long Beach, CA: Steward & Sons, 1977.
- Moschytz, G. S. *Linear Integrated Network—Design*. New York: Van Nostrand Reinhold, 1975.
- Van Valkenburg, M. E. *Analog Filter Design*. New York: Holt, Rinehart and Winston, 1982.
- Schaumann, R., Ghausi, M. S., and Laker, K. R. *Design of Analog Filters: Passive, Active RC and Switched Capacitor*. Englewood Cliffs, NJ: Prentice Hall, 1990.
- Schaumann, R., and Van Valkenburg, M. E. *Design of Analog Filters*. New York: Oxford University Press, 2001.
- Chen, W. K., Ed. *The Circuits and Filters Handbook*. Boca Raton, FL/New York: CRC Press/IEEE Press, 1995.
- Schaumann, R., Soderstrand, M. S., and Lakcr, R., Eds. *Modern Active Filter Design*. IEEE Press Selected Reprint Series. New York: IEEE Press, 1981.
- Allen, P. E., and Sanchez-Sinencio, E. *Switched-Capacitor Circuits*. New York: Van Nostrand Reinhold, 1984.
- Tsividis, Y. P., and Voorman, J. O., Eds. *Integrated Continuous-Time Filters—Principles, Design, and Applications*. IEEE Press Selected Reprint Series. New York: IEEE Press, 1993.
- Schaumann, R. "The Design of Continuous-Time Fully Integrated Filters: A Review," *IEE Proceedings*, Vol. 136, Pt. G, Electronic Circuits and Systems, pp. 184–190, Aug. 1989.
- Schaumann, R. "Continuous-Time Integrated Filter—A Tutorial," Paper 1.1 in Reference [11].
- Voorman, J. A. "Continuous-Time Analog Integrated Filters," Paper 1.2 in Reference [11].
- Sedra, A. S., Zharab, M. A., and Martin, K. "Optimum Configuration of Single-Amplifier Biquadratic Filters," *IEEE Transactions on Circuits and Systems*, Vol. CAS-27, pp. 1155–1163, 1980.
- De Veirman, G. A., and Yamasaki, R. G. "Monolithic 10–30 MHz Tunable Bipolar Bessel Low-pass Filter," *IEEE Int. Symp. on Circuits and Systems*, pp. 1444–1447, May 1991.
- Zverev, A. I. *Handbook of Filter Synthesis*. New York: Wiley, 1967.
- Saal R. *Handbook of Filter Design*. Berlin: AEG-Telefunken, 1979.

* References 1, 6, and 24–27.

† Reference 27.

19. Wu, P., and Schaumann, R. "A 200 MHz Elliptic OTA-C Filter in GaAs Technology," *Proc. IEEE Int. Symp. on Circuits and Systems*, pp. 1745–1748, June 1991.
20. Tan, M. A., and Schaumann, R. "Simulating General-Parameter LC-Ladder Filters for Monolithic Realizations with Only Transconductance Elements and Grounded Capacitors," *IEEE Transactions on Circuits and Systems*, Vol. CAS-36, No. 2, pp. 299–307, Feb. 1989.
21. Snelgrove, W. M., and Sedra, A. S., "Optimization of Dynamic Range in Cascade Active Filters," *Proc. IEEE Int. Symp. on Circuits and Systems*, pp. 151–155, 1978.
22. Chiou, C. F., and Schaumann, R. "Comparison of Dynamic Range Properties of High-Order Active Bandpass Filters," *Proc. IEE*, Vol. 127, Pt. G, Electronic Circuits and Systems, No. 3, pp. 101–108, 1980.
23. Khoury, J. M. "Design of a 15-MHz CMOS Continuous-Time Filter with On-Chip Tuning," *IEEE Journal on Solid State Circuits*, Vol. SC-26, No. 12, pp. 1988–1997, 1991.
24. Gray, P. R., Hodges, D. A., and Broderon, R. W., Eds. *Analog MOS Integrated Circuits*, IEEE Press Selected Reprint Series, 1980.
25. Martin, K., and Sedra, A. S. "Effects of Op Amp Finite Gain and Bandwidth on the Performance of Switched Capacitor Filters," *IEEE Transactions on Circuits and Systems*, Vol. CAS 28, pp. 822–829, Aug. 1981.
26. Fleischer, P. E., and Laker, K. R. "A Family of Active Switched-Capacitor Biquad Building Blocks," *The Bell System Technical Journal*, Vol. 58, No. 10, Dec. 1979. (Reprinted in [9])
27. Choi, T. C., et al. "High-Frequency CMOS Switched-Capacitor Filters for Communications Applications," *IEEE J. Solid-State Circuits*, Vol. SC-18, pp. 652–664, Dec. 1983.
28. Li, D., and Tsividis, Y. "Active LC Filters on Silicon," *IEE Proceedings-Circuits, Devices and Systems*, Vol. 147, No. 1, pp. 49–56, Feb. 2000.

11 Attenuators

*Revised by
Bruno O. Weinschel*

Definitions	11-2
Typical Designs of Resistive Attenuators	11-2
Resistance Networks for Attenuators	
Power Dissipation Within a Tee Pad	
Connectors	11-4
Measurement of Attenuation	11-4
Fixed-Frequency	
Broadband	
Swept or Stepped in Frequency	

DEFINITIONS

An attenuator is a network that reduces the input power by a predetermined ratio. The ratio of input power to output power is expressed in logarithmic terms such as decibels (dB).

$$\text{Attenuation in dB} = 10 \log_{10} P_{in}/P_{out} \\ \equiv 20 \log_{10} E_{in}/E_{out}$$

NOTE: $Z_{source} = Z_{load} = Z_{attenuator}$
All resistive, matched

Examples:

- 1) $P_{in}/P_{out} = 13.18 = 10 \times 1.318$
 $10 \log_{10} P_{in}/P_{out} \text{ dB} = 10 (\log_{10} 10 + \log_{10} 1.318) \text{ dB}$
 $= 10(1 + 0.1199) \text{ dB} \approx 11.2 \text{ dB}$
- 2) $E_{in}/E_{out} = 3.630$
 $20 \log_{10} E_{in}/E_{out} \text{ dB} = 20 \log_{10} 3.630 \text{ dB}$
 $= 20 \times 0.560 \text{ dB} = 11.2 \text{ dB}$

To convert attenuation in decibels into power or voltage ratio:

$$P_{in}/P_{out} = \log_{10}^{-1} \text{ dB}/10 = 10^{\text{dB}/10}$$

$$E_{in}/E_{out} = \log_{10}^{-1} \text{ dB}/20 = 10^{\text{dB}/20}$$

Examples:

- 3) $11.2 \text{ dB } P_{in}/P_{out} = 10^{11.2/10} = 10^{1.12} = 13.18$
- 4) $11.2 \text{ dB } E_{in}/E_{out} = 10^{11.2/20} = 10^{0.56} = 3.630$

Table 1 lists a few decibel values together with the corresponding power and voltage ratios.

The power ratio used in the *characteristic insertion loss* relation (page 82, IEEE Standard 100—1972) is shown in Fig. 1. The ratio specified is equal to the characteristic insertion loss since source and load are reflectionless. Usually, one is interested in resistive attenuators of nominally equal impedance as the source and load impedance. Resistive devices are most common. Attenuators include absorptive and reflective devices having at least two ports.

TABLE 1. DECIBELS VERSUS POWER AND VOLTAGE RATIO

dB	P_{in}/P_{out}	E_{in}/E_{out}
10	10	—
20	10^2	10
30	10^3	—
40	10^4	10^2
50	10^5	—
60	10^6	10^3
70	10^7	—
80	10^8	10^4
90	10^9	—
100	10^{10}	10^5
110	10^{11}	—
120	10^{12}	10^6

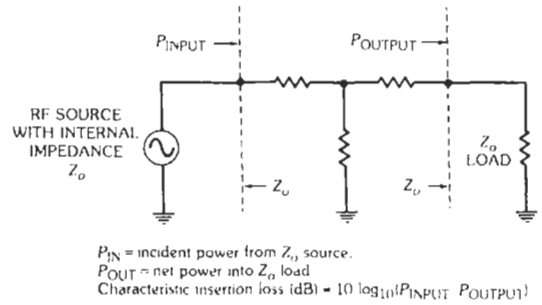


Fig. 1. Definition of characteristic insertion loss.

There are fixed attenuators, continuously variable attenuators, and step attenuators. This chapter treats principally Z_0 matched, two-port attenuators, for use in a Z_0 system. They are reversible except for a high-power attenuator designed for unilateral power flow. Minimum-loss L-pads for matching between two different impedances are also treated.

TYPICAL DESIGNS OF RESISTIVE ATTENUATORS

Simple wirewound resistors are used to upper audio frequencies. Nonreactive wirewound resistors (mica card, Ayrton-Perry winding, woven resistors) can be used to higher frequencies for low characteristic impedance values of, for example, 50 ohms. Preferred US impedance values are 600 ohms for audio and 50 or 75 ohms for coaxial video, RF, and microwave transmissions. For coaxial applications to over 26.5 GHz, thin-film resistors on low-loss insulators are used as attenuating elements. Up to a few gigahertz, lumped rod and disc resistors can be used as series and shunt elements in coaxial lines. The rods should be shorter than $\lambda/8$ at the maximum frequency (Figs. 2 and 3). For

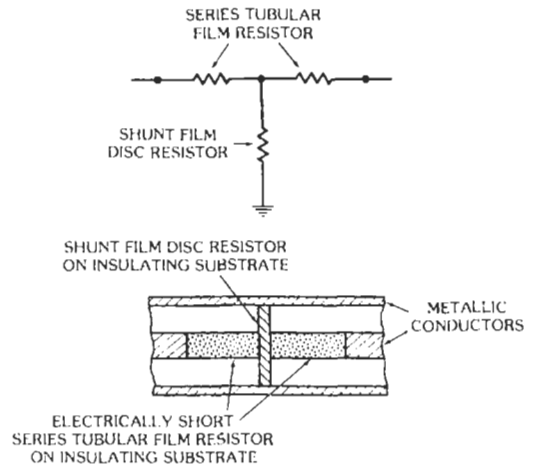


Fig. 2. Coaxial tee section.

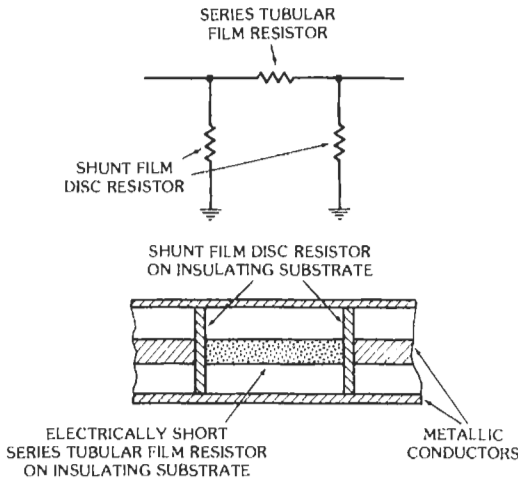


Fig. 3. Coaxial pi section.

higher frequencies, distributed resistive films on suspended substrates are used (Fig. 4).

Resistance Networks for Attenuators

Four types of symmetrical pads are shown in Fig. 5. The formulas for the resistance values in ohms for these pads when $R_o = 1$ ohm are given below. If R_o has a

value other than 1 ohm, multiply each of the resistance values (a , b , c , $1/a$, $1/b$, and $1/c$) by R_o .

$$a = (10^{dB/20} - 1)/(10^{dB/20} + 1)$$

$$b = 2 \cdot 10^{dB/20} / (10^{dB/10} - 1)$$

$$c = 10^{dB/20} - 1$$

An unsymmetrical matching pad is shown in Fig. 6. The formulas for this pad are:

$$j = R_1 - kR_2 / (k + R_2)$$

$$k = [R_1 R_2^2 / (R_1 - R_2)]^{1/2}, \text{ where } R_1 > R_2$$

$$\text{Minimum loss in dB} = 20 \log_{10} \{ [(R_1 - R_2) / R_2]^{1/2} + (R_1 / R_2)^{1/2} \}$$

Typical values for the pads in Fig. 5 are listed in Table 2. Typical values for the pad in Fig. 6 are listed in Table 3.

Balanced pads are used in balanced transmission lines. For voice, Z_o is typically 600 ohms; for VHF TV, Z_o is typically 300 ohms. Bridged tee pads are common in low-frequency step attenuators. Lumped-element pi pads are short and easily matched up to about 1.5 GHz, with 75- or 50-ohm impedance. Coaxial tee pads of 50-ohm impedance can work above 10 GHz if the series elements are short in comparison with the wavelength. For low frequency sensitivity, thin resistive films (Ni-

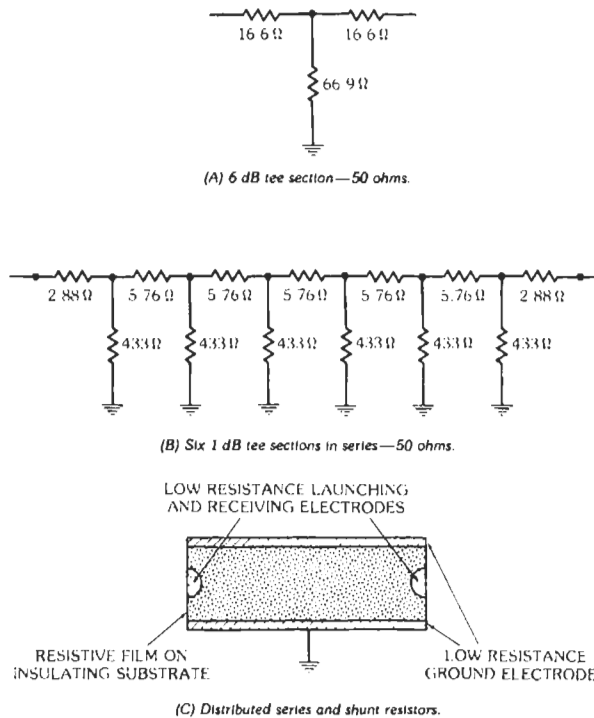
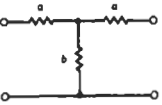
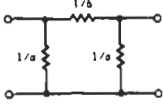

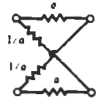


Fig. 4. Distributed series and shunt resistive attenuator element on suspended substrate in coaxial line.

TABLE 2. RESISTANCE VALUES FOR ATTENUATOR PADS WHEN $R_o = 1 \text{ OHM}^1$

dB	Tee Pad		Pi Pad		Bridged Tee Pad		Balanced Pad	
								
	a	b	$1/b$	$1/a$	c	$1/c$	a	$1/a$
0.1	0.005 7567	86.853	0.011 514	173.71	0.011 580	86.356	0.005 7567	173.71
0.2	0.011 513	43.424	0.023 029	86.859	0.023 294	42.930	0.011 513	86.859
0.3	0.017 268	28.947	0.034 546	57.910	0.035 143	28.455	0.017 268	57.910
0.4	0.023 022	21.707	0.046 068	43.438	0.047 128	21.219	0.023 022	43.438
0.5	0.028 775	17.362	0.057 597	34.753	0.059 254	16.877	0.028 775	34.753
0.6	0.034 525	14.465	0.069 132	28.965	0.071 519	13.982	0.034 525	28.965
0.7	0.040 274	12.395	0.080 678	24.830	0.083 927	11.915	0.040 274	24.830
0.8	0.046 019	10.842	0.092 234	21.730	0.096 478	10.365	0.046 019	21.730
0.9	0.051 762	9.6337	0.103 80	19.319	0.109 18	9.1596	0.051 762	19.319
1.0	0.057 501	8.6668	0.115 38	17.391	0.122 02	8.1954	0.057 501	17.391
2.0	0.114 62	4.3048	0.232 30	8.7242	0.258 93	3.8621	0.114 62	8.7242
3.0	0.171 00	2.8385	0.352 30	5.8481	0.412 54	2.4240	0.171 00	5.8481
4.0	0.226 27	2.0966	0.476 97	4.4194	0.584 89	1.7097	0.226 27	4.4194
5.0	0.280 13	1.6448	0.607 97	3.5698	0.778 28	1.2849	0.280 13	3.5698
6.0	0.332 28	1.3386	0.747 04	3.0095	0.995 26	1.0048	0.332 28	3.0095
7.0	0.382 48	1.1160	0.896 04	2.6145	1.2387	0.807 27	0.382 48	2.6145
8.0	0.430 51	0.946 17	1.0569	2.3229	1.5119	0.661 43	0.430 51	2.3229
9.0	0.476 22	0.811 83	1.2318	2.0999	1.8184	0.549 94	0.476 22	2.0999
10.0	0.519 49	70 273*	1.4230	1.9250	2.1623	46 248*	0.519 49	1.9250
20.0	0.818 18	20 202*	4.9500	1.2222	9.0000	11 111*	0.818 18	1.2222
30.0	0.938 69	6330.9*	15.796	1.0653	30.623	3265.5*	0.938 69	1.0653
40.0	0.980 198	2000.2*	49.995	1.0202	99.000	1010.1*	0.980 198	1.0202
50.0	0.993 70	632.46*	158.11	1.0063	315.23	317.23*	0.993 70	1.0063
60.0	0.998 00	200.00*	500.00	1.0020	999.00	100.10*	0.998 00	1.0020
70.0	0.999 37	63.246*	1581.1	1.0006	3161.3	31.633*	0.999 37	1.0006
80.0	0.999 80	20.000*	5000.0	1.0002	9999.0	10.001*	0.999 80	1.0002
90.0	0.999 94	6.3246*	15 811	1.0001	31 622	3.1633*	0.999 94	1.0001
100.0	1.0000	2.0000*	50 000	1.0000	99 999	1.0000*	1.0000	1.0000

* These values have been multiplied by 10^5 .

NOTES:

1. If $R_o \neq 1 \text{ ohm}$, multiply all values by R_o .
2. For other decibel values, use formulas in text.

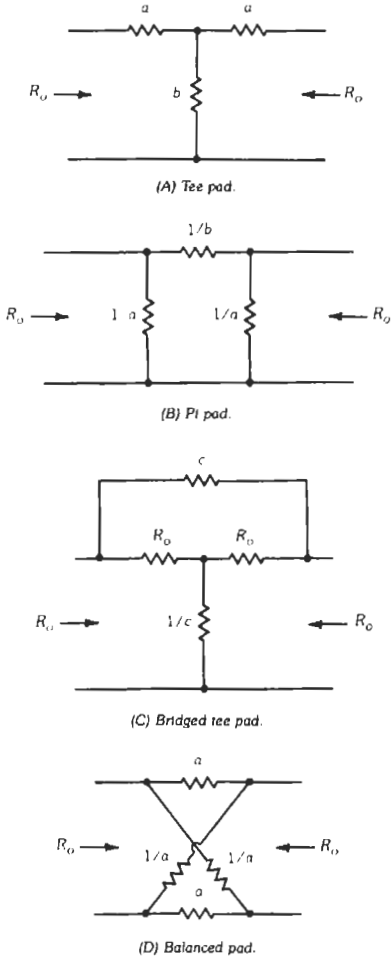


Fig. 5. Symmetrical pads for use between equal characteristic impedances.

chrome, cracked carbon, tantalum nitride) on an insulating substrate (mica, ceramic, quartz) are used. Figs. 2, 3, and 4 illustrate coaxial cross sections.

For a broadband match between impedances R_1 and R_2 , use the minimum-loss L pad (Fig. 6).

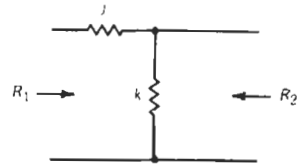


Fig. 6. Unsymmetrical matching pad between R_1 and R_2 .

Power Dissipation Within a Tee Pad

Table 4 lists values of power dissipation within a tee pad. The values in the table are for an input of 1 watt; for other input powers, multiply the values in the table by the input power.

CONNECTORS

Each connector of an attenuator must mate nondestructively with connectors that comply with the governing standard, typically Mil-C-39012, IEEE Std. 287-1968, or Mil-F-3922, as shown in Table 5.

MEASUREMENT OF ATTENUATION

The methods covered are fixed-frequency, broadband, and swept or stepped in frequency. All accuracies exclude the effects of mismatch uncertainty and connector nonrepeatability. In Fig. 7, the most commonly used attenuation-measurement systems are compared as to dynamic range and typical accuracy.

Fixed-Frequency

Measurement of Small Attenuation (< 5 dB)

Dual-channel audio substitution (Fig. 8)

Accuracy: 0.01 dB

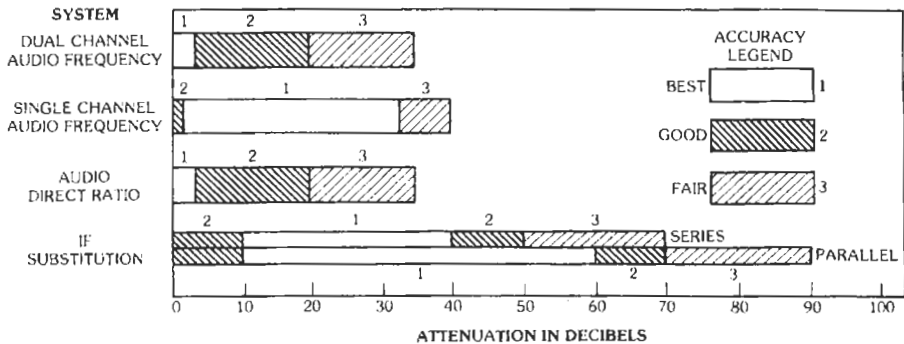
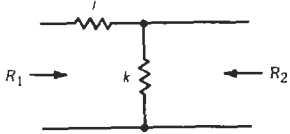


Fig. 7. Comparison of most commonly used attenuation measurement systems.

TABLE 3. RESISTANCE VALUES AND ATTENUATION FOR MINIMUM-LOSS L PAD*



R_1/R_2	j	k	dB
20.0	19.49	1.026	18.92
16.0	15.49	1.033	17.92
12.0	11.49	1.044	16.63
10.0	9.486	1.054	15.79
8.0	7.484	1.069	14.77
6.0	5.478	1.095	13.42
5.0	4.472	1.118	12.54
4.0	3.469	1.155	11.44
3.0	2.449	1.225	9.96
2.4	1.833	1.310	8.73
2.0	1.414	1.414	7.66
1.6	0.9798	1.633	6.19
1.2	0.4898	2.449	3.77
1.0	0	∞	0

* For $R_2 = 1$ ohm and $R_1 > R_2$. If $R_2 \neq 1$ ohm, multiply values by R_2 . For ratios not in the table, use the formulas in the text.

Examples of use of table:

If $R_1 = 50$ ohms and $R_2 = 25$ ohms, then $R_1/R_2 = 2.0$, and $j = k = 1.414 \times 25$ ohms = 35.35 ohms.

If $R_1/R_2 = 1.0$, minimum loss = 0 dB.

For $R_1/R_2 = 2.0$, the insertion loss with the use of j and k for matching is 7.66 dB above that for $R_1/R_2 = 1.0$.

Advantage: Single radio-frequency source

Attenuation standards: Audio attenuator (resistive)
Ratio transformer (inductive)

Measurement of Moderate Attenuation (1 to 30 dB)

Audio substitution

- (A) Single channel (Fig. 9)
Accuracy: 0.1 dB
- (B) Dual channel (Fig. 8)
Accuracy: 0.05 dB

The major advantage of systems (A) and (B) is simplicity.

Measurement of High Attenuation (30 to > 100 dB)

Intermediate-frequency substitution

- (A) Series substitution (Fig. 10)

Accuracy: 0.1 to 1 or 2 dB

Disadvantage: Minimum loss of standard attenuator (typically 30 dB for a waveguide-below-cutoff attenuator) in series with the mixer reduces dynamic range.

- (B) Parallel substitution (Fig. 11)

Accuracy: 0.1 to 1 or 2 dB

Advantage: Minimum loss of standard intermediate-frequency attenuator does not affect system sensitivity; gain stability of intermediate-frequency amplifier does not influence accuracy.

Disadvantage: More complex

All intermediate-frequency systems have an inherent limitation due to the range of power conversion linearity of the mixer. Beyond 90 dB, corrections may have to be made for linearity deviation. Partial radio-frequency substitution is often used to extend range, but mismatch and leakage errors are prevalent.

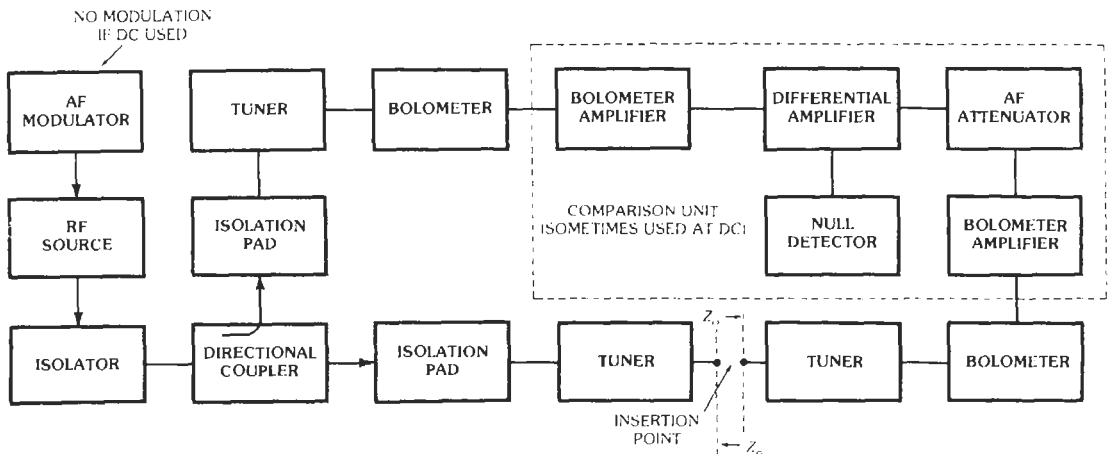
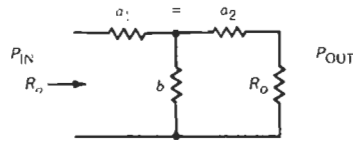


Fig. 8. Dual-channel audio substitution system (Reference 1).

TABLE 4. POWER DISSIPATION WITHIN TEE PAD*



dB	Watts Input Series Resistor	Watts Shunt Resistor	Watts Output Series Resistor
0.1	0.005 76	0.0112	0.005 625
0.3	0.0173	0.0334	0.016 113
0.5	0.0288	0.0543	0.025 643
0.7	0.0403	0.0743	0.034 279
0.9	0.0518	0.0933	0.0421
1.0	0.0575	0.1023	0.0456
1.2	0.0690	0.120	0.0523
1.4	0.0804	0.1368	0.0582
1.6	0.0918	0.1525	0.0635
1.8	0.103	0.1672	0.0679
2.0	0.114	0.1808	0.0718
2.2	0.126	0.1953	0.0758
2.4	0.137	0.2075	0.0787
2.6	0.149	0.2205	0.0818
2.8	0.160	0.232	0.0839
3.0	0.170 998	0.242 114	0.085 698
3.2	0.182	0.2515	0.0870
3.4	0.193	0.2605	0.0882
3.6	0.204	0.2695	0.0890
3.8	0.215	0.2775	0.0897
4.0	0.226	0.285	0.0898
5	0.280	0.3145	0.0884
6	0.332	0.332	0.0833
7	0.382	0.341	0.0761
8	0.430	0.343	0.0681
9	0.476 218	0.337 94	0.059 9527
10	0.519	0.328	0.519
12	0.598	0.3005	0.0377
14	0.667	0.266	0.0266
16	0.726 386	0.230 36	0.018 2460
18	0.776	0.1955	0.0123
20	0.818	0.1635	0.0100
30	0.938	0.0593	0.0010
40	0.980	0.0196	0.0001

*For 1-watt input and matched termination. If input \neq 1 watt, multiply values by P_{in} .

Broadband

Measurement of Small Attenuation (Less Than 5 dB)

Dual-channel audio substitution (Fig. 8)

Accuracy: 0.01 dB

Measurement of Attenuation Above 1 dB

Single- or dual-channel audio (Figs. 8 and 9)

Accuracy: 0.01 to 0.05 dB

Advantage: Single radio-frequency source

TABLE 5. LIST OF POPULAR COAXIAL CONNECTORS

Connector Type	Contact Type	Specifications
N	Female	Mil-C-39012/2
N	Male	Mil-C-39012/1
C	Male	Mil-C-39012/6
C	Female	Mil-C-39012/7
SC	Male	Mil-C-39012/39
SC	Female	Mil-C-39012/36
TNC	Male	Mil-C-39012/26
TNC	Female	Mil-C-39012/27
BNC	Male	Mil-C-39012/16
BNC	Female	Mil-C-39012/17
SMA	Male	Mil-C-39012/55
SMA	Female	Mil-C-39012/57
7-mm GPC	Sexless	IEEE Std. 287-1968
14-mm GPC	Sexless	IEEE Std. 287-1968
Waveguide	Choke or cover flange	Mil-F-3922

References—Mil-C-39012 on Coaxial Connectors; Mil-F-3922 on Waveguide Flanges, General Purpose; IEEE Std. 287-1968 on Precision Coaxial Connectors

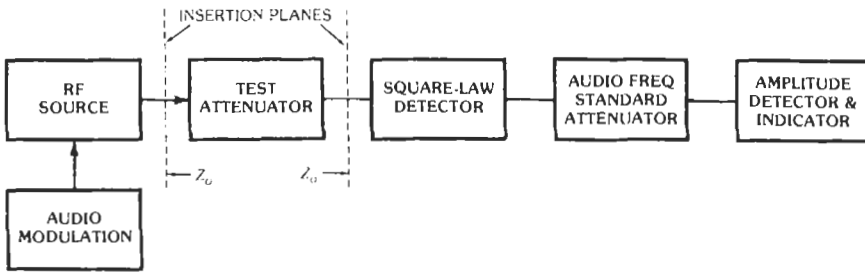
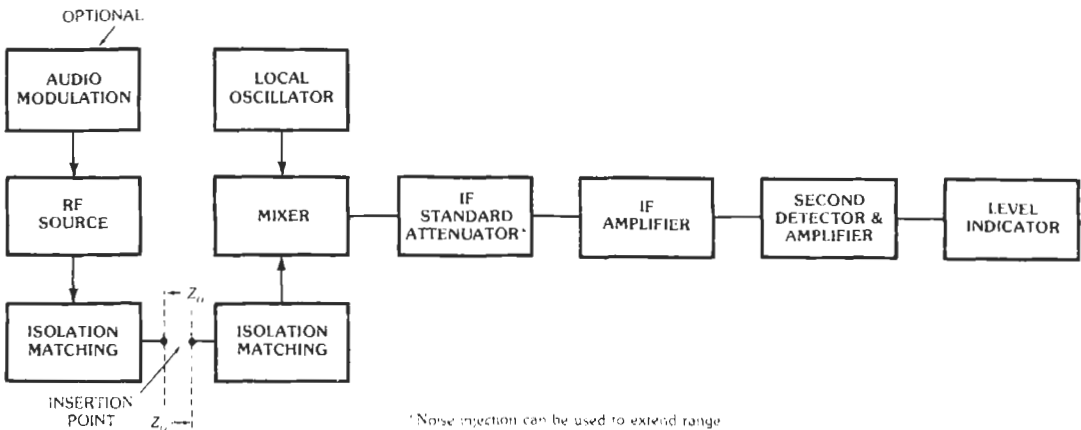


Fig. 9. Single-channel audio substitution system (Reference 1).



* Noise injection can be used to extend range

Fig. 10. Series type intermediate-frequency substitution system (References 3 and 4).

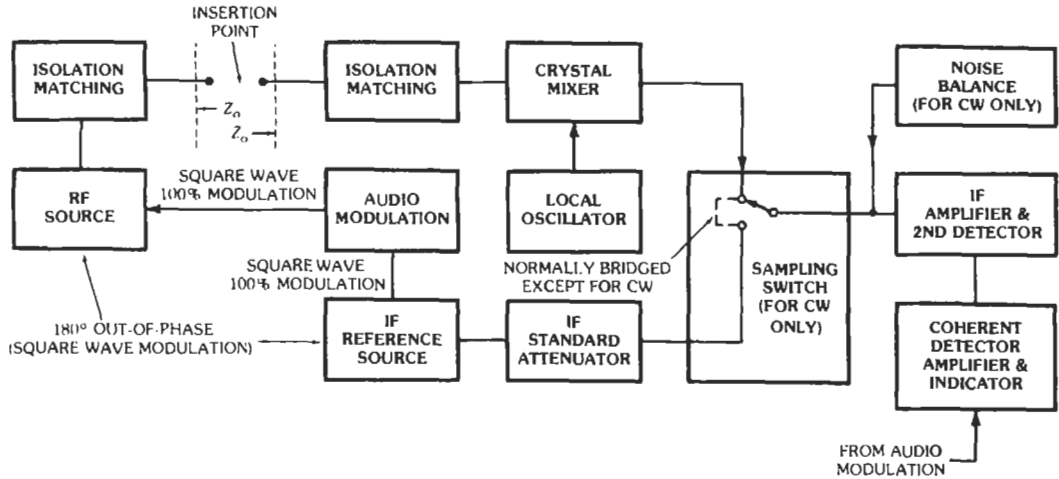


Fig. 11. Parallel IF substitution system with noise injection and coherent detection (References 5 and 6).

Disadvantage: Limited dynamic range (about 30 dB) without partial radio-frequency substitution

The major difficulty in broadband techniques is maintaining a Z_0 -matched source and terminating impedance. The increased mismatch uncertainty due to impedance variations of source and load with frequency may degrade the above accuracy figures by an order of magnitude. A technique producing an equivalent Z_0 -matched source impedance is discussed in Reference 9.

Swept or Stepped in Frequency

Measurement of Attenuation Less than 30 dB

- (1) Audio-substitution dual channel (Fig. 8 with addition of level-stabilized swept radio-frequency source)

Accuracy: 0.05 to 0.1 dB

Advantage: Potentially most accurate system

Disadvantage: Impedance variation

- (2) Audio ratio technique (Fig. 12)

Accuracy: 0.1 to 0.3 dB

Advantage: Rapid measurement, simplicity

Disadvantage: Curvilinear reference lines which can be straightened through the use of memory circuits

- (3) Sampling (Fig. 13)

Accuracy: 0.1 to 0.5 dB dependent on uncertainty of system precalibration

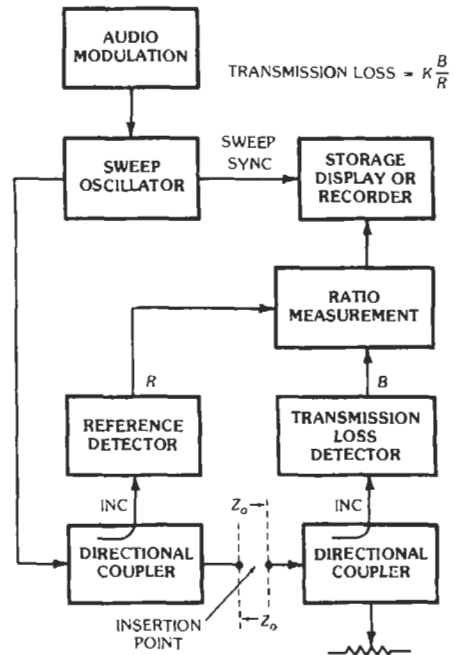


Fig. 12. Diagram of broadband direct-ratio system (References 2 and 9).

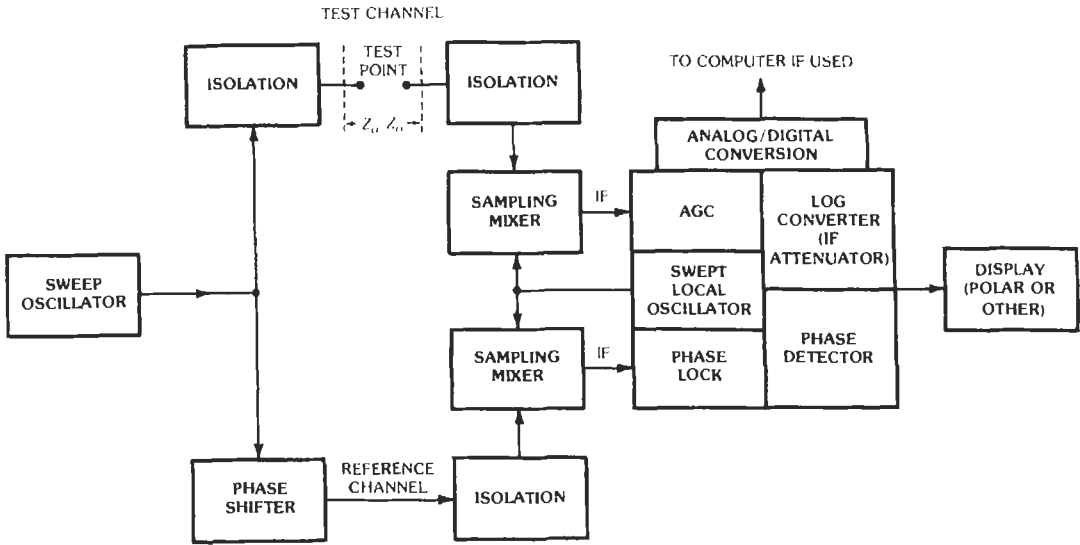


Fig. 13. Sampling conversion method of attenuation measurement (Reference 8).

- Advantages: (A) Digital technique permits use of logic circuitry for automatic correction of measurement errors
 (B) Extremely rapid measurement technique

Disadvantage: Very complex

Measurement of Attenuation Greater than 30 dB

Sampling (Fig. 13)

Accuracy: 0.5–2 dB to 50 or 60 dB dependent on uncertainty of system precalibration

- Advantages: (A) Digital technique permits use of logic circuitry for automatic correction of measurement errors
 (B) Extremely rapid measurement technique

Disadvantage: Very complex

The aforementioned measurement techniques are not intended to be all-inclusive. These are representative of current practices only. No mention has been made of specialized power ratio, modulated subcarrier, radio-frequency substitution, homodyne, impedance, or self-calibrating methods, or range extension by using a supervised change of incident power level.

REFERENCES ON ATTENUATION MEASUREMENT

1. Sorger, G. U., and Weinschel, B. O., "Comparison of Deviations from Square Law for RF Crystal Diodes and Barretters," *IRE Transactions on Instrumentation*, Vol. I-8, Dec. 1959, pp. 103–111.
2. Weinschel, B. O., "Measurement of Microwave Parameters by the Ratio Method," *Microwave Journal*, Aug. 1969, pp. 69–73.
3. Terman, F. E., "Linear Detection of Heterodyne Signals," *Electronics*, Nov. 1930, pp. 386–387.
4. Harnett, D. E., and Case, N. P., "The Design and Testing of Multi-Range Receivers," *Proceedings of the IRE*, Vol. 23, June 1935, pp. 578–593.
5. Gainsborough, G. F., "A Method of Calibrating Standard Signal Generators and Radio Frequency Attenuators," *Journal of the IEE*, Vol. 94, Pt. III, May 1947, pp. 203–210.
6. Weinschel, B. O., Sorger, G. U., and Hedrich, A. L., "Relative Voltmeter for VHF/UHF Signal Generator Attenuator Calibration," *IRE Transactions on Instrumentation*, Vol. I-8, Mar. 1959, pp. 22–31.
7. Schafer, G. E., and Rumpfelt, A. Y., "Mismatch Errors in Cascade-Connected Variable Attenuators," *IRE Transactions on Microwave Theory and Techniques*, Vol. MTT-7, Oct. 1959, pp. 447–453.

8. Adam, S. F., "A New Precision Automatic Microwave Measurement System," *IEEE Transactions on Instrumentation and Measurement*, Vol. IM-17, Dec. 1968, pp. 308-313.
9. Engen, G. F., "Amplitude Stabilization of a Microwave Signal Source," *IRE Transactions on Microwave Theory and Techniques*, Vol. MTT-6, Apr. 1958, pp. 202-206.

OTHER GENERAL REFERENCES

1. MIL-A-3933, General Specification for Fixed Attenuators.
2. Section 14A on Attenuators of MIL-HDBK-216 on RF Transmission Lines and Fittings.

References 1 and 2 can be obtained through:

Commanding Officer
Naval Publications and Forms Center
5801 Tabor Avenue
Philadelphia, PA 19120

3. IEEE STD-474-1973, Specifications and Test Methods for Fixed and Variable Attenuators, DC to 40 GHz.
4. IEEE STD-100-1977, Dictionary of Electrical & Electronics Terms

References 3 and 4 can be obtained through:

IEEE Service Center
445 Hoes Lane
Piscataway, NJ 08854

12 Measurements and Analysis

*Arthur Fong, Robert Coackley, John J. Dupre, Stephen J. Greer,
Siegfried Linkwitz, John R. Page, Ronald E. Pratt,
Douglas K. Rytting, and Gary D. Sasaki*

- Impedance Bridges 12-3
- Fundamental Alternating-Current or Wheatstone Bridge
 - Wagner Ground Connection
 - Capacitor Balance
 - Series-Resistance-Capacitance Bridge
 - Wien Bridge
 - Owen Bridge
 - Resonance Bridge
 - Maxwell Bridge
 - Hay Bridge
 - Schering Bridge
 - Substitution Method for High Impedances
 - Measurement With Capacitor in Series With Unknown
 - Measurement of Direct Capacitance
 - Felici Mutual-Inductance Balance
 - Mutual-Inductance Capacitance Balance
 - Hybrid-Coil Method
 - Q of Resonant Circuit by Bandwidth
 - Q Meter (Hewlett-Packard 4342A)
 - Twin-T Admittance-Measuring Circuit (General Radio Type 821-A)
 - Ratio-Arm Bridges (Wayne Kerr)
 - Automatic Impedance Meters

- Network Analysis *12-13*
 - Sources
 - Receivers
 - Test Sets
 - Displays and Output
 - Six-Port Network
 - Error Correction
 - Frequency and Time-Domain Relationships
 - Large-Signal Measurements

- Signal Analysis *12-17*
 - Amplitude Measurement Range
 - Signal-Analysis Characteristics as Determined by the IF Filter, Detector, and Video Amplifier

- Time and Frequency Measurement *12-20*
 - Time Measurement
 - Frequency Measurement
 - Frequency and Time-Interval Analysis

- RF and Microwave Power Measurements *12-24*
 - Thermistor Sensors
 - Thermal Converters
 - Thermocouple Sensors
 - Diode Sensors
 - Power-Measurement Definitions

- Microwave-Link Analysis *12-26*
 - Insertion Loss or Gain
 - Amplitude Response
 - Envelope-Delay Distortion (EDD)
 - Measurement of Group-Delay Distortion
 - Return Loss
 - Measurement of Return Loss
 - Baseband Measurements
 - Carrier/Noise Measurement
 - Diagnostic Measurements and System Performance

- Computer Control of Instruments *12-33*
 - IEEE 488.1 or IEC 625-1 General Purpose Interface Bus (GPIB)
 - IEEE 488.2 or IEC 625-2
 - Standard Commands for Programmable Instruments (SCPI)

- Electromagnetic Compatibility, Interference, and Susceptibility *12-35*
 - EMC/EMI/EMS
 - Regulations
 - EMI/EMS Measurements

The aim of this chapter is to provide a condensed description of methods of electronic measurements. Science and technology would indeed be vague without the ability to measure. Lord Kelvin cautioned that "Knowledge not expressible in numbers is of a meager and unsatisfactory kind"; he was identifying an essential aspect of scientific knowledge.

In the past, measurements of voltage, frequency, impedance, or power were made by using bridges or substitution methods with prototype standards derived from primary standards. Now, using digital techniques, one can create signals with precise voltage, frequency, and phase and can measure signals with the same precision. Many measurements are now made by stimulating the network or system with a precise signal, measuring the response, and using a microprocessor to compute the impedance, gain, phase, or other required results.

Another method for determination of measured quantities uses the technique of ratio measurements. For example, the complex ratio of E/I can be determined very precisely for a given frequency; then a microprocessor instantaneously calculates the impedance or response. In addition, the sequential steps of measurements, computation, and final display can also be controlled by the same microprocessor. Touch and read instruments were created by these technologies.

These measurements topics will be covered:

- Impedance
- Networks
- Signals
- Time and Frequency
- Power
- Microwave Links
- IEEE 488 Data Bus
- EMI/EMC/EMS

It is assumed that the reader is already familiar with such fundamental topics as voltage, current, gain, phase, distortion, etc. Specialized measurements for linear systems, digital signals, time domains, fields, magnetics, etc., are left for specialized publications. (Some chapters in this book covering specialized topics do contain information on the associated measurements.)

IMPEDANCE BRIDGES

In the diagrams of bridges in this section, the source (generator) and the detector (headphones) may be interchanged as dictated by the location of grounds. For all but the lowest frequencies, a shielded transformer is required at either the input or output (but not usually both) terminals of the bridge. The detector is chosen according to the frequency of the source. When insensitivity of the ear makes direct use of headphones impractical, a simple radio receiver or its equivalent is essential. Some selectivity is desirable to discriminate against harmonics, for the bridge is often frequency sensitive. The source may be modulated to obtain an

audible signal, but greater sensitivity and discrimination against interference are obtained by the use of a continuous-wave source and a heterodyne detector. An oscilloscope is sometimes preferred for observing nulls. In this case, it is convenient to have an audible output signal available for the preliminary setup and for locating trouble, since much can be deduced from the quality of the audible signal that would not be apparent from observation of amplitude only.

Fundamental Alternating-Current or Wheatstone Bridge

Refer to Fig. 1. The balance condition is $Z_x = Z_s Z_a / Z_b$. Maximum sensitivity exists when Z_d is the conjugate of the bridge output impedance and Z_x is the conjugate of its input impedance. Greatest sensitivity exists when the bridge arms are equal; for example, for resistive arms

$$Z_d = Z_a = Z_b = Z_x = Z_s = Z_g$$

Wagner Ground Connection

None of the bridge elements (Fig. 2) is grounded directly. First balance the bridge with the switch to B . Throw the switch to G , and rebalance by means of R and C . Recheck the bridge balance and repeat as required. The capacitor balance C is necessary only when the frequency is above the audio range. The transformer may have only a single shield as shown, with the capacitance of the secondary to the shield kept to a minimum.

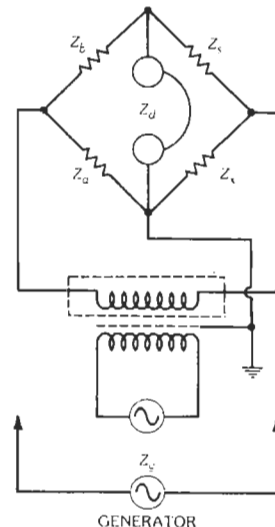


Fig. 1. Fundamental ac bridge.

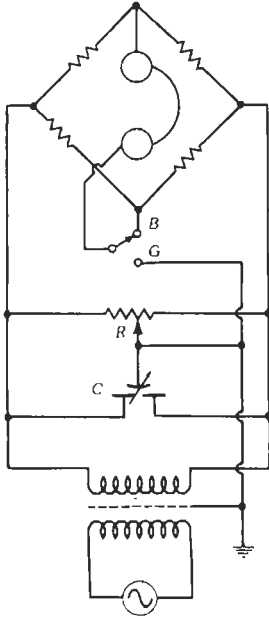


Fig. 2. Wagner ground connection.

Capacitor Balance

A capacitor balance is useful when one point of the bridge must be grounded directly and only a simple shielded transformer is used (Fig. 3). Balance the bridge, then open the two arms at *P* and *Q*. Rebalance by auxiliary capacitor *C*. Close *P* and *Q* and check the balance.

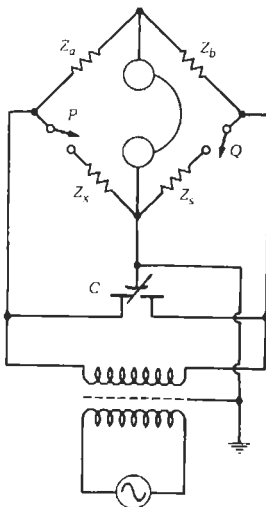


Fig. 3. Capacitor balance.

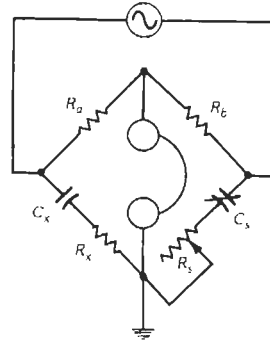


Fig. 4. Series-resistance-capacitance bridge.

Series-Resistance-Capacitance Bridge

In the bridge of Fig. 4:

$$C_x = C_s R_b / R_a$$

$$R_x = R_s R_a / R_b$$

Wien Bridge

In the bridge of Fig. 5:

$$C_x C_s = (R_b / R_a) - (R_s / R_x)$$

$$C_s / C_x = 1 / \omega^2 R_s R_x$$

For measurement of frequency, or in a frequency-selective application, if we make $C_x = C_s$, $R_x = R_s$, and $R_b = 2R_a$, then

$$f = (2\pi C_s R_s)^{-1}$$

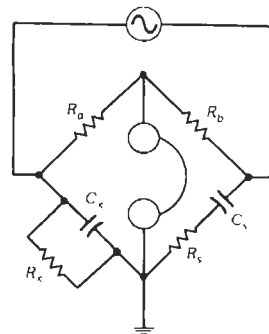


Fig. 5. Wien bridge.

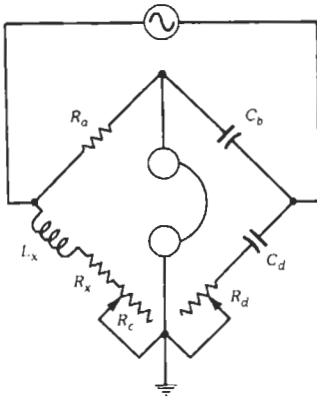


Fig. 6. Owen bridge.

Owen Bridge

In the bridge of Fig. 6:

$$L_x = C_b R_a R_d$$

$$R_x = (C_b R_a / C_d) - R_c$$

Resonance Bridge

In the bridge of Fig. 7:

$$\omega^2 LC = 1$$

$$R_x = R_s R_a / R_b$$

Maxwell Bridge

In the bridge of Fig. 8:

$$L_x = R_a R_b C_s$$

$$R_x = R_a R_b / R_s$$

$$Q_x = \omega(L_x / R_x) = \omega C_s R_s$$

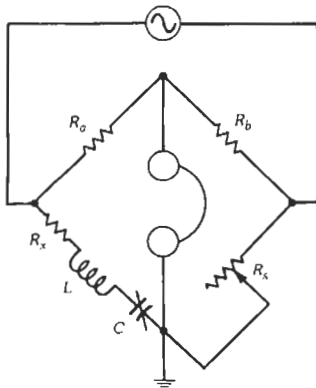


Fig. 7. Resonance bridge.

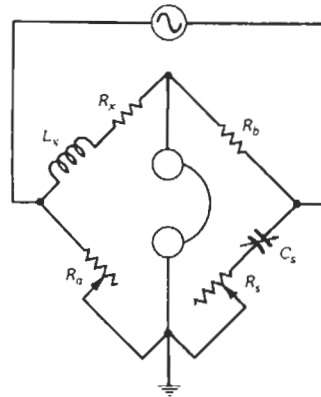


Fig. 8. Maxwell bridge.

Hay Bridge

The bridge of Fig. 9 is for the measurement of large inductance.

$$L_x = R_a R_b C_s / (1 + \omega^2 C_s^2 R_s^2)$$

$$Q_x = \omega L_x / R_x = (\omega C_s R_s)^{-1}$$

Schering Bridge

In the bridge of Fig. 10:

$$C_x = C_s R_b / R_a$$

$$1/Q_x = \omega C_x R_x = \omega C_b R_b$$

Substitution Method for High Impedances

Refer to Fig. 11.

Initial balance (unknown terminals $x - x$ open):

C_s' and R_s'

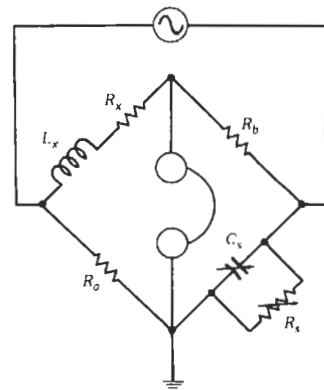


Fig. 9. Hay bridge.

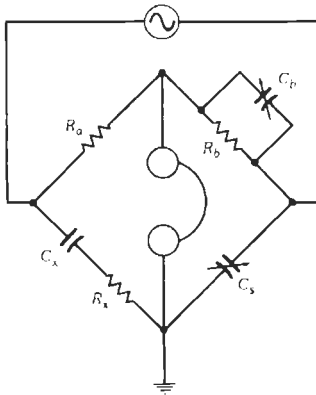


Fig. 10. Schering bridge.

Final balance (unknown connected to $x - x$):

$$C_s'' \text{ and } R_s''$$

Then when $R_x > 10/\omega C_s'$, there results, with error < 1 percent

$$C_x = C_s' - C_s''$$

The parallel resistance is

$$R_x = [\omega^2 C_s'^2 (R_s' - R_s'')]^{-1}$$

If unknown is an inductor

$$L_x = -(\omega^2 C_x)^{-1} = [\omega^2 (C_s'' - C_s')]^{-1}$$

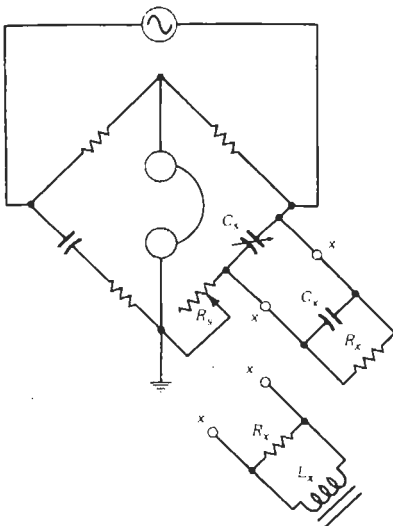


Fig. 11. Substitution method.

Measurement With Capacitor in Series With Unknown

Refer to Fig. 12.

Initial balance (unknown terminals $x - x$ short-circuited):

$$C_s' \text{ and } R_s'$$

Final balance ($x - x$ unshorted):

$$C_s'' \text{ and } R_s''$$

Then

$$R_x = (R_s'' - R_s')R_a/R_b$$

$$C_x = \frac{R_b C_s' C_s''}{R_a (C_s' - C_s'')}$$

$$= \frac{R_b}{R_a} C_s' \left(\frac{C_s'}{C_s' - C_s''} - 1 \right)$$

When $C_s'' > C_s'$

$$L_x = \frac{1}{\omega^2} \frac{R_a}{R_b C_s'} \left(1 - \frac{C_s'}{C_s''} \right)$$

Measurement of Direct Capacitance

Refer to Fig. 13.

Connection of N to N' places C_{nq} across the detector and C_{np} across R_b , which requires only a small readjustment of R_s .

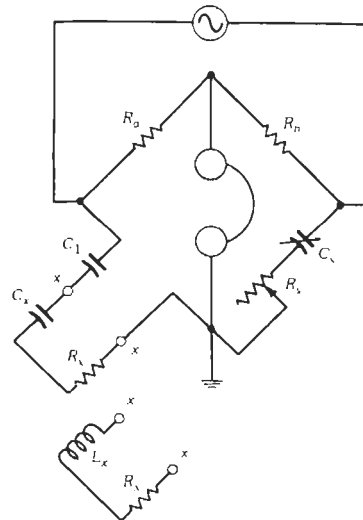


Fig. 12. Measurement with capacitor in series with unknown.

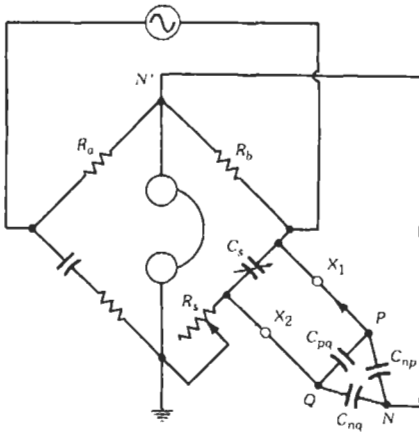


Fig. 13. Measurement of direct capacitance.

Initial balance: Lead from *P* disconnected from *X*₁ but lying as close to the connected position as practical.
 Final balance: Lead connected to *X*₁.
 By the substitution method above

$$C_{pq} = C_s' - C_s''$$

Felici Mutual-Inductance Balance

At the null (Fig. 14):

$$M_x = -M_s$$

This is useful at lower frequencies where capacitive reactances associated with windings are negligibly small.

Mutual-Inductance Capacitance Balance

With a low-loss capacitor (Fig. 15), at the null:

$$M_x = 1/\omega^2 C_s$$

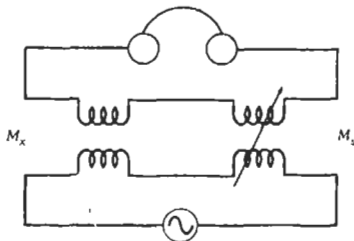


Fig. 14. Felici mutual-inductance balance.

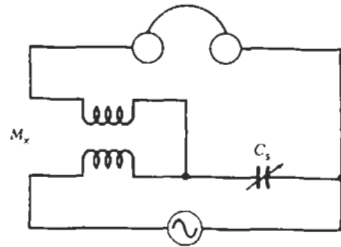


Fig. 15. Mutual-inductance capacitance balance.

Hybrid-Coil Method

At the null (Fig. 16):

$$Z_1 = Z_2$$

The transformer secondaries must be accurately matched and balanced to ground. This is useful at audio and carrier frequencies.

Q of Resonant Circuit by Bandwidth

The method of Fig. 17 may be used to evaluate *Q* by finding the 3-dB, or half-power, points. The source should be loosely coupled to the circuit. Adjust the frequency to each side of resonance, noting the bandwidth between the points where *V* = 0.71 × (*V* at resonance). Then

$$Q = (\text{resonance frequency})/(\text{bandwidth})$$

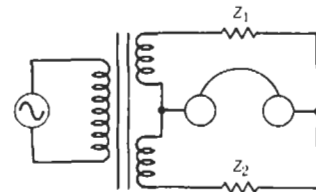


Fig. 16. Hybrid-coil method.

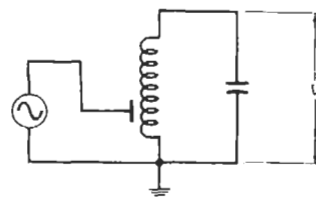


Fig. 17. Measurement of *Q* by bandwidth.

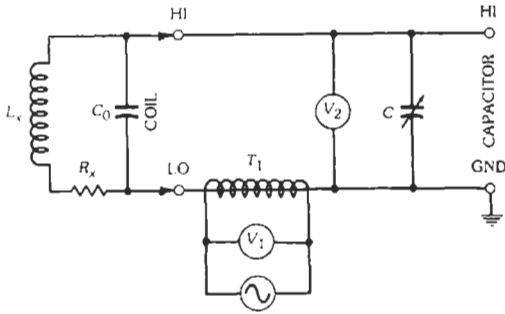


Fig. 18. Q meter.

Q Meter (Hewlett-Packard 4342A)

Refer to Fig. 18. In this circuit, T_1 is a wideband transformer with n turns in the primary and one turn in the secondary. The secondary impedance is 0.001 ohm. The combination $L_x R_x C_0$ represents an unknown coil plugged into the COIL terminals for measurement; V_2 is a very high impedance voltmeter. With this arrangement,

$$Q = nV_2/V_1$$

Correction of Q Reading—The value of Q corrected for distributed capacitance C_0 of the coil is given by

$$Q_{true} = Q[(C + C_0)/C]$$

where,

- Q = reading of Q-meter (corrected for internal resistors R_1 and R_2 if necessary),
- C = capacitance reading of Q-meter.

Measurement of C_0 and True L_x —The plot of $1/f^2$ against C is a straight line (Fig. 19).

$$L_x = \text{true inductance} = \frac{1/f_2^2 - 1/f_1^2}{4\pi^2(C_2 - C_1)}$$

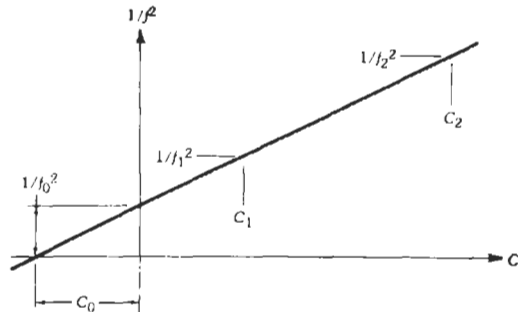


Fig. 19. Plot of C and $1/f^2$.

C_0 = negative intercept

f_0 = natural frequency of coil

When only two readings are taken and $f_1/f_2 = 2.00$

$$C_0 = (C_2 - 4C_1)/3$$

With values in microhenrys; megahertz, and picofarads

$$L_x = 19\,000/f_2^2(C_2 - C_1)$$

Measurement of Admittance—An initial reading, $C'Q'$, is taken as in Fig. 20A ($L R_p$ is any suitable coil). For the final reading, $C''Q''$ (Fig. 20B):

$$1/Z = Y = G + jB = 1/R_p + j\omega C$$

Then

$$C = C' - C''$$

$$1/Q = G/\omega C$$

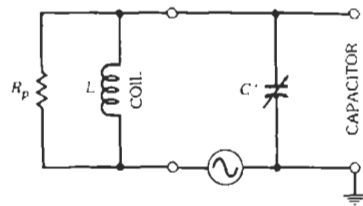
$$= C'/C(1000/Q'' - 1000/Q') \times 10^{-3}$$

If Z is inductive

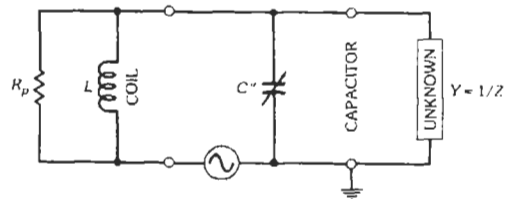
$$C'' > C'$$

Measurement of Impedances Lower Than Those Directly Measurable—For the initial reading, $C'Q'$, the CAPACITOR terminals are open (Fig. 21A).

On the second reading, $C''Q''$, a capacitive divider, $C_a C_b$ (Fig. 21B), is connected to the CAPACITOR terminals.



(A) Initial reading.



(B) Final reading.

Fig. 20. Measurement of admittance.

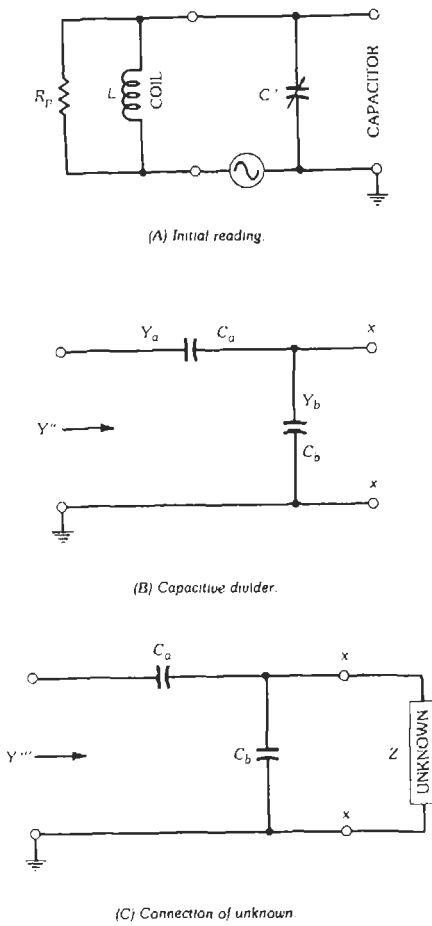


Fig. 21. Measurement of low impedances.

For the final reading, $C''Q''$, the unknown is connected to $x - x$ (Fig. 21C). Admittances Y_a and Y_b are

$$Y_a = G_a + j\omega C_a \quad Y_b = G_b + j\omega C_b$$

with G_a and G_b not shown in the diagrams. Then the unknown impedance is

$$Z = [Y_a / (Y_a + Y_b)]^2 (Y''' - Y'')^{-1} - (Y_a + Y_b)^{-1} \text{ ohms}$$

where, with capacitance in picofarads and $\omega = 2\pi \times$ frequency in megahertz

$$(Y''' - Y'')^{-1} = \frac{10^6 / \omega}{C'(1000/Q''' - 1000/Q'') \times 10^{-3} + j(C'' - C''')}$$

Usually G_a and G_b may be neglected; then there results

$$Z = \left(\frac{1}{1 + C_b/C_a} \right)^2 (Y''' - Y'')^{-1} + j \frac{10^6}{\omega(C_a + C_b)} \text{ ohms}$$

For many measurements, C_a may be 100 picofarads. Capacitance $C_b = 0$ for very low values of Z and for highly reactive values of Z . For unknowns that are principally resistive and of low or medium value, C_b may take sizes up to 300 to 500 picofarads. When $C_b = 0$

$$Z = (Y''' - Y'')^{-1} + j(10^6 / \omega C_a) \text{ ohms}$$

and the "second" reading above becomes the "initial," with $C' = C''$ in the equations.

Measurement of Coupling Coefficient of Loosely Coupled Coils—The coefficient of coupling

$$k = M / (L_1 L_2)^{1/2}$$

between two high- Q coils can be obtained by measuring the inductance L with S_1 closed and again with S_1 open (Fig. 22). From these two measurements

$$k = (1 - L_{\text{closed}} / L_{\text{open}})^{1/2}$$

When the coil self-inductances are known, a measurement of L_a and L_b (Fig. 23) yields

$$k = (L_a - L_b) / 4(L_1 L_2)^{1/2}$$

If $L_1 = L_2$

$$k = (L_a - L_b) / (L_a + L_b)$$

Neither of the above methods provides adequate precision when two high- Q coils are only about critically coupled. In that case, the Q of each coil is measured with the other coil open-circuited. Then the coupled coils ($L_1 R_1$ and $L_2 R_2$) and a low-loss adjustable capacitor (C_2) are connected to the Q -meter as shown in Fig. 24, C_2 is disconnected, and C is adjusted to maximize

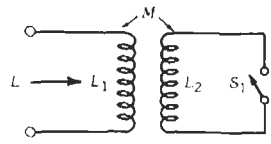


Fig. 22. Method of determining coefficient of coupling between two coils.

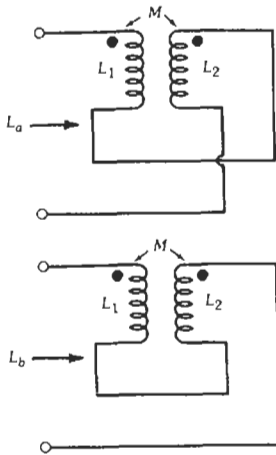


Fig. 23. Method of determining coefficient of coupling when self-inductances are known.

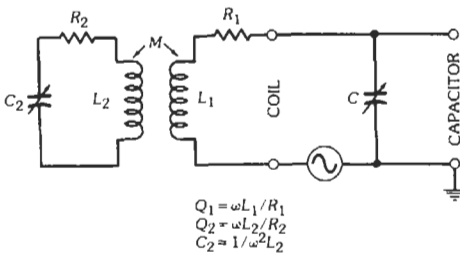


Fig. 24. Method of determining coefficient of coupling when coils are only about critically coupled.

the Q -meter reading; C_2 is then connected and adjusted to minimize the reading. If the final reading is Q_0

$$K = [1/Q_2(1/Q_0 - 1/Q_1)]^{1/2}$$

If the final reading is too small to be read accurately, Q_2 may be reduced by inserting a small resistance in series with L_2R_2 .

Twin-T Admittance-Measuring Circuit (General Radio Type 821-A)

The circuit in Fig. 25 may be used for measuring admittances in a range somewhat exceeding 400 kilohertz to 40 megahertz. It is applicable to the special measuring techniques described above for the Q -meter.

Conditions for a null in the output are

$$G + G_1 = R\omega^2 C_1 C_2 (1 + C_g/C_3)$$

$$C + C_b = 1/\omega^2 L - C_1 C_2 (1/C_1 + 1/C_2 + 1/C_3)$$

With the unknown disconnected, call the initial balance C_b' and C_g' . With the unknown connected, the final balance is C_b'' and C_g'' . Then the components of the unknown, $Y = G + j\omega C$, are

$$C = C_b' - C_b''$$

$$G = (R\omega^2 C_1 C_2 / C_3)(C_g'' - C_g')$$

Ratio-Arm Bridges (Wayne Kerr)*

Transformer ratio-arm bridges can be designed to operate at radio frequencies up to about 250 MHz. Beyond that point, other forms of bridges based on transmission lines become practicable.

Fig. 26 illustrates a practical circuit for a bridge

* From Calvert, R. *The Transformer Ratio-Arm Bridge*. Bognor Regis, Sussex, England: Wayne Kerr Co. Ltd.

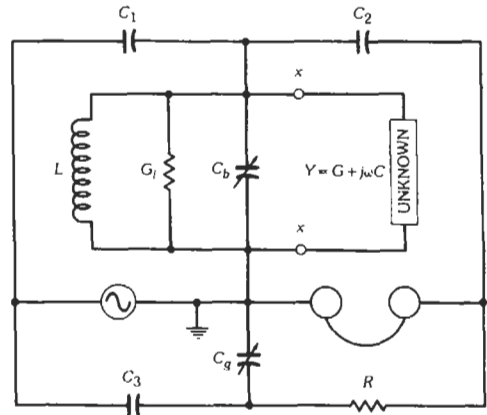


Fig. 25. Twin-T admittance-measuring circuit.

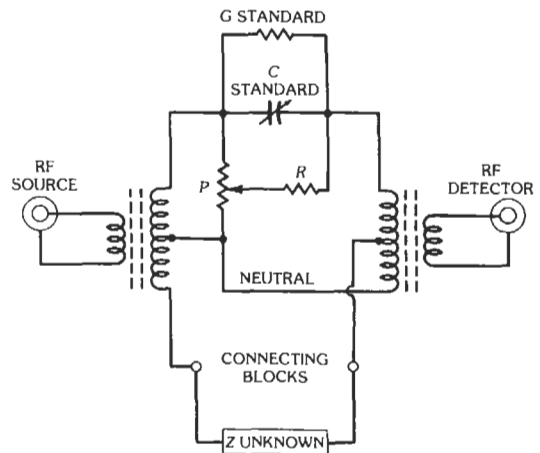


Fig. 26. Ratio-arm bridge.

capable of operating at frequencies up to 100 MHz. The transformers are formed by winding thin silver tapes onto ferrite or ferrous-dust ring cores, which are mounted inside individual screening cans. Drums of low-inductance resistors forming fixed conductance standards are arranged to engage with spring contacts. A variable conductance for interpolation is formed by means of resistor R , which is fed with a voltage derived from a resistive potential divider, P .

Radio-frequency bridge measurements require that considerable care should be taken in setting up the apparatus. Any leakage of power from the source to the detector that bypasses the bridge network will give errors.

Automatic Impedance Meters

The availability of digital voltmeter technology has created a generation of touch-and-read impedance meters. One basic technique is to impress a known current or voltage upon an impedance and measure the magnitude and phase of the voltage or current resulting from it (Fig. 27). In its simplest form, the ratio of the complex voltage to the complex current is the value of the impedance:

$$Z \angle \theta = E \angle \theta_1 / I \angle \theta_2$$

Automation can also be achieved by using feedback to balance a bridge with an unknown impedance and a known standard impedance. The complex voltage required for balance is a measure of the unknown impedance.

At frequencies above 10 MHz, it is convenient to compare an impedance to a precise 50-ohm coaxial transmission line and measure the reflection coefficient. The impedance is calculated using microprocessor techniques.

An example of each technique is shown below.

Four-Terminal-Pair Bridge, 5 Hz-13 MHz—Refer to Fig. 28. A four-terminal-pair configuration is

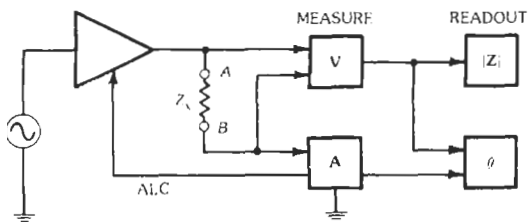


Fig. 27. Vector impedance meter. (From Alonzo, G. J., et al. "Direct-Reading Vector Impedance Meters." Hewlett-Packard Journal, January 1967. © 1967 Hewlett-Packard, used with permission.)

used to avoid errors caused by mutual coupling between leads. The bridge provides a complex voltage, V_x , across the device under test (DUT) and another complex voltage V_r proportional to the current through the DUT at bridge balance. The bridge is balanced by a heterodyne method through the frequency range of 5 Hz to 13 MHz.*

A vector ratio detector is used to measure V_x/V_r . Since

$$V_x/Z_x = i_x = i_r = -V_r/R_r$$

therefore,

$$Z_x = -R_r(V_x/V_r)$$

Automatic LCR Meters, 1 MHz—Refer to Fig. 29. A reference voltage, e_r , is applied to the unknown, and the feedback voltage, e_f , feeds the standard resistor for C - G measurements. The voltages are reversed for L - R measurements. If the bridge is not balanced, an unbalanced current, i_d , flows into the current detector, which produces an error voltage, e_d . This voltage is amplified, phase-detected, and rectified to produce dc voltages E_1 proportional to the real part of e_d and E_2 proportional to the imaginary part of e_d .†

Voltages E_1 and E_2 are integrated and used to modulate reference voltage e_r and je_r , thus creating e_1 and e_2 , which correspond to the real and imaginary parts of voltage e_d . At balance, the vector sum e_1 reduces i_d to zero.

At balance

$$G_x + j\omega C_x = -e_1/e_r R_r$$

or

$$R_x + j\omega L_x = -(e_2/e_r)R_r$$

A digital voltmeter is used to measure the vector ratio of e_1 and e_r .

Vector Impedance Analyzer, to 1000 MHz—Refer to Fig. 30. Outputs from a directional RF bridge channel and a compensated reference channel go to a receiver for synchronous detection and processing by a dual-slope integrator, thus providing a vector ratio that is a measure of the reflection coefficient, Γ . With a microprocessor, the impedance, admittance, inductance, capacitance, or Q of the device under test can be calculated.

The key component of the RF bridge is the balun transformer. It is wound with fine semirigid coaxial cable on a ferrite core. The construction must maintain

* Y. Narimatsu, et al, "A Versatile LF Impedance Analyzer," *H-P Journal*, Sept. 1981.

† Maeda, K., et al. "An Automatic, Precision 1-MHz Digital LCR Meter." *H-P Journal*, March 1974.

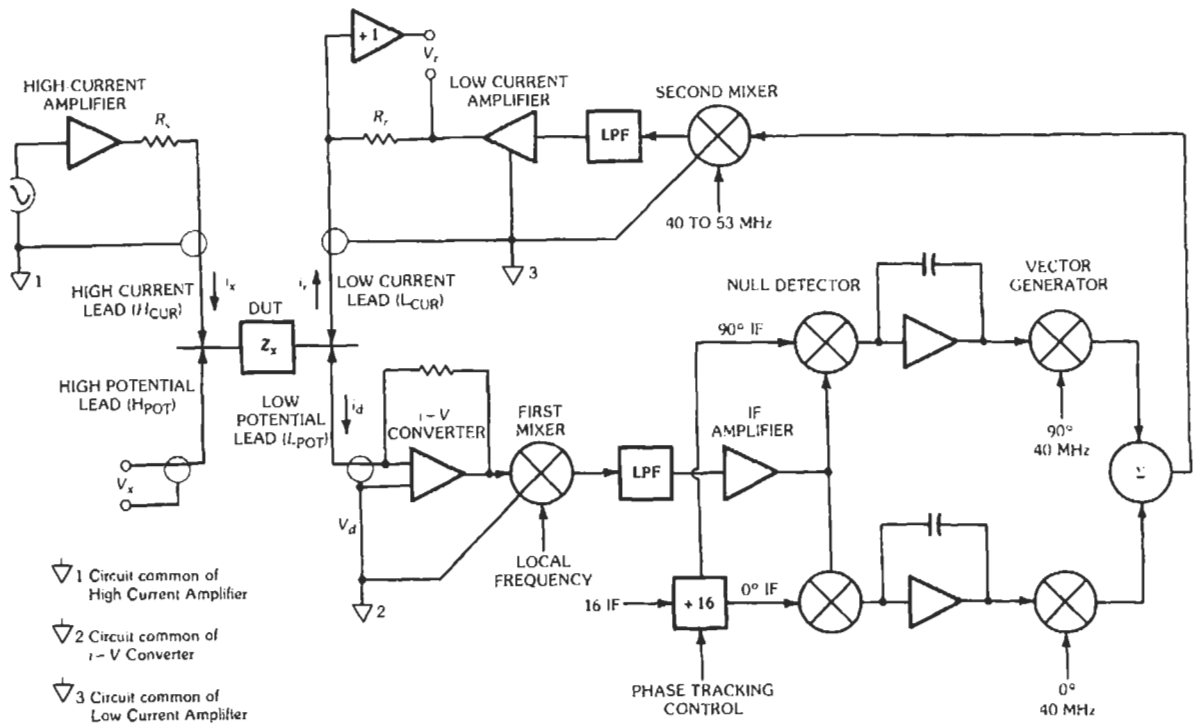


Fig. 28. 5-Hz to 13-MHz bridge. (From Narimatsu, Y., et al. "A Versatile LF Impedance Analyzer." Hewlett-Packard Journal, September 1981. © 1981 Hewlett-Packard, used with permission.)

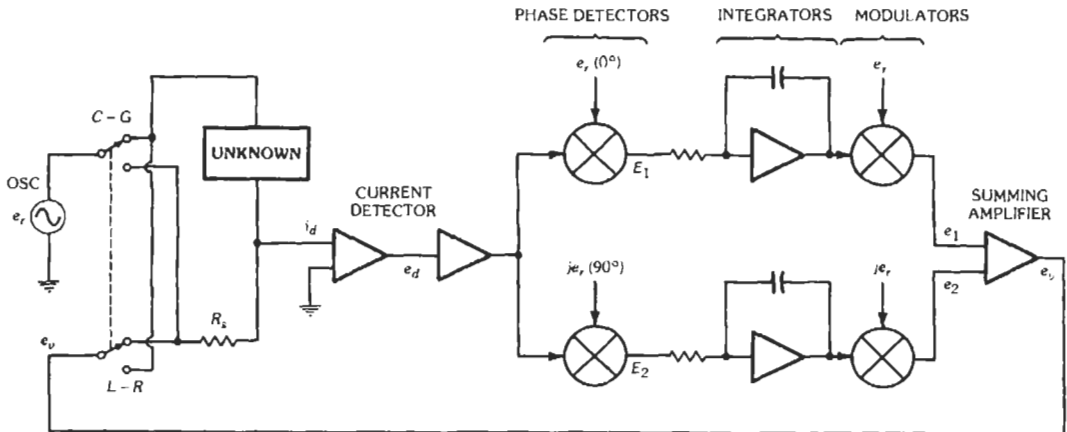


Fig. 29. Automatic LCR meter, 1 MHz. (From Maeda, K., et al. "An Automatic, Precision 1 MHz Digital LCR Meter." Hewlett-Packard Journal, March 1974. © 1974 Hewlett-Packard, used with permission.)

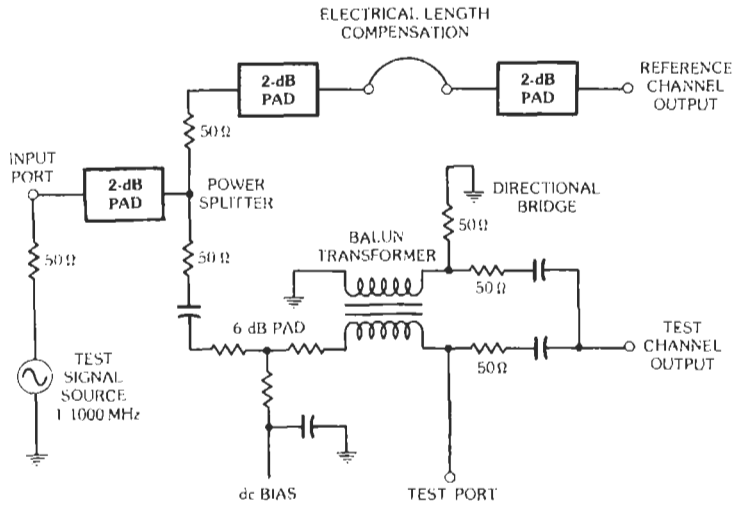


Fig. 30. RF impedance analyzer. (From Ichino, T., et al. "Vector Impedance Analysis at 1000 MHz." Hewlett-Packard Journal, January 1980. © 1980 Hewlett-Packard, used with permission.)

electrical balance up to 1000 MHz and must be temperature controlled.*

For a 50-ohm coaxial line, calculations are:

Impedance

$$(R + jX) = 50 \times (1 + \Gamma)/(1 - \Gamma)$$

Admittance

$$(G + jB) = 0.02 \times (1 - \Gamma)/(1 + \Gamma)$$

Inductance

$$L = X/2\pi f$$

Capacitance

$$C = B/2\pi f$$

Quality factor

$$Q = |X|/R$$

NETWORK ANALYSIS

A network analyzer consists of the five basic blocks shown in Fig. 31. At low frequencies, voltage and current can be measured by probes to determine the Z, Y, or h parameters (Fig. 32). Typically, a network analyzer is used to characterize small signal parameters. For higher frequencies (up through microwave frequencies), network analyzers primarily characterize the

magnitude and phase of reflection coefficients, transmission coefficients, and S parameters of networks (Fig. 33). (See Chapter 31 for a definition of S parameters.)

At microwave frequencies, it is difficult to build perfect structures that have no frequency-response, directivity, or port-match errors. It is widely accepted practice to reduce these errors with a mathematical error-correction procedure implemented with a computer.

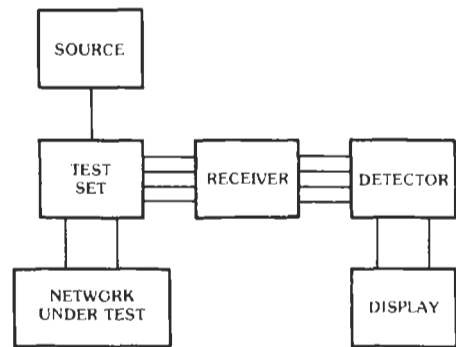
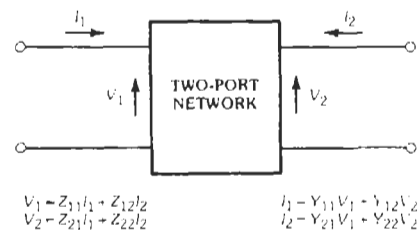


Fig. 31. Block diagram of a network analyzer.



$$\begin{aligned} V_1 &= Z_{11}I_1 - Z_{12}I_2 & I_1 &= Y_{11}V_1 - Y_{12}V_2 \\ V_2 &= -Z_{21}I_1 + Z_{22}I_2 & I_2 &= -Y_{21}V_1 + Y_{22}V_2 \end{aligned}$$

Fig. 32. Z and Y parameters typically used at low frequency.

* Ichino, T., et al. "Vector Impedance Analysis at 1000 MHz." H-P Journal, Jan. 1980.

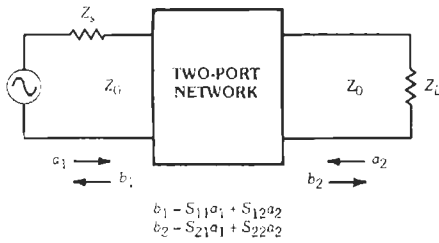


Fig. 33. Parameters used at higher frequencies.

Sources

The most common source used for network analysis is a sweep oscillator. This provides a real-time frequency response of the network characteristics. When greater frequency accuracy is required, synthesizers are used. However, the switching times can be long (10 ms to 50 ms). These synthesizers provide the best phase accuracy and repeatability, but at the expense of fast real-time measurements.

Receivers

There are three common receiver techniques, the diode detector, homodyne or self-mixing, and the heterodyne approach (both fundamental and harmonic), usually with multiple inputs.

Test Sets

The test-set portion provides connection of the source and receiver to the network under test. At low frequencies, voltage and current probes are sufficient, but at microwave frequencies, couplers, directional bridges, power splitters, and slotted lines are required.

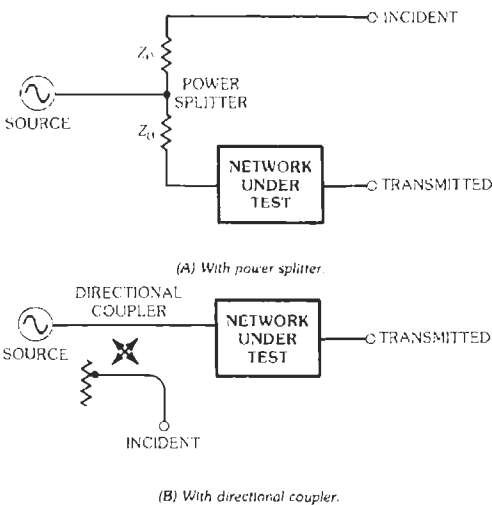


Fig. 34. Transmission measurements.

Test configurations are shown as follows: for transmission measurements, Fig. 34; for reflection measurements, Fig. 35; for S parameters, Fig. 36; and for transmission and reflection measurements, Fig. 37.

Displays and Output

The most useful display is the digital storage CRT system (Fig. 38). Data can be written into the digital display at any speed from real time to very slow (digital sweep for example) and then read out at a flicker-free rate. The display can also store traces, and trace math can be calculated to remove frequency-response errors, etc. Almost all commercial network analyzers have a digital I/O port to an external controller to provide for data transfer and operational commands.

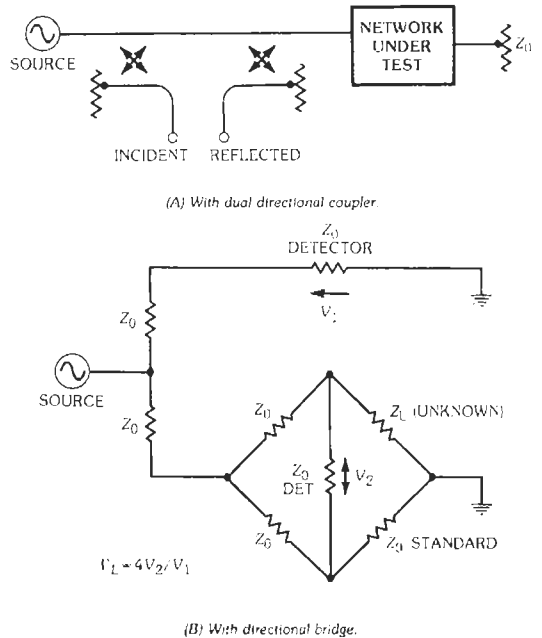


Fig. 35. Reflection measurements.

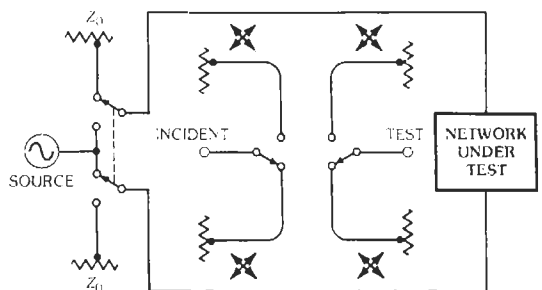


Fig. 36. Measurement of S parameters, using directional couplers with switching.

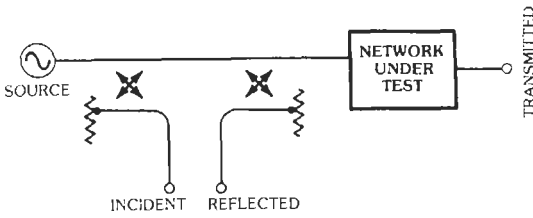


Fig. 37. Transmission and reflection measurements.

Six-Port Network

A six-port structure can be used to make magnitude and phase measurements using magnitude-only detectors such as the diode, thermistor, or bolometer. The general solution is shown below (refer to Fig. 39).

$$\rho_i = |b_i|^2 = |A_i a + B_i b|^2 \quad i = 3 \text{ to } 6$$

where A_i and B_i are functions of the S parameters of the six-port network, which can be determined by calibration. Rearranging yields

$$\rho_i = |B_i|^2 |b_i|^2 (A_i/B_i) \Gamma + 1|^2 \quad i = 3 \text{ to } 6$$

The above four equations can be solved for Γ , the reflection coefficient of the device under test:

$$\Gamma = \frac{\sum C_i \rho_i + j \sum S_i \rho_i}{\sum \beta_i \rho_i} \quad i = 3 \text{ to } 6$$

where C_i , S_i , and β_i are functions of A_i and B_i . Note that Γ consists of a sum of noncomplex terms.

A dual six-port structure can be used to measure the S parameters of a two-port network provided that the restriction $S_{12} = S_{21}$ is applied (Fig. 40). The mathematical approach of calibrating and measuring with the six-port structure is divided into two steps. Step one is the calibration of the six-port structure to determine

phase from the magnitude-only measurements. Step two is the same as the error-correction procedure for the four-port network analyzers. There has been some excellent work done on six-port networks by Glen Engen and Cletus Hoer of the National Bureau of Standards.

The approach can be realized in all frequency bands from audio to greater than 100 GHz. Great stability of the measurement system is required to achieve accurate results. There are also limitations in measuring active devices where power setting flexibility is needed and when S_{21} does not equal S_{12} . Typically, data gathering time is slow but accuracy is very good.

Error Correction

There are error-correction procedures by which the linear time invariant errors of the network analyzer can be characterized and then, by an inverse mathematical procedure, be removed from the measured data. See references 1 through 9 (listed at the end of this chapter) for descriptions of these error-correction procedures.

Frequency and Time-Domain Relationships

With the advent of the modern computer, measuring data in the frequency domain and then converting to the time domain by means of the inverse-Fourier transform has become practical. This ability to observe measurements in the time domain adds additional insight to microwave measurements. Many times in microwave measurements, circuit discontinuities are separated by lengths of transmission lines. Taking the Fourier transform of such a circuit allows us to isolate the various impulse responses to the circuit discontinuities (see Fig. 41C).

Also, by mathematically generating a step stimulus, a traditional time-domain reflectometer (TDR) display with high resolution and stability is achieved (see Fig.

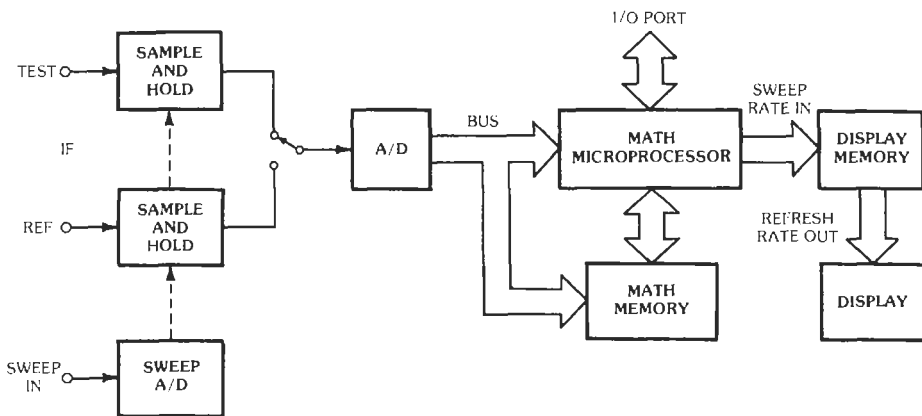


Fig. 38. Digital storage CRT system.

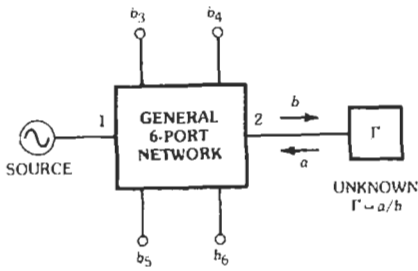
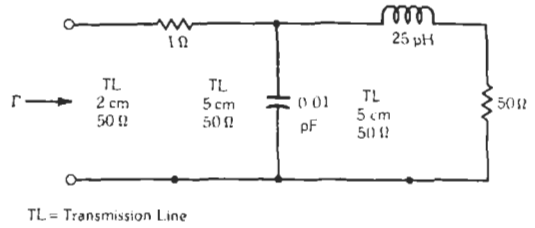


Fig. 39. Measurements with six-port system.



(A) Circuit.

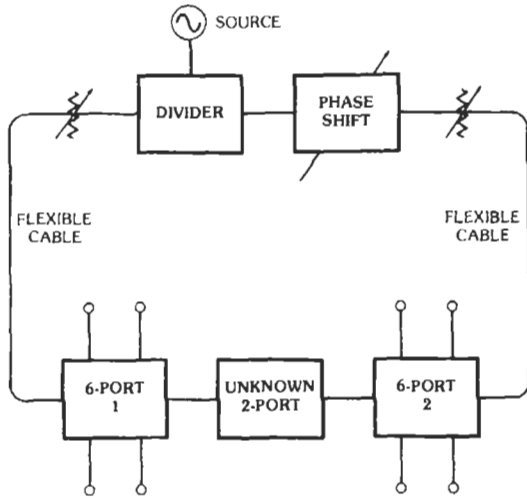
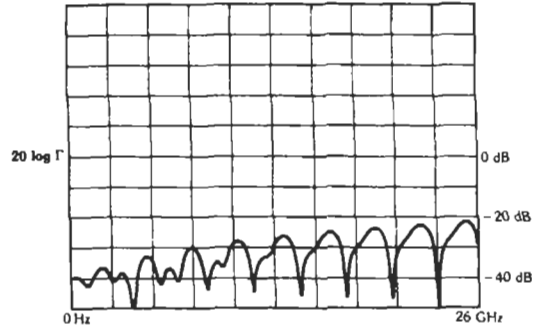


Fig. 40. A dual six-port structure used to measure a two-port network.

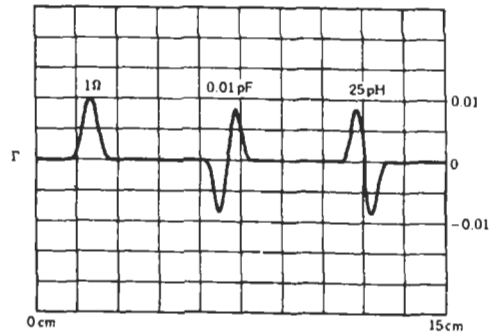
41D). With the time data isolated, it is possible to filter out unwanted time-domain responses and then transform back to the frequency domain. By means of this technique, connector discontinuities, launches, and other unwanted responses can be filtered (or deconvolved) from the original data (Fig. 42). This technique can greatly enhance measurement accuracy.

Large-Signal Measurements

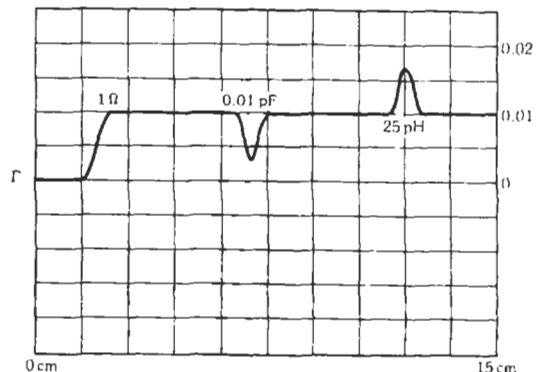
Devices and amplifiers are often characterized in other than class-A linear operation. A technique called load-pull is one way to observe the effect of changing device parameters as a function of load impedance and delivered power. The load-pull measurement system is shown in Fig. 43. The output tuner realizes most positive real load conditions. The power delivered to the device and load can be measured by power meters connected to the coupler side arms. The impedance of the load can be measured with a network analyzer. Fig. 44 shows contours of constant delivered power as a function of load impedance. These plots can be generated for the fundamental or the harmonics to give good insight into the nonlinear behavior. With these data,



(B) Frequency response.

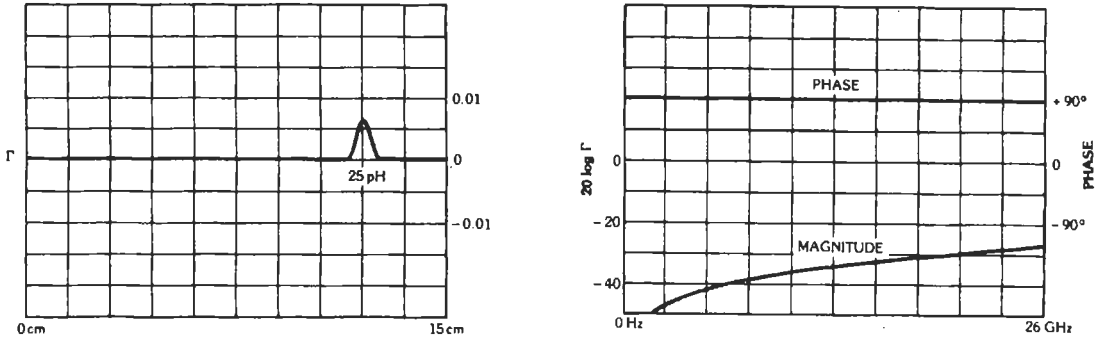


(C) Impulse time domain.



(D) Step time domain or time-domain reflectometer.

Fig. 41. Frequency and time relationships.



(A) Time-gated Inductor.

(B) Frequency response after gating.

Fig. 42. Time-domain filtering.

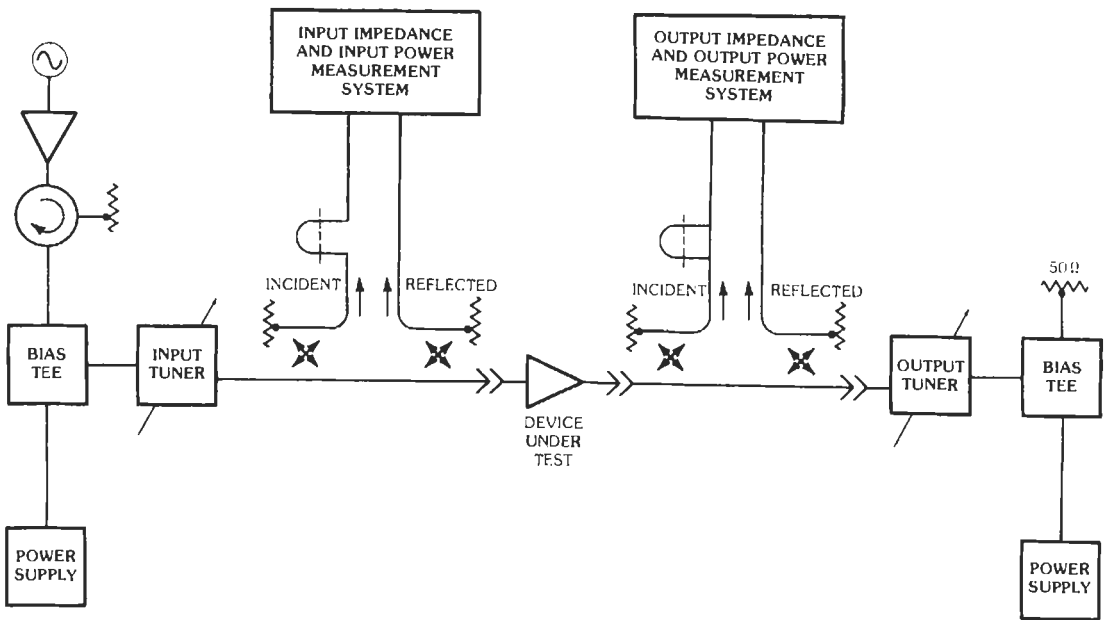


Fig. 43. Large-signal measurement system utilizing load-pull technique.

matching networks can be designed that will optimize the delivered power under nonlinear operating conditions.

SIGNAL ANALYSIS

Analysis of signals in the frequency domain is an important technique to determine the characteristics of audio, RF, and microwave signals. While the oscilloscope is the basic tool for time-domain analysis, the spectrum analyzer provides a display of signal amplitude as a function of frequency and is the basic tool for frequency-domain analysis. The relationship between the time and frequency domains is illustrated in Fig. 45 for the case of two sine waves. Signal amplitude vs time and frequency is shown in the three-dimensional representation of Fig. 45A. In the time-domain representa-

tion of Fig. 45B, the two signal components add at each instant of time, producing the composite as would be viewed on an oscilloscope. Fig. 45C is the frequency-domain view. Each frequency is represented by a vertical line, with the height representing the amplitude and the horizontal position representing the frequency. Although the information about the signal is the same in either domain, the frequency domain provides a tool to observe small but important signal characteristics such as a small amount of harmonic distortion.

Conceptually, a spectrum analyzer could consist of a parallel bank of rectangular bandpass filters each followed by an envelope detector as in Fig. 46. The filters cover adjacent frequency bands and do not overlap. The detector outputs are sequentially scanned to provide an amplitude-vs-frequency display on a cathode-ray tube. The resolution of the spectrum analyzer, that is, its

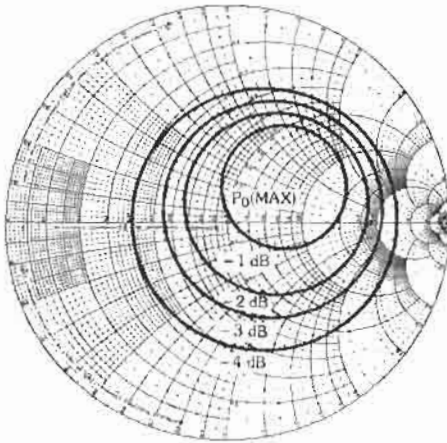


Fig. 44. Contours of constant delivered power vs load impedance under large-signal conditions.

ability to resolve closely spaced spectral lines, is determined by the width of the filter bandpass. Obviously, a very large number of filters would be required to obtain wide frequency spans with narrow resolution. Practical filters would have long time constants, thereby requiring a long time for their outputs to reach a steady-state value.

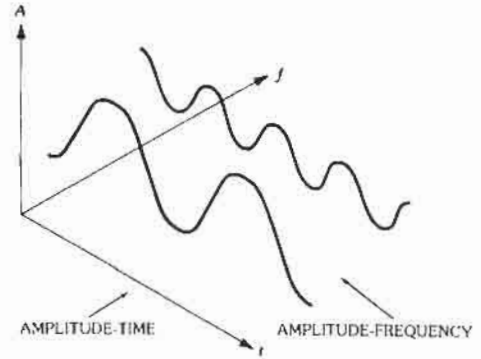
Another type of spectrum analyzer uses digital signal-processing techniques. Called a Fourier analyzer or dynamic signal analyzer, it is shown in Fig. 47. The input waveform is sampled, and each sample is digitized. The fast Fourier transform (FFT) is used to compute the spectrum from the sampled data. This type of analyzer solves the complexity and response-time problems of the parallel filter analyzer; however, ADC and computation speeds limit the input frequency to about 10 MHz.

The most common form of spectrum analyzer for RF and microwave frequencies is the swept superheterodyne receiver shown in simplified form in Fig. 48. A sweep ramp provides the tuning voltage for a voltage-tuned oscillator and is simultaneously applied to the horizontal-deflection system of a cathode-ray-tube display. The local-oscillator frequency is mixed with the input signal to produce an intermediate frequency (IF). The IF is detected, filtered, and applied to the vertical-deflection plates of the CRT. To achieve broad frequency coverage with a single local oscillator, harmonic mixing is often used so that

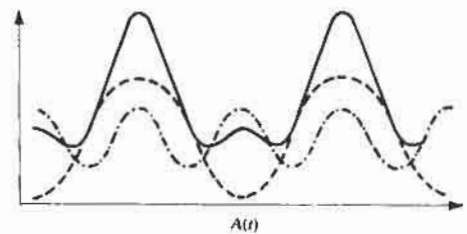
$$f_s = n f_{LO} \pm f_{IF}$$

where,

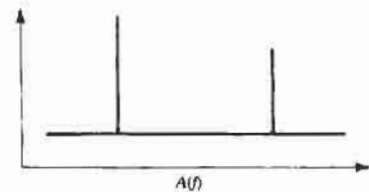
- f_s = signal frequency,
- f_{LO} = local-oscillator frequency,
- f_{IF} = intermediate frequency,
- $n = 1, 2, 3, \dots$



(A) Three dimensional coordinates showing time, frequency, and amplitude.



(B) View seen in the time domain



(C) View seen in the frequency domain.

Fig. 45. Relationship between time and frequency domains. (From Hewlett-Packard Application Note No. 150 on Spectrum Analysis. © Hewlett-Packard, used with permission.)

A tracking preselector filter at the input selects the desired response.

Amplitude Measurement Range

The range of amplitudes that may be measured on a spectrum analyzer depends on several factors and is diagrammed in Fig. 49. The minimum signal level is limited by the displayed noise that is a function of the noise figure and the IF bandwidth:

$$N = k T_0 B$$

where,

- N = displayed noise level, in watts,
- T_0 = absolute temperature, in kelvins,
- k = Boltzmann's constant,
- F = spectrum-analyzer noise figure,
- B = IF bandwidth.

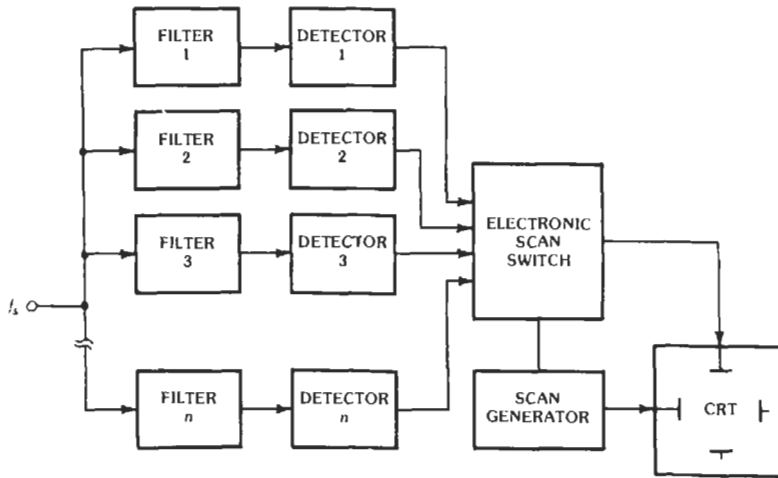


Fig. 46. A multiple-filter spectrum analyzer. (From Hewlett-Packard Application Note No. 150 on Spectrum Analysis. © Hewlett-Packard, used with permission.)

For example, a noise factor of 100 (20 dB) and an IF bandwidth of 10 kHz give a displayed noise level of 4×10^{-15} watt (-114 dBm).

The maximum signal level is the damage level of the input circuitry, which is typically +13 dBm for an input mixer and +30 dBm for an attenuator. Below that level, the input mixer compresses the input signal, causing an

amplitude inaccuracy. The maximum signal at the input mixer for less than 1 dB of gain compression is typically 0 dBm (1 milliwatt).

Distortion products produced within the spectrum analyzer pose yet another constraint on maximum signal amplitude. For a specified spurious level, the signal level to the mixer must not exceed a given level. The dynamic range is the ratio of the largest to the smallest signal amplitudes that can be displayed simultaneously with no internally generated spurious products being present. Dynamic range depends on the display range, the distortion characteristics, and the noise level.

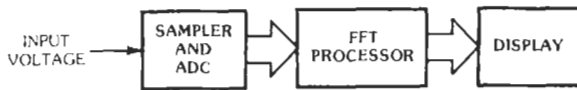


Fig. 47. Block diagram of a Fourier analyzer.

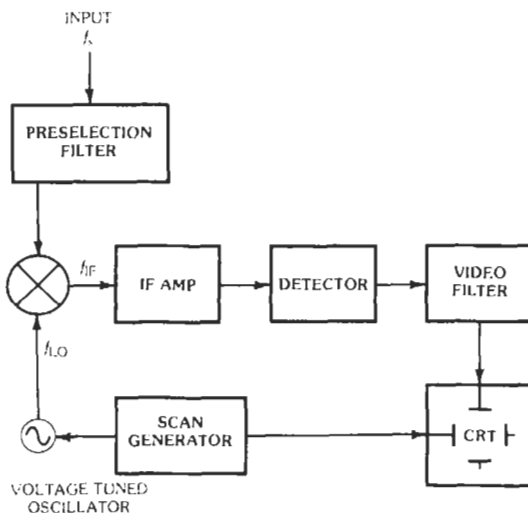


Fig. 48. Swept superheterodyne spectrum analyzer.

Signal-Analysis Characteristics as Determined by the IF Filter, Detector, and Video Amplifier

The input mixer and local oscillator translate the input frequency to a convenient intermediate frequency. Typically, several conversions are used to reduce spurious responses. The final IF filter, detector, video amplifier, and display-processing circuitry largely determine the signal-analysis characteristics of a spectrum analyzer.

The IF-filter bandwidth and shape determine the frequency resolution of the analyzer, assuming the local oscillators are sufficiently stable. Displayed CW signals have the shape of the IF filter, as shown in Fig. 50. To resolve two equal-amplitude signals, the IF 3-dB bandwidth should be less than the signal separation. Resolving unequal-amplitude signals requires narrower bandwidths or rectangularly shaped filters. Filters approximating a rectangular shape, however, have poor pulse response and limit the scan time. Filters approximating a Gaussian shape allow the shortest scan times without introducing amplitude errors due to overshoot

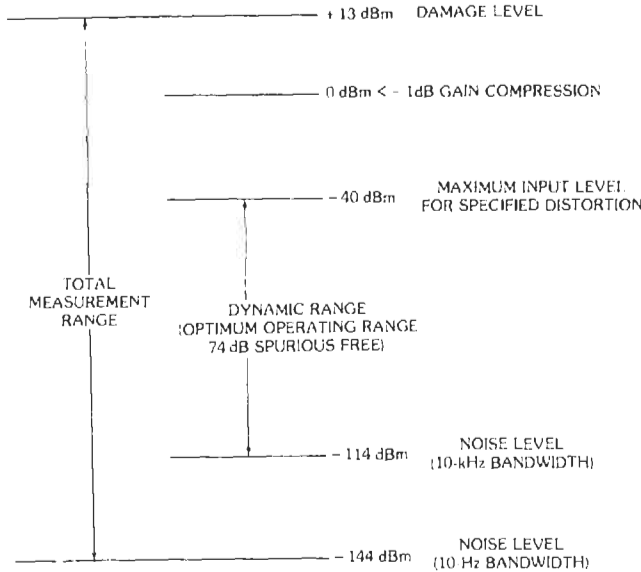


Fig. 49. Typical spectrum-analyzer input signal range.

or ringing. The fastest scan rate for Gaussian-shaped filters is approximately given by:

$$R = BW^2$$

where,

- R = maximum scan rate in Hz/s,
- BW = 3-dB IF filter bandwidth.

The detector following the IF amplifier is an envelope detector that provides a video signal proportional to the input RF or microwave signal. The bandwidth of the detector should exceed the widest resolution bandwidth to preserve the fidelity of signal modulation. However, post-detection low-pass filtering is often useful to provide signal integration—for example, to recover low-level signals in the presence of random noise.

While analog spectrum analyzers apply the video signal directly to the vertical-deflection system of the

display, newer instrumentation using digital displays first performs an A/D conversion. Having the video data in digital form allows further opportunities for signal processing. Trace-to-trace maximums or minimums can be stored and displayed. For instance, the occupied spectrum of the FM signal shown in Fig. 51 is easily determined from the “max hold” trace.

Trace-to-trace signal averaging is possible without the sweep-time limitation associated with video filtering. Fig. 52 shows the result of digital averaging on a noise-modulated signal. The signal power distribution and average frequency are much more apparent.

Digital storage of multiple traces allows arithmetic operations among traces. Fig. 53 shows how trace arithmetic is used to correct for amplitude errors in a measurement system. Trace B may contain measurement-system correction factors, antenna factors, etc.

TIME AND FREQUENCY MEASUREMENT

The measurements of time and frequency were among the first to be made in a digital form. The world standard definition of 1 second is the length of time it takes for a cesium 131 atom to “vibrate” 9 192 631 770 times, itself a digital representation. Most simply, frequency can be measured by counting the number of phase cycles taking place within a second, again a digital representation. Frequency may be measured over any interval of time, however, and not just over 1 second.

A class of instruments called counters can be used to measure time, frequency, phase, the number of events, or several other parameters. They offer a precise measurement in a relatively short measurement time, afford-

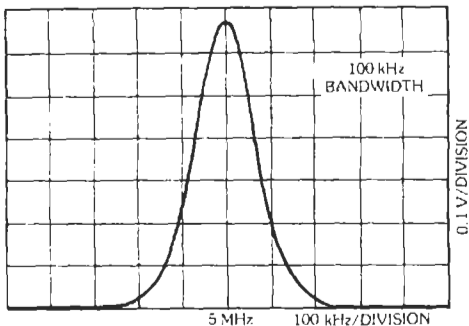
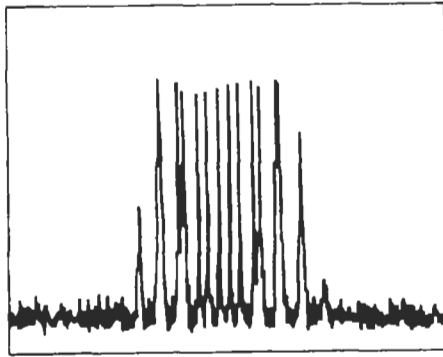
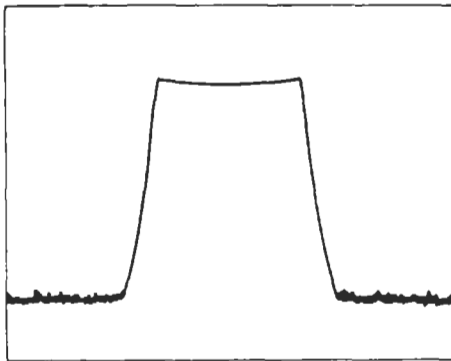


Fig. 50. Displayed CW signal has shape corresponding to IF filter bandpass.



(A) Normal spectrum-analyzer display.



(B) Digitally processed "max hold" trace.

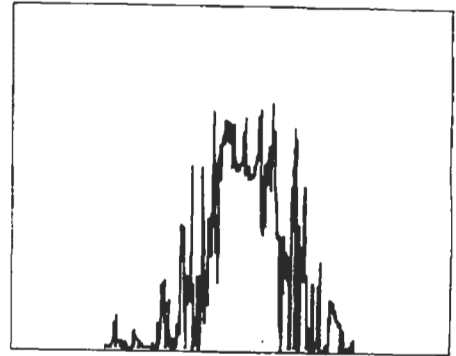
Fig. 51. Occupied spectrum of an FM signal.

ing good measurement throughput. Most of these measurements are accomplished by first converting the signal to be measured into a pulse train. This is done by sending the signal through a buffer amplifier and then a comparator. Once this is done, the pulse train can be counted with digital logic. The method of counting depends on the desired measurement. Also found in all counters is a time base. This is usually a signal, derived from a calibrated 10-MHz quartz crystal, which can also be counted.

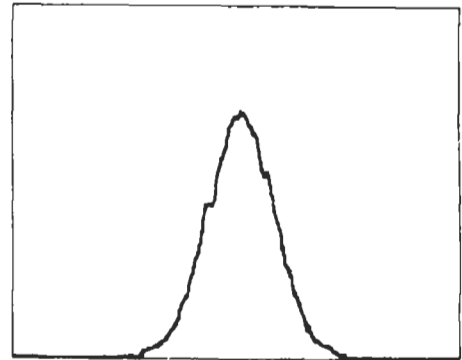
Time Measurement

To measure time, the counter simply counts the number of 10-MHz time-base cycles that take place during the time to be measured. The quantization of time is the inverse of 10 MHz, or 100 ns. The start and stop signals come from the user's input signal. The length of a period can be measured by gating the time-base counter open at the start of the period and closing the gate at the end of the period. (This is shown in the shaded portion of Fig. 54.)

A time interval measurement is the time between two separate signals. The gate is opened upon receiving a signal from one input channel, and it is closed upon receiving a signal from a second input channel. The basic measurement architecture is the same as for a



(A) Normal trace.



(B) 100 digitally averaged traces.

Fig. 52. Noise-modulated FM signal.

period measurement, except that there are now two inputs, each with its own comparator, feeding into the gating logic. Slope switches for both input channels permit use of the rising or falling edge of a signal to open or close the gate. Pulse width, rise and fall times, and relative phase also can be measured.

Resolution extension can be done in several ways. One way is to use a higher frequency time base. A 500-MHz time base gives 2-ns resolution. If the signal is repetitive, a second resolution extension method is to average. Time interval averaging is accomplished by accumulating several measurements and taking their average. Resolution is extended by the square root of the number of measurements averaged. A third method is called interpolation. This uses interpolation circuitry to measure the time between the edge of the signal pulse and the next time-base pulse edge at both the opening and closing of the gate. Extensions of 10 to 1000 times the resolution are found. By employing combinations of these techniques, the resolution can be extended to better than 1 ps.

Frequency Measurement

Simple counters measure frequency by counting input pulses during a decade value of time, such as 0.1, 1, or 10 seconds. This makes the process of dividing

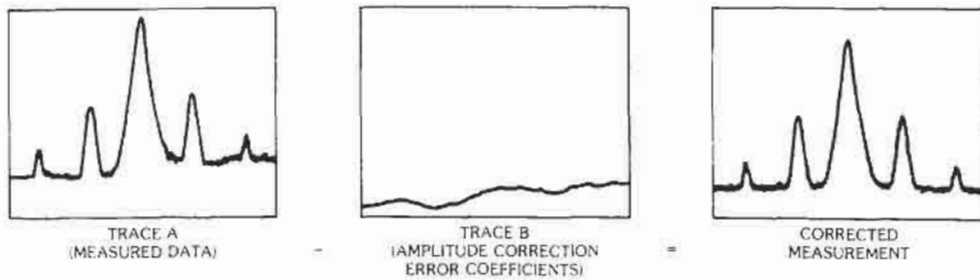


Fig. 53. Error correction by arithmetic operations on multiple traces.

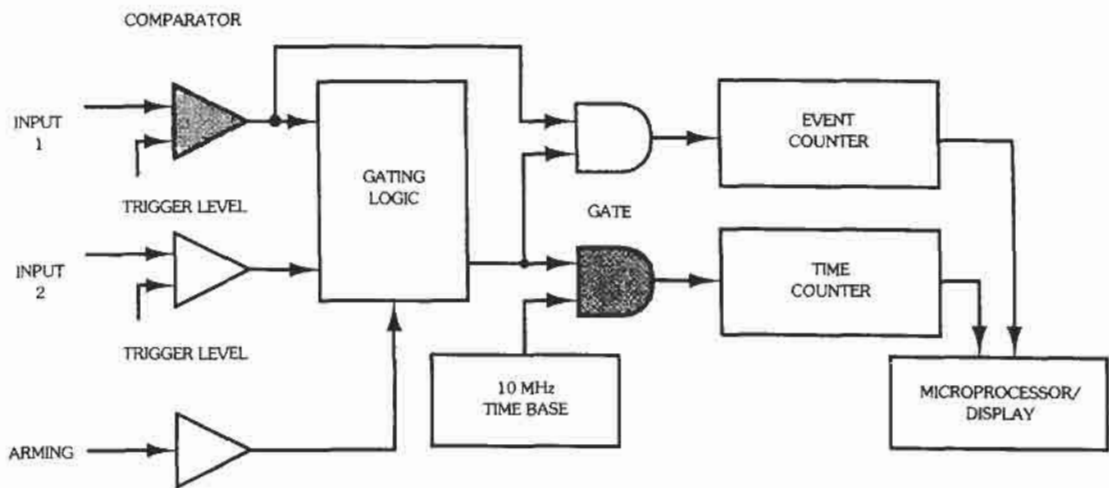


Fig. 54. Basic reciprocal counter. To make a simple period measurement, the shaded portion is used.

event pulses by time a trivial matter, and one not needing a microprocessor. Modern counters use a microprocessor to perform the division, so the gate can be open over any arbitrary time. To make the measurement, the input signal pulses are counted with one circuit while the length of time the gate was open is measured with another circuit. Additionally, an arming signal can be used to control the timing of the measurement if a particular portion of the signal is being measured. The arrangement described is called a reciprocal counter (Fig. 54).

The resolution of the frequency measurement in a reciprocal counter is determined by the resolution of a time measurement. If a counter can resolve time to 1 ns, 9 digits of frequency resolution can be displayed given a 1-second measuring time. This same counter will provide 3 digits given a 1- μ s measuring time. A more recent resolution extension technique employs continuous count signal processing. A truer average of the frequency can be calculated if the signal is sampled multiple times during the measurement period; up to 12 digits of resolution in a 1-second measurement time can be obtained with 1-ns time resolution.

Microwave frequencies can be measured with the same counter architectures, but usually some form of frequency conversion needs to be used, because the counter logic typically is limited to operation below 500 MHz. One approach is to use a prescaler. This is simply a very fast counter placed between the comparator and the gating logic. A 3-GHz divide-by-32 prescaler (a five-bit counter) brings the count rate down to a maximum of below 100 MHz, and the microprocessor multiplies the count by 32 for proper display.

A sampler can also be used to mix down the signal to a more countable rate. This is a harmonic heterodyne down-conversion technique. A relatively low-frequency local-oscillator signal is sent through a step recovery diode to generate microwave-frequency harmonics for the sampler to use. The microwave input signal is mixed down, and the resulting intermediate frequency is counted. The microprocessor is used to dither the local-oscillator frequency to determine which harmonic came into play so that it can add the proper frequency offset for display. This technique can be used for microwave and millimeter-wave frequencies. It has the advantage of added frequency resolution because the

phase is preserved, but the act of finding the proper harmonic number slows the measurement. Also, confusion can be introduced if multiple signals are present, so sometimes a tuned YIG filter is used to preselect the signal. For pulsed microwave signals, a detector is added to form a signal that the counter uses to open and close the gate properly.

Frequency and Time-Interval Analysis

A new class of instruments allows further analysis of time and frequency. These instruments share similarities with the counter, but have several significant differences. First, the instrument does not have a gate, so the counter circuits are constantly counting the input and time-base signals. The contents of the counters are read synchronously and stored into memory. From two

sets of these readings, a frequency measurement can be effectively constructed. This architecture is called *continuous count* (Fig. 55).

The first advantage of continuous-count architecture is that a frequency measurement can now be constructed from several *contiguous* samples, and by using a least-squares fit or other algorithm, a truer estimate of frequency can be made.

The second advantage is the ability to store these samples into a deep memory. This allows for the analysis of time-varying signals, such as the measurement of the step response of a voltage-controlled oscillator (Fig. 56A). If the signal is being frequency or phase modulated (either to carry data or to make the signal spread-spectrum), the modulation pattern can be displayed in a frequency-vs-time or phase-vs-time manner (Figs. 56B and 56C). If required, the modulation can be displayed in an eye diagram form. The maximum

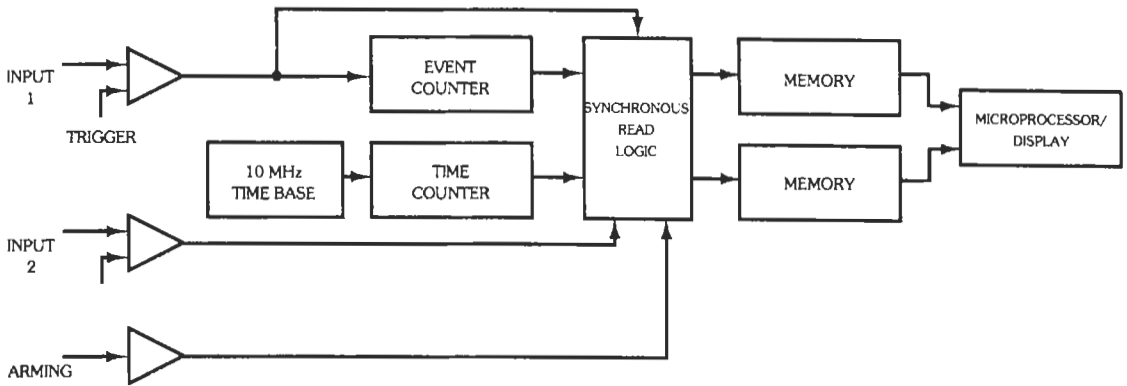
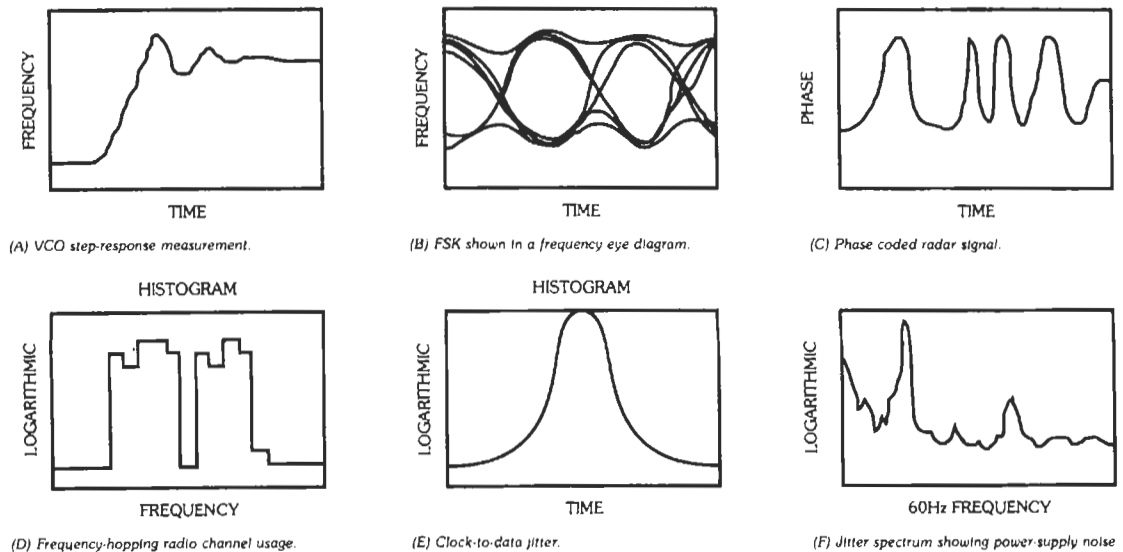


Fig. 55. Continuous-count architecture.



(A) VCO step-response measurement.

(B) FSK shown in a frequency eye diagram.

(C) Phase coded radar signal.

(D) Frequency-hopping radio channel usage.

(E) Clock-to-data jitter.

(F) Jitter spectrum showing power supply noise

Fig. 56. Continuous-count measurements.

modulation rate that can be analyzed is governed by the sample rate of the counter, or how often the counters are read into memory. Note that phase-modulation measurement is possible only because the signal is sampled continuously.

The measurements can also be stored and analyzed in histogram form (Figs. 56D and 56E). This can be useful for analyzing a frequency-agile signal. Note that these measurements cannot be easily done with a spectrum analyzer.

If the parameter being measured is time or time interval, the histogram can show information on data-to-data or clock-to-data jitter. The information in the memory can also be passed through an FFT to show jitter spectrum (Fig. 56F). These measurements are useful in the design of high-speed computers, telecommunications systems, and magnetic storage systems.

RF AND MICROWAVE POWER MEASUREMENTS

A common technique for measuring power at high frequencies is to employ a sensing element that converts the RF power to a measurable dc or low-frequency signal. The sensing element is often designed to form a termination that is matched to the characteristic impedance (Z_0) of the input transmission line. Various types of sensing elements are used.

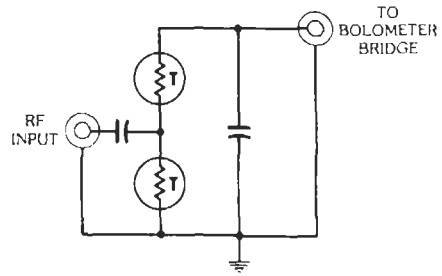
Thermistor Sensors

Thermistor sensors provide a change of resistance. The typical power range is $1 \mu\text{W}$ to 10 mW ; the maximum frequency is greater than 100 GHz . Fig. 57A shows a typical power sensor employing thermistors. The thermistors form the termination for the RF input, and dc or audio power from the self-balancing bridge in Fig. 57B raises the temperature of the thermistors until they each have a resistance of $2Z_0$. The RF impedance then becomes equal to Z_0 .

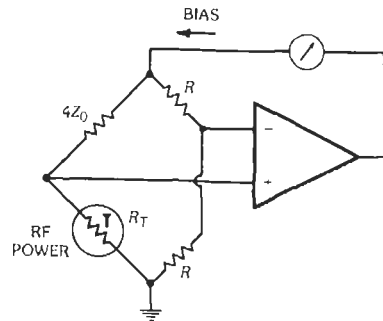
Since the bridge keeps the thermistor resistance constant, any heat added by the RF power causes a corresponding reduction in bias power. The RF power level is determined by measuring this change in bias power.

Thermal Converters

Thermal converters provide a dc voltage (less than 10 mV). The typical power range is 0.1 to 100 mW ; the maximum frequency is less than 1 GHz . Thermal converters employ a number of thermocouples (thermopile) mounted with good thermal contact to the RF termination (Fig. 58). The RF power heats the termination, and the thermopile output voltage is proportional to the amount of power dissipated. The converter can be calibrated by applying a precisely known dc or RF power level at the input. This calibration yields the value of C .



(A) Thermistor sensor



(B) Self-balancing bridge.

Fig. 57. Power measurement with thermistor sensor.

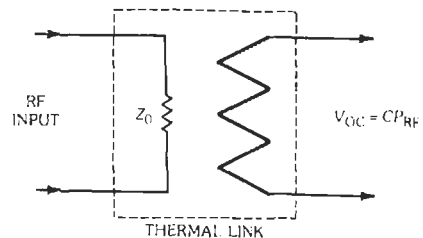


Fig. 58. Thermal converter.

Thermocouple Sensors

Thermocouple sensors provide a dc voltage (less than 10 mV). The typical power range is $0.1 \mu\text{W}$ to 100 mW ; the maximum frequency is greater than 100 GHz . Thermoelectric sensors differ from thermal converters in that the thermocouples are used as the terminating resistors (Fig. 59). This type of sensor must be calibrated with a precise RF power level to determine the value of C .

Diode Sensors

Diode sensors provide a dc voltage (approximately 1 V at 10 mW). The typical power range is 0.1 nW to 100 mW ; the maximum frequency is greater than 18 GHz .

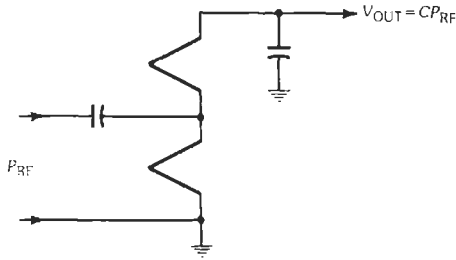


Fig. 59. Thermocouple sensor.

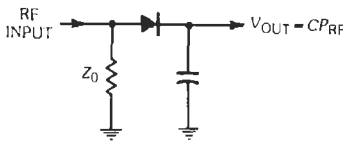


Fig. 60. Diode power sensor.

Diode power sensors (Fig. 60) use point-contact or Schottky barrier diodes to detect the RF signal. If the RF voltage is less than 20 mV, the diode output follows the square of the applied voltage, so the dc voltage is a function of RF power. At higher levels, the rectified output gradually changes to the more familiar peak-detection mode, and harmonics in the signal can cause errors in the power reading. Diode sensors must be calibrated with a precise RF signal level to determine the value of C .

Power-Measurement Definitions

Maximum Available Power, P_{avs} —The power obtainable from a source when it is terminated in a load whose impedance is the complex conjugate of the source impedance. This condition is usually obtained by installing a tuning device between the source and the power sensor and adjusting the tuner for a peak in the power reading.

Z_0 Available Power, P_{z_0} —The power obtainable from a source when it is terminated with a load matched to the Z_0 of the transmission line. The Z_0 available power is related to the maximum available power by:

$$P_{z_0} = P_{avs} (1 - \rho_g^2)$$

where $\rho_g = |\Gamma_g|$

Power Incident Upon the Load, P_i —The amount of power the source transmits toward the load is usually measured by use of a directional coupler connected such that its coupled output is governed by the wave

emerging from the source. The incident power is related to the Z_0 available power by:

$$P_i = P_{z_0} / |1 - \Gamma_g \Gamma_L|^2$$

Power Reflected by the Load, P_r —This power is usually measured by using a directional coupler connected such that its coupled output is governed by the amount of power reflected by the load. The reflected power is related to the incident power by:

$$P_r = \rho_L^2 P_i$$

Power Dissipated in a Load, P_L —The power dissipated in the load is related to the other power levels by:

$$\begin{aligned} P_L &= P_i - P_r = P_i(1 - \rho_L^2) \\ &= P_{avs} \frac{(1 - \rho_g^2)(1 - \rho_L^2)}{|1 - \Gamma_g \Gamma_L|^2} \end{aligned}$$

Source and Load Reflection Coefficients, Γ_g and Γ_L —Some of the foregoing equations refer to the match of the source and load impedance to the characteristic impedance of the transmission line. Values of reflection coefficients Γ_g and Γ_L are given by:

$$\Gamma_g = (Z_g - Z_0) / (Z_g + Z_0)$$

$$\Gamma_L = (Z_L - Z_0) / (Z_L + Z_0)$$

$$\rho_g = |\Gamma_g|$$

$$\rho_L = |\Gamma_L|$$

Effective Efficiency, η —Efficiency of a power sensor is defined as

$$\eta = P_{sub} / P_0$$

where,

P_{sub} is the dc substituted or RF calibration power.

P_0 is the amount of power dissipated in the sensor.

Calibration Factor, K_B —Calibration factor is related to efficiency by:

$$K_B = P_{sub} / P_i = \eta(1 - \rho^2)$$

where ρ is the magnitude of the reflection coefficient of the sensor.

Multiple Mismatch Error, M_u —The precise amount of power dissipated in the sensor is a function of many variables, the largest one being the interaction of the source and load reflections. Fig. 61 can be used to estimate the worst-case uncertainty for a simple source-to-load power measurement. In the example shown, an

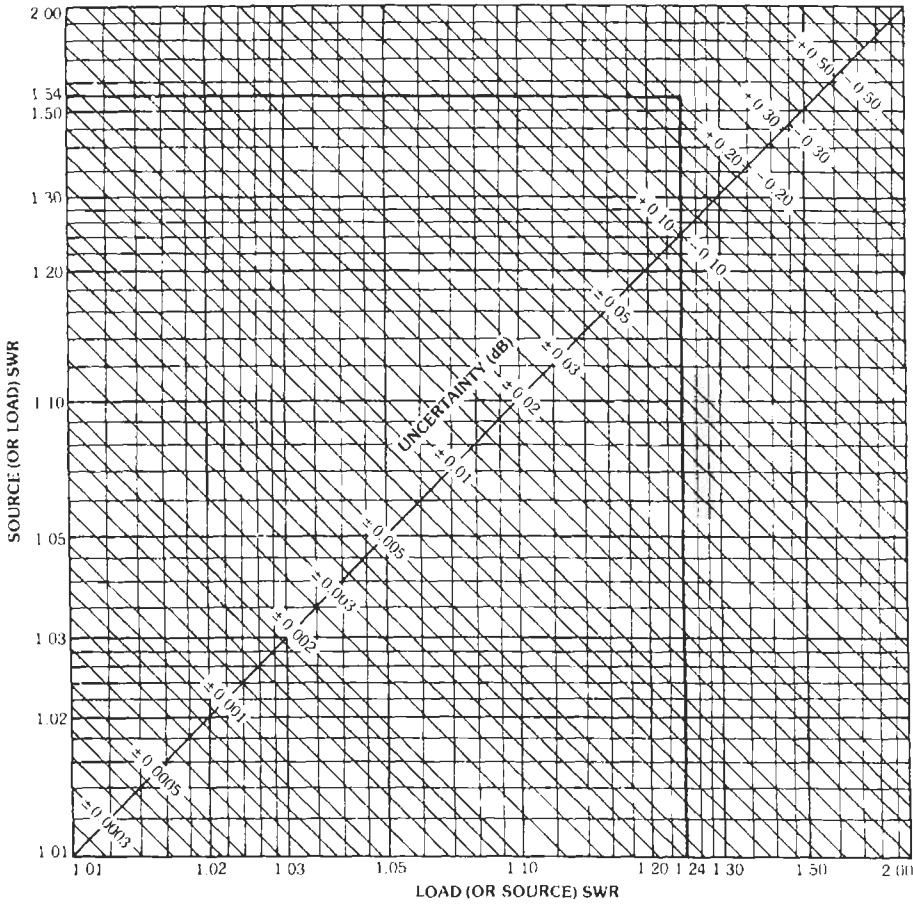


Fig. 61. Mismatch uncertainty limits.

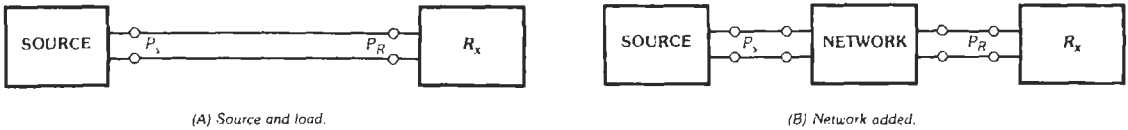


Fig. 62. Insertion loss or gain.

SWR of 1.54 interacting with an SWR of 1.24 introduces about a 0.2-dB uncertainty.

$$M_u = (1 \pm \rho_g \rho_l)^{-2}$$

MICROWAVE-LINK ANALYSIS

Insertion Loss or Gain

Insertion loss or gain is defined as the loss or gain that is apparent upon inserting the network to be measured between a given source and a given receiver. In Fig. 62, P_s is the transmitted power of the source, and

P_R is the received power. In Fig. 62A, it is clear that $P_R = P_s$, but in Fig. 62B, the value of P_R is modified by the insertion of the network. This change in power is the insertion loss or gain and is usually quoted in decibels:

$$\text{Insertion Loss or Gain} = 10 \log_{10}(P_R/P_s)$$

Amplitude Response

Amplitude response is sometimes referred to as "flatness" or frequency response and is the variation in gain or loss with change in frequency over a defined frequency band. While it may be thought that this

parameter is not important for frequency-modulated radio links, it generally must be well controlled, especially for high-capacity systems. Amplitude response is also very important for the RF and IF sections of digital radio links.

In order to measure amplitude response accurately, it is important that harmonics of the transmitted signal not be included in the measurement of received power. Thus the ideal measurement technique would be one using a tunable receiver. In practice, however, the measurement of radio links using intermediate frequencies of 70 MHz or 140 MHz requires a swept frequency range that allows a simple low-pass filter to remove harmonics. Additionally, errors due to the use of a broadband detector are generally low. This, then, allows a simple measurement principle to be used that lends itself to remote or "straightaway" measurements (Fig. 63).

Measurement receivers usually incorporate a frequency-tracking loop that reduces the effective sweep range in the analyzer where the benefits arising from a reduced sweep may be used for group-delay and/or linearity measurements. The frequency-tracking loop, however, also gives a recovered sweep signal that is used to drive the X deflection of the display and makes the "straightaway" measurement possible (Fig. 64).

The swept range of the input signal is reduced by the loop gain of the arrangement shown in Fig. 64, where the derived low-frequency signal used to drive the VCO is a facsimile of the frequency variation against time characteristic of the source signal.

These techniques are quite satisfactory when used to measure amplitude response of IF or RF sections of radio links. They are, however, generally inadequate for the measurement of baseband amplitude response, where a selective receiver should be used. This implies that the receiver should be capable of following the frequency changes of the input signal; i.e., it must "track" the input frequency. Various techniques have been developed for this, even for those measurements where frequency "stepping" synthesizers are used. These methods usually require that a receiver be programmed to tune to its next frequency only after having successfully completed a measurement at the current frequency. The method therefore usually relies on the generator stepping at a slower rate than the receiver.

Envelope-Delay Distortion (EDD)

Envelope-delay distortion is often termed group-delay distortion and is effectively the variation in the derivative of the phase-vs-frequency response. For zero group-delay distortion, the phase/frequency characteristic must be linear. This means that all frequencies are then transmitted through the system with equal time delay. Thus group delay is defined in units of time, but it is the variation that distorts broadband signals and is therefore of interest in measurement.

$$\text{Group Delay} = d\phi/df$$

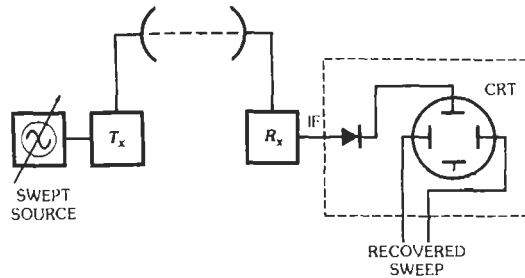


Fig. 63. Amplitude-response measurement.

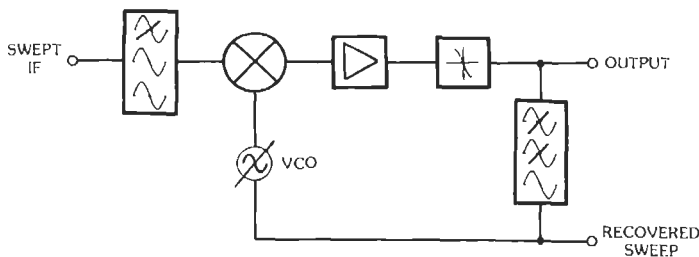


Fig. 64. Frequency-tracking loop.

The nonlinear phase characteristic of Fig. 65 will give rise to distortion of a broadband signal that will manifest itself in the form of noise and intermodulation for analog systems. The group-delay response derived from Fig. 65 will be of the form shown in Fig. 66. As can be seen, there is a mean level of group delay that is not very important; it is the variation of distortion that must be equalized.

Equalization is achieved by means of networks that give inverse group-delay variations. The most common types required are those that compensate for linear group-delay slope or parabolic group-delay distortion. It is important that equalizers be connected as close as possible to the source of group-delay distortion or at least that no nonlinear networks be connected before the equalizer. This is because nonlinearities that intro-

duce amplitude-to-phase-modulation conversion will produce effects that cannot be removed by group-delay equalization.

Measurement of Group-Delay Distortion

There are several methods for the measurement of group-delay distortion, but the basis is usually that of comparing the phase of an output modulation envelope with the phase of a reference signal. The usual method employs a frequency-modulated signal that is swept over the frequency band of interest. Use of a phase-locked loop at the receiver makes possible the recovery of the frequency modulation, which then is used to give the variation in phase as the input signal is swept (Fig. 67).

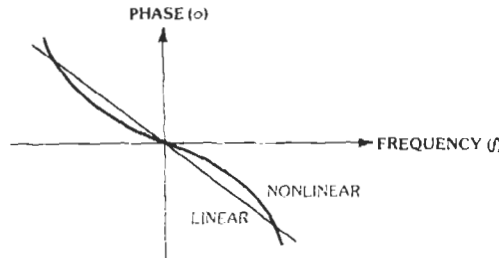


Fig. 65. Phase versus frequency slope.

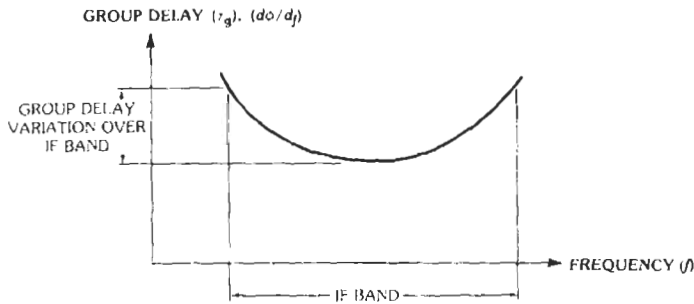


Fig. 66. Group-delay distortion.

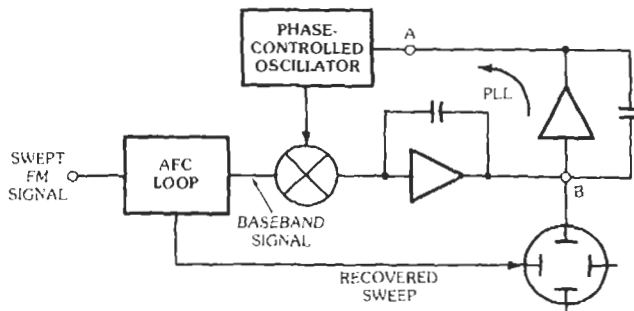


Fig. 67. Measurement of group-delay distortion.

It was shown earlier that a frequency-tracking loop could be used to recover the sweeping signal for the X deflection of the CRT. Additionally, the tracking loop reduces the swept range of the received IF signal so that distortions arising from the nonideal nature of networks in the receivers are minimized. For example, a signal that is used to test a radio link may sweep over a band of more than 30 MHz but will be sweep reduced in the analyzer to yield a sweep range of less than 100 kHz. Consequently, the errors introduced by the analyzer are minimized. The phase-locked loop in the group-delay measurement then uses the demodulated baseband frequency to control the phase of its reference oscillator, which will be held at the mean phase of the modulating frequency. Variations in phase are then available at point B in Fig. 67.

Another way of looking at this is to say that if point A represents the reference phase of the modulating signal, then going backward around the double integrating loop means that the signal at point B is the derivative of this phase, or $d\phi/dt$. In this system, however, time is the same as frequency, because as time proceeds, the IF is swept; hence the instantaneous value of the dc voltage at B represents the group delay ($d\phi/df$) at the instantaneous value of the IF.

The reason for using a phase-tracking detector is that it is always working very close to zero phase difference between the inputs. This means that the phase detector is always working on the central part of its dc/phase conversion slope, and nonlinearity of this response does not affect the measurement.

The measurement of group-delay distortion involves the careful selection of several test parameters:

- Sweep range
- Sweep rate
- Modulation frequency (test tone)
- Modulation index
- Post-detection bandwidth

It is important to select a sweep range appropriate for the device or system under test. Some consideration should be given to the spectrum of the modulating signal, since the device under test will be subjected to the total spectrum. This is usually important only where a high-frequency test tone (> 1 MHz) is used, but it can be important for testing components such as narrow-band filters with lower test-tone frequencies.

The sweep rate employed is often in the range of 50 to 100 Hz, but this range may not be suitable for systems such as satellite communication links where use of a lower sweep rate will allow a narrow bandwidth to be selected for the post-detection bandwidth. This will then enhance the measurement resolution by reducing noise power.

The modulation frequency used is a compromise between two conflicting effects:

1. Use of too high a frequency will tend to conceal rapid fluctuations in group delay such as the ripple produced by imperfect impedance matching.

2. Use of too low a frequency will produce a low voltage at the output of the group-delay detector, and the signal-to-noise ratio of the display will be poor.

It is therefore usual to select frequencies between 50 kHz and 500 kHz to give an appropriate compromise. There are cases, however, where much higher modulation frequencies can be used. This is usually true when either television or broadband telephone systems are carried and the test signal has to reveal low values of nonlinearity that produce intermodulation. In these cases, the measurement sensitivity is enhanced by using a modulation frequency between 1 MHz and 12 MHz. In using these high frequencies, it is not appropriate to refer to the measurement as a group-delay measurement, but rather as a measurement of "differential phase."

In practice, this measurement is not as important as the measurement of "differential gain," which can be very useful in tracing problems of amplitude-to-phase modulation conversion.

Television systems are often measured with a test-tone frequency equal to the television color subcarrier frequency. The measurement of differential phase is then used to define system performance.

Differential gain is an extension of the measurement of linearity that uses techniques similar to those for group delay, but, whereas this reveals how the phase for the modulation envelope varies, linearity shows how the amplitude of the envelope varies. The same measurement principle is applied to modulators and demodulators to measure nonlinearity of the voltage/frequency characteristic. It is interesting to note that for a nonlinearity occurring in a discriminator or modulator, the measured nonlinearity will be relatively independent of the modulation frequency employed. This contrasts with the case for IF nonlinearities where the distortion value changes with the square of the change in test-tone frequency. This property may then be used in analyzing system deficiencies.

It can be valuable to be able to separate key characteristics that contribute significantly to intermodulation distortion. Table 1 shows the relative effects of distortion styles to give an approximate indication of which shapes are most serious.

Table 1 does not show the effects of the "coupled" responses where amplitude-to-phase modulation conversion interacts with distortions in other parameters, but these interactions can have a serious effect in high-capacity radio systems. Additionally, nonflatness can be important, especially for digital radio systems.

Return Loss

In the alignment of microwave radio links, it is important that the impedance match of the various sections be well maintained. This is especially important where cabling is used between a source and a load, since any mismatch will produce time-delayed reflec-

TABLE 1. DISTORTION FOR 100-pW INTERMODULATION NOISE (1800-CHANNEL SYSTEM WITH PREAMPHASIS)

Parameter	Distortion Over 10 MHz	Test Tone
Differential-Gain Slope	1%	2.4 MHz
Group-Delay Slope	2 ns	—
Differential-Phase Slope	3% rad	2.4 MHz
Group Delay, Cubic	3.5 ns	—
Group Delay, Parabolic	14 ns	—

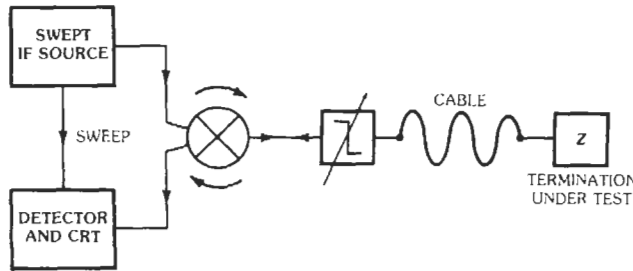


Fig. 68. Return-loss measurement by long-cable method.

tions that may impair link performance. The normal way of describing the mismatch of a source and load is by using the term "return loss." Return loss is the measure of the ratio between the transmitted and reflected signals:

$$\text{Return loss} = 20 \log_{10} |E_i/E_r|$$

where,

- E_i is the incident signal,
- E_r is the reflected signal.

Return loss is a measure of the magnitudes of incident and reflected signals and does not take account of phase relationships.

Measurement of Return Loss

Long-Cable Method—In this method, a "long cable" is connected to the termination under test as shown in Fig. 68. When the swept IF signal is applied to the long cable and its termination, a series of ripples will appear on the CRT display. If the test termination is removed, the open circuit produces a large amplitude of ripple that is then adjusted by the attenuator to equal the level when the termination is connected. The return loss is equal to 2 times the attenuation inserted.

Note that the cable must be long enough to produce at least one ripple over the swept range, but more than ten are required to observe variations across the band. Thus, for normal IF measurements, a length in excess of 20 meters is required.

Standard-Mismatch Method—This method (Fig. 69) relies on the measurement of power passed from a

hybrid when it is terminated by a known mismatch and then by the test item. A typical mismatch of 17 dB is used to calibrate the power meter/detector, after which the actual return loss may be measured directly.

Baseband Measurements

Most of the above measurements are made at the intermediate frequency (IF) of the radio system or at the microwave carrier frequency. Obviously, some measurements such as modulator linearity also involve the direct use of baseband inputs, but there is an additional class of measurements that define system performance and are carried out at baseband level. The most obvious

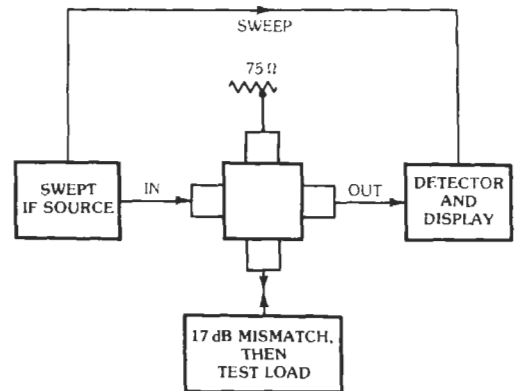


Fig. 69. Measurement of return loss by standard-mismatch method.

of these is the measurement of baseband amplitude/frequency response, but checks for spurious frequencies should also be made. These measurements can be made with special instruments that provide spectrum analysis.

For multichannel telephony systems, the prime concern is to minimize the level of noise in each channel. There can be many sources of noise, but they are mainly either thermal or caused by cross talk or intermodulation. Thermal noise is not affected by the traffic level (loading) of the system, whereas intermodulation distortion is sensitive to loading. The system designer therefore attempts to define a loading level at which intermodulation noise is low and yet the ratio of signal level to thermal noise is also satisfactory.

In order to verify that the correct balance of thermal and intermodulation noise is obtained, it is necessary to be able to simulate the traffic load and then measure the effects in channels throughout the baseband. This is achieved by using a band of thermal noise to represent the traffic.

Noise Loading—Traffic simulation using a band of thermal noise is possible because the amplitude distribution of such noise is Gaussian and the power-versus-frequency spectrum is uniform. This latter characteristic is analogous to white light, so the band of noise used to simulate traffic is often referred to as "white noise."

In order to simulate a given traffic load, the formula normally used to calculate the noise power required is

$$\text{Mean noise power} = -15 + 10 \log_{10} N$$

where,

Mean noise power is in dBm0 (power in dBm referred to a point of zero relative transmission level).

N is the number of telephone channels.

For this formula to be valid, N must be greater than 240, because as the noise bandwidth is reduced to simulate high-capacity systems, the signal departs significantly from a Gaussian distribution. The formula then used is

$$\text{Mean noise power} = -1 + 4 \log_{10} N$$

where,

Mean noise power is in dBm0,
 N is less than 240.

Modern systems may, of course, be loaded with data. In this case, a modified formula is used:

$$\text{Mean noise power} = -10 + 10 \log_{10} N$$

where,

Mean noise power is in dBm0,
 Data loading N is greater than 12 channels.

It is necessary to limit the noise bandwidth to represent the bandwidth of the system. Filters for

controlling noise bandwidth are prescribed by the various telecommunications authorities.

Noise Power Ratio—When the system has been loaded with noise to simulate traffic, it is then necessary to measure intermodulation noise. This is done by introducing a quiet "slot" into the noise band. This slot is then inspected to determine the level of noise introduced by intermodulation. This is known as the measurement of noise power ratio (NPR), which may be defined as the ratio of the noise power in a measurement channel with the baseband fully loaded to the noise power in that channel when the baseband is noise-loaded except for the measurement channel. In Fig. 70, P_1 (Fig. 70A) represents the power in the measurement channel due to the loading signal, P_2 (Fig. 70B) would equal zero for the ideal case where there is no thermal noise or intermodulation noise, and P_3 (Fig. 70C) shows the effect that these noise components will have. Thus,

$$\text{NPR} = P_1/P_3$$

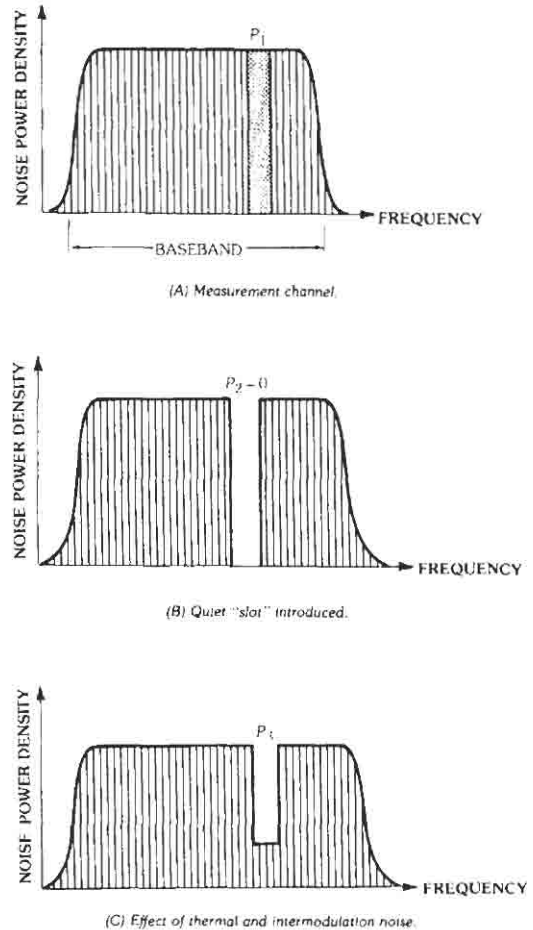


Fig. 70. Principle of NPR measurement.

If the noise load is removed, the thermal noise component in the measurement channel can be measured. (Thermal noise is sometimes referred to as basic noise or intrinsic noise.) Then if the noise load is added, the additional noise introduced by intermodulation and cross talk can be measured.

By varying the level of the baseband noise loading signal, it is possible to explore the combined effects of thermal noise and intermodulation noise to determine the optimum operating point for the system.

The NPR curve is often referred to as a "V" curve because of the characteristic shape (Fig. 71). At low loading levels, thermal noise will predominate, and NPR will simply be proportional to the loading power level. Thus this part of the characteristic is linear. As intermodulation noise rises, however, the noise power ratio will fall rapidly and give the shape of Fig. 71. It is possible to gain troubleshooting information by inspecting the V curve for slots placed at different frequencies throughout the baseband. The method employed is to select a slot and observe the change in NPR for a 1-dB change in noise loading. In the linear region of the curve, a 1-dB change in load will give a 1-dB change in NPR. For the nonlinear segment of the V curve, the change in NPR can then be used to determine the order of the nonlinearity involved.* For example, a 2-dB change in NPR will indicate that third-order nonlinearity is the cause.

The conventional form of test set† for carrying out NPR measurements is shown in Fig. 72. Fixed filters are used at both transmit and receive points in the system. This form of test set has the advantage that

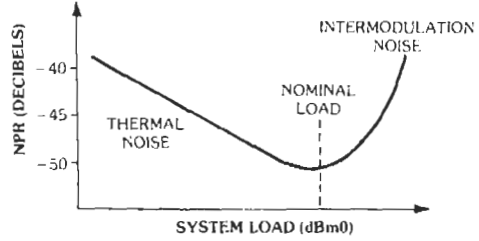


Fig. 71. NPR curve for a multichannel radio system.

intermodulation in its receive mixer is not important since bandpass filters reduce the noise band presented to the mixer. Unfortunately, it is necessary to have a range of bandpass filters and oscillators to cope with all the slots and systems that may have to be measured.

The modern design of highly linear mixers, however, allows a selective receiver to make the NPR measurement. Such receivers use a synthesizer as the local oscillator, so the requirement for separate filters and oscillators is avoided.

Conversions—It is sometimes necessary to convert from NPR to other measures, such as signal-to-noise ratio (SNR). The following relations are used.

1. Weighted signal-to-noise ratio in decibels:

$$\text{Weighted SNR} = \text{NPR} + 10 \log_{10}(f/3.1) - P + W$$

where,

- f = bandwidth of baseband (kHz),
- P = nominal load (dBm0),
- W = weighting factor (dB).

* References 24 and 25.

† Reference 24.

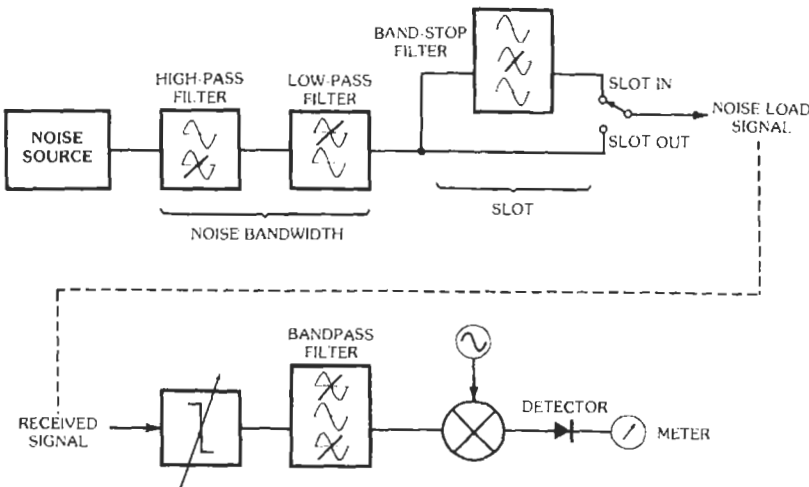


Fig. 72. NPR measurement.

When applied to systems having greater than 240 channels and using *C*-message weighting, this simplifies to

$$\text{SNR} = \text{NPR} + 17.8$$

2. The notation dBnC0 means decibels of noise with respect to reference noise of -90 dBm and with *C* message weighting, corrected to a point of zero transmission level. Conversion to dBnC0 for systems that have more than 240 channels is given by

$$\text{dBnC0} = 71.7 - \text{NPR}$$

Note that the European psophometric weighting and *C* message differ by 1 dB, but in practice the above formula is often used for psophometric weighting.

Carrier/Noise Measurement

For digital radio systems, the normal overall performance measurement is based on "bit error rate," where a pseudo-random-sequence bit pattern is applied to the system and errors are indicative of system imperfections. A more detailed analysis is made possible by inspection of the so-called "eye" diagram at the sampling stage of digital signal recovery, but measurement of error rate as the carrier-to-noise ratio is varied can be a more revealing method of analyzing system performance.

Diagnostic Measurements and System Performance

The relationships between diagnostic tests and system performance are often very complex. Thus it is difficult to relate measurements of envelope-delay distortion to the intermodulation that will be produced. It is possible to make such predictions,* but this generally requires the use of a computer. It is generally more useful to derive a clear understanding of how these measurements can best be used to optimize system performance, or at least to be able to separate those measurement responses that have a significant effect on system performance from those that have a lesser effect.

COMPUTER CONTROL OF INSTRUMENTS

A series of standards provides mechanisms for computers to communicate with instruments. IEEE 488.1 or IEC 625-1 describes a physical interface including electrical and mechanical constraints. IEEE 488.2 or IEC 625-2 describes how complete messages are formatted and exchanged. SCPI defines the semantics of the messages.

IEEE 488.1 or IEC 625-1 General Purpose Interface Bus (GPIB)*

The GPIB provides a means by which system components communicate with each other. It is analogous to the telephone system that allows one person to exchange verbal information with another person—or, in the case of a conference call, with several other people.

Every participating device in a GPIB system performs at least one of three roles: talker, listener, or controller. A talker transmits data to other devices via the bus. A listener receives data from other devices. A controller manages communications on the bus, primarily by designating which devices are to send or receive data during each measurement sequence. The controller may also interrupt and command specific actions within devices.

Many devices are both talkers and listeners. For example, a programmable multimeter or electronic counter listens when receiving its program instructions, and it talks when sending its measurements to another system component such as a printer or computer. There can be several listeners, but to avoid confusion, only one talker can be active on the bus at any one time.

In its simplest form, a GPIB system can consist of only one talker and one listener—for example, an electronic counter linked to a printer. In this application, there is no need for a controller; the counter "talks" and the printer "listens."

The full versatility of the GPIB becomes apparent when there is more stimulative interaction among the interconnected devices. To achieve this requires a controller that can schedule measurement tasks, set up instruments to perform specified tests and measurements, monitor processes on-line, and process data and analyze and interpret the results. A desktop computer, or "computing controller," is often used as a system manager to control the bidirectional data flow of the GPIB.

Interface circuitry is part of each interfaced instrument or computing controller—the GPIB cable itself is passive. Inside the cable, 16 active signal lines are grouped into three sets, according to function (Fig. 73).

The eight data lines carry coded messages—such as addresses, program data, measurements, and status bytes—among as many as 15 devices interconnected with a single bus. The same data lines are used for both input and output, and the messages are in bit-parallel, byte-serial form. Data are exchanged asynchronously, for maximum compatibility among a wide variety of devices.

For unambiguous, intelligible communication between instrument and computer devices, some rules or protocol must apply to the communication process itself. Thus, the exchange of data is controlled by the

* Reference 17.

* From Hewlett-Packard publication 5021-1927. © Hewlett-Packard, used with permission.

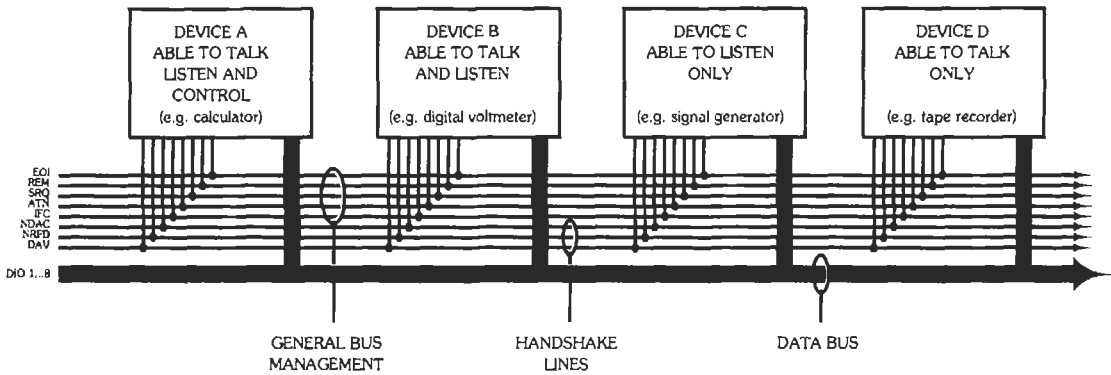


Fig. 73. General purpose interface bus. (From Hewlett-Packard publication 5021-1927. © Hewlett-Packard, used with permission.)

second set of signal lines, the three data byte transfer control lines. All listeners addressed on the GPIB send a signal when they are "ready for data." The talker then initiates data-byte transfer by signaling that "valid data" are on the data lines. When each listening device successfully receives the message, it acknowledges "data accepted."

When two or more devices are listening, all must acknowledge "data accepted" before new data can be transmitted on the bus. Thus, data are transferred at the rate of the slowest listener participating in the particular conversation.

High-speed and slow devices can be connected on the same bus. Data transfer between high-speed devices is not adversely affected, as long as the slow devices are not addressed or participating in the conversation.

The remaining five general interface management lines are used for such things as activating all the connected devices at once, clearing the interface, and so forth.

Up to 15 GPIB devices can be linked together on one bus. The GPIB cable connects them in parallel functionally, in either star or linear fashion physically. Total cable transmission path length on any one bus should not exceed 20 meters or 2 meters per device, whichever is less, except where the distance is extended by devices such as a common carrier interface.

A single computing controller can manage more than one bus instrument cluster at once. This is particularly useful for achieving concurrent actions at different test stations in a production test measurement application.

Interface circuitry for controllers and instruments is commercially available in integrated-circuit packages.

IEEE 488.2 or IEC 625-2

As illustrated in Fig. 74, IEEE 488.2 is a higher level standard than GPIB. It has also been used on top of IEEE 1155, VME Extensions for Instrumentation (VXI). While the physical interface provides a connection between system components, IEEE 488.2 provides rules on how the communication on the interface will take place.

The major areas covered by IEEE 488.2 are:

1. Reliable transfer of complete messages between a computer and an instrument.
2. A precise description of the syntax in those messages.
3. The form and meaning of a set of commands useful in all instruments.
4. Common status reporting capability.
5. Techniques for guaranteeing synchronization of application programs with instrument functions.

IEEE 488.1 provides a reliable means of transferring bytes between a talker and a listener. In addition, a reliable means of transferring programming commands and measurement results is needed.

A response is generated only after a query message has been received. An instrument responds only with the data asked for by the controller. An instrument is not allowed to generate arbitrary data when it is addressed to talk. Any violations of the message exchange protocol are reported as query errors in the event-status register described by the status reporting model.

IEEE 488.2 requires that devices listen in a "forgiving" manner. The level of forgiveness is, however, very precisely defined. The syntactic meaning of every data byte is well defined, even when multiple characters have the same meaning.

While the listening syntax may be very forgiving, an instrument is required to talk very precisely. As a result, the chances of a message being accepted by any controller are greatly increased.

Certain functions are useful across all types of instruments. IEEE 488.2 defines a set of commands which all instruments must implement exactly as described. Some of the more generic commands are resetting and identifying the instrument. Other commands are described which may be useful in some instruments. They are included in the standard so that if the function is implemented, implementation will be done in a standard way.

IEEE 488.2 describes a status-reporting model that is hierarchical. The contents of other data structures are

BUS

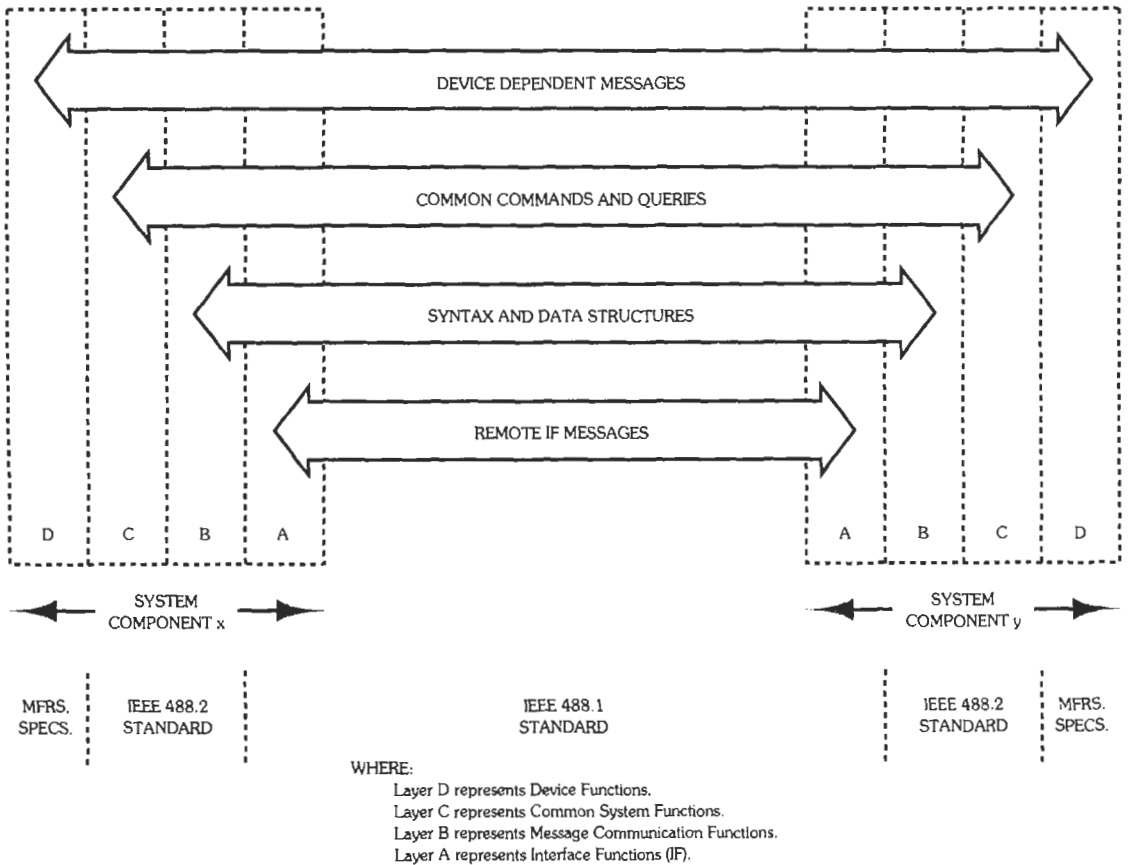


Fig. 74. IEEE 488.2 layers. (From Hewlett-Packard publication 5021-1927. © Hewlett-Packard, used with permission.)

summarized in the status byte. Some of the conditions required are whether a message is available and whether an error has occurred.

Instruments are often able to accept commands faster than they can be executed. While this feature can improve system throughput, it also creates a need to know when the instrument has actually completed all its commands. IEEE 488.2 requires commands which enable the application program to detect when all pending operations are complete.

Standard Commands for Programmable Instruments (SCPI)

SCPI fills in the top layer of Fig. 74 by describing the semantics of the messages. The most important objective of SCPI is to be consistently usable with a large variety of instruments. This horizontal compatibility means that a power supply can be programmed similarly to a microwave network analyzer. Even though these two instruments have very different capabilities, there is no reason to use two different programming languages.

Although having a consistent style is convenient, the ability to use the exact same command for the same function (even with different hardware) is more important. Thus, a voltmeter uses the same command as an oscilloscope for measuring dc voltage.

This compatibility is derived from using a generalized model of an instrument. This model, along with the underlying standards, enforces a discipline on the language designers. In an effort to meet the needs of advancing instruments needs, a consortium of instrument manufacturers manages the growth of the language.

ELECTROMAGNETIC COMPATIBILITY, INTERFERENCE, AND SUSCEPTIBILITY

EMC/EMI/EMS

Electromagnetic compatibility (EMC) is the ability of items of electronic equipment to function properly together in the electronic environment. Most electronic

CHART 1. TYPICAL EMI/EMS MEASUREMENTS

Unwanted Emissions	Unwanted Susceptibility
Conducted	Conducted
Power Lines	Power Lines
Signal lines	Continuous wave
Control lines	Dropouts, surges
Antenna leads	Fast transients
Radiated	Signal lines
Electric fields	Control lines
Magnetic fields	Receiving antenna
	Radiated
	Electric Field
	Magnetic field
	Electrostatic Discharge (ESD)

equipment is designed to perform an assigned task. However, when improperly designed or poorly maintained, a device may emit unintentional radiation to interfere with the proper functioning of other equipment. Sometimes equipment is designed to emit radiation at designated frequencies and with prescribed bandwidths but may inadvertently emit at other frequencies or harmonics. These situations create electromagnetic interference (EMI) to other equipment.

Engineers would prefer to design equipment to survive under the most severe electromagnetic interference,^{*} but they are often constrained by size, weight, and cost requirements. Designing equipment to work at prescribed interference levels is designing for acceptable susceptibility (EMS).

Regulations

In the United States, the Federal Communications Commission is the body responsible for the control of EMI. Part 15 (for radio-frequency devices) and Part 18 (for medical, scientific, and industrial equipment) of the FCC Rules and Regulations[†] contain sections on the control of interference. Devices intended for commercial and scientific applications with limited usage are governed by class A limits. Devices intended for home or consumer use are covered by class B limits, which are up to 22 dB lower in level of conducted emission and 10 dB lower in level of radiated emission. Certification from the FCC is required before equipment may be sold in the USA.

Another source of EMI regulations, particularly for builders of military equipment, is MIL-STD-461. Test methods for the frequency range 30 Hz to 40 GHz are covered in MIL-STD-462.

European Community countries use unified EMC standards (EN)[‡] that regulate the level of permissible electromagnetic emission and also susceptibility. These standards are based on IEC[§] (International Electrotechnical Commission) and CISPR[§] (International Special Committee on Radio Interference) recommendations. Certification to these regulations must be performed before the equipment is offered for sale.

Most countries in Asia and other parts of the world follow the CISPR recommendations in the development of their national EMI standards.

Other international bodies, such as the ITU (International Telecommunication Union), provide detailed limits of EMI for specialized equipment. In general, compliance with their standards is voluntary, but compliance will often mean better chances of success in interfacing with existing equipment.

Typical emission limits are shown in Fig. 75.

EMI/EMS Measurements

Chart 1 lists typical EMI/EMS measurements. An open area test site^{**} (OATS) with little or no electromagnetic ambient noise is ideal for making emission measurements. An anechoic chamber, a shielded room with absorbent material lining all walls, is ideal for susceptibility measurements. Emission measurements in a semianechoic chamber with a conductive floor as ground plane may substitute for OATS measurements within limitations. Unlined screen rooms are afflicted with multiple reflections, which add to or subtract from the true magnitude of the radiation, and are useful only for conducted-EMC measurements.

Inexpensive techniques for the testing of radiated EMC are offered by TEM cells^{††} (transmission lines) if

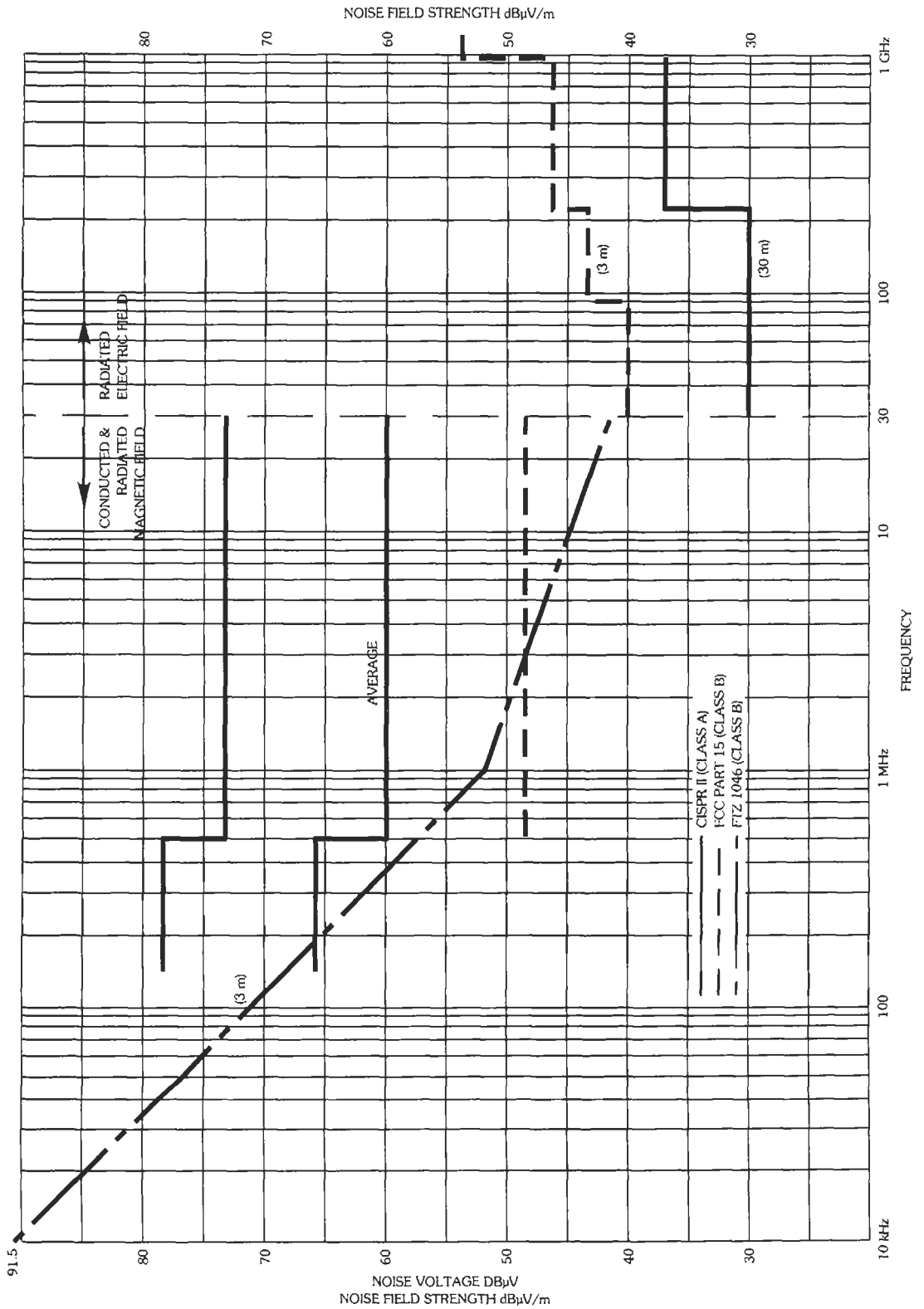
^{*} Ott, Henry W., *Note Reduction Techniques in Electronic Systems*. New York: Wiley-Interscience.

[†] *Code of Federal Regulations, Title 47*. Washington, DC 20402: US Government Printing Office.

[‡] EN, IEC, and CISPR standards may be purchased from: Sales Department, American National Standards Institute, 11 West 42nd Street, New York, NY 10036.

^{**} *American National Standard Guide for Construction of Open Area Test Sites for Performing Radiated Emission Measurements*. ANSI C63.7-1988. (IEEE, Inc., 345 East 47th Street, New York, NY 10017.)

^{††} Crawford, M. L., Workman, J. L. *Using a TEM Cell for EMC Measurements of Electronic Equipment*, NBS Technical Note 1013, 1981. Boulder, CO 80303: NIST.



the equipment is small relative to the dimensions of the cell. For EMI measurements, a receiver is connected as the load to the cell, and for EMS measurements, a signal generator is connected as the source. The TEM cell is completely shielded. As the cell dimensions become comparable to the wavelength, higher-order modes will affect the calibration thus creating a high-frequency usage limit known as the cutoff frequency. The cutoff frequency is raised considerably in a special form of the TEM cell* where the voltage and current on the transmission line are internally terminated and the associated electromagnetic field is absorbed.

Strip-line cells† are sometimes used for testing small electronic modules. They are inexpensive, easy to build, and economical of bench space. Unfortunately, they are not shielded, and their use is therefore limited to preliminary testing and troubleshooting.

Emission sources may be located during troubleshooting with inexpensive, small hand-held probes‡ for magnetic fields. Such probes may also be used to inject fields for localized susceptibility tests.

REFERENCES

Error Correction

1. Kruppa, W., and Sodomsy, K. F. "An explicit solution for the scattering parameters of a linear two-port measured with an imperfect test set." *IEEE Trans. Microwave Theory Tech.*, Vol. MTT-19, Jan. 1971, pp. 122-123.
2. Speciale, R. A. "A generalization of the TSD network-analyzer calibration procedure, covering n-port scattering-parameter measurements, affected by leakage errors." *IEEE Trans. Microwave Theory Tech.*, Vol. MTT-25, Dec. 1977, pp. 1100-1115.
3. Fitzpatrick, J. "Error models for systems measurements." *Microwave J.*, No. 5, 1978, pp. 63-66.
4. Engen, G. F., and Hoer, C. A. "Thru-reflect-line: An improved technique for calibrating the dual six-port automatic network analyzer." *IEEE Trans. Microwave Theory Tech.*, Vol. MTT-27, Apr. 1979, pp. 987-993.
5. Soares, R. A., et al. "A unified mathematical approach to two-port calibration techniques and some applications." *IEEE Trans. Microwave Theory Tech.*, Vol. 37, Nov. 1989, pp. 1669-1674; also Vol. 38, Aug. 1990, pp. 1144-1145.
6. Eul, H. J., and Schick, B. "A generalized theory and new calibration procedures for network analyzer self-calibration." *IEEE Trans. Microwave Theory Tech.*, Vol. 39, Apr. 1991, pp. 724-731.
7. Eul, H. J., and Schiek, B. "Reducing the number of calibration standards for network analyzer calibration." *IEEE Trans. Instrum. Meas.*, Vol. 40, Aug. 1991, pp. 732-735.
8. Butler, John V.; Rytting, Douglas K.; Iskander, Magdy F.; Pollard, Roger D.; and Vanden Bossche, Marc. "16-term error model and calibration procedure for on-wafer network analysis measurements." *IEEE Trans. Microwave Theory Tech.*, Vol. 39, No. 12, Dec. 1991, pp. 2211-2217.
9. Silvonon, Kimmo J. "A general approach to network analyzer calibration." *IEEE Trans. Microwave Theory Tech.*, Vol. 40, No. 4, Apr. 1992.
10. Barnes, J. A., et al. "Characterization of Frequency Stability." *IEEE Trans. Instrum. Meas.*, Vol. IM-20, No. 2, May 1971, pp. 105-120.
11. "Standard Frequencies and Time Signals." Vol. VII, CCIR XIII Plenary Assembly. Geneva: International Telecommunication Union, 1975, p. 18.
12. Fischer, Michael C. "Frequency Stability Measurement Procedures." Proceedings, Eighth Annual PTTI Meeting, Goddard Space Flight Center, Code 250, Greenbelt, Md., 1976, pp. 575-618.
13. Fischer, Michael C. "Analyze Noise Spectra With Tailored Test Gear." *Microwaves*, Vol. 18, No. 7, July 1979, pp. 66-75.
14. *Fundamentals of the Electronic Counters*, Application Note 200. Palo Alto, Calif.: Hewlett-Packard, 1978.
15. *Fundamentals of Time and Frequency Standards*, Application Note 52-1. Palo Alto, Calif.: Hewlett-Packard, 1974.
16. Kartaschoff, Peter. *Frequency and Time*. New York: Academic Press, Inc., 1978.
17. Manassewitsch, Vadim. *Frequency Synthesizers Theory and Design*. 2nd ed. New York: John Wiley & Sons, Inc., 1980.
18. Oliver, B. M., and Cage, J. M., eds. *Electronic Measurements and Instrumentation*. New York: McGraw-Hill Book Co., 1971, Ch. 6.
19. Rutman, Jacques. "Characterization of Phase and Frequency Instabilities in Precision Frequency Sources: Fifteen Years of Progress." *Proceedings of the IEEE*, Vol. 66, No. 9, Sept. 1978, pp. 1048-1075.
20. *Timekeeping and Frequency Calibration*, Application Note 52-2. Palo Alto, Calif.: Hewlett-Packard, 1979.
21. Vessot, Robert F. C., et al. "Research With a Cold Atomic Hydrogen Maser." *Proceedings 33rd Frequency Control Symposium*. Washington, D.C.:

Time and Frequency Measurement

* Wilson, P., Hansen, D., Koenigstein, D. "Simulating Open Area Test Site Emission Measurements Based on Data Obtained in a Novel TEM Cell," IEEE 1989 International Symposium on EMC, IEEE Cat. No. 89CH2736-7, pp. 171-177.

† *Engineering Design Handbook—EMC*, DARCOM-P 706-410 (US Army; pp. 7-49).

‡ Terrien, M. "EMI Troubleshooting Techniques Utilizing a Broadband Magnetic Field Sensor." International Zurich Symposium on EMC, 1987.

Electronic Industries Association, 1979, pp. 511–514.

22. Winkler, Gernot M. R. "A Brief Review of Frequency Stability Measures." Proceedings, Eighth Annual PTTI Meeting, Goddard Space Flight Center, Code 250, Greenbelt, Md., 1976, pp. 489–527.

Microwave-Link Analysis

23. *Differential Phase & Gain at Work*, Application Note 175-1. Palo Alto, Calif.: Hewlett-Packard, Nov. 1975.

24. Tant, M. J. *The White Noise Book*. White Crescent Press Ltd., 1974.

25. "An Integrated Test Set for Microwave Radio Link Baseband Analysis." *Hewlett-Packard Journal*, April 1982.

26. Smith, Emerson C. *Glossary of Communications*. Telephony Publishing Corp., 1971.

27. Preston, C. L. *Broadband Transmission System*, a volume from the series Noise Analysis Theory and Technique. New York: American Telephone and Telegraph Corp., 1972.

13 **Magnetic-Core Transformers and Reactors**

Revised by

*Charles F. Hempstead, Durwood R. Kressler, Richard W. Avery,
Fred J. Banzi, Bernard D. Carniglia, Arthur Olsen, Jr.,
Kent W. Sternstrom, and Richard C. Walker*

Introduction *13-3*

Definition of Transformer and Inductor
Transformer Types and Frequency Ranges
Inductor Types and Frequency Ranges
Generalized Equivalent Circuit for a Transformer

Power Transformers *13-3*

Types of Magnetic Cores
Design of Power Transformers for Rectifiers
Effect of Duty Cycle on Design
Methods of Winding Transformers
Dielectric Insulation and Corona
Temperature and Humidity

Ferroresonant Transformers *13-10*

Converter/Inverter Transformers *13-11*

Audio-Frequency Transformers *13-12*

Types of Magnetic Cores
Design of Audio-Frequency Transformers

Wideband Transformers	<i>13-15</i>
Core-Material Considerations	
Design Example for Carrier Frequencies	
High-Frequency Wideband Transformers	
Pulse Transformers	<i>13-17</i>
Magnetic-Core Reactors	<i>13-18</i>
Rectifier-Filter Reactors	
AC-Filter Reactors	
Magnetic-Core Inductors	<i>13-19</i>
Audio-Frequency Inductors	
Precision Adjustable Inductors	

INTRODUCTION

Definition of Transformer and Inductor

Magnetic-core transformers are static devices containing magnetically coupled windings. They are used in power systems to change values of voltage and current at a single frequency. In communications circuits, often over a wide band of frequencies, they are used to provide direct-current isolation, signal splitting and combining functions, specific current or voltage ratios, impedance matching, and phase inversion.

The Institute of Electrical and Electronics Engineers, Inc. (IEEE) has defined a transformer as follows: "A static device consisting of a winding, or two or more coupled windings, with or without a magnetic core, for introducing mutual coupling between circuits. Note: Transformers are extensively used in electric power systems to transfer power by electromagnetic induction between circuits at the same frequency, usually with changed values of voltage and current."*

Magnetic-core inductors and reactors are static devices containing one or more windings to introduce inductance into an electric circuit. Reactors are used in power circuits primarily to filter alternating current from direct current. Inductors are used in communications systems primarily in frequency-selective circuits.

In this chapter, only those devices having magnetic cores will be considered. The type of core material is known as soft magnetic material, which is defined as ferromagnetic material which, once having been magnetized, is very easily demagnetized (i.e., requires only a slight coercive force to remove the resultant magnetism). A ferromagnetic material usually has relatively high values of specific permeability, and it exhibits hysteresis. The principal ferromagnetic materials are iron, nickel, cobalt, and certain of their alloys.†

Transformer Types and Frequency Ranges

The major types of transformers for both power and communications applications are listed below, along with the general operating frequencies for each type.

<i>Power</i>	
Power transformers	50, 60, and 400 Hz
Ferroresonant transformers	50, 60, and 400 Hz
Converter transformers	100 Hz to 150 kHz
<i>Communications</i>	
Audio-frequency transformers	20 Hz to 20 kHz
Carrier-frequency transformers	20 kHz to 20 MHz

High-frequency transformers	20 MHz to 1000 MHz
Pulse transformers	Repetition rates to 4 MHz

Inductor Types and Frequency Ranges

The major types of reactors or inductors are rectifier-filter reactors, alternating-current reactors, audio-frequency inductors, and precision adjustable inductors for filters. Inductors are used from 20 Hz to 1000 MHz or higher.

Generalized Equivalent Circuit for a Transformer

Fig. 1 shows the equivalent circuit for a generalized transformer having two windings. Commonly accepted nomenclature* is as follows:

- a = turns ratio = N_p/N_s
- C_p = primary equivalent shunt capacitance
- C_s = secondary equivalent shunt capacitance
- E_g = root-mean-square generator voltage
- E_{out} = root-mean-square output voltage
- k = coefficient of coupling
- L_p = primary inductance
- l'_p = primary leakage inductance
- l'_s = secondary leakage inductance
- R_c = core-loss equivalent shunt resistance
- R_g = generator impedance
- R_l = load impedance
- R_p = primary-winding resistance
- R_s = secondary-winding resistance

POWER TRANSFORMERS

Power transformers operate from a low source impedance at a low frequency. Depending on the source of power, the frequency may vary as much as $\pm 5.5\%$ at 50 or 60 Hz, and as much as $\pm 20\%$ at 400 Hz.

Types of Magnetic Cores

The magnetic cores used for power transformers are usually E and I laminations stamped from silicon-iron sheet when low cost is of primary importance. When

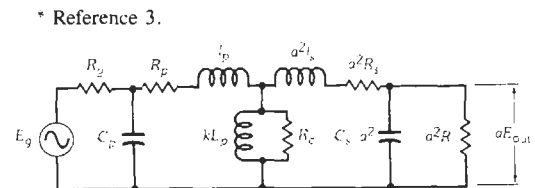


Fig. 1. Equivalent network of a transformer.

* Reference 1.
† Reference 2.

* Reference 3.

minimum size or low loss is of greatest concern, wound-cut "C" cores of oriented silicon steel or the more expensive Supermendur could be considered.

Table 1 lists basic properties of soft magnetic materials. Table 2 gives the maximum-flux-density operating conditions for various core materials at 60 and 400 hertz. Two types of laminations are listed as typical, although there are more types that are of both a higher and lower grade. The 0.014 M-6 material is of a higher grade than the 0.0185 M-19 material, but the cost of the former may be 50% more per pound than the cost of the latter, depending on the lamination size chosen. On the other hand, less of the better-grade material is required for the same performance. In a recent study,* it was shown that transformers and reactors constructed of the higher-grade materials were always smaller, as expected, but that in some cases the total material cost of the structures having the better-grade laminations was the lowest.

Wound-cut "C" cores are wound in tape form on a rectangular mandrel, impregnated, and cut into halves. These halves are then banded around the wound transformer coil.

The sheet-form magnetic materials may also be made in thin strip form and supplied as toroidal tape-wound bobbin cores for high frequencies. The 79 permalloy in 0.000 125-inch tape thickness on a bobbin core may be

used at frequencies as high as 1 MHz. These cores are not cut into halves.

Tables 1 and 2 make reference to "Metallic Glass." METGLAS[†] is Allied Corporation's registered trademark for amorphous alloys of metals. This material is new to the list of available core materials. Its advantage is low core loss and exciting volt-amperes. It is currently manufactured in ribbon form with a nominal thickness of 0.0011 inch and widths of 1.0 to 4.0 inches. The material can be stamped into laminations, but the form lends itself best to tape-wound cores. METGLAS material was originally developed for low-frequency power transformers at a commercially competitive price compared with the standard Si-Fe material. The general-purpose alloy with the formulation $Fe_{81} B_{13.5} Si_{3.5} C_2$, known as Alloy 2605SC, is specified in Tables 1 and 2. Other alloys are available.

The approximate number of exciting volt-amperes is largely dependent on the air gap or equivalent air gap in series with the magnetic path. The values specified in Table 2 are for the material exclusive of the air gap. Information on core loss and exciting volt-amperes for specific lamination sizes may be obtained from Thomas and Skinner.**

Similar information for cut-cores may be obtained from catalogs of the Arnold Engineering Co.

* Reference 4.

† Reference 5.

** Reference 6.

TABLE 1. PROPERTIES OF SOFT MAGNETIC MATERIALS

Material	Initial Permeability	B_S (kG)	Specific Gravity (g/cm ³)	Curie Temp. (°C)	Resistivity (Ω-cm)	Operating Frequency (Hz)
Sheet Form						
SiFe (Unoriented)	400	20	7.65	740	47×10^{-6}	60 to 1000
SiFe (Oriented)	1500	20	7.65	740	50×10^{-6}	60 to 1000
50-50 NiFe (Oriented)	2000	16	8.25	360	40×10^{-6}	60 to 1000
79 Permalloy	12 000	8	8.74	450	55×10^{-6}	1k to 75k
	to	to				
	100 000	11				
Supermendur (Co ₄₉ Fe ₄₉ V ₂)	800	23	8.15	980	26×10^{-6}	60 to 5000
Metallic Glass (Fe ₈₁ B _{13.5} Si _{3.5} C ₂)	2000	16	7.32	370	125×10^{-6}	60 to 100k
Bonded Powder Form						
Permalloy Powder	14	3		450	1.0	10k to 200k
	to					
	550					
Iron Powder	5	10		770	10^4	100k to 100M
	to					
	80					
Ferrite—MnZn	750	3	4.5	100	10	10k to 2M
	to	to	to	to	to	
	15 000	5	5.2	300	100	
Ferrite—NiZn	10	3	3.7	150	10^6	200k to 100M
	to	to	to	to		
	1500	5	5.3	450		

Sources: Allegheny Ludlum, Allied Corp., Armco, Arnold Engineering, Indiana General, Magnetics

TABLE 2. MAXIMUM-FLUX-DENSITY OPERATING CONDITIONS FOR CORE MATERIALS AT 60 AND 400 HZ

Freq. (Hz)	Material Thickness (Inches)	Core Material	Core Flux Density B_m (kG)	Approx. Core Loss (W/lb)	Approx. Exciting VA* (VA/lb)
60	0.0185	M-19 Nonoriented Laminations	12	1.0	3.3
60	0.014	M-6 Oriented Laminations	15	0.65	0.83
60	0.012	Silectron Wound-Cut Core	16	1.1	4.0
60	0.004	Supermendur Wound-Cut Core	21	1.1	4.0
60	0.0011	Metallic-Glass Wound-Cut Core	14	0.14	0.16
400	0.004	Silectron Wound-Cut Core	15	10.0	15.0
400	0.004	Supermendur Wound-Cut Core	20	7.0	30.0
400	0.0011	Metallic-Glass Wound-Cut Core	14	1.6	1.82

Sources: Allied Corp., Arnold Engineering, U.S. Steel
 * Exciting VA values are for the material exclusive of the air gap.

Design of Power Transformers for Rectifiers

The design of transformers for rectifiers was chosen as an example because it represents one of the more difficult transformer design problems. The reason it is difficult is that all parts of all windings are not operating continuously with a sinusoidal current. The result is that the capability of the transformer structure is not as fully utilized as it would be if the transformer supplied a nonrectified resistive load. An early paper by R. W. Armstrong* gave the effective, or rms, values of voltages and currents for the primary and secondary windings of a transformer that provides the power for various types of single-phase and polyphase rectification followed by a reactor-input filter. Table 3 lists single-phase rectifier transformer ratings for both inductive (*L*) and resistive (*R*) filter inputs. A more comprehensive listing of all the circuit parameters of these single-phase as well as polyphase rectifier circuits is found in Chapter 14 of this book. Should the input to the filter be capacitive, higher rms values of current appear in the secondary winding, resulting in a higher volt-ampere rating. A detailed analysis of the effect of the capacitor-input filter on the transformer-winding voltage and current requirements was made by O. H. Schade in 1943.† More recent analyses have been made by R. Lee** and N. R. Grossner.‡

The first step in designing a transformer is to determine its volt-ampere (VA) rating. In the case of the transformer for a rectifier,

$$VA = VA_T E_{dc}' I_{dc}$$

where,

VA_T may be obtained from Table 3, E_{dc}' is the sum of the output dc voltage (E_{dc}), the diode drop (E_{dd}), and the reactor drop (E_{rd}).

In the case of the bridge rectifier, the current passes through two diodes on each half cycle, and it is necessary to use the value of two diode drops in the calculation of E_{dc}' .

At this point, it is handy to have a table of VA ratings versus lamination sizes, such as Table 4. Note the constraints of the 60-Hz frequency and the 75 °C winding temperature rise in an ambient of 75 °C. If a lower frequency or temperature rise is required, then the VA average will be lower for a given structure, with other parameters also being affected. It might also be noted that if a bobbin type of construction is used on the smaller structures instead of a layer-wound construction with interlayer insulation, then up to 30% more copper may be used in a given structure, resulting in up to 30% more VA average at the same frequency and temperature rise. The above calculated VA rating may be used to select a suitable structure from Table 4.

In order to determine the number of turns and wire size of each winding, it is first necessary to determine the voltage and current for each winding. Table 3 gives the relationships between the rms voltages and currents and the dc voltages and currents. The value of E_{dc}'

* Reference 7.
 ‡ Reference 8.
 † Reference 9.
 ** Reference 10.

TABLE 3. * SINGLE-PHASE RECTIFIER TRANSFORMER RATINGS**

	Half-Wave	Full-Wave, Center-Tapped		Bridge	
	R Input	R Input	L Input	R Input	L Input
Ripple frequency . . .	f	$2f$	$2f$	$2f$	$2f$
Ripple V_{rms}	1.11	0.471	0.471	0.471	0.471
Ripple V_{rms} total† . . .	1.21	0.482	0.482	0.482	0.482
Primary E_p	2.22	1.11	1.11	1.11	1.11
Primary I_p	1.21	1.11	1.00	1.11	1.00
Primary VA	2.69	1.23	1.11	1.23	1.11
Secondary E_s	2.22	1.11†	1.11†	1.11	1.11
Secondary I_s	1.57	0.785	0.707	1.11	1.00
Secondary VA	3.49	1.74	1.57	1.23	1.11
Average VA_T	3.09	1.49	1.34	1.23	1.11
PIV per Diode	3.14	3.14	3.14	1.57	1.57
I_{pk} per Diode	3.14	1.57	1.00	1.57	1.00
I_{av} per Diode	1.00	0.50	0.50	0.50	0.50
I_{rms} per Diode	1.57	0.785	0.707	0.785	0.707

* From N. R. Grossner, *Transformers and Electronic Circuits* (New York: McGraw-Hill Book Co., 1967).

** Voltages, currents, and VA are based on unity dc output voltage and current. Voltage drops in transformer, diode, and inductor are neglected.

† Secondary voltage at each side of center tap.

‡ Includes second and third harmonics of ripple frequency.

calculated above is used for the dc voltage. The primary voltage, E_p , is determined by the specified input voltage. The primary rms current, I_p , is determined as follows:

$$I_p = [(VA) (Primary VA)] / [(Efficiency)(E_p)]$$

The quantities "primary VA" and "efficiency" are obtained from Tables 3 and 4, respectively. The secondary rms voltage, E_s , and the secondary rms current, I_s , are obtained from the following relationships:

$$E_s = (Secondary E_s)(E_{dc'})$$

$$I_s = (Secondary I_s)(I_{dc})$$

The quantities "Secondary E_s " and "Secondary I_s " are obtained from Table 3.

Now the number of turns in each winding can be calculated. The number of primary-winding turns, N_p , is obtained as follows:

$$N_p = (E_p \times 10^5) / (4.44 \times f \times A_c \times B_m)$$

where f is the frequency of the applied voltage, which is assumed to be sinusoidal. If the applied voltage has other than a sinusoidal waveform, then the 4.44 constant must be changed. An ASTM standard explains this.* The quantity A_c is the area of the core, and B_m is

the core flux density. If the metric unit the kilogauss (kG) is used, as listed in Table 4, then A_c must be in square centimeters. If B_m is in kilolines per square inch, then the units of A_c must be square inches. The A_c of the structure selected from Table 4 may be computed by squaring the specified tongue width (in the appropriate units) and multiplying by the core manufacturer's recommended stacking factor. Although Table 4 structures are all square stacks (i.e., the stack height equals the tongue width), the stack heights are often increased or decreased for a larger or smaller VA rating.

The number of secondary-winding turns, N_s , is obtained as follows:

$$N_s = N_p(E_s/E_p)[1 + (Copper Regulation / 100)]$$

where "Copper Regulation" (in percent) is obtained from Table 4.

The wire sizes for each winding are determined by dividing the current calculated for the winding by the current density specified in Table 4. This determines the cross-sectional area needed for the conductor. A wire table, such as Table 5, will give the bare area for each wire size. This table also gives the thickness of insulation necessary to support each layer of a given wire size in the layer-wound construction, as well as the maximum diameter of the wire over its film insulation.

The next step is to design the coil that will result in the configuration shown in Fig. 2. This wound coil will then be laminated to form the finished product as shown in Fig. 3.

* Reference 11.

TABLE 4.* VA RATINGS VS LAMINATION SIZES—60 HZ, 75 °C TEMPERATURE RISE**

Lamination E-I Type	Tongue Width, in	Area Product	VA Average	Flux Density, Kilogauss	Current Density, A/in ²	Effi- ciency, %	Core Loss, Watts	Copper Loss, Watts	Copper Regulation, %	Weight, lb	
										Iron	Copper
625	5/8	0.114	9.1	14.0	4 060	63.5	0.5	4.7	52	0.37	0.098
75	3/4	0.237	20.3	14.2	3 480	73.2	0.88	6.6	32.6	0.63	0.182
87	7/8	0.441	40.0	14.4	3 040	80.9	1.47	8.7	22	1.00	0.32
100	1	0.750	72.5	14.6	2 580	84.5	2.3	11.0	15.1	1.49	0.46
125	1 1/4	1.825	163	14.8	2 220	88.6	4.7	16.2	9.9	2.91	1.0
138	1 3/8	2.66	229	14.8	2 130	90.2	6.2	19.0	8.3	3.88	1.44
150	1 1/2	3.80	298	14.8	2 000	91.0	8.1	21.8	7.3	5.05	1.75
175	1 3/4	7.0	524	14.8	1 845	92.8	12.8	27.8	5.3	8.0	2.86
212	2 1/8	15.3	1 050	14.8	1 550	94.7	22.6	37.4	3.6	14.2	5.18
250	2 1/2	29.3	1 823	14.8	1 335	95.7	35.5	47.5	2.6	22.2	8.5
251	2 1/2	68.8	3 551	14.8	935	96.5	49.8	79.2	2.23	31.1	26.6

* *Electro-Technol.*, Vol. 67, No. 1, p. 61, January, 1961. Copyright C-M Technical Publications Corp., 1961.

** Table based on 29-gauge, grain-oriented (M6) silicon steel, square stack. Exciting VA/input VA is 23.5 percent (EI625) to 12.7 percent (EI251). Operating temperature = 75 °C (amb) + 75 °C (rise) = 150 °C. Copper weight will ordinarily be less than values in the table.

TABLE 5. WIRE TABLE FOR TRANSFORMER DESIGN*

AWG Size†	Nom. Bare Area (in ² × 10 ⁻³)	Maximum Overall Diameter		Layer Factor	Nominal Resist. ** (Ω/1000 ft)	Nominal Weight (lb/1000 ft)	Minimum Margin <i>m</i> (inches)	Kraft Layer Insulation (inches)
		Single Build (inches)	Heavy Build (inches)					
14	3.227	0.0666	0.0682	0.90	2.52	12.44	0.1875	0.0100
15	2.561	0.0594	0.0609	0.90	3.18	9.87	0.1562	0.0100
16	2.027	0.0531	0.0545	0.90	4.02	7.812	0.1562	0.0100
17	1.612	0.0475	0.0488	0.90	5.05	6.213	0.1562	0.0070
18	1.276	0.0424	0.0437	0.90	6.39	4.914	0.125	0.0070
19	1.012	0.0379	0.0391	0.90	8.05	3.900	0.125	0.0070
20	0.804	0.0339	0.0351	0.90	10.13	3.099	0.125	0.0050
21	0.638	0.0303	0.0314	0.90	12.77	2.459	0.125	0.0050
22	0.503	0.0270	0.0281	0.90	16.20	1.937	0.125	0.0050
23	0.401	0.0243	0.0253	0.90	20.30	1.546	0.125	0.0050
24	0.317	0.0217	0.0227	0.90	25.67	1.223	0.125	0.0020
25	0.252	0.0194	0.0203	0.90	32.37	0.970	0.125	0.0020
26	0.1986	0.0173	0.0182	0.89	41.02	0.7650	0.125	0.0020
27	0.1584	0.0156	0.0164	0.89	51.44	0.6101	0.125	0.0020
28	0.1247	0.0140	0.0147	0.89	65.31	0.4806	0.125	0.0015
29	0.1003	0.0126	0.0133	0.89	81.21	0.3866	0.125	0.0015
30	0.0785	0.0112	0.0119	0.89	103.71	0.3025	0.0938	0.0015
31	0.0622	0.0100	0.0108	0.88	130.9	0.2398	0.0938	0.0015
32	0.0503	0.0091	0.0098	0.88	162.0	0.1937	0.0938	0.0013
33	0.0396	0.0081	0.0088	0.88	205.7	0.1526	0.0938	0.0013
34	0.0312	0.0072	0.0078	0.88	261.3	0.1201	0.0938	0.0010
35	0.0246	0.0064	0.0070	0.88	330.7	0.0949	0.0938	0.0010
36	0.0196	0.0058	0.0063	0.87	414.8	0.07569	0.0938	0.0010
37	0.0159	0.0052	0.0057	0.87	512.1	0.06128	0.0938	0.0010
38	0.01257	0.0047	0.0051	0.87	648.2	0.04844	0.0625	0.0010
39	0.00962	0.0041	0.0045	0.86	846.6	0.03708	0.0625	0.0007
40	0.00755	0.0037	0.0040	0.86	1079.2	0.02910	0.0625	0.0007
41	0.00616	0.0033	0.0036	0.85	1323.	0.02374	0.0625	0.0007
42	0.00491	0.0030	0.0032	0.85	1659.	0.01892	0.0625	0.0005
43	0.00380	0.0026	0.0029	0.85	2143.	0.01465	0.0625	0.0005
44	0.00314	0.0024	0.0027	0.85	2593.	0.01210	0.0625	0.0005

* Data for this table are courtesy of NEMA (ref. 12) and Phelps Dodge (ref. 13).

† Square or rectangular wire is recommended for wire sizes heavier than AWG 14.

Single build is not recommended for these sizes.

** Resistance values are at 20 °C.

In designing the coil, one first selects the core tube thickness (J in Fig. 2). This thickness may vary from 0.025 to 0.050 inch for lamination stack heights (p in Fig. 3) of 0.5 to 2 inches. Next, calculate the coil-build (a in Fig. 3) for each winding, using the number of turns and wire sizes previously calculated.

$$a = 1.1[n_i(D + t) - t + t_c]$$

where,

n_i is the number of layers in the winding,

D is the diameter of the insulated wire,

t_c is the thickness of the insulation under and over the winding.

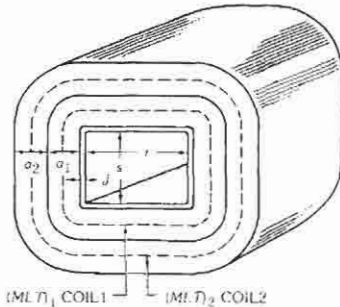
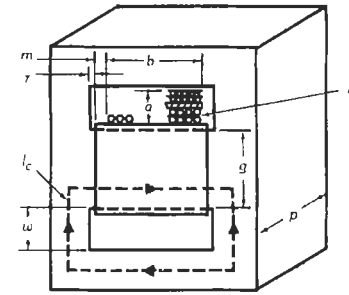


Fig. 2. Dimensions relating to coil mean length of turn (MLT).



- A_c = core area = $(gp)k$
- u = height of coil
= coil-build
- b = coil width
- g = width of lamination tongue
- l_c = average length of magnetic-flux path
- k = stacking factor
- = 0.90 for 14-mil lamination
- = 0.80 for 2-mil lamination or ribbon-wound core
- m = marginal space given in Table 5
- p = height of lamination stack
- t = thickness of interlayer insulation
- u = width of core window
- r = window length tolerance
- = 1/16 in. total

Fig. 3. Dimensions relating to the design of a transformer coil-build and core.

The numeric 1.1 allows for a 10% bulge factor. This factor may be reduced, depending on the winding equipment and methods of winding used. The total coil-build should not exceed 85–90% of the window width.

Compute the mean length per turn (MLT) of each winding from the geometry of the core and windings (Fig. 2).

$$(MLT)_1 = 2(r + J) + 2(s + J) + \pi a_1$$

$$(MLT)_2 = 2(r + J) + 2(s + J) + \pi(2a_1 + a_2)$$

where,

- a_1 is the build of the first winding,
- a_2 is the build of the second winding,
- J is the thickness of the winding form,
- r and s are the winding-form dimensions.

Calculate the total length and resistance of each winding, and determine the IR drop and I^2R loss for each winding.

Make corrections, if required, in the number of turns of the windings to allow for the IR drops, so as to have the required E_s .

$$E_s = (E_p - I_p R_p)(N_s/N_p) - I_s R_s$$

Compute the core losses from the weight of the core and Table 2 or reference 6.

Determine the percent efficiency η and voltage regulation (VR) from

$$\eta = \frac{W_{out} \times 100}{W_{out} + (\text{core loss}) + (\text{copper loss})}$$

$$(VR) = \frac{I_s [R_s + (N_s/N_p)^2 R_p]}{E_s}$$

For a more accurate evaluation of voltage regulation, determine leakage-reactance drop = $I_{dc} \omega l_{sc} / 2\pi$, and add to the above (VR) the value of $I_{dc} \omega l_{sc} / 2\pi E_{dc}$. Here, l_{sc} = leakage inductance viewed from the secondary; see "Methods of Winding Transformers" (in this chapter) to evaluate l_{sc} .

Bring out all terminal leads. Use the wire of the coil, insulated with suitable sleeveings, for all sizes of wire heavier than 21. Use 7-30 stranded and insulated wire for smaller sizes.

Effect of Duty Cycle on Design

If a transformer is operated at different loads according to a regular duty cycle, the equivalent volt-ampere (VA) rating is

$$(VA)_{eq} = \left[\frac{(VA)_1^2 t_1 + (VA)_2^2 t_2 + (VA)_3^2 t_3 + \dots + (VA)_n^2 t_n}{t_1 + t_2 + t_3 + \dots + t_n} \right]^{1/2}$$

where $(VA)_1$ = output during time t_1 , etc.

Example: 5 kilovolt-ampere output, 1 minute on, 1 minute off.

$$(VA)_{eq} = \left[\frac{(5000)^2(1) + (0)^2(1)}{1 + 1} \right]^{1/2} = \left[\frac{(5000)^2}{2} \right]^{1/2}$$

$$= 5000/(2)^{1/2} = 3535 \text{ volt-amperes}$$

Methods of Winding Transformers

The most common methods of winding transformers are shown in Fig. 4. Leakage inductance is reduced by interleaving, i.e., by dividing the primary or secondary coil into two sections and placing the other winding between the two sections. Interleaving may be accomplished by concentric and by coaxial windings, as shown in Figs. 4B and C. Reduction of leakage inductance is computed from

$$I_{sc} = \frac{10.6 N^2 (MLT) (2nc + a)}{n^2 b \times 10^9}$$

where l_{sc} is the leakage inductance (referred to the winding having N turns) in henrys and the dimensions are in inches and to be the same for Figs. 4B and C.

Means of reducing leakage inductance are:

- (A) Minimize turns by using high-permeability core.
- (B) Reduce build of coil.

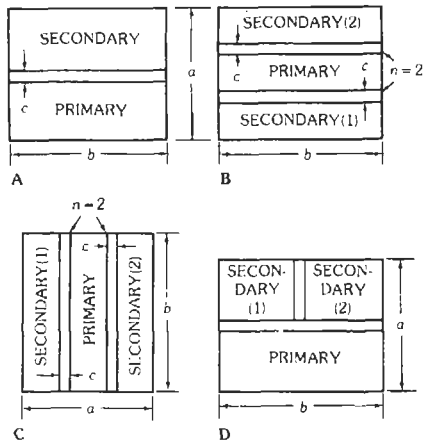


Fig. 4. Methods of winding transformers.

- (C) Increase winding width.
 (D) Minimize spacing between windings.
 (E) Use bifilar windings.

Means of minimizing capacitance are:

- (A) Increase dielectric thickness t .
 (B) Reduce winding width b and thus area A .
 (C) Increase number of layers.
 (D) Avoid large potential differences between winding sections, as the effect of capacitance is proportional to applied potential squared.

Note: Leakage inductance and capacitance requirements must be compromised in practice since corrective measures are opposites.

Effective interlayer capacitance of a winding may be reduced by sectionalizing it as shown in Fig. 4D. This can be seen from

$$C_e = (4C_l/3n_l)(1 - 1/n_l)$$

where,

- C_e = effective capacitance in picofarads,
 n_l = number of layers,
 C_l = capacitance of one layer to another
 = $0.225 A\epsilon/t$ picofarads,
 A = area of winding layer = $(MLT)b$ inches²,
 t = thickness of interlayer insulation in inches,
 ϵ = dielectric constant ≈ 3 for paper.

Dielectric Insulation and Corona

For class A (Table 6), a maximum dielectric strength of 40 volts/mil is considered safe for small thicknesses of insulation. At high operating voltages, due regard must be paid to corona that occurs before dielectric breakdown and will in time deteriorate insulation and cause dielectric failure. Best practice is to operate

insulation at least 25 percent below the corona starting voltage. Approximate 60-hertz rms corona voltage V is

$$\log \frac{V \text{ (in volts)}}{800} = (2/3) \log (100t)$$

where t = total insulation thickness in inches. This may be used as a guide in determining the thickness of insulation. With the use of varnishes that require no solvents, but solidify by polymerization, the bubbles present in the usual varnishes are eliminated, and much higher operating voltages and, hence, reduction in the size of high-voltage units may be obtained. Epoxy resins and some polyesters belong in this group. In the design of high-voltage transformers, the creepage distance required between wire and core may necessitate the use of insulating channels covering the high-voltage coil, or taping of the latter. For units operating at 10 kilovolts or higher, oil insulation will greatly reduce creepage and, hence, the size of the transformer.

Temperature and Humidity

Table 6 lists the standard classes of insulating materials and their limiting operating temperatures. Table 7 compares the properties of four high-temperature wire-insulating coatings.

Open-type constructions generally permit greater cooling than enclosed types, thus allowing smaller sizes for the same power ratings. Moderate humidity protection may be obtained by impregnating and dip-coating or molding transformers in polyester or epoxy resins; these units provide good heat dissipation but are not as good in this respect as completely open transformers.

Protection against the detrimental effects of humidity is commonly obtained by enclosing transformers in hermetically sealed metal cases. This is particularly important if very fine wire, high output voltage, or direct-current potentials are involved. Heat conductivity to the case exterior may be improved by the use of asphalt or thermosetting resins as filling materials. Best conductivity is obtained with high-melting-point silica-filled asphalts or resins of the polyester or epoxy types. Coils impregnated with these resins dissipate heat best, since voids in the heat path may be eliminated.

Immersion in oil is an excellent means of removing heat from transformers. An air space or bellows must be provided to accommodate expansion of oil when heated.

FERRORESONANT TRANSFORMERS

Ferroresonant transformers, also known as constant-voltage transformers or ferroresonant voltage regulators, make use of an alternating-current phenomenon involving at least one magnetic component in combination with a suitable capacitor to provide stabilized voltage to a load when the source voltage fluctuates. A

TABLE 6. CLASSIFICATION OF ELECTRICAL INSULATING MATERIALS†

Class	Insulating Material	Limiting Insulation Temperature (Hottest Spot) in °C
O	Materials or combinations of materials such as cotton, silk, and paper without impregnation*	90
A	Materials or combinations of materials such as cotton, silk, and paper when suitably impregnated or coated or when immersed in a dielectric liquid*	105
B	Materials or combinations of materials such as mica, glass fiber, asbestos, etc., with suitable bonding substances*	130
F	Same as for Class B	155
H	Materials or combinations of materials such as silicone elastomer, mica, glass fiber, etc., with suitable bonding substances such as appropriate silicone resins*	180
C	*	220
Over C	Materials consisting entirely of mica, porcelain, glass, quartz, and similar inorganic materials*	Over 220

* (Other) materials or combinations of materials may be included in this class if by experience or accepted tests they can be shown to have comparable thermal life at the temperature given in the right-hand column.

These temperatures are, and have been in most cases over a long period of time, benchmarks descriptive of the various classes of insulating materials, and various accepted test procedures have been or are being developed for use in their identification. They should not be confused with the actual temperatures at which these same classes of insulating materials may be used in the various specific types of equipment nor with the temperatures on which specified temperature rise in equipment standards are based.

In the above definitions the words "accepted tests" are intended to refer to recognized test procedures established for the thermal evaluation of materials by themselves or in simple combinations. Experience or test data, used in classifying insulating materials, are distinct from the experience or test data derived for the use of materials in complete insulation systems. The thermal endurance of complete systems may be determined by suitable test procedures.

A material that is classified as suitable for a given temperature may be found suitable for a different temperature, either higher or lower, by an insulation system test procedure. For example, it has been found that some materials suitable for operation at one temperature in air may be suitable for a higher temperature when used in a system operated in an inert gas atmosphere. Likewise some insulating materials when operated in dielectric liquids will have lower or higher thermal endurance than in air.

It is important to recognize that other characteristics (in addition to thermal endurance) such as mechanical strength, moisture resistance, and corona endurance are required in varying degrees in different applications for the successful use of insulating materials.

†From "Insulation class ratings," *IEEE Standard Dictionary of Electrical and Electronics Terms*, IEEE Std. 100-1972

magnetic component must have a magnetization characteristic that is sharply nonlinear, and it must operate above the "knee" of the characteristic, where the change in voltage across it is small compared to the change in current through it. A set of design equations for commonly used ferroresonant regulators and regulated rectifiers are derived in a paper by Hart and Kakalec.* These equations are for the basic open-loop type of ferroresonant voltage regulator.

A more recent development is the closed-loop, or feedback-controlled, ferroresonant voltage regulator.† Whereas the regulating function of the open-loop type depends on the stability of the frequency of the source and the magnetic characteristics of the core, the regulating function of the closed-loop type is dependent on the

voltage stability of the reference in the control circuit, and also on the gain of the feedback loop. Regulation better than 0.5% for line, load, frequency, and temperature changes can be attained. The circuit of reference 15 requires the use of a magnetic component, an inductor, in the control circuit. A later development** eliminates the need for this inductor as a separate component by the use of a double-shunt feedback ferroresonant transformer.

CONVERTER/INVERTER TRANSFORMERS

The IEEE‡ defines a converter as a machine that changes alternating-current power to direct-current

* Reference 14.
† Reference 15.

** Reference 16.
‡ Reference 1.

TABLE 7. COMPARISON OF FOUR WIRE INSULATING COATINGS*

NEMA Std. Designation	MW 15C	MW 28C	MW 35C	MW 16C
Thermal Class	105 °C	130 °C	200 °C	220 °C
Insulating Material	Polyvinyl formal, modified	Polyurethane, followed by Nylon	Modified Polyester, followed by Polyamide Imide	Aromatic Polyimide
Min. Dielectric Strength at Rated Temperature of Single-Film Coated 36 AWG	1900 V	1725 V	1900 V	1900 V
Scrape Resistance (Min. Grams-to-Fail of Single-Film Coated 30 AWG)	250	250	250	160
Min. Thermo-plastic Flow Temperature of Single-Coated 36 AWG	180 °C	170 °C	300 °C	400 °C
Trade Names				
Anaconda	Formvar	Nylac	AP 2000	ML
Phelps Dodge	Formvareze	Nyleze	Armored Polythermaleze	ML
Rea	Formvar	Nysol	Thermamid	Pyre-ML
Westinghouse	Formvar	Nythane	Omegaklad	ML

* From NEMA Standards Publication No. MW 1000-1981, Revised July 1982.

power or vice versa, or from one frequency to another. It also defines an inverter as a machine, device, or system that changes direct-current power to alternating-current power.

Transformers are used in static dc-to-dc converters and dc-to-ac inverters. The core materials in these types of transformers may or may not be driven to a saturation flux density, depending on the control circuits involved. Also, the core may operate in a double-ended (bipolar) or single-ended (unipolar) mode, again dependent on the external circuitry. Bipolar is also known and defined by ASTM* as a symmetrically cyclically magnetized condition (SCM); unipolar is defined as a cyclically magnetized condition (CM). The latter has flux-current loops that are not symmetrical with respect to the origin of the axes.

Details of how to design the circuitry and the magnetics (transformer and inductor) may be found in reference 17. Details on the selection of magnetic materials for static converter and inverter transformers may be found in reference 18.

AUDIO-FREQUENCY TRANSFORMERS

Audio-frequency transformers are used mainly for matching impedances and transmitting audio frequencies. They also provide isolation from direct currents and present balanced impedances to lines or circuits.

Types of Magnetic Cores

The magnetic core for this type of transformer is usually an EI or EE type using either audio-grade silicon steel or nickel-alloy steel (refer to Tables 1 and 10). High-permeability nickel-alloy tape cores in toroidal form are used for extreme bandwidths. High-permeability ferrite cores with highly polished mating surfaces are also used to obtain wide bandwidths.

If there is no direct current in the windings, it is possible to design a small, high-quality audio transformer with a high-permeability ferrite core. Material with initial permeability above 10 000 is available. To minimize the air gap in the assembled core, the mating surfaces are ground flat and polished to a mirror finish. Because of the small core volume, when these transformers are used at low frequencies care must be taken

* Reference 11.

to avoid core saturation. The maximum flux density in a core is given by:

$$B_{\max} = (3.49 \times 10^6) E/fNA_c$$

where,

- B_{\max} = maximum flux density (gauss),
- E = rms volts,
- f = frequency (hertz),
- N = number of winding turns,
- A_c = cross section of core (square inches).

Design of Audio-Frequency Transformers

Important Parameters—Important parameters are generator and load impedances, R_g and R_l , respectively; generator voltage E_g ; frequency band to be transmitted; harmonic distortion; and operating voltages (for adequate insulation). See Fig. 1. Refer to the section on power-transformer design for details about physical design, cores, winding, and so forth.

Midband Frequencies—The relative low- and high-frequency responses are taken with reference to midband frequencies where

$$aE_{\text{out}}/E_g = [(1 + R_s/R_l) + R_l/a^2R_1]^{-1}$$

Low Frequencies—At low frequencies, the equivalent unity-ratio network of a transformer becomes approximately as shown in Fig. 5.

$$\text{Amplitude} = [1 + (R'_{\text{par}}/X_m)^2]^{-1/2}$$

$$\text{Phase angle} = \tan^{-1}(R'_{\text{par}}/X_m)$$

where,

$$\begin{aligned} R'_{\text{par}} &= (R_1R_2a^2)/(R_1 + R_2a^2) \\ R_1 &= R_g + R_p \\ R_2 &= R_l + R_s \\ X_m &= 2\pi fL_p \end{aligned}$$

In a good output transformer, R_p , R_s , and R_c may be neglected. In input or interstage transformers, R_c may be omitted.

High Frequencies—At high frequencies, neglecting the effect of winding and other capacitances, the equivalent unity-ratio network becomes approximately as in Fig. 6.

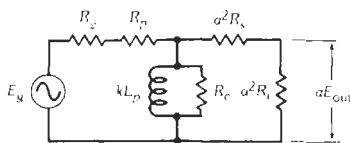


Fig. 5. Equivalent network of an audio-frequency transformer at low frequencies.

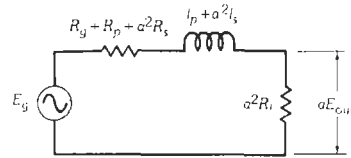


Fig. 6. Equivalent network of an audio-frequency transformer at high frequencies, neglecting the effect of the winding shunt capacitances.

$$\text{Amplitude} = [1 + (X_l/R'_{\text{sc}})^2]^{-1/2}$$

$$\text{Phase angle} = \tan^{-1}(X_l/R'_{\text{sc}})$$

where,

$$\begin{aligned} R'_{\text{sc}} &= R_1 + R_2a^2, \\ X_l &= 2\pi f l_{\text{scp}}, \\ l_{\text{scp}} &= \text{inductance measured across primary with secondary short-circuited} = l_p + a^2l_s. \end{aligned}$$

The low- and high-frequency responses are shown by the curves of Fig. 7.

If at high frequencies the effect of winding and other capacitances is appreciable, the equivalent network on a 1:1-turns-ratio basis becomes as shown in Fig. 8. In a step-up transformer, C_2 = equivalent shunt capacitances of both windings. In a step-down transformer, C_2 shunts both leakage inductances and R_2 . The relative high-frequency response of this network is given by

$$\frac{(R_1 + R_2)/R_2}{[(R_1/X_c + X_l/R_l)^2 + (X_l/X_c - R_g/R_l - 1)^2]^{1/2}}$$

This high-frequency response is plotted in Fig. 9 for $R_2 = R_1$ (matched impedances), based on simplified equivalent networks as indicated. At frequency f_r , $X_l = X_c$ and $B = X_c/R_l$.

Harmonic Distortion—Harmonic distortion requirements may constitute a deciding factor in the design of transformers. Such distortion is caused by either variations in load impedance or nonlinearity of magnetizing current. The percent harmonic voltage appearing in the output of a loaded transformer is given by

$$\begin{aligned} &(\text{Percent harmonics}) \\ &= 100E_h/E_f = 100I_h/I_f(R'_{\text{par}}/X_m)[1 - (R'_{\text{par}}/4X_m)] \end{aligned}$$

where $100I_h/I_f$ = percent of harmonic current measured with a zero-impedance source (values in Table 8 are for a 4-percent silicon-steel core).*

* N. Partridge, "Harmonic Distortion in Audio-Frequency Transformers," *Wireless Engineer*, Vol. 19; September, October, and November 1942.

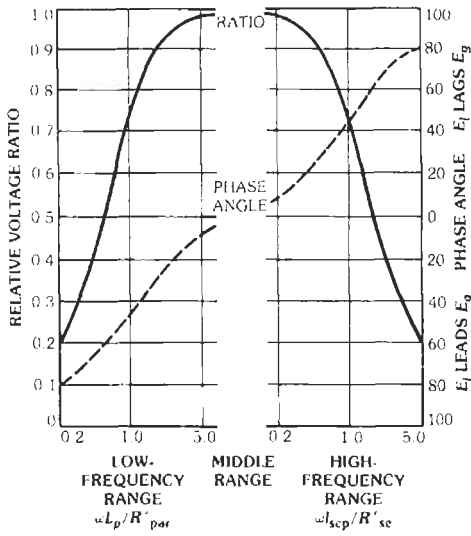


Fig. 7. Universal frequency and phase response characteristics of output transformers. (Courtesy McGraw-Hill Book Co.)

Insertion Loss—Insertion loss is the loss introduced into the circuit by addition of the transformer. At midband, the loss is caused by winding resistance and

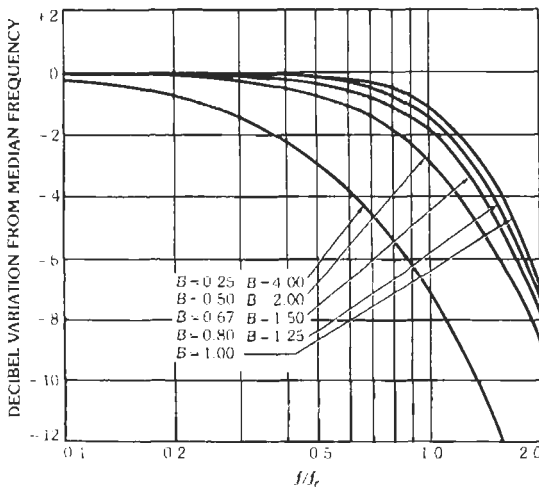
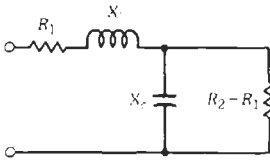


Fig. 9. Transformer characteristics at high frequencies for matched impedances. At frequency f_1 , $X_l = X_c$ and $B = X_c/R_1$. (Reprinted from *Electronic Transformers and Circuits*, by R. Lee, 2nd ed., p. 151, 1955; by permission, John Wiley & Sons, New York.)

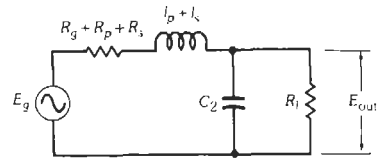


Fig. 8. Equivalent network of a 1:1-turns-ratio audio-frequency transformer at high frequencies when effect of winding shunt capacitances is appreciable.

core loss. Frequency discrimination adds to this at low and high frequencies. Insertion loss is input divided by output expressed in decibels, or (in terms of measured voltages and impedance)

$$(\text{dB insertion loss}) = 10 \log[(E_g^2 R_l)/(4E_o^2 R_g)]$$

Impedance Match—For maximum power transfer, the reflected load impedance should equal the generator impedance. Winding resistance should be included in this calculation: For matching

$$R_g = a^2(R_l + R_s) + R_p$$

Also, in a properly matched transformer

$$R_g = a^2 R_l = (Z_{oc} \times Z_{sc})^{1/2}$$

where,

Z_{oc} = transformer primary open-circuit impedance,
 Z_{sc} = transformer primary impedance with secondary winding short-circuited.

If more than one secondary is used, the turns ratio to match impedances properly depends on the power delivered from each winding.

$$N_s/N_p = [(R_n/R_g)(W_n/W_p)]^{1/2}$$

Example: Using Fig. 10

$$N_2/N_p = [(10/600)(10/16)]^{1/2} = 0.102$$

$$N_3/N_p = [(50/600)(5/16)]^{1/2} = 0.161$$

$$N_4/N_p = [(100/600)(1/16)]^{1/2} = 0.102$$

TABLE 8. HARMONICS PRODUCED BY VARIOUS FLUX DENSITIES B_m IN A 4-PERCENT SILICON-STEEL-CORE AUDIO TRANSFORMER

B_m	Percent 3rd Harmonic	Percent 5th Harmonic
100	4	1.0
500	7	1.5
1000	9	2.0
3000	15	2.5
5000	20	3.0
10 000	30	5.0

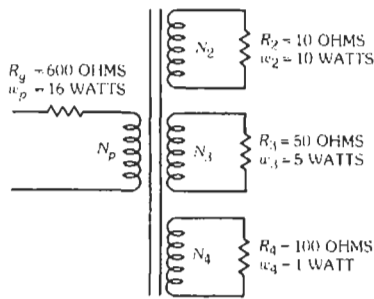


Fig. 10. Multisecundary audio transformer.

WIDEBAND TRANSFORMERS

Wideband transformers operate over a range of one to three or more frequency decades without significant attenuation. The frequency band of interest can be considered to lie between 20 kHz and 1 GHz. This frequency range can be broken into two subranges: carrier frequency (20 kHz to 20 MHz) and high frequency (20 MHz to 1 GHz). Wideband transformers are used to provide impedance matching, voltage or current ratios, dc isolation, connection of balanced to unbalanced circuits, power splitting, and phase inversion.

Maximum bandwidth is achieved when the ratio of shunt inductance to leakage inductance is maximized, neglecting the effect of shunt capacitance. To increase this ratio, we must use the highest μ_e possible and make the ratio of winding width to core length approach unity. The latter occurs for a toroidal core, but is smaller for spool-wound structures. Modern wideband transformers utilize ferrite cores. Some commonly used core shapes are EE, EI, EP, RM, pot, Q, toroid, and X.

Core-Material Considerations

Carrier-frequency transformers often use cores of MnZn ferrite having material permeability (μ) in the 1000 to 20 000 range. High-frequency transformers sometimes use NiZn ferrite cores having material permeability as low as 15, especially when low core loss is important. Most high-frequency transformers operate at relatively low impedance levels. To obtain effective permeability above about 10 000, either toroidal cores or split cores having extremely small air gaps are needed. Highly polished, clean, and coplanar core mating surfaces can achieve air gaps as small as 0.2 μm . Even at such small air gaps, μ_e can be considerably reduced from μ and is given by:

$$\mu_e = l_e \mu / [l_c + (\mu - 1)l_g]$$

where,

- μ_e = effective permeability,
- μ = material permeability,
- l_e = effective core magnetic path length,
- l_g = length of air gap.

Electrical loss in ferrite material is usually expressed as a relative loss coefficient, $\tan \delta / \mu = \omega L_p / R_p$. It has three components: hysteresis loss, eddy-current loss, and residual loss. The effect in transformer cores can be more conveniently expressed as a parallel resistance per turn squared in ohms (R_c / N^2) and plotted as a function of frequency. Ferrite core loss is a frequency-dependent nonlinear resistance; hence it cannot be equalized easily by a simple reactive network, should its effects be significant.

The Curie temperature of the core material can be a consideration when transformers are operated at high ambient temperatures. High-permeability materials tend to have lower Curie temperatures. Core temperature and disaccommodation factors are usually not as critical in wideband transformers as in stable filter-class inductors.

Design Example for Carrier Frequencies

Assume the following problem: Design a 75-to-150-ohm unbalanced transformer to operate over a frequency range of 50 kHz to 10 MHz. In-band loss shall be less than 0.1 dB, and loss shaping at band ends shall be less than 0.1 dB. The transformer will not carry any direct current and should be small in size (less than 0.75 in² in area and less than 0.60 in high).

The initial step is to determine the minimum inductance for low-frequency loss shaping.

$$\text{Loss in dB} = 10 \log_{10} [1 + (R/\omega L_p)^2]$$

where,

- $R = (R_g R') / (R_g + R')$,
- $R_g = 75 \text{ ohms}$,
- $R_l = 150 \text{ ohms}$,
- $R' = a^2 R_l$,
- $a = \text{turns ratio} = N_p / N_s$

For a loss of 0.1 dB maximum

$$1 + (37.5/\omega L_p)^2 = 10^{0.1/10}$$

From this, it is found that $\omega L_p = 245.7 \text{ ohms}$ minimum at 50 kHz, and L_p therefore must be greater than 0.782 millihenries.

A minimum value for L_p could also be determined from a return-loss requirement:

$$\text{Return loss in dB} = 10 \log_{10} [1 + (2\omega L_p / R')^2]$$

The minimum inductance previously selected for loss shaping would provide a return loss relative to R' of only 16.4 dB. If a greater return loss were specified, L_p minimum would have to be larger. Return-loss requirements for wideband transformers frequently cause more severe constraints than transformer loss shaping. Transformer loss is related to return loss by the expression:

$$\text{Transformer loss} = -10 \log_{10} (1 - 10^{-RL/10})$$

The expression is applicable when the equivalent circuit is considered to contain pure reactances.

The EP13 core is a small, efficient core that meets the size constraints of this example. It is available with an A_L value of 3500 nH/N². With no direct current, no air gap is required between the core halves, and the full 3500 nH/N² may be utilized.

$$N_p = (782/3.5)^{1/2} = 14.9 \text{ turns minimum}$$

$$N_s/N_p = (150/75)^{1/2} = 1.41$$

Primary turns of 17 and secondary turns of 24 will closely approximate the 1.41 ratio. This number of primary turns results in a minimum shunt inductance of 1.012 millihenrys ($17^2 \times 3.5 \mu\text{H}/N^2$), providing an inductance margin to allow for variations in assembly.

High-frequency loss shaping is controlled by both leakage inductance and effective shunt capacitance. Interleaved windings reduce leakage inductance but cause increased shunt capacitance. For a low-impedance transformer, the leakage inductance dominates, while in a high-impedance transformer the shunt capacitance is controlling. Windings usually can be oriented to balance and minimize the total contributions to shaping loss from leakage inductance and shunt capacitance. As a starting point, allocate half the requirement to each effect; thus, assume 0.05 dB caused by the leakage inductance.

The leakage inductance of the primary winding with the secondary winding shorted (l_{sc}) is that of both primary and secondary windings referred to the primary.

$$\text{Loss in dB} = 10 \log_{10} [1 + (\omega l_{sc})^2 / (R_s + R')^2]$$

$$1 + (\omega l_{sc}/150)^2 = 10^{0.05/10}$$

$$\omega l_{sc} = 16.14 \text{ ohms maximum}$$

$$l_{sc} = 0.257 \mu\text{H max. at 10 MHz}$$

By singly interleaving the secondary around the primary, leakage can be reduced to a sufficiently low value. The primary is a single layer of 34 AWG wire wound as a parallel pair to achieve a smooth full single layer. Each secondary section is wound as a single layer of 31 AWG wire. Two layers of 2.5-mil Mylar tape are applied between windings. From the winding geometry, the leakage inductance referred to the primary is calculated to be 0.18 microhenry.

The loss contribution from effective shunt capacitance is also allowed to be 0.05 dB.

$$\text{Loss in dB} = 10 \log_{10} [1 + (\omega RC)^2]$$

$$1 + (\omega RC)^2 = 10^{0.05/10}$$

$$\omega RC = 0.1076$$

$$C = 45.7 \text{ pF max. at 10 MHz}$$

The direct capacitance between each secondary interleaved winding and the primary is calculated to be 32.8 pF. Each is reflected to the primary as $(C \text{ direct})/3$ with

appropriate corrections made for turns ratio, relative winding direction, and reversing or nonreversing circuit connection. The effective shunt capacitance is calculated to be less than 40 pF.*

Midband loss results from both core and copper losses. Core loss (R_c) is a minimum of 22 ohms/N². Therefore, R_c is ($17^2 \times 22$), or 6358 ohms minimum.

$$\begin{aligned} \text{Loss in dB} &= 20 \log_{10} (1 + R/R_c) \\ &= 0.051 \text{ dB max.} \end{aligned}$$

The series ac resistances of the primary and secondary windings are R_p and $a^2 R_s$ when referred to the primary winding. The value of R_p is calculated to be 0.177 ohm maximum, and the value of $a^2 R_s$ is calculated to be 0.062 ohm maximum. These values lead to a copper loss of

$$\begin{aligned} \text{Loss in dB} &= 20 \log_{10} [1 + (R_p + a^2 R_s) / (R_s + R')] \\ &= 0.014 \text{ dB max.} \end{aligned}$$

The total midband loss should be less than 0.065 decibel.

This transformer design was realized. The loss-frequency characteristic is shown in Fig. 11B. Loss shaping at the band ends, 50 kHz and 10 MHz, is 0.04 dB and 0.07 dB, respectively. Midband loss is 0.035 dB.

High-Frequency Wideband Transformers

Circuit impedances are usually low at high frequencies (20 MHz to 1 GHz). Transformer designs typically require only a few turns because shunt inductances are small and leakage inductances must be kept low. Transformers are small in size, which is compatible with low leakage inductance and low effective shunt capacitance. Winding dc resistances are usually negligible, but core loss can be important for low turns and some core materials. Small toroidal cores are often used because of their ability to achieve low leakage inductance.

Transmission-line techniques can be applied to transformer windings at high frequencies by tightly coupling the primary and secondary windings to form a transmission line having a particular characteristic impedance. The overall characteristic impedance (Z) of a transformer is approximately equal to $(Z_p Z_s)^{1/2}$. This impedance is matched as nearly as possible to a transmission-line winding whose characteristic impedance (Z_0) is equal to $(L/C)^{1/2}$. In this expression, L is the inductance of the transmission line when the far ends are shorted, and C is the effective capacitance of the transmission line with the far ends open. Windings are typically wound tightly coupled as parallel or twisted groups of

* Reference 19.

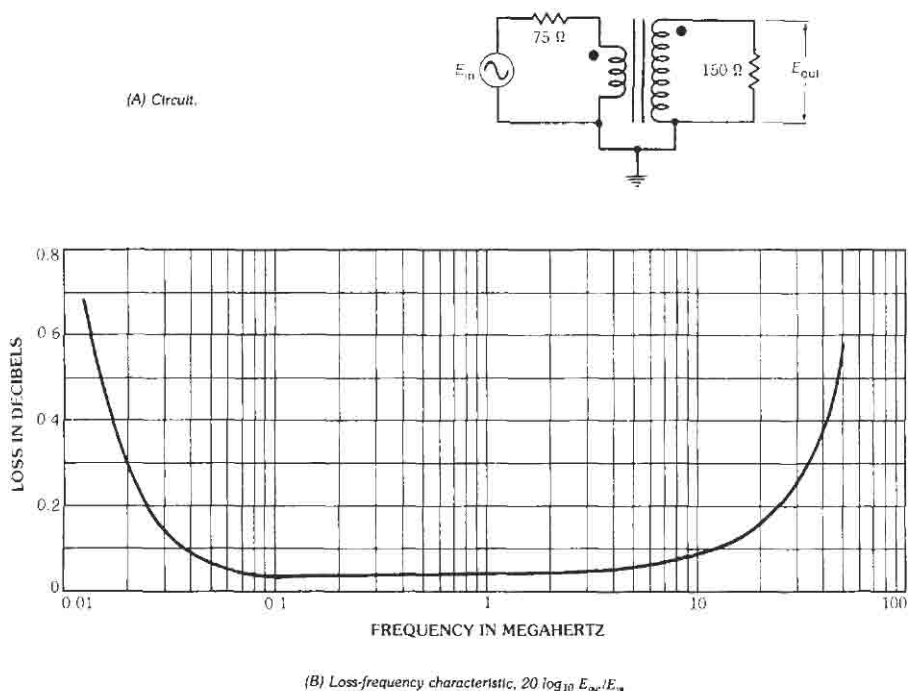


Fig. 11. Response of transformer in example.

two, three, or four wires to obtain a particular ratio. Only a few integer turns ratios are easily obtainable; ratios of 1:1, 1:2, and 1:3 are the most common. Tightly coupled windings produce low leakage inductance. By controlling wire size, wire insulation, and degree of twisting, various uniform characteristic impedances can be obtained. Values of Z_0 between 50 and 100 ohms are most easily made.

Autotransformers provide very low leakage inductances and lend themselves well to a distributed transmission-line analysis. They are preferred in applications where direct-current isolation of windings is not a factor.

Use of NiZn cores, which have lower permeability than MnZn cores, will restrict low-end response considerably. However, NiZn cores usually have lower core loss at high frequencies, which is important when only a few turns are required.

PULSE TRANSFORMERS

Pulse transformers are designed to transmit rectangular waves or trains of pulses while maintaining as closely as possible the original shape. Functions performed by pulse transformers are similar to those of broadband analog transformers. Examples include impedance conversion, dc isolation, coupling between balanced and unbalanced circuits, voltage transformation, and phase inversion. Analytical design is usually done in the time domain instead of the frequency domain, and transient considerations become important.

Pulse waveshapes contain a wide range of frequencies. Lower frequencies relate to pulse duration and repetition rate, while higher frequencies determine the shape of pulse edges. Core and winding parameters limit both extremes of the frequency response of a transformer. Fig. 12 portrays a typical transformer output pulse compared with the corresponding input pulse. (In the strictest sense, pulse rise and decay times are measured between the 10- and 90-percent values; width is measured between the 50-percent values.) Pulse transformers can be analyzed by considering the leading edge, top, and trailing edge of the pulse separately. Fig. 13 illustrates simplified equivalent circuits applicable to each time interval.

As shown in Fig. 13A, leading-edge reproduction is controlled by leakage inductance, I_{scp} , winding capacitances C_p and C_s , and external impedances. Analysis for step-up and step-down transformers varies slightly, as shown in Figs. 13B and 13C. Leakage inductance and winding capacitance must be minimized to achieve a sharp rise; however, output voltage may overshoot input voltage (Fig. 12), and oscillation may be encountered where very abrupt rise times are invoked. Rise time T_r may be related to the high-frequency 3-dB loss point f_h of a wideband analog transformer as follows:*

$$f_h = 0.382/T_r$$

* Reference 20.

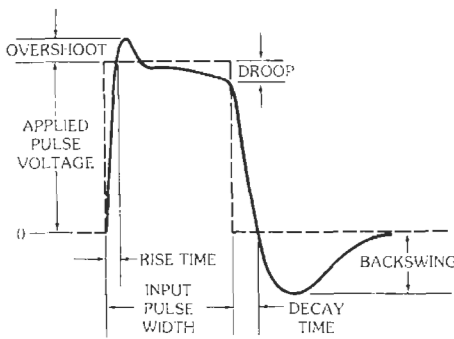


Fig. 12. Output pulse shape.

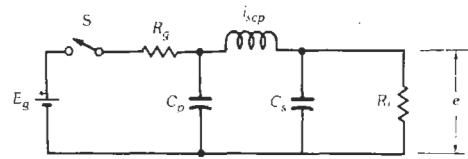
A graphical solution of rise time may be found in additional references.*

Pulse-top response (droop, Fig. 12) depends on the magnitude of open-circuit inductance L_p and external impedances as shown in Fig. 13D. A greater magnitude of L_p results in reduced droop. The circuit is similar to the low-frequency equivalent circuit for wideband analog transformers. However, the value of open-circuit inductance is proportional to the average slope of the part of the $B-H$ loop traversed and may differ when measured under continuous-wave and pulse conditions. A CW input signal generally is sinusoidal, and the core $B-H$ loop traversed is symmetrical about the origin. As shown in Fig. 14A, a unipolar pulse traverses a smaller, asymmetrical portion of the $B-H$ loop, resulting in a lower average slope. Fig. 14B shows the $B-H$ loop generated with bipolar pulse excitation (see reference 20).

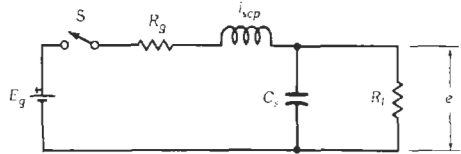
Trailing-edge response is controlled by the dissipation of energy stored in the transformer after the supporting pulse voltage has been removed. It depends on the open-circuit inductance, secondary winding capacitance, and external impedances as shown in Fig. 13E. A lower capacitance results in a faster rate of voltage decay. Negative backswing is proportional to the magnitude of the transformer magnetizing current. A graphical analysis of trailing-edge response can be found in reference 10.

Choice of a pulse-transformer core is usually determined by the Et product, which refers to the pulse voltage E and the duration t in microseconds. The core should not saturate, so Et is proportional to the flux swing traversed on the $B-H$ loop just below saturation and to core cross-sectional area.†

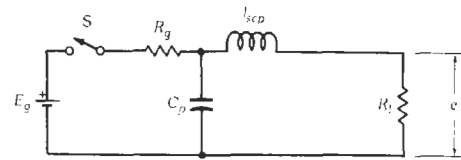
Ferrite materials, specified in terms of pulse excitation, are typically used for pulse-transformer cores. Gaps may be used to provide greater flux swing for unipolar pulse applications. Core configurations may be toroidal where self-shielding and low leakage induc-



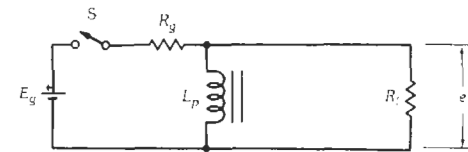
(A) Leading-edge equivalent circuit.



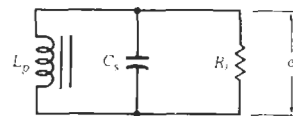
(B) Leading-edge equivalent circuit for step-up-ratio transformer.



(C) Leading-edge equivalent circuit for step-down-ratio transformer.



(D) Top-of-pulse equivalent circuit.



(E) Trailing-edge equivalent circuit.

Fig. 13. Pulse-transformer equivalent circuits.

tance are paramount, "E" shape for typically lower-cost bobbin winding techniques, and cup-core or RM type for greater magnetic efficiency with bobbin windings.

MAGNETIC-CORE REACTORS

The purpose of a reactor is to introduce reactance into a circuit. Inasmuch as this reactance is inductive, this device is also called an inductor. Reactors consist of one or more windings and may or may not have a

* Reference 10.

† Reference 10.

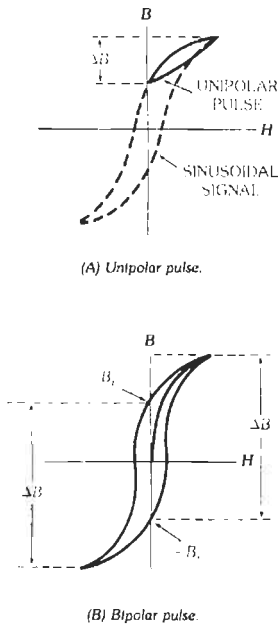


Fig. 14. Traverses of B - H loop.

magnetic core. This section is concerned with those reactors that have magnetic cores. A filter reactor is defined* as “. . . a reactor used to reduce harmonic voltage in alternating-current or direct-current circuits.”

Rectifier-Filter Reactors

The rectifier-filter reactor is used mainly in direct-current power supplies to smooth the output ripple voltage. It carries all the direct current of the rectifier filter and must be designed not to saturate with direct current in the reactor winding.

Optimum design data may be obtained from Hanna curves, Fig. 15. These curves relate direct-current energy stored in the core per unit volume, LI_{dc}^2/V , to magnetizing field NI_{dc}/l_c (where l_c = average length of flux path in core), for an appropriate air gap. Heating is seldom a factor, but direct-current-resistance requirements affect the design; however, the transformer equivalent volt-ampere ratings of chokes given in Table 9 should be useful in determining their sizes. This is based on the empirical relationship $(VA)_{eq} = 188LI_{dc}^2$.

As an example, take the design of a choke that is to have an inductance of 10 henrys with a superimposed direct current of 0.225 ampere and a direct-current resistance ≤ 125 ohms. This reactor is to be used for suppressing harmonics of 60 hertz, where the alternating-current ripple voltage (second harmonic) is about 35 volts.

(A) $LI^2 = 0.51$. Based on the data of Table 9, try a 4% silicon-steel core, type EI-125 lamination, with a core buildup of 1.5 inches.

(B) From Table 9: $V = (11.4/1.25) \times 1.5 = 13.7$ in³, $l_c = 7.50$, $LI^2/V = 0.51/13.7 = 0.037$. From Fig. 15: $NI/l_c = 88$ ampere-turns per inch, $N = 88l_c/I = (88 \times 7.5)/0.225 = 2930$ turns, $l_g/l_c = 0.0032$, the length of air gap $l_g = 0.0032 \times 7.5 = 0.024$ inch.

(C) From Table 9, coil (MLT) = $(7.21 + 0.5)/12 = 0.643$ foot, and length of coil = $N \times (MLT) = 2930 \times 0.643 = 1884$ feet. Since the maximum resistance is 125 ohms, the maximum ohms/ft = $125/1884 = 0.0663$, or 66.3 ohms/1000 ft. From Table 5, the nearest size of wire is No. 28.

(D) Now see if 2930 turns of No. 28 single-insulated wire will fit in the window space of the core. (Determine turns per layer, number of layers, and coil-build, as explained in the design of power transformers.)

(E) This is an actual coil design; in case the lamination window space is too small (or too large), change stack of laminations, or size of laminations, so that the coil meets the electrical requirements and the total coil-build ≈ 0.85 to $0.90 \times$ (window width).

Note: To allow for manufacturing variations in permeability of cores and resistance of wires, use at least 10-percent tolerance.

Swinging reactors are used where the direct current in the rectifier circuit varies. These reactors are designed to saturate under full load current while providing adequate inductance for filtering. At light load current, higher inductance is available to perform proper filtering and prevent “capacitor effect.” The equivalent reactor size is determined from:

$$LI^2 = (L_{max} \times L_{min})^{1/2} I_{dc(max)}^2$$

The design is similar to that of a normal reactor and is based on meeting both L and I_{dc} extremes. The typical swing in inductance is 4:1 for a current swing of 10:1.

AC-Filter Reactors

In the rectifier-filter reactor, the amount of ac flux in the magnetic core is small in comparison to the dc flux. With the ac-filter reactor, there is no dc flux, and the ac flux is large. In the rectifier-filter reactor, the direct current through the coil causes the core to operate in a unipolar, or CM, mode, which necessitates the use of an air gap as calculated from the Hanna curves. In the ac-filter reactor, the core operates in a bipolar, or SCM, mode, but the peak values of the alternating current are large, which also necessitates the use of an air gap in the core to prevent it from saturating.

A method for designing ac-filter reactors may be found in Chapter VIII of reference 18.

MAGNETIC-CORE INDUCTORS

Magnetic-core inductors used in frequency-selective filters must provide a specific inductance value, a high

* Reference 1.

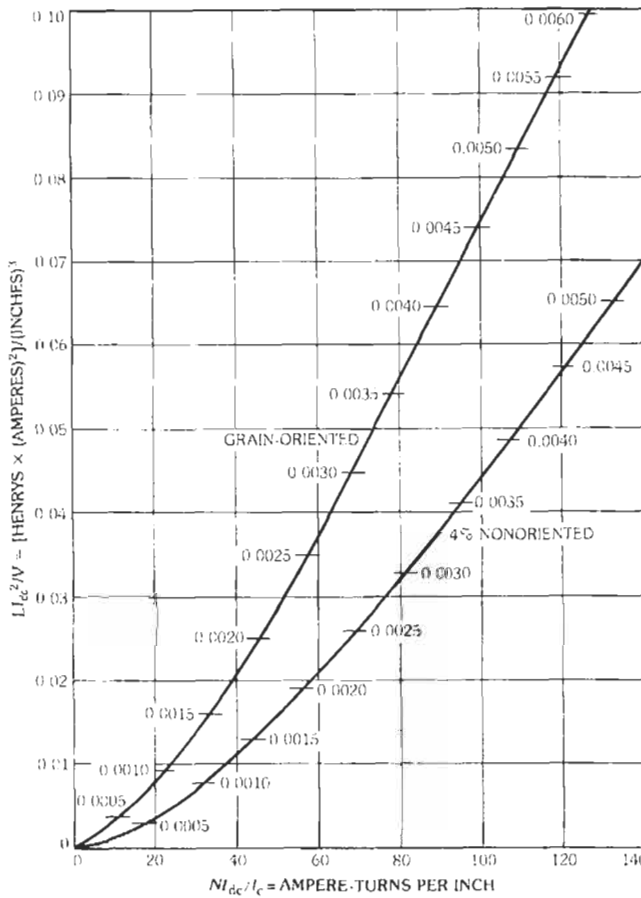


Fig. 15. Hanna curves for silicon steel. The numbers on the curves represent length of air gap l_g /length of flux path l_c .

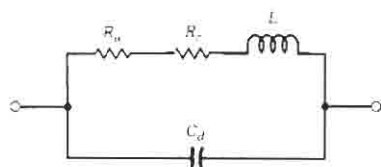
TABLE 9. EQUIVALENT LI^2 RATING OF FILTER REACTOR FOR RECTIFIERS*

LI^2 †	Current Density† (A/in ²)	EI-Type Lamination‡	Stack Height ρ (in)	Core Volume V (in ³)	Magnetic Path Length l_c (in)	Average Copper Mean Length per Turn (MLT)(in)
0.0195	3200	EI-21	0.5	0.80	3.25	3.12
0.0288	2700	EI-625	0.625	1.45	3.75	3.62
0.067	2560	EI-75	0.75	2.51	4.50	4.33
0.088	2560	EI-75	1.00	3.35	4.50	4.83
0.111	2330	EI-11	0.875	3.88	5.25	5.04
0.200	2130	EI-12	1.00	5.74	6.00	5.71
0.300	2030	EI-12	1.50	8.61	6.00	6.71
0.480	1800	EI-125	1.25	11.4	7.50	7.21
0.675	1770	EI-125	1.75	16.0	7.50	8.21
0.850	1600	EI-13	1.50	19.8	9.00	8.63
1.37	1500	EI-13	2.00	26.4	9.00	9.63
3.70	1200	EI-19	1.75	39.4	13.0	12.8

* L = inductance in henrys, I = direct current in amperes. The rating is based on power-supply frequencies up to 400 hertz and 50 °C temperature rise above ambient. The LI^2 values should be reduced for lower temperature rises and high-voltage operation.

† From *Radio Components Handbook* (Cheltenham, PA.: Technical Advertising Associates, 1948; page 92).

‡ Lamination designation and constants per Allegheny Ludlum Corp., Pittsburgh, Pa.



L = calculated inductance based on turns and magnetic-core constant
 C_d = distributed capacitance of winding
 R_w = copper losses in winding
 R_c = magnetic-core losses reflected in series with winding
 Q = quality factor $\omega L / (R_w + R_c)$, $\omega = 2\pi f$

Fig. 16. Equivalent network of an audio filter coil.

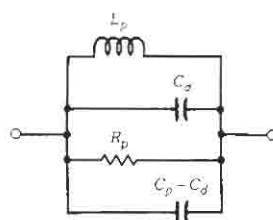
Q value (ratio of reactance to effective resistance), and stable parameter values in the operating environment. The inductance value depends on the core inductance coefficient, $A_L = L/N^2$, and the number of turns in the winding. The A_L values for specific cores are a function of the effective permeability, μ_e , and the core geometry; core manufacturers tabulate these data in their catalogs. Achieving high Q is more complex. The inductor is modeled with a set of independent loss tangent ($1/Q$) values that together describe all the loss mechanisms for the inductor. An attempt is made to minimize the major contributors at a given frequency. At the same time, the designer must also attempt to minimize distributed capacitance, μ_e , and mechanical stresses on the cores in order to promote stable performance. The completed structure must be mechanically sound to minimize the effects of shock and vibration on the parameter values.

Audio-Frequency Inductors

Audio-frequency inductors operate in the frequency band of about 200 Hz to 20 kHz. At the lowest frequencies, laminated cores are generally used, although recent designs have used ferrite. Where low-cost, nonadjustable inductors are needed in the range 300 Hz to 10 kHz, permalloy powder toroids are often used. For frequencies above 10 kHz, or for applications requiring adjustable inductors, ferrite pot cores are preferred (see Table 10). Pot-core inductors are discussed in the section on precision adjustable inductors.

To design stable, high- Q inductors, it is necessary to understand the equivalent circuit of an inductor (Fig. 16) and its application to the parallel-tuned mesh (Fig. 17) and series-tuned mesh (Fig. 18). The distributed capacitance, C_d , is very important in series-tuned meshes, because it affects both the inductance and the Q . If possible, it is desirable to keep the distributed capacitance to one tenth of the tuning capacitance, C_s . Means of minimizing C_d are discussed in succeeding paragraphs.

The power loss in the coil winding, R_w , is made up of the dc resistance of the wire and the eddy-current losses generated in the winding by stray flux from the core that cuts the winding. This eddy-current loss can be minimized at high frequencies by dividing the wire



L_p = equivalent parallel inductance
 R_p = equivalent parallel resistance
 C_s = distributed capacitance of winding
 C_p = parallel tuning capacitance
 Q_p = effective parallel Q
 $L_p \sim L$, for $Q > 10$
 $R_p \sim (R_w + R_c)(Q^2 + 1)$, for $Q_p \sim Q$

Fig. 17. Parallel-tuned mesh.

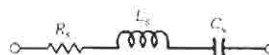
into many small strands. This type of wire is called litz wire.

The loss due to the magnetic core reflected in series with the coil winding is represented by R_c . The calculation of this factor is discussed in succeeding paragraphs.

Laminated-Core Audio-Filter Inductors—Externally, a laminated-core audio-filter inductor resembles a transformer. The laminations are usually ‘‘F’’ in shape to create an air gap in the center leg of the core. For high Q , nickel-steel laminations either 14 or 6 mils thick are used (Table 10). The air gap in the core reduces the effective permeability, μ_e , which reduces the effective core losses and stabilizes the inductance.

The winding is usually layer wound as described in the section on methods of winding transformers. The techniques for reducing capacitance described in that section can be used to reduce the distributed capacitance of the winding. A detailed method for designing this type of inductor can be found in *High Q Reactors for Low Frequencies*, Bulletin A10, Magnetic Metals Co., Camden, N.J.

Toroidal Core Types—Toroidal coils are doughnut shaped, with the winding covering the entire core. The core is usually made of pressed molybdenum-permalloy powder, although some cores are made of carbonyl powdered iron for very-high-frequency applications (Table 10).



L_s = equivalent series inductance
 R_s = equivalent series resistance
 C_s = series tuning capacitance
 Q_s = effective series Q

For $Q > 10$
 $L_s = L(1 - \omega^2 L C_d)$
 $R_s = (R_w + R_c) / (1 - \omega^2 L C_d)^2$
 $Q_s = Q(1 - \omega^2 L C_d)$
 $L = L_s / (1 + \omega^2 L_s C_d)$

Fig. 18. Series-tuned mesh.

TABLE 10. CHARACTERISTICS OF SOME CORE MATERIALS FOR AUDIO-FILTER COILS

Material or Alloy	Initial Permeability (μ_0)	Resistivity (ohm-cm)	Hysteresis Coefficient ($a \times 10^6$)	Residual Coefficient ($c \times 10^6$)	Eddy-Current Coefficient ($e \times 10^9$)	Gauge (mils)	Application and Frequency Range (kilohertz)
4% silicon steel	400	60×10^{-6}	120	75	870	14	Rectifier filters
Low nickel	3500 to 10 000	44×10^{-6}	0.4	14	1550	14	Audio filters up to 0.2
	284				6	Audio filters up to 10	
High nickel	10 000 to 20 000	57×10^{-6}	0.05	0.05	950	14	Audio filters up to 0.2
	175				6	Audio filters up to 10	
Molybdenum permalloy powder	550†	1.0	1.5	88	27	—	Audio filters 0.1–6
	200†	1.0	0.7	21	25	—	Audio filters 0.1–7
	160†	1.0	0.9	25	17	—	Audio filters 0.1–10
	125†	1.0	0.9	32	15	—	Audio filters 0.2–20
	60†	1.0	1.5	50	7.5	—	Audio filters 5–50
	25†	1.0	4.0	96	7.0	—	Audio filters 15–60
	26†	1.0	4.0	96	7.0	—	Audio filters 15–60
	14†	1.0	7.0	143	6.5	—	Audio filters 40–150
Carbonyl types:							
C	55	—	9	80	7	—	High-frequency filters
P	26	—	3.4	220	27	—	High-frequency filters
Th	16	—	2.5	80	8	—	High-frequency filters
Ferrites‡:							
3B7	2300	100	§	§	§	—	Audio filters 0.2–300
3B9	1800	100	§	§	§	—	Audio filters 0.2–300
3D3	750	10^5	§	§	§	—	HF filters 200–2 500
4C4	125	10^5	§	§	§	—	HF filters 1000–20 000

$R_c/(\mu_0 L t) = aB_m + c + ef$, where R_c = series resistance in ohms due to core loss.*

* Data and coefficients a , c , and e are from V. E. Legg and F. J. Given, "Compressed Powdered Molybdenum Permalloy for High Quality Inductance Coils," *Bell System Technical Journal*, Vol. 19, No. 3, July 1940; pp. 385–406.

† Data from Catalog PC303T, Magnetics, Inc., Butler, Pa.

‡ Data from Bulletin 220-C, Ferroxcube Corporation of America, Saugerties, N.Y.

§ See Fig. 19.

Permalloy-powder toroidal cores are made in various sizes from 0.140 to 2.25 inches OD and with ten material compositions to change the effective permeability, μ_e , from 14 to 550 for use at different frequency ranges.

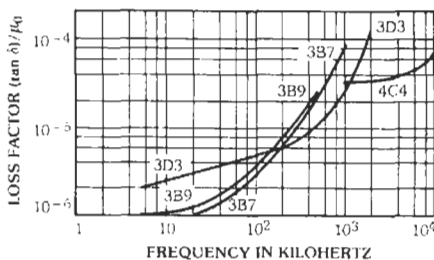


Fig. 19. Loss factor as a function of frequency. (From Ferrite Pot Cores, Bulletin 220-C, Ferroxcube Corp. of America, Saugerties, N.Y.)

Fig. 20 gives the Q -vs-frequency characteristics of four sizes of toroids with five different effective permeabilities. In addition, the inductance in millihenrys per 1000 turns, A_L , is given for each size of core.

Fig. 21 gives the number of turns of various sizes of enameled wire that can be placed on the toroid as well as the mean length of turn (MLT) for the respective toroidal cores of Fig. 20. These curves are based on a single winding wound with 180° traverse using commercially available toroidal winding machines.

Fig. 22 gives the core-loss factor in ohms per millihenry for various effective permeabilities and frequencies. This factor is multiplied by the inductance in millihenrys to determine R_c . The curves are based on a magnetic induction of 100 gauss. For other levels, R_c can be calculated from the Legg coefficients given in Table 10.

Table 11 gives an approximate value of distributed capacitance C_d that can be expected for the different sizes of cores and winding methods.

When maximum temperature stability of inductance is required, most manufacturers of permalloy toroidal cores can provide various types of stabilized cores. The stabilizations most used are listed in Table 12.

Since the distributed capacitance changes rapidly with temperature, it should be kept to a minimum to avoid changing the inductance.

Fig. 23 illustrates the most common methods of winding toroids. In Fig. 23A the toroid is rotated over a 360° arc for every layer of the winding (called 360° traverse winding). In Fig. 23B, the toroid is rotated over a 180° arc for each layer until half the coil is wound. The other half of the winding is similarly wound. This is called 180° traverse winding. In Fig. 23C the toroid is rotated over only 90° until one quarter is wound. The other three quarters are wound in the same manner. This winding method is called 90° traverse or quadrature winding.

Since most toroid cores are made with a tolerance of ±8% on the A_L (millihenrys/1000 turns), it is usually necessary to adjust the inductance after winding. This is done by winding about 5% more turns than calculated on the core and removing turns until the inductance reaches the desired value. Coils for series-tuned meshes should be adjusted to the resonance frequency of the mesh with the tuning capacitor in series to eliminate the effect of distributed capacitance. Coils for parallel tuned meshes should be adjusted at low frequency in such a way that the tuning capacitance is 1000 times the distributed capacitance for 0.1% accuracy.

Toroidal Core Design Example—It is desired to design an inductor of 100 millihenrys for a series-tuned mesh that resonates at 10 kilohertz. The Q must be 150 minimum and the size as small as possible.

(A) Consulting the Q curves of Fig. 20B shows that a 1.06-inch OD toroid core with $\mu_r = 125$ is the smallest core that will meet the Q requirements. This has an A_L value of 157 millihenrys per 1000 turns.

(B) From Table 11, $C_d \approx 200$ picofarads. Calculate L from Fig. 18:

$$L = \frac{L_s}{1 + \omega^2 L_s C_d}$$

$$= \frac{100}{1 + (2\pi 10\,000)^2 \times 0.1 \times 200 \times 10^{-12}}$$

$$= 92.68 \text{ millihenrys}$$

(C) Compute the number of turns required from

$$N = 1000(L_{\text{millihenrys}}/A_L)^{1/2}$$

$$= 1000(92.68/157)^{1/2}$$

$$= 768 \text{ turns}$$

(D) Fig. 21 gives the maximum size of wire and mean length of turn. Use No. 30 heavy enameled wire for 768 turns. On a 1.06-inch OD core the (MLT) is 1.57 inches.

(E) Calculate R_w from (MLT), N , and Table 5.

$$R_w = \frac{(\text{MLT}) \times N \times \text{ohms}/1000 \text{ ft}}{12\,000}$$

$$= \frac{1.57 \times 768 \times 103.2}{12\,000}$$

$$= 10.36 \text{ ohms}$$

(F) Calculate R_c from R/L values of Fig. 22 for $\mu_r = 125$ and $f = 10$ kilohertz. $R/L = 0.23$ ohm/millihenry. $R_c = (R/L)L = 0.23 \times 92.68 = 21.32$ ohms.

(G) Calculate Q per Fig. 16.

$$Q = \frac{\omega L}{R_w + R_c} = \frac{2\pi 10 \times 92.68}{10.36 + 21.32} = 183.8$$

(H) Calculate Q_s per Fig. 18.

$$Q_s = Q(1 - \omega^2 L C_d)$$

$$= 183.8[1 - (2\pi 10\,000)^2 \times 0.09268 \times 200 \times 10^{-12}]$$

$$= 170.4$$

This should approximate the measured Q .

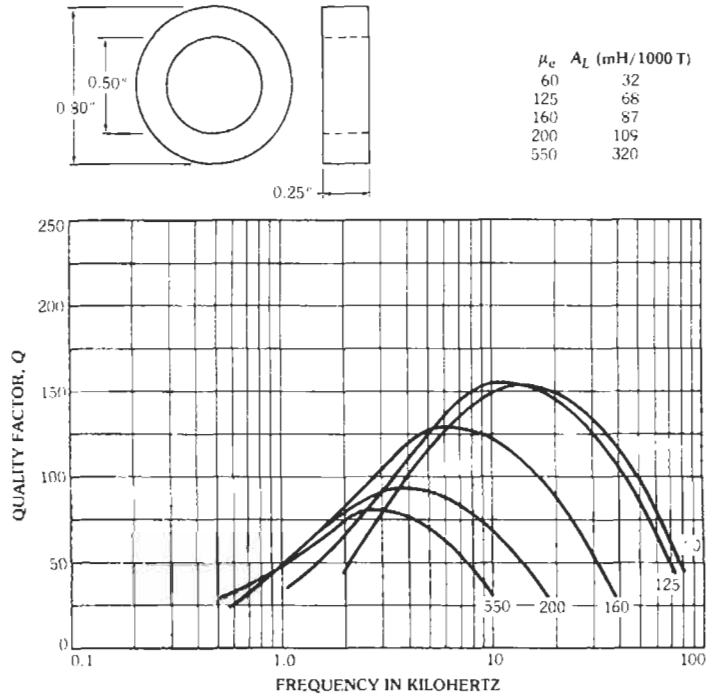
(I) The coil is wound using 360° traverse winding, and, since no special temperature stability is required, a standard unstabilized core is used.

Precision Adjustable Inductors

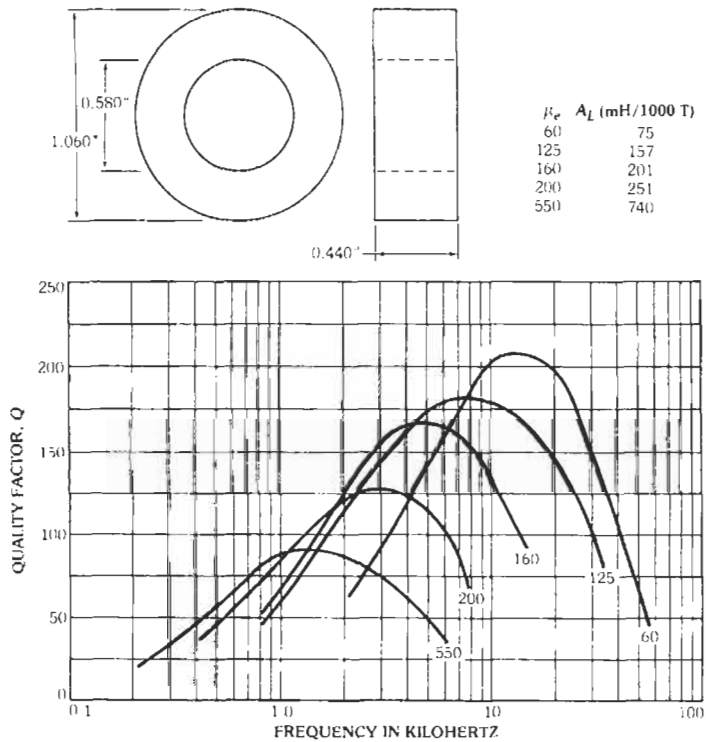
Pot Cores—Filter-class inductors are precision, adjustable devices designed to provide high Q values at a particular frequency. Although Q is often the greatest concern, the temperature and time stability are usually also important. While round pot-core structures are very popular, other core shapes such as RM, X, and Q are used to improve packing densities or magnetic performance.

A cross-section of a typical pot-core inductor is shown in Fig. 24. The two core halves are clamped or cemented around a bobbin that contains the winding. The ceramic ferrite cores are chosen specifically for particular performance characteristics. A wide range of permeabilities, loss characteristics, and stability performance exists. The information is available in ferrite-supplier catalogs.

The cores have a recessed center post so that when they are mated, an air gap exists. Although the gap reduces the effective permeability, μ_r , more importantly it reduces the core loss and temperature coefficient of inductance. The gap is used to control the A_L value (in nanohenrys per turn squared), a most important inductor design parameter. Inductors operating at low frequencies have relatively high values of μ_r and A_L . As the operating frequency is increased, lower values of these parameters usually are adequate.

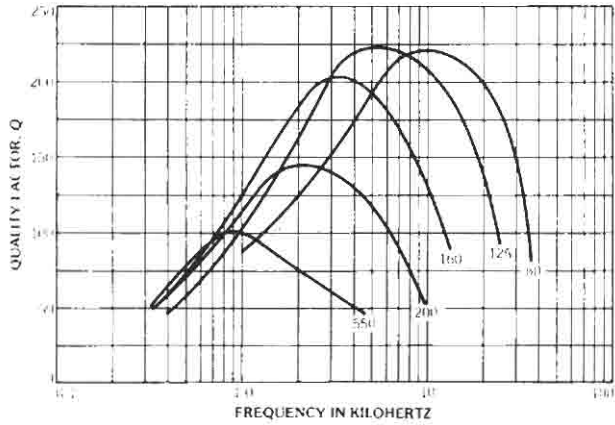
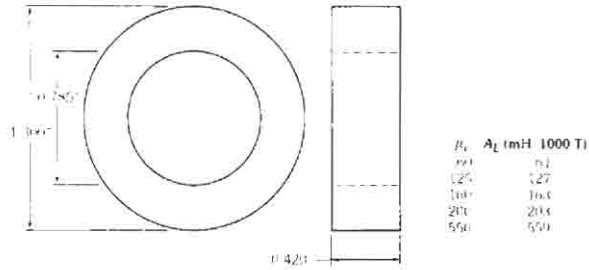


(A) 0.8 inch OD.

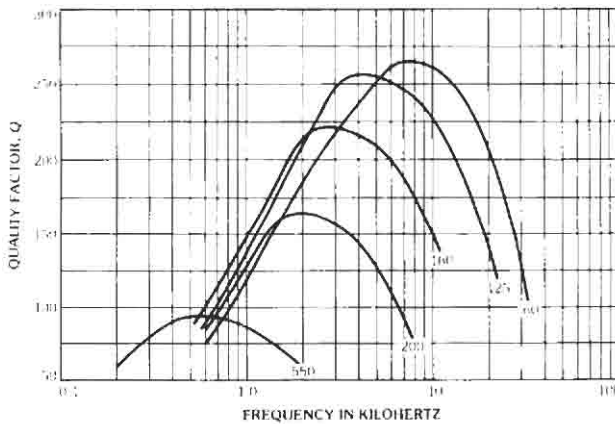
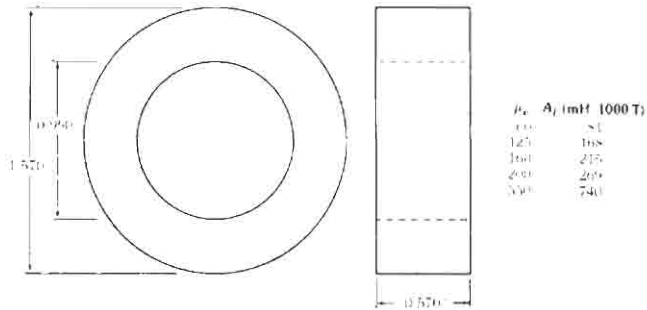


(B) 1.06 inch OD.

Fig. 20. Quality factor Q as a function of frequency for several sizes of permalloy-dust toroids. (Data replotted from Permalloy



(C) 1.3 inch OD.



(D) 1.57 inch OD.

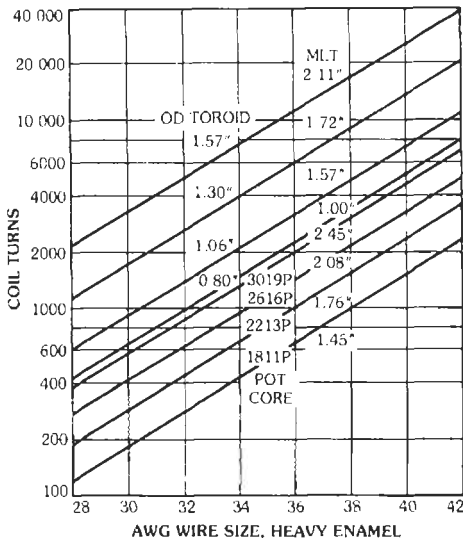


Fig. 21. Coil turns as a function of wire gauge for various core sizes and types.

Loss Mechanisms and Q —Inductor Q , the ratio of reactance to effective resistance at a given frequency, is controlled by four loss factors: core loss, winding resistance loss, capacitance loss, and proximity-effect loss. These are expressed in terms of loss tangents. Inductor Q is the reciprocal of the sum of these four loss

tangents. Many inductor-loss equations are given by Snelling,* but only four equations are needed in practice. Optimization equations are given by Banzi.†

Most manufacturers of inductor cores list the core-material loss, $\tan \delta_m$, which has been normalized with respect to the core-material permeability μ_i . To calculate the effective core-loss tangent, it is necessary only to multiply the advertised value by the effective permeability, μ_e , listed separately in the core literature for specific geometries and inductance factors (A_L values). The effective magnetic core-loss equation is

$$\tan \delta_{m_e} = (\tan \delta_m)(\mu_e/\mu_i) = \mu_e/\mu_i Q_m$$

where Q_m is the core-material Q . While it is apparent that a low ratio of μ_e/μ_i (implying a large air gap) would reduce the effective core loss, a large gap results in higher winding losses because of flux fringing at the gap.

Although coils generally have associated capacitances that are distributed in nature, a single lumped capacitance shunting the winding may be considered a sufficiently accurate model for most Q calculations. The capacitive loss is the product of a ratio and a constant, as shown in the equation below. The ratio is the parallel winding capacitance, C_p , divided by the capacitance,

* Reference 19.

† Reference 21.

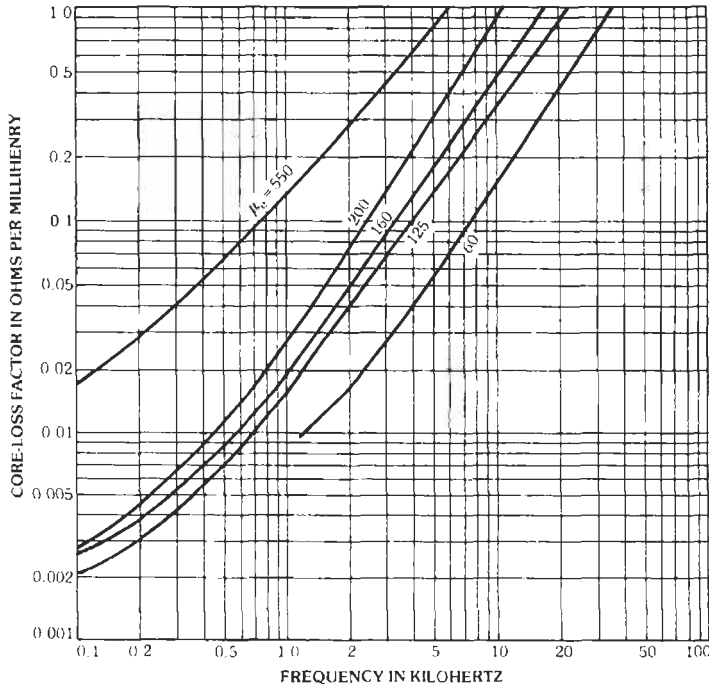


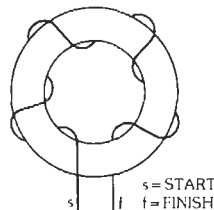
Fig. 22. Molybdenum permalloy-dust-core loss characteristics. $B_m = 100$ gauss. (Data replotted from Permalloy Powder Cores, Catalog MPP-303T, Magnetics, Inc., Butler, Pa.)

TABLE 11. ESTIMATION OF DISTRIBUTED CAPACITANCE FOR CORE SIZE AND WINDING METHOD

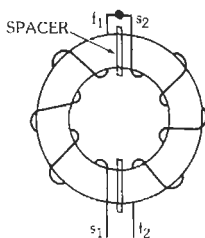
Toroid Core Size	Distributed Capacitance C_d (picofarads)		
	360° Traverse	180° Traverse	90° Traverse
0.8" OD	120	60	30
1.06" OD	200	100	50
1.30" OD	300	150	75
1.57" OD	360	180	90

C_r , necessary to resonate the inductor at the frequency of interest. The value of C_p can be measured, or it can be calculated based on the winding geometry.† Constant K_D depends on the dielectric constant of the wire insulating material. The capacitive loss-tangent formula is

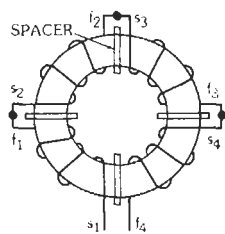
* Reference 8.



(A) 360° traverse.



(B) 180° traverse.



(C) 90° traverse.

Fig. 23. Methods of winding toroid cores.

TABLE 12. STABILIZATION OF TOROIDAL CORES

Identification	Inductance Temperature Stability	Temperature Range (°C)
	(%)	
B	±0.1	+13 to +35
D	±0.1	0 to +55
W	±0.25	-55 to +85

$$\tan \delta_c = (C_p/C_r)K_D$$

Experimental work has shown the value of K_D to be 0.02 for nylon-served litz wire (groups of individually insulated wires twisted into a bundle that is then wrapped in nylon yarn) and 0.01 for polyurethane-coated solid wire. It is not necessarily advantageous to use solid wire to minimize capacitive losses. For multiple-layer windings, litz wire usually results in lower capacitance than does solid wire because litz-wire insulation provides more separation between layers. The high loss associated with the nylon serving, then, tends to be offset by the lower value of C_p . For single-layer windings, though, there may be very little difference in capacitance between the two wire types, depending on the number of turns used, so that here solid wire presents an advantage in minimizing capacitive loss.

The winding-resistance loss tangent, $\tan \delta_r$, is the ratio of the ac resistance of the winding, R_{ac} , to the coil reactance, ωL , at the frequency of interest:

$$\tan \delta_r = R_{ac}/\omega L$$

Because the dc resistance can be determined, it would be convenient to calculate the ac resistance in terms of the dc resistance. The skin-effect factor, SEF, is the ratio of ac to dc resistance. The SEF versus d/Δ (wire diameter/skin depth) characteristic is used in conjunction with the skin depth versus frequency characteristic. These characteristics are plotted by Snelling.* The ac resistance is thus the dc resistance times SEF. This can be rewritten in more useful terms:

$$\tan \delta_r = [K_R (N/n)/\omega L](SEF)$$

† Reference 1.

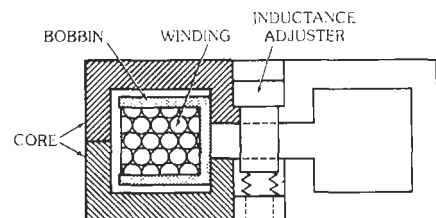


Fig. 24. Typical pot-core coil.

where,

- K_R = dc resistance per strand turn of wire,
- N = number of turns of wire,
- n = number of wire strands,
- SEF = skin-effect factor.

The proximity-effect loss-tangent formula is

$$\tan \delta_{pe} = K_E f N n d^4 (\text{PEF}) / A_L$$

where,

- K_E = a constant that is geometry dependent,
- f = frequency of operation in hertz,
- N = number of winding turns,
- n = number of strands of wire,
- d = bare diameter of an individual strand of wire in millimeters,
- PEF = proximity-effect factor, which is frequency dependent,
- A_L = inductance factor in nH/N^2

The value of K_E depends on the winding configuration and core type and must be determined experimentally; it varies slowly with frequency. Snelling gives typical values of K_E as well as means for obtaining PEF.*

Combining these effects:

$$Q = 1 / (\tan \delta_{mc} + \tan \delta_c + \tan \delta_r + \tan \delta_{pe})$$

For details see reference 21.

Inductor Q values as high as 1000 are obtainable with proper design and at frequencies in the range around 100 kHz. Modern designs use cores with a high A_L value and low mean length of turn to achieve high Q ; unfortunately, this choice leads to higher effective permeability, which adversely affects stability of inductance.

Stability Considerations—Temperature stability of inductance may be critical for a precision adjustable inductor. The inductor is part of a tuned LC circuit in which the temperature coefficient (TC) of the inductor may be compensated by that of the capacitor to give a net TC near zero. The TC of an inductor is roughly proportional to the effective permeability of its core. Core temperature factors as low and tight as 0.7 ± 0.3 ppm/°C near room temperature are available. With an effective permeability of 100 for its core, an inductor would have a TC of 70 ± 30 ppm/°C. Ferrite cores exhibit magnetostrictive effects, so care must be taken to avoid stress on the core halves when bonding or clamping them together. A properly designed adjuster assembly is also critical for good temperature stability.

A slow decrease of inductance value with time, calling aging, is characteristic of ferrite-cored inductors. Described in terms of a disaccommodation factor, this aging usually manifests itself as a linear decrease in inductance per decade of time, on a semilog plot. Manufacturers usually measure disaccommodation factor between 10 and 100 minutes, but an extrapolation to

years of life is risky, since the longer-term slope may increase above that measured over a few hours. Inductor aging of less than 0.5% over a 20-year life is achievable.

Inductance variation is usually obtained by the use of a threaded ferrite adjustment mechanism in the center hole of the core. As the ferrite slug bridges the air gap in the center leg of the core (Fig. 24), μ_e is increased, and hence inductance increases. In many applications, an adjustment range of greater than $\pm 8\%$, with a sensitivity of 0.1%, is possible. Most ferrite inductor cores have preadjusted air gaps that result in an A_L tolerance of $\pm 3\%$; thus guaranteed inductor adjustment ranges of about $\pm 2\%$ are practical, even with a small number of turns.

REFERENCES

1. *IEEE Standard Dictionary of Electrical and Electronics Terms*, IEEE Std., 100-1972.
2. Graf, Rudolf F. *Modern Dictionary of Electronics*. Indianapolis: Howard W. Sams & Co., Inc., 1984.
3. *IEEE Standard for High-Power Wide-Band Transformers (100 Watts and Above)*, IEEE Std. 264-1977.
4. Workman, T. J., of Thomas & Skinner, Inc. "Transformer Lamination Study." *Proceedings of Coil Winding/Electrical Manufacturing Expo '81*.
5. *Catalogue 15M-10/81*. Parsippany, N.J.: Allied Corp., 1981.
6. *Mini-Log IV*. Indianapolis: Thomas & Skinner, Inc.
7. Armstrong, R. W. "Polyphase Rectification Special Connections." *Proceedings of the Institute of Radio Engineers*, Vol. 19, No. 1, January 1931.
8. Grossner, N. R. *Transformers for Electronic Circuits*. New York: McGraw-Hill Book Co., 1967. (Second edition published in 1983)
9. Schade, O. H. "Analysis of Rectifier Operation." *Proceedings of the IRE*, July 1943.
10. Lee, R. *Electronic Transformers and Circuits*. 2nd ed. New York: John Wiley & Sons, Inc., 1961.
11. *Standard Definitions of Terms, Symbols and Conversion Factors Relating to Magnetic Testing*. ANSI/ASTM A 340-77.
12. NEMA Standards Publication No. MW-1000-1981. Washington, D.C.: National Electrical Manufacturers Association, October 1981.
13. *Engineering Data for Film Insulations*. Fort Wayne, Ind.: Phelps Dodge Copper Products, January 1, 1961.
14. Hart, H. P., and Kakalec, R. J. "The Derivation and Application of Design Equations for Ferroresonant Voltage Regulators and Regulated Rectifiers." *IEEE Transactions on Magnetics*, Vol. Mag-7, No. 1, March 1971.
15. Kakalec, R. J. "A Feedback-Controlled Ferroresonant Voltage Regulator." *IEEE Transactions on Magnetics*, Vol. Mag-6, No. 1, March 1970.

16. Hart, H. P., and Kakalec, R. J. "A New Feedback Controlled Ferroresonant Regulator Employing a Unique Magnetic Component." *IEEE Transactions on Magnetics*, Vol. Mag-7, No. 3, September 1971.
17. Pressman, A. I. *Switching and Linear Power Supply, Power Converter Design*. Rochelle Park, New Jersey: Hayden Book Co., Inc., 1977.
18. McLyman, W. T. *Transformer and Inductor Design Handbook*. New York: Dekker Inc., 1978.
19. Snelling, E. C. *Soft Ferrites*. The Chemical Rubber Co., 1969.
20. *IEEE Standard for Low-Power Pulse Transformers*, IEEE Std., 390-1975.
21. Banzi, Fred J., Jr., "Higher Q From Pot Core Inductors." *IEEE Trans. Parts, Hybrids, and Packaging*, Vol. PHP-13, No. 4, December 1977.

14 Power Electronics— Rectifiers, Filters, and Power Supplies

Fred G. Turnbull
Revised by Ondrej Pauk

Characteristics of Power Semiconductor Devices 14-3

- Rectifiers
- Schottky Rectifiers
- Zener Diodes
- Power Transistors
- Insulated-Gate Bipolar Transistors
- Field-Effect Transistors
- Thyristors (Silicon Controlled Rectifiers)
- MOS-Controlled Thyristors
- Triacs
- Transient-Voltage Suppressors

AC–DC Converter Circuits 14-12

- Rectifier Circuits
- Single-Phase Voltage Multipliers
- Phase-Controlled Thyristor Circuits
- Linear Transistor Circuits
- Switch-Mode Power Supplies
- Basic Converters
- Forward-Mode Converter Basics
- Flyback or Boost-Mode Converters Basics
- Power Semiconductors in Switch-Mode Power Supplies
- Power MOSFETs
- Bipolar Power Transistors
- IGBTs
- Rectifiers
- Power Losses in Power Semiconductors
- Turn-on Losses

Turn-off Losses
Snubbers and Clamps
Switch-Mode Converter Control
Voltage-Mode Control
Current-Mode Control
Resonant Transistor Circuits
Power Factor Correction Circuits
Power Factor Definition
International Regulatory Standards
Active Power Factor Correction

Filter Circuits and Design 14-39
 Inductor-Input Filter Design
 Capacitor-Input Filter Design
 Harmonic Currents and Power Factor
 Phase-Controlled-Thyristor Filter Design
 Resonant-Filter Design

DC-Output Power Supplies 14-43
 Unregulated Power Supplies
 Linear Regulated Power Supplies
 Phase-Control Power Supplies
 Switching Power Supplies
 Distributed Power Supplies

AC-Output Power Supplies 14-46
 AC Voltage Regulators
 DC-AC Inverters
 Cycloconverters

CHARACTERISTICS OF POWER SEMICONDUCTOR DEVICES

This section briefly describes the terminal properties of various types of power semiconductor devices to aid in understanding their operation in power conversion equipment. Power semiconductor devices discussed are rectifiers, including fast recovery and Schottky; bipolar, field-effect, and insulated-gate transistors; thyristors, MOS-controlled thyristors, and triacs; and voltage transient clippers.

Rectifiers

Voltage and Current Ratings—Silicon-rectifier ratings* are generally expressed in terms of reverse-voltage ratings and of mean-forward-current ratings in a half-wave circuit operating from a 60-hertz sinusoidal supply and into a purely resistive load.

There are three reverse-voltage ratings of importance:

Peak transient reverse voltage	V_{RM}
Maximum repetitive reverse voltage	$V_{RM(rep)}$
Working peak reverse voltage	$V_{RM(wkg)}$

Peak transient reverse voltage (V_{RM}) is the rated maximum value of any nonrecurrent surge voltage, and this value must not be exceeded under any circumstances, even for a microsecond. Maximum repetitive reverse voltage ($V_{RM(rep)}$) is the maximum value of reverse voltage that may be applied recurrently, e.g., in every cycle, and will include any circuit oscillatory voltage that may appear on the sinusoidal supply voltage. Working peak reverse voltage ($V_{RM(wkg)}$) is the crest value of the sinusoidal voltage of the supply at its maximum limit. The manufacturer generally recommends a $V_{RM(wkg)}$ that has an appreciable safety margin in relation to the V_{RM} to allow for the commonly experienced transient overvoltages on power mains.

Three forward-current ratings are similarly of importance:

Nonrecurrent surge current	$I_{FM(surg)}$
Repetitive peak forward current	$I_{FM(rep)}$
Average forward current	$I_{F(av)}$

Silicon diodes have comparatively small thermal mass, and care must be taken to ensure that short-term overload currents are limited. The nonrecurrent surge current is sometimes given as a single value that must not be exceeded at any time, but it is more generally given in the form of a graph of permissible surge current versus time. It is important to observe whether the surge-

current scale is marked in peak, rms, or average value for the data to be correctly interpreted. The repetitive peak forward current is the peak value of the forward current reached in every cycle and excludes random peaks caused by transients. Its relation to the average forward current depends on the circuit used and on the load that is applied. For example, the repetitive peak is about three times the average for a half-wave or bridge circuit working into a resistive load; it may be many times greater when the same circuits work into capacitive loads.

Forward Characteristics—The manufacturer generally supplies curves of instantaneous forward voltage versus instantaneous forward current at one or more operating temperatures; a typical characteristic is shown by the solid curve in Fig. 1. Such curves are not exact for all rectifiers of a given type but are subject to normal production spreads. They are of particular importance in determining the power dissipated by the rectifier under given working conditions.

Calculation of power dissipation from the voltage-current curves need not be done in every instance, since the manufacturer gives curves of power dissipation versus forward current for a limited number of commonly used circuits. However, cases do arise for which the particular form of circuit or load is not covered, and it is then necessary to calculate the dissipation for these particular conditions. The calculation can be greatly simplified, with little loss of accuracy in most cases, by approximating the actual V - I character-

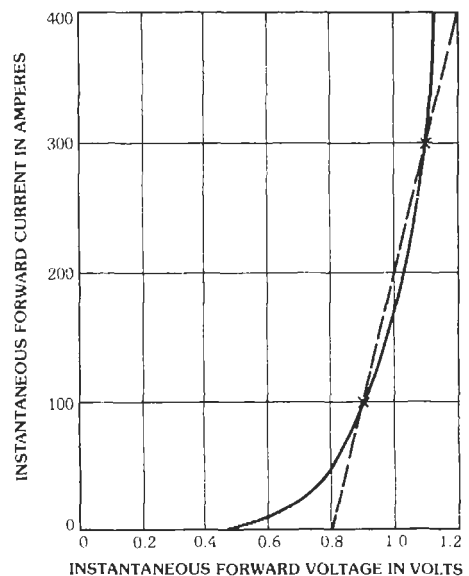


Fig. 1. Instantaneous forward voltage-current characteristic for a typical 100-ampere diode operated at 100°C junction temperature. The ideal threshold forward voltage V_T is the value where the broken line intercepts zero forward current.

* For a complete list of silicon-rectifier ratings, refer to the EIA-JEDEC *Recommendations for Letter Symbols, Abbreviations, Terms and Definitions for Semiconductor Device Data Sheets and Specifications*, published by Electronic Industries Association, 2001 I Street, NW, Washington, DC 20006.

istic curve to a straight line, as shown by the broken line in Fig. 1. The approximate characteristic corresponds to that of a fixed voltage (the threshold voltage) plus a fixed resistance (the slope resistance). For any shape of current waveform, the power dissipated at constant voltage is the product of the average current and this fixed voltage, while the power dissipated at constant resistance is the product of the square of the rms current and this fixed resistance. Thus, the following simple equation can be used:

$$P = I_{F(\text{av})} \times V_T + I_{F(\text{rms})}^2 \times R_S \quad (\text{Eq. 1})$$

where,

P is the forward power dissipation,

$I_{F(\text{av})}$ is the average forward current through the rectifier, averaged over one complete cycle,

$I_{F(\text{rms})}$ is the rms value of the forward current through the rectifier,

V_T is the threshold voltage, and

R_S is the slope resistance.

For the best accuracy, the straight-line approximation should be drawn through points on the current curve corresponding to 50% and 150% of the peak current at which the rectifier is to be used. Thus, in Fig. 1 the broken line would correspond to a peak working current of 200 amperes.

Carrier Storage—On switching from forward conduction to reverse blocking, a silicon diode cannot immediately revert to its blocking state because of the presence of the stored carriers in the junction. These have the effect of allowing current to flow in reverse, as through a forward-biased junction, when reverse voltage is applied. The current is limited only by the external voltage and circuit parameters. However, the carriers are rapidly removed from the junction both by internal recombination and by the sweep-out effect of the reverse current, and when this has happened, the diode reverts to its blocking condition in which only a low leakage current flows. This sudden cessation of a large reverse current can cause objectionable voltage transients if there is appreciable circuit inductance and surge-suppression components have not been included. The reverse current due to carrier storage is not excessive in normal operation of power rectifier circuits and does not in itself constitute a hazard; however, its effect can sometimes lead to complications in switching arrangements. For example, in an inductively loaded circuit, the current will “free-wheel” through the diodes after the supply has been removed until the inductive energy has been discharged. Should the supply be reapplied while this process is going on, some of the diodes will be required to conduct in a forward direction, but others will be required to block; while the latter are recovering from the carrier storage injected by the free-wheeling current, the short-circuit across the supply can cause a damaging surge current to flow.

A technique for reducing this problem is to use “fast-recovery” rectifiers. These rectifiers are tailored to operate in high-frequency circuits with reduced and specified amounts of recovered charge. The devices change from conducting reverse current to leakage current in a “nonabrupt” manner. The finite rate of change in current reduces the transient voltages that appear on the power semiconductors and other circuit elements.

Schottky Rectifiers

The Schottky rectifier has the same terminal characteristics as do conventional and fast-recovery rectifiers; that is, it conducts current in one direction and blocks voltage of one polarity. Rather than rely on a pn junction for rectification, the Schottky diode uses a metal-to-semiconductor contact. Depending on the metal, the forward voltage drop is typically one-half that of a conventional silicon pn rectifier. The reverse voltage rating is lower than that of a pn rectifier, being on the order of 50 to 200 volts. The rectifying action of the Schottky rectifier depends upon the majority carriers, so there is very little reverse current caused by minority-carrier recombination. Therefore, these devices are characterized by very low recovered charge and as such are suitable for high-frequency rectification, up to typically 1 to 5 MHz. The devices are used extensively as rectifiers on the low-voltage windings of high frequency switching mode power supplies. The output voltage of these systems is typically 5 volts dc, and they operate typically at 20 kHz to 1 MHz. The Schottky rectifiers operate in these systems with considerably lower conduction and switching losses than silicon pn rectifiers.

Zener Diodes

“Zener” is the name given to a class of silicon diodes having a sharp turnover characteristic at a particular reverse voltage, as shown in Fig. 2. If such a diode is operated on this part of its characteristic, no breakdown (in the sense of dielectric breakdown) occurs, and the process is reversible without damage. The steepness of the reverse part of the current-voltage characteristic in the turnover region makes these diodes excellent elements for voltage reference and voltage regulation.

The temperature coefficient for a typical range of zener diodes is shown in Fig. 3. It will be seen that the coefficient changes from negative to positive in the region of 5 volts. Use is sometimes made of this phenomenon to match diodes of opposite coefficient to produce a series pair having a low effective temperature coefficient in combination.

Zener diodes are used to provide stable reference voltages for electronic control circuits and as a voltage reference for closed-loop regulating systems in which the actual voltage is compared to the reference voltage in order to develop an error signal. They are available

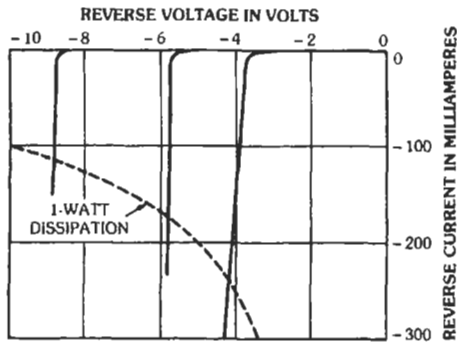


Fig. 2. Typical reverse characteristics for three low-voltage zener diodes of 1-watt rating. (From J. M. Waddell and D. R. Coleman, "Zener Diodes—Their Properties and Applications," *Wireless World*, Vol. 66, No. 1, p. 18, Fig. 2; January 1960. © 1959, Iliffe Electrical Publications, Ltd., London, England.)

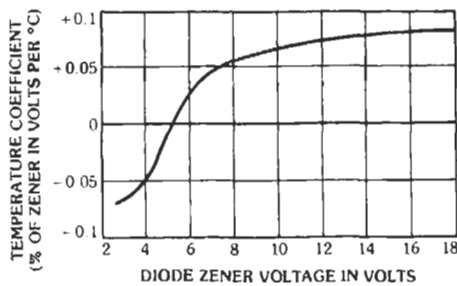


Fig. 3. Temperature coefficient for a typical range of low-voltage zener diodes. (From J. M. Waddell and D. R. Coleman, "Zener Diodes—Their Properties and Applications," *Wireless World*, Vol. 66, No. 1, p. 18, Fig. 4; January 1960. © 1959, Iliffe Electrical Publications, Ltd., London, England.)

in voltage ratings from 3 to 200 volts and power ratings from less than a watt to 50 watts. The voltage tolerance can be specified as $\pm 20\%$, $\pm 10\%$, $\pm 5\%$, and, with further selection, to closer tolerances. Special units designed for transient voltage clipping are discussed in the section on transient-voltage suppressors.

Power Transistors

Power transistors are three-terminal semiconductor devices that are widely used in power supplies and other power conversion equipment. The addition of a third terminal allows the device to have the capability of electronic control of its impedance to the flow of current. The device impedance can range from a very low forward voltage drop, called *saturation voltage*, with the current limited by the external load to a very low leakage current during its off condition. The transistor can operate in its linear region with a simultaneous high voltage across the device and a high current

flowing through the device. This linear mode is used in the linear regulated power supplies described in a later section. The power transistor can also be operated in the switching mode, where it is either on or off with the shortest possible times to accomplish the turn-on and turn-off. To increase the gain defined as the ratio of collector current to base current, two transistors are arranged in a Darlington connection in which the collector current of one transistor provides the base current to the output transistor. This connection increases the collector-to-emitter saturation voltage. Silicon high-voltage transistors useful for high-frequency switching circuits are characterized by a set of voltage, current-gain, and switching characteristics that are briefly discussed below.

Breakdown Sustaining Voltages—Power transistors during conditions of turn-off are required to withstand immediately a forward collector-to-emitter voltage. Because of the rapid change from forward conduction to forward blocking, the sustaining voltages are less than the steady-state forward blocking voltage ratings.

The collector-to-emitter sustaining-voltage rating is specified as a function of the base-to-emitter voltage conditions:

$V_{CE0(SUS)}$ = collector-emitter sustaining voltage, base open

$V_{CER(SUS)}$ = collector-emitter sustaining voltage, base-emitter resistor

$V_{CES(SUS)}$ = collector-emitter sustaining voltage, base-emitter reverse-bias voltage

The $V_{CE0(SUS)}$ rating is the lowest of the three ratings and is specified at a low level of collector current. It is important in switching inductive loads, and this voltage should not be exceeded during the switching interval.

DC Current Gain and Saturation Voltage—The collector-emitter saturation voltage is the voltage drop that occurs when the transistor is carrying current. It is specified at a given collector current and a given base current or as a gain that is the ratio of collector current to base current. The saturation voltage is also a function of junction temperature. The gain has a peak value at nominal current levels and decreases at both lower and higher values of collector current. The typical value of saturation voltage for single transistors is one volt. This voltage drop times the collector current represents a major part of the conduction losses in the transistor. The collector-emitter saturation-voltage versus collector-current characteristics of a Darlington power transistor are plotted in Fig. 4 for a 500-volt, 50-ampere NPN power transistor.

Forward- and Reverse-Biased Second Breakdown—When transistors are being turned on and off, the possibility for simultaneous high currents and high voltages exists. This results in localized high power dissipation in the transistor, leading to increasing tem-

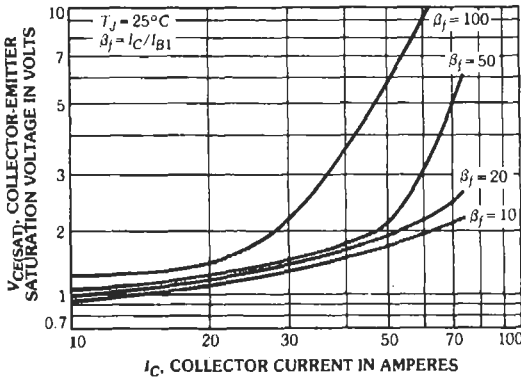


Fig. 4. Plots of V_{CESAT} versus I_C . $T_J = 25^\circ\text{C}$. (From *Transistors-Diodes*. Auburn, NY: General Electric Co., Semiconductor Products Dept., 1982.)

perature rise and potential device failure. This failure is described as "second breakdown." Power-transistor manufacturers present a family of curves outlining a safe operating area for the device. The actual device switching locus should remain within the designated area. Two sets of curves are provided, one with forward bias on the base and the other with reverse bias on the base (turn-off). The forward-biased safe operating area for a 500-volt, 100-ampere Darlington transistor is shown in Fig. 5. The curves are bounded by peak current limits, peak voltage limits, thermal limits, and second breakdown limits. The safe operating area

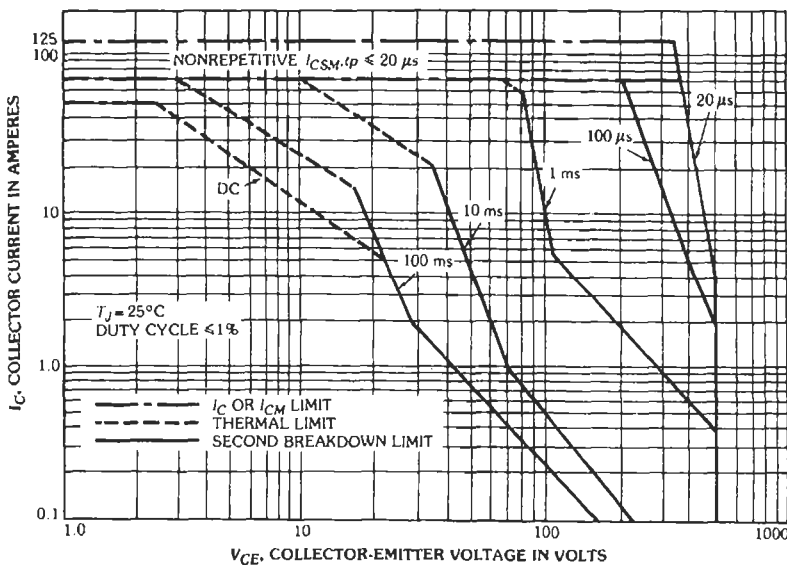


Fig. 5. Forward-bias safe operating area for a 100-ampere power transistor. (From *Transistors-Diodes*. Auburn, NY: General Electric Co., Semiconductor Products Dept., 1982.)

increases as the switching times decrease. The reverse-biased safe operating area for a 500-volt, 100-ampere Darlington transistor is shown in Fig. 6. The curve shown is for turn-off with a voltage clamp on the maximum collector-to-emitter voltage.

Switching Times—The switching times, both on and off, are important during high-frequency switching because both voltage and current are high, resulting in high peak power dissipation. Reduction in these switching times results in less average power dissipation due to switching losses. The turn-on time is composed of two parts, delay time and rise time. With inductive loads, these times are quite short, typically less than 1 microsecond. The turn-off time is also composed of two parts, a storage time and a fall time. Fall time is the most important parameter because the peak device power dissipation is high and the heat generated is nonuniform owing to current crowding. The switching times are functions of turn-on base current, turn-off base current, base-to-emitter reverse voltage, junction temperature, and device design. Typical values are from 0.1 to 15 microseconds.

Reverse Voltage Operation—If the emitter is biased positive with respect to the collector and current is supplied to the base, then the device will conduct in the "inverted mode." The gain is low, and there exists a possibility that the current capability of the device is low. In normal practice, a rectifier is placed in inverse parallel with the collector and emitter of the transistor so that the rectifier instead of the transistor carries the reverse current.

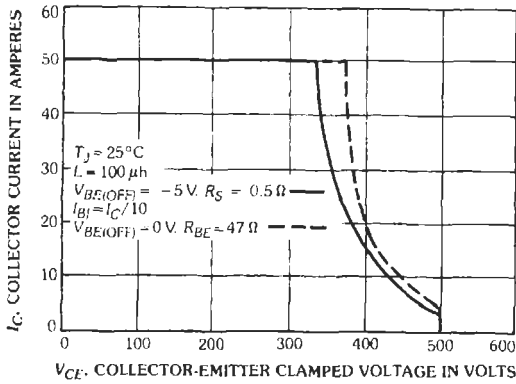


Fig. 6. Reverse-bias safe operating area (clamped) for a 100-ampere power transistor. (From *Transistors-Diodes*. Auburn, NY: General Electric Co., Semiconductor Products Dept., 1982.)

Insulated-Gate Bipolar Transistors

The insulated-gate bipolar transistor (IGBT) is a power semiconductor device used as a controlled switch. A bipolar transistor section carries the load current, and a field-effect transistor is connected to the base of the bipolar section; the field-effect transistor is an integral part of the bipolar transistor. The equivalent circuit is shown in Fig. 7. The addition of the field-effect transistor allows the IGBT to be turned on and off with a voltage signal rather than by supplying turn-on and turn-off base current as in a conventional power transistor. The device can be driven from logic-level signals. The power portion of the device, shown as a PNP transistor in Fig. 7, has all of the features described in the previous section on power transistors.

The saturation characteristic of an IGBT rated at 600 volts and 100 amperes is shown in Fig. 8. The safe-operating-area curve for the same device is shown

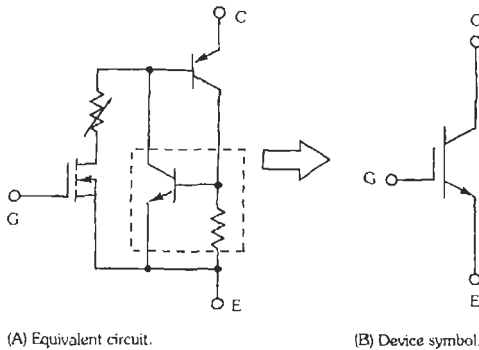


Fig. 7. Insulated-gate bipolar transistor. (From *Toshiba IGBT Application Notes*, Toshiba Corporation, 1991.)

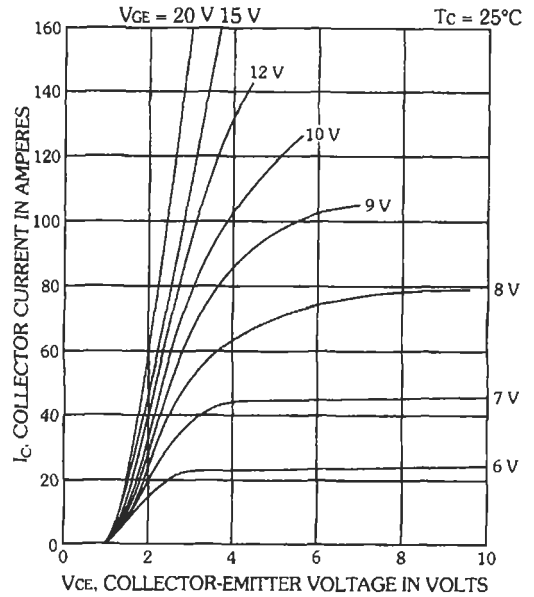


Fig. 8. On-state collector current versus collector-emitter voltage. (From *Fuji IGBT Application Manual*. Fuji Electric, 1989.)

in Fig. 9. The IGBT switching times can be controlled by selecting an appropriate value of gate resistor. The IGBT gate-to-emitter junction appears as a capacitance; decreasing the series gate resistor decreases the

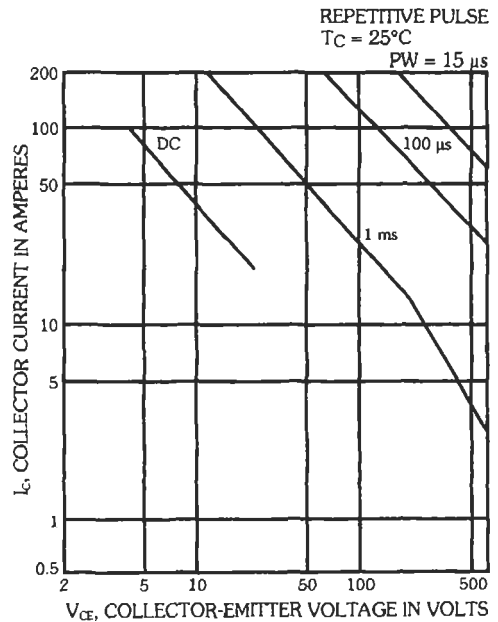


Fig. 9. Safe operating area for a 600-volt, 100-ampere IGBT. (From *Fuji IGBT Application Manual*. Fuji Electric, 1989.)

time to reach the switching voltage. The change in switching times for the 600-volt, 100-ampere IGBT is shown in Fig. 10.

Modules with current ratings of 400 amperes are currently available. Voltage ratings of 1200 volts are available. Growth in both current and voltage ratings is anticipated. Switching speeds range from less than a microsecond to 10 microseconds, resulting in operating frequencies up to 20 kHz.

Field-Effect Transistors

The field-effect transistor is also used as a controlled switch in high-voltage and high-frequency power circuits. The three terminals, drain, gate, and source, in an n-channel device bear the same relationship as the collector, base, and emitter in an NPN bipolar transistor. That is, a positive signal from gate to source causes the device to conduct a positive drain current.

The advantages of power field-effect transistors over bipolar transistors are:

- (1) Field-effect transistors have faster switching speeds with reduced delay, rise, storage, and fall times.
- (2) Devices are voltage controlled rather than current controlled and can be driven from logic-level signals.
- (3) The second-breakdown failure mechanism of bipolar transistors is absent in field-effect transistors.

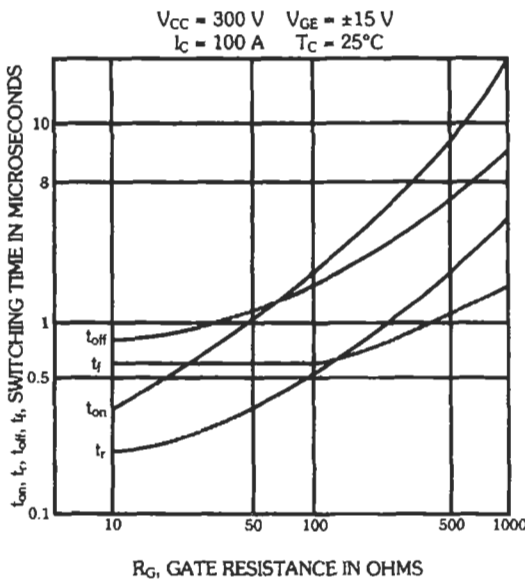


Fig. 10. Switching time versus gate resistance for an IGBT. (From *Fuji IGBT Application Manual*. Fuji Electric, 1989.)

- (4) Field-effect transistors, because of the conduction voltage drop versus temperature characteristic, tend to share current when operated directly in parallel.
- (5) The device does not block reverse voltage but has a "built-in" rectifier that has a current rating equivalent to the drain current rating.

A disadvantage of present field-effect transistors is the higher conduction voltage drop when compared with a bipolar transistor of the same current rating. The value of "on" resistance is a function of the drain-source voltage rating of the device. Higher-voltage devices have higher on resistances and therefore lower drain currents for the same temperature rise. The voltage drop can be comparable with the voltage drop of a Darlington transistor.

At present, the voltage and current ratings of field-effect transistors are not as high as those available in bipolar transistors. Field-effect transistors have replaced some bipolar transistors in switching power supplies at generally higher operating frequencies, typically over 50 kHz. Fig. 11 shows the drain-to-source voltage, during saturation, versus drain current characteristic for a 500-volt, 10-ampere n-channel field-effect transistor. The safe operating region for the same field-effect transistor is plotted in Fig. 12.

Thyristors (Silicon Controlled Rectifiers)

The thyristor* is much like a normal rectifier that has been modified to "block" in the forward direction

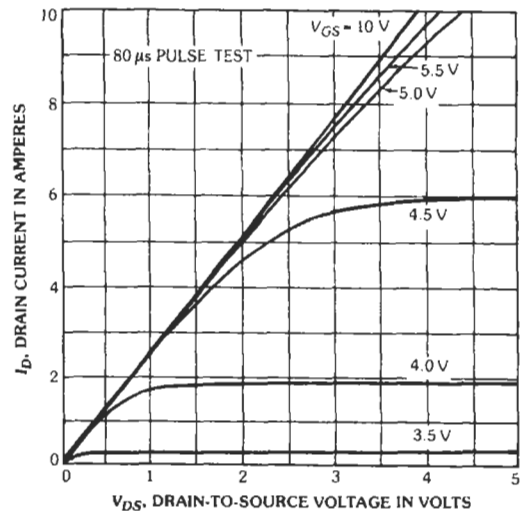


Fig. 11. Typical saturation characteristics for a 10-ampere FET. (From *HEXFET Databook*. El Segundo, CA: International Rectifier Corp., Semiconductor Div., 1981.)

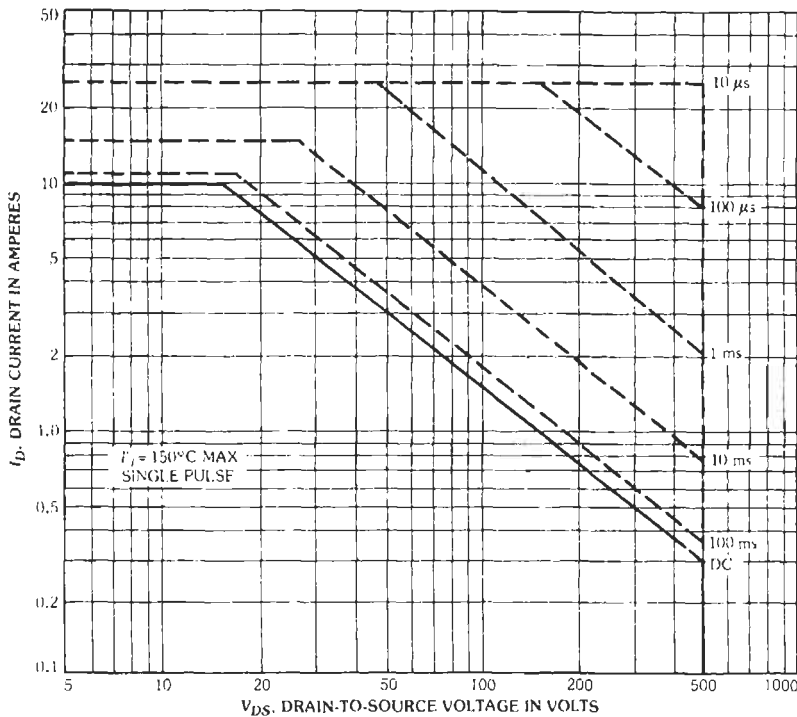


Fig. 12. Maximum safe operating area for a 10-ampere FET. (From *HEXFET Databook*, El Segundo, CA: International Rectifier Corp., Semiconductor Div., 1981.)

until a small signal is applied to the control (gate) electrode. After the signal is applied, the device conducts in the forward direction with a forward characteristic very similar to that of a normal silicon rectifier. It continues to conduct even after the control signal has been removed.

A small pulse is required at the gate electrode to switch a thyristor on, and the anode supply must be removed, reduced, or reversed to switch the thyristor off. If proportional control is required, means must be provided for adjusting the phase of the trigger pulse with respect to the supply to control the proportion of the cycle during which the thyristor is permitted to conduct. In ac circuit applications, the turn-off is obtained by the natural reversal of the supply voltage every half-cycle; in dc circuits it is usual to charge a "commutating" capacitor during the "on" period and to apply this charge in negative polarity between the anode and cathode when it is desired to turn the thyristor off.

* This term, which indicates a general class of solid-state controlled rectifiers, is used throughout this section instead of the term "silicon controlled rectifier" (a four-layer PNP device that is the most common member of the class).

Ratings—Voltage and current ratings are generally expressed in terms similar to those for silicon rectifiers, as discussed in a previous section. It is necessary, however, to add the following ratings.

Peak forward blocking voltage is the maximum safe value that may be applied, under recurrent or nonrecurrent transient conditions, while the thyristor is in the blocking state. The thyristor may break over into a conducting state regardless of gate drive if either (A) too high a positive voltage is applied between anode and cathode or (B) a positive anode-cathode voltage is applied too quickly (dv/dt firing). Even small voltage pulses, if their leading edges are sufficiently steep at the anode, can turn the thyristor on. Firing by condition (A) is generally avoided by making the peak forward blocking voltage lower than the breakover voltage of any thyristor of a particular type and by seeing that this voltage is not exceeded in practice. Trouble from dv/dt effects may be minimized by locating the gate wires to avoid stray coupling between anode and gate, by use of negative bias on the gate during blocking, and by use of R-C damping circuits between anode and cathode to slow down the rate of change of applied voltage.

Continuous forward blocking voltage covers operation under dc conditions, with forward voltage from anode to cathode.

Peak forward gate voltage is quoted for the anode positive with respect to the cathode and for the anode negative with respect to the cathode. The voltage rating is quite low in the latter case (typically 0.25 V), since the reverse voltage rating is reduced by forward gate current.

Peak reverse gate voltage is generally the same whether the anode is positive or negative with respect to the cathode.

Peak forward gate current involves several factors. Forward gate impedance is a finite value subject to quite large variations between samples and over a temperature range. It is usually necessary to plot a load line on the gate current-voltage characteristics to determine the gate current that may flow due to a given external gate voltage and source resistance. Care must be taken that the rating is not exceeded with all known spreads of gate-cathode characteristic and temperature.

Gate dissipation is generally given in terms both of average rating and of peak rating.

Characteristics—Characteristics of the thyristor important for circuit design are as follows:

Leakage currents are specified for both forward and reverse blocking, at maximum applied voltage and at maximum rated temperature. Although these currents are low in comparison with forward conducting currents and can be neglected in assessing power losses, they must be taken into account in certain circumstances. An example would be a circuit in which a capacitor is slowly charged from an external source and then suddenly discharged through a thyristor into a second circuit (as in pulse modulators); the capacitor charging operation may be affected by the amount of forward leakage current conducted by the thyristor in its blocking state.

Holding current is the minimum anode-cathode current that will keep the thyristor conducting after it has been switched on. In some applications, a thyristor with a high holding current is wanted so that it can be turned off easily without the need for reducing the anode current to a very low level. In other applications, where a low load current is normal, it might be desirable to have a low holding current to ensure that the thyristor latches on reliably with light loading.

Forward voltage drop is important in assessing power loss. The same methods of assessing power loss in terms of the forward current-voltage characteristic apply as in the case of silicon rectifiers. It is impracticable to measure the junction temperature under working conditions, and therefore the manufacturers list maximum values of stud or case temperature related to the forward current. This relationship is expressed in the form of a graph, as shown in Fig. 13.

Note that since the ratio of rms to average forward current varies with the angle of conduction, the power dissipation for any average current also varies with this angle. Fig. 13 is drawn for fractional sine waves, as would apply to the cases of single-phase half-wave or bridge rectifier circuits working into a resistive load.

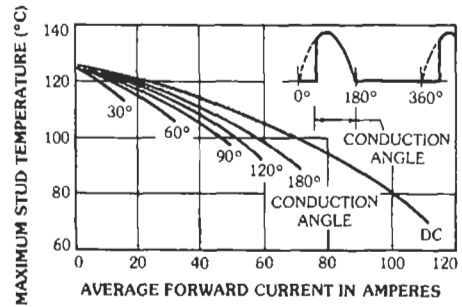


Fig. 13. Maximum permitted stud temperature vs average thyristor forward current in single-phase half-wave or bridge circuit feeding a resistive load. (70-ampere average rated thyristor.)

Fig. 14 shows a typical relationship between power dissipation and average current for a 70-ampere thyristor under the same circuit conditions. The two types of graphs illustrated by Figs. 13 and 14 together enable one to calculate the thermal resistance of heat sink required to keep the thyristor below its maximum temperature ratings when it is used under given working conditions of current and ambient temperature.

Gate trigger sensitivity is specified in terms of a minimum voltage and/or current that must be applied to ensure that all samples of a particular type of thyristor will be triggered into conduction. The minimum voltage is not temperature sensitive, but the minimum trigger current varies considerably with temperature, more current being required to turn on at low temperature than at high. The basic requirements of a gate drive circuit are therefore that the driving voltage and source resistance must be such that either the mini-

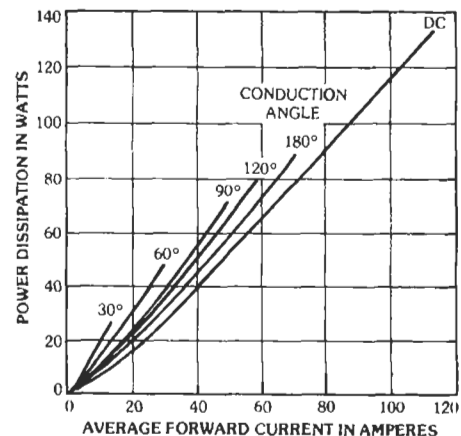


Fig. 14. Power dissipation vs average thyristor forward current in single-phase half-wave or bridge circuit feeding a resistive load. (70-ampere average rated thyristor.)

imum voltage or the minimum current (or both) is exceeded but that the rated gate dissipation is not exceeded.

Switching times of importance are the turn-on and turn-off times, the latter generally being at least one order of magnitude greater than the former. When a gate signal is applied to the thyristor, there is a finite delay time during which the anode current remains at its normal blocking level; this is followed by a "rise time" during which the anode current increases from its blocking level to a value determined by the external load circuit. Turn-on time is the sum of these two times. For a given thyristor, the turn-on time is influenced by the magnitude of gate drive, the load current to be achieved, and, to a lesser extent, the applied anode supply voltage. The time is reduced by high gate drive, low load current, and high anode supply voltage. Turn-off time is similarly composed of two individual periods; the first is a storage time, analogous to that obtained with a saturated transistor, and the second is a recovery time. Forward voltage may not be reapplied before the completion of both phases of the turn-off process, or the thyristor may conduct load current again. After this period, however, forward voltage may be applied, and the thyristor will remain in its blocking state provided that the rate of rise of anode voltage is not allowed to exceed the specified maximum dv/dt , as already discussed.

Gate Turn-Off Thyristors—These devices are fabricated to be able to turn off and interrupt the anode current by extracting current from the gate-to-cathode terminals. The ratio of controllable anode current to gate current is from 3/1 to 5/1. However, the reverse gate voltage is low, yielding a high power gain. These thyristors are rated up to 4500 volts and 2500 amperes.

MOS-Controlled Thyristors

The MOS-controlled thyristor is a four-layer, three-junction semiconductor similar to the thyristor described in the preceding section, except that an integral field-effect transistor has been added in parallel with the anode and the anode gate. The equivalent circuit is shown in Fig. 15. The addition of the field-effect structures allows the device to be turned on and off by a logic-level voltage signal. Reference 1 describes the device structure and features. Devices in a TO-218 plastic package are available with ratings of 600 volts and 50 amperes. Growth in voltage and current ratings is anticipated. Higher-current modules have been fabricated using selected lower-current devices operating in parallel. Switching times are currently in the range from 0.5 to 3 microseconds.

Triacs

A triac is a three-terminal ac semiconductor switch that is triggered into conduction by a gate signal much as a thyristor is. The triac was developed to provide a

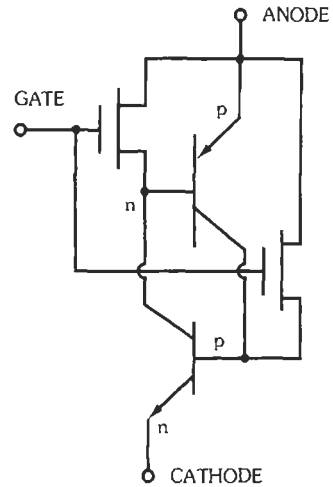


Fig. 15. Equivalent circuit of a MOS-controlled thyristor.

single device that could control current flow in both directions with the application of a gate signal. This allowed for the replacement of two inverse-parallel thyristors and the complex gate-drive circuit they require with a single device that has a single gate.

The terms "anode" and "cathode" are not used; instead, the power terminals are numbered MT_1 and MT_2 . The gate terminal is associated with MT_1 . The triggering characteristics are such that with MT_2 positive with respect to MT_1 , either positive or negative gate current will trigger the device into conduction. With MT_1 positive with respect to MT_2 , again either positive or negative gate current will trigger the device into conduction. However, positive gate current should be used only if needed, because the gain is lower.

The other major difference from two inverse thyristors is in the commutating dv/dt rating. When two thyristors are connected in inverse parallel and operating with a lagging-power-factor load, one device ceases to conduct, and the voltage rises on the previously conducting device in the reverse or non-turn-on direction. The other device sees a positive dv/dt in the turn-on direction, but it has had an entire half-cycle for turn-off or commutation. In a triac, the only period for turn-off occurs when the current goes through zero and a pulse of recovery current is conducted through the device. The voltage rises to the circuit voltage level in the turn-on direction very rapidly, and the device does not have an entire half-cycle to regain its forward blocking capability. A suitable R-C snubber is generally provided.

Triacs are available for standard ac voltages (120, 240, and some higher voltages) and up to 40–100 amperes operating at 50, 60, or 400 hertz. The remainder of their ratings are consistent with thyristor ratings.

Transient-Voltage Suppressors

Transient-voltage suppressors have a terminal characteristic such that, above a certain voltage level, the voltage drop across the device increases slightly for large changes in current magnitude. Some of these devices exhibit voltage clamping characteristics for both voltage polarities, whereas others exhibit clamping characteristics for only one voltage polarity. Two types are currently available. One is similar to a silicon zener diode, selected for a high-pulse-power, low-average-power duty cycle. The second type is composed of a polycrystalline nonlinear resistive material. Both of these categories of devices come in a range of voltage and power ratings. They are used as transient-voltage clippers on incoming ac or dc utility lines and as transient-voltage limiters for semiconductors and other voltage-sensitive equipment. Fig. 16 shows a comparison of several types and sizes of devices as a function of standardized pulse current waveforms and supply voltages.

AC-DC CONVERTER CIRCUITS

This section describes ac-to-dc converter circuits, both unregulated and electronically regulated, using rectifiers, phase-controlled thyristors, linear transistors, and switching-mode transistors. The basic design equations together with some of the technical advantages, disadvantages, and common applications are given and discussed.

Rectifier Circuits

Table 1 shows seven of the most commonly used power-rectifier circuits and general design information for each type. Their advantages, disadvantages, and common applications follow.

Single-Phase Half-Wave Rectifier—Since only half of the input wave is used, the efficiency is low and the regulation is relatively poor. Capacitors are commonly used in half-wave circuits to increase the output voltage and decrease the voltage ripple. The output voltage and degree of filtering are determined by the value of capacitance used in relation to the load current. Transformer design is complicated, and the unidirectional secondary current causes core saturation and poor regulation. Most half-wave circuits operate either directly from ac lines or at a high voltage with a relatively low current.

Single-Phase Full-Wave Center-Tap Rectifier—The efficiency is good, but the transformer ac voltage is approximately 2.2 times the dc output voltage. The circuit requires a larger transformer than an equivalent bridge rectifier, with the added complication of a center tap. Each arm of the center-tap circuit must block the full terminal voltage of the transformer. Because of this, center-tap connections are economical only in voltage ranges where not more than one rectifier per arm is required. If series units must be used to obtain the required output voltage, a bridge circuit is preferable. This circuit is used with low-voltage logic-level

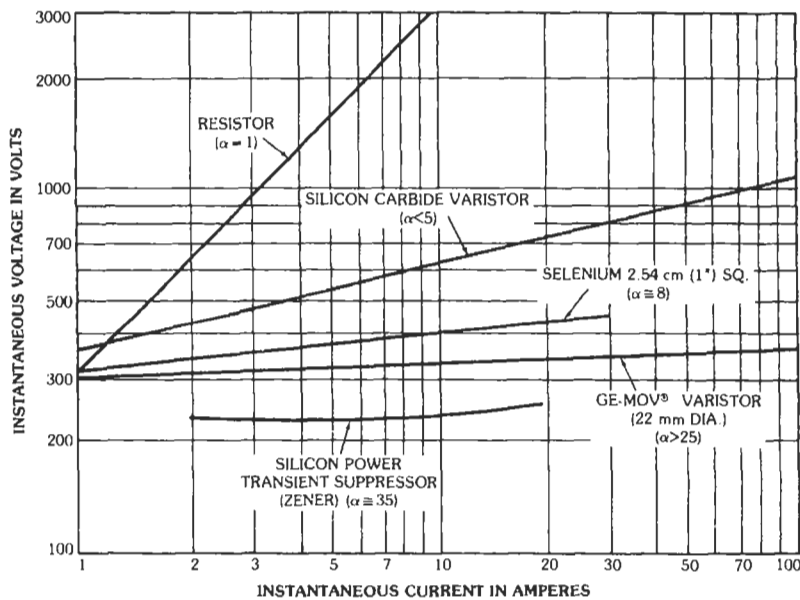


Fig. 16. Voltage-current characteristics of four transient-suppressor devices. (From *Transient Voltage Suppression*, 3rd ed. Auburn, NY: General Electric Co., Semiconductor Products Dept., 1982.)

power supplies where the efficiency is improved because there is only one diode in series with the output.

Single-Phase Full-Wave Bridge Rectifier—If single-phase full-wave output is required, a bridge circuit is commonly used. Efficiency is good and transformer design is easy. Filtering is simplified because the ripple frequency is twice the input frequency.

Three-Phase Wye (or Star) Half-Wave Rectifier—This circuit is commonly used if dc output-voltage requirements are relatively low and current requirements are moderately large. The dc output voltage is approximately equal to the phase voltage. However, each of the three arms must block the line-to-line voltage, which is approximately 2.5 times the phase voltage. For this reason, it is desirable to use a three-phase half-wave connection only where one series unit per arm will provide the required dc output. The transformer design and utilization are somewhat complicated because there is a tendency to saturate the core with unidirectional current in each winding.

Three-Phase Full-Wave Bridge Rectifier—This circuit is commonly used if high dc power is required and if efficiency must be considered. The ripple component in the load is 4.2% at a frequency six times the input frequency, and so additional filtering is not required in most applications. The dc output voltage is approximately 25% higher than the phase voltage, and each arm must block only the phase voltage. Transformer utilization is good. This is the most common three-phase rectifier connection.

Three-Phase Diametric Half-Wave Rectifier—The characteristics of this circuit approximate those of the three-phase double-wye circuit without an interphase transformer. Popular applications include requirements for very high dc load currents in low-to-medium voltage ranges (approximately 6 to 125 volts dc).

Three-Phase Double-Wye Half-Wave Rectifier—A three-phase double-wye connection is recommended if a very high direct current is required at a relatively low dc voltage. Each arm is required to block the full phase voltage of the secondary windings. The dc output current rating is double that of a three-phase bridge or half-wave connection. However, the output voltage is only 75% of the phase voltage. The transformer design is complicated by additional connections and extra insulation, and an interphase transformer (or balance coil) is required.

Single-Phase Voltage Multipliers

These circuits use the principle of charging capacitors in parallel from the ac input and adding the voltages across them in series to obtain dc voltages higher

than the source voltage. Filtering must be of the capacitor-input type.

Conventional and Cascade Voltage Doublers—In the conventional circuit (Fig. 17), capacitors C_1 and C_2 are each charged, during alternate half-cycles, to the peak value of the alternating input voltage. The capacitors are discharged in series into load R_L , thus producing an output across the load of approximately twice the ac peak voltage.

In the cascade circuit (Fig. 18), C_1 is charged to the peak value of the ac input voltage through rectifier CR_2 during one half-cycle, and during the other half-cycle it discharges in series with the ac source through CR_1 to charge C_2 to twice the ac peak voltage.

The “conventional” circuit has slightly better regulation, and, since the ripple frequency is twice the supply frequency, the output is easier to filter, the percentage ripple being approximately the same in both circuits. In addition, both capacitors are rated at the peak ac voltage, whereas C_2 in the cascade circuit must be rated at twice this value. With both circuits, the peak inverse voltage across each rectifier is twice the ac peak. The cascade circuit, however, has the advantage of a common input and output terminal and, therefore, permits the combination of units to give higher-order voltage multiplications. The regulation of both circuits is poor, so that only small load currents can be drawn.

Bridge Rectifier or Voltage Doubler—If SW in Fig. 19 is open, the circuit is a bridge rectifier; if the peak ac applied voltage is $2E_p$, the dc output voltage is $2E_p$. If SW is closed, the circuit is a voltage doubler; an ac input of E_p gives a dc output of $2E_p$. This circuit is used if the same dc voltage is desired and the input voltage can be either 115 or 230 volts ac.

Further Voltage Multiplication—The cascade voltage doubler shown in Fig. 18 can be combined several times to obtain higher dc voltages, as shown in Fig. 20. The voltage ratings of all the capacitors and rectifiers are twice the ac peak voltage, but the capacitors must have the values shown. The value of C will be the same as that for the cascade voltage doubler (Fig. 18), which is the basic unit for the circuit in Fig. 20. The load current must be small. The increasing size of capacitors and the deterioration in regulation limit the voltages that can be obtained from this type of circuit.

Phase-Controlled Thyristor Circuits

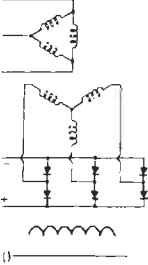
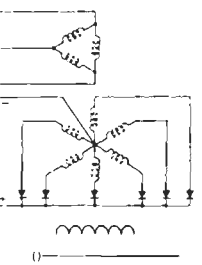
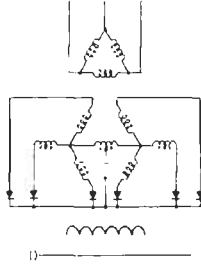
As described earlier, the thyristor can block forward voltage until it is triggered into conduction. The delay angle of gating the thyristor into conduction is designated as alpha (α). With $\alpha = 0$, the thyristor behaves like a rectifier. As the delay angle is increased, the average dc output voltage decreases, and therefore

TABLE I. RECTIFIER

Type of Circuit→		Single-Phase Half-Wave	Single-Phase Center Tap	Single-Phase Bridge	Three-Phase Star (Wye)
Primary→					
Secondary→					
One Cycle Wave of Rectifier Output Voltage (No Overlap)					
Number of rectifier elements	=	1	2	4	3
Rms dc volts output	=	1.57	1.11	1.11	1.02
Peak dc volts output	=	3.14	1.57	1.57	1.21
Peak reverse volts per rectifier element	=	3.14	3.14	1.57	2.09
	=	1.41	2.82	1.41	2.45
	=	1.41	1.41	1.41	1.41
Average dc output current	=	1.00	1.00	1.00	1.00
Average dc output current per rectifier element	=	1.00	0.500	0.500	0.333
Rms current per rectifier element:					
Resistive load	=	1.57	0.785	0.785	0.587
Inductive load	=	—	0.707	0.707	0.578
Peak current per rectifier element:					
Resistive load	=	3.14	1.57	1.57	1.21
Inductive load	=	—	1.00	1.00	1.00
Ratio of peak to average current per element:					
Resistive load	=	3.14	3.14	3.14	3.63
Inductive load	=	—	2.00	2.00	3.00
% Ripple (rms of ripple/ average output voltage)	=	121%	48%	48%	18.3%
Ripple Frequency	=	1	2	2	3
		Resistive Load		Inductive Load or Large Choke Input Filter	
Transformer secondary rms volts per leg	=	2.22	1.11 (to center tap)	1.11 (total)	0.855 (to neutral)
Transformer secondary rms volts line-to-line	=	2.22	2.22	1.11	1.48
Secondary line current	=	1.57	0.707	1.00	0.578
Transformer secondary volt-amperes	=	3.49	1.57	1.11	1.48
Transformer primary rms amperes per leg	=	1.21	1.00	1.00	0.471
Transformer primary volt-amperes	=	2.69	1.11	1.11	1.21
Average of primary and secondary volt-amperes	=	3.09	1.34	1.11	1.35
Primary line current	=	1.21	1.00	1.00	0.817
Line power factor	=	—	0.900	0.900	0.826

* The data assume zero forward drop and zero reverse current in rectifiers and no alternating-current line or source reactance.

CIRCUIT CHART²³

Three-Phase Bridge	Six-Phase Star (Three-Phase Diametric)	Three-Phase Double Wye With Interphase Transformer		
				To Determine Actual Value of Parameter in Any Column, Multiply Factor Shown by Value of:
6	6	6		
1.00	1.00	1.00	×	Average dc voltage output
1.05	1.05	1.05	×	Average dc voltage output
1.05	2.09	2.42	×	Average dc voltage output
2.45	2.83	2.83	×	Rms secondary volts per transformer leg
1.41	1.41	1.41	×	Rms secondary volts line-to-line
1.00	1.00	1.00	×	Average dc output current
0.333	0.167	0.167	×	Average dc output current
0.579	0.409	0.293	×	Average dc output current
0.578	0.408	0.289	×	Average dc output current
1.05	1.05	0.525	×	Average dc output current
1.00	1.00	0.500	×	Average dc output current
3.15	6.30	3.15		
3.00	6.00	3.00		
4.2%	4.2%	4.2%		
6	6	6	×	Line frequency, <i>f</i>
Inductive Load or Large Choke Input Filter				
0.428 (to neutral)	0.740 (to neutral)	0.855 (to neutral)	×	Average dc voltage output
0.740	1.48 (max)	1.71 (max-no load)	×	Average dc voltage output
0.816	0.408	0.289	×	Average dc output current
1.05	1.81	1.48	×	Dc watts output
0.816	0.577	0.408	×	Average dc output current
1.05	1.28	1.05	×	Dc watts output
1.05	1.55	1.26	×	Dc watts output
1.41	0.817	0.707	×	(Avg. load current × sec. leg voltage)/primary line voltage
0.955	0.955	0.955		

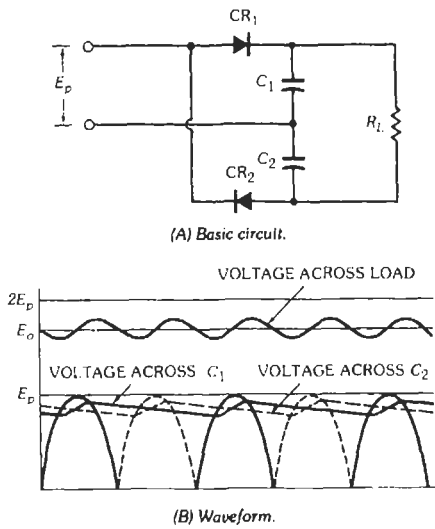


Fig. 17. Conventional voltage doubler.

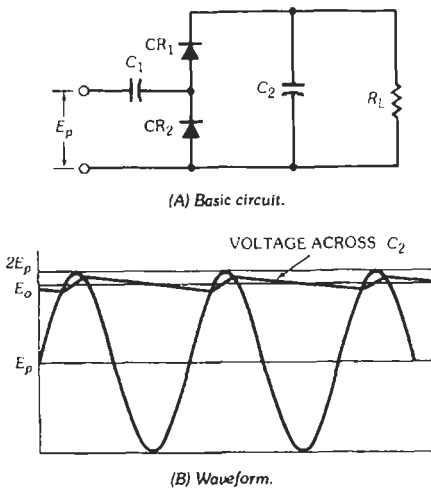


Fig. 18. Cascade voltage doubler.

these circuits provide electronic regulation of the dc output voltage.

Figs. 21 through 26 are basic circuits of thyristors used as controlled rectifiers. In many applications, it is not necessary for all rectifier elements to be controllable, and it is common to find bridges composed of thyristors and ordinary diodes in equal numbers; such circuits are called *half-controlled*, whereas those containing only thyristors are termed *full-controlled*.

Figs. 21 and 22 show a full-controlled single-phase thyristor bridge, the thyristor circuit being the same whether a resistive or inductive load is being driven,

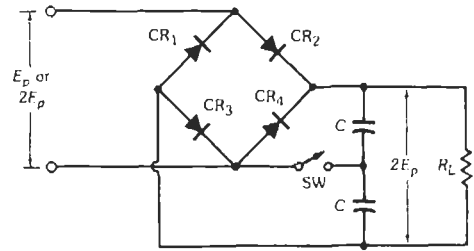


Fig. 19. Bridge rectifier or voltage doubler.

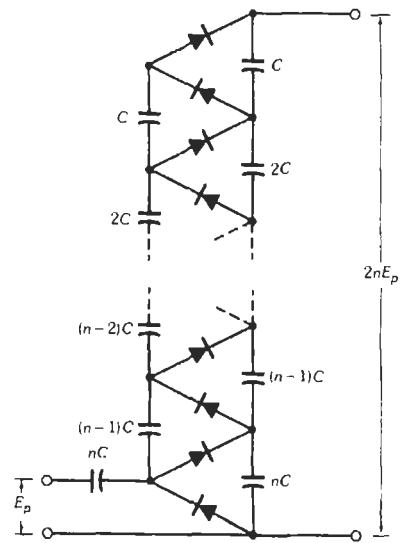


Fig. 20. Circuit for high-order voltage multiplication.

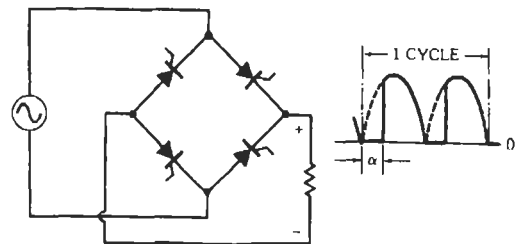


Fig. 21. Full-controlled single-phase thyristor bridge. The broken lines indicate the path of the normal waveform at full conduction, i.e., where $\alpha = 0^\circ$ and the circuit behaves like a diode rectifier. The quantity α is the firing angle delay. Delay range required for this circuit is 0° to 180° .

The voltage waveforms are different, however, in the two cases. The principal difference of practical importance is that the range of firing-pulse phase control is required to be different in the two cases. For the resis-

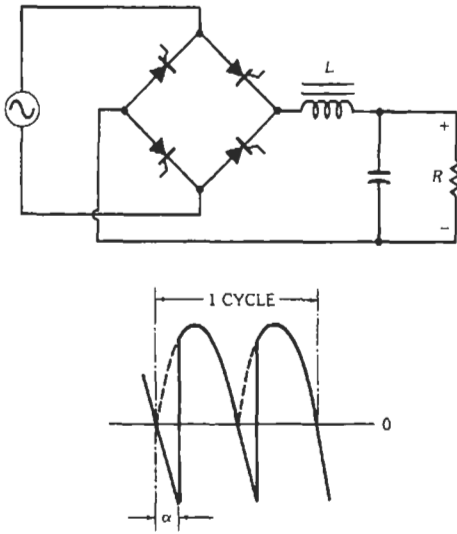


Fig. 22. Full-controlled single-phase thyristor bridge working into an inductive load where $\omega L \gg R$. Delay range required is 0° to 90° .

tive load, phase control over the range 0° to 180° is necessary to obtain full control from maximum output voltage down to zero; for an almost pure inductive load (*i.e.*, a very high $\omega L/R$ ratio), full-phase control is obtained with a range of only 0° to 90° .

Fig. 23 shows a half-controlled single-phase rectifier driving an inductive load. For such loads a bypass diode,* CR_3 , must be added at the output: at the end of a voltage half-cycle, current still flows in the choke, but in this circuit the current is transferred at the end of the voltage half-cycle to the bypass diode. The two important effects of this diode are (A) that the output

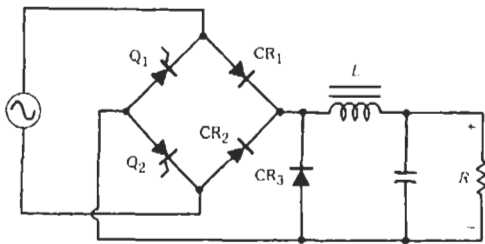


Fig. 23. Half-controlled single-phase thyristor bridge working into inductive load. Common-anode thyristor connection. Waveform and delay range as for Fig. 21.

* Sometimes called a "commutating," "freewheel," or "flyback" diode.

voltage is clamped to zero while this inductively maintained current flows, so that the output waveform is the same as that of Fig. 21 (having the characteristic flat portion), and (B) that the transfer of load current from the thyristor to the diode turns the thyristor off. If CR_3 were not present, the waveform would still be the same, since a zero voltage clamp would exist through a series combination CR_1-Q_1 or CR_2-Q_2 , depending on which thyristor was conducting during the previous half-cycle. However, this flow of uncontrolled current through the thyristor, bypassing the supply, is undesirable; for example, should it be required that the output voltage be turned off by removal of gate pulses, this action may prove to be impossible. The thyristor can be held on continuously through the negative half-cycle by inductive circulating current and is then ready to conduct on the next positive half-cycle. Thus the gate has lost control, and a continuous half-wave output is produced. The bypass diode overcomes these difficulties by ensuring thyristor turn-off at the end of each voltage half-cycle.

Fig. 24 is a push-pull controlled rectifier circuit. Fig. 25 is a three-phase full-controlled rectifier circuit. Fig. 26 shows the circuit of a three-phase half-controlled rectifier, the bypass diode being necessary only for inductive loads.

In all the cases illustrated in which a bypass diode is used, this diode must be rated continuously for a maximum average current equal to the load current if the full load current is to be drawn when the average output voltage is reduced almost to zero. In practice, a larger diode is often used in this position than in the bridge arms, or several diodes may be used in parallel.

For each of the circuits illustrated in Figs. 21 through 26, Table 2 gives equations for the average dc output voltage, V_{do} , at any angle α in terms of the maximum average dc output voltage, V_{do} , obtained at $\alpha = 0^\circ$. The table also shows the value of V_{do} for each circuit in terms of the peak sinusoidal input voltage, E_p . For the single-phase push-pull circuit, the peak input voltage is E_p-0-E_p , and for all the three-phase circuits E_p is defined as the peak value of the line-to-line voltage.

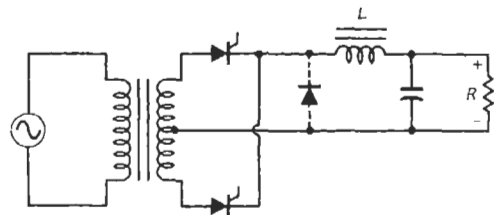


Fig. 24. Full-controlled single-phase push-pull thyristor rectifier working into inductive load. With diode connected, waveform and delay range are as for Fig. 21. Without diode, waveform and delay range are as for Fig. 22 (if $\omega L \gg R$). Addition of diode reduces critical inductance of L for continuous current.

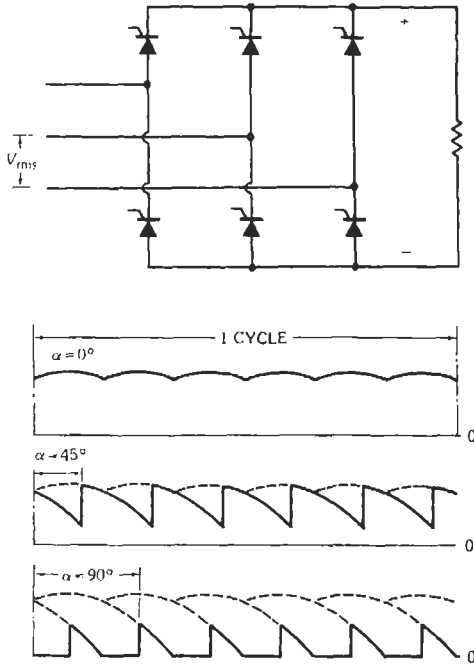


Fig. 25. Full-controlled three-phase thyristor bridge working into resistive load; waveforms of voltage output are shown for three values of α . Waveform develops characteristic flat portion when $\alpha \geq 60^\circ$. Delay range required is 0° to 120° .

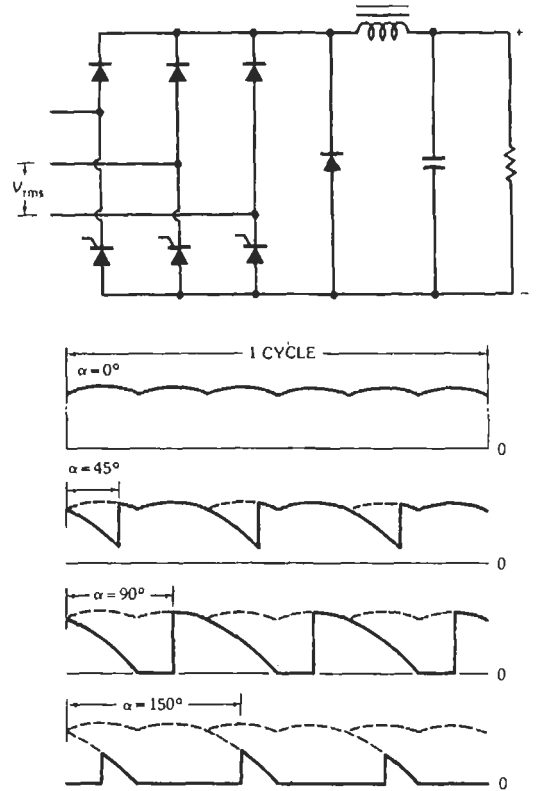


Fig. 26. Half-controlled three-phase thyristor bridge. The bypass diode is necessary only when feeding an inductive load. The waveforms are the same for resistive or inductive load. Delay range required is 0° to 180° .

Linear Transistor Circuits

A block diagram of a series voltage regulator is shown in Fig. 27. The regulator is composed of three basic parts, a comparator, an amplifier, and a series pass transistor. The purpose of the comparator is to compare the actual load voltage with a reference voltage and develop an error signal proportional to the difference between the two. This error voltage is

amplified and level shifted to provide base current to a series transistor operating in its linear region. The advantages of this system are its simplicity, its current availability as a complete integrated circuit from several manufacturers, its fast response, and its low switching noise. Its disadvantage is its power dissipa-

TABLE 2. MEAN DC OUTPUT VOLTAGE FOR THYRISTOR-CONTROLLED RECTIFIERS

Circuit	V_{do}	$V_{d\alpha}$
Fig. 21	$2E_p/\pi$	$V_{do}(1/2)(1+\cos\alpha)$
Fig. 22	$2E_p/\pi$	$V_{do} \cos\alpha$
Fig. 23	$2E_p/\pi$	$V_{do}(1/2)(1+\cos\alpha)$
Fig. 24, no diode and inductive load	$2E_p/\pi$	$V_{do} \cos\alpha$
Fig. 24, all other cases	$2E_p/\pi$	$V_{do}(1/2)(1+\cos\alpha)$
Fig. 25	$3E_p/\pi$	$V_{do} \cos\alpha$, for $\alpha=0^\circ$ to 60° $V_{do}[1+\cos(\alpha+60^\circ)]$, for $\alpha=60^\circ$ to 120°
Fig. 26	$3E_p/\pi$	$V_{do}(1/2)(1+\cos\alpha)$

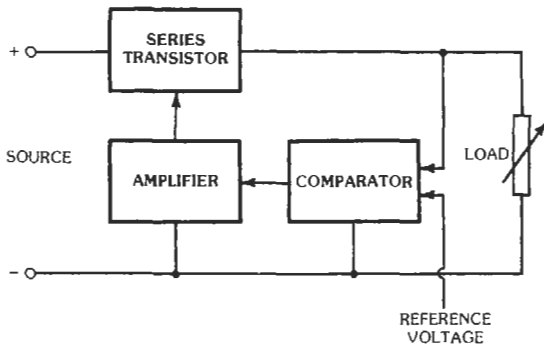


Fig. 27. Block diagram of linear series transistor voltage regulator.

tion, since the series transistor operates as a linear resistive element. Series pass transistors can be provided external to available integrated circuits for higher-voltage and higher-current applications. These integrated control circuits can be provided with load-current limit, remote start-up and shut-down, and over-temperature limits.

Switch-Mode Power Supplies

To achieve higher output power capability, transistor switch-mode power circuits have been developed. The transistor is operated either in its saturated state or its cutoff state. Switching frequencies can range from 2 kHz up to several megahertz. Load-voltage control is provided by adjustment of the repetition rate, pulse width, or on-to-off ratio. A functional block diagram of a switching regulator is shown in Fig. 28. An input and output filter are shown to reduce the ripple current in the source and load, respectively. The transistor either is on, connecting the source and the load, or is cut off, disconnecting the source and the load.

Basic Converters

There are three basic dc-to-dc switching converter configurations: the step-down (buck), the step-up

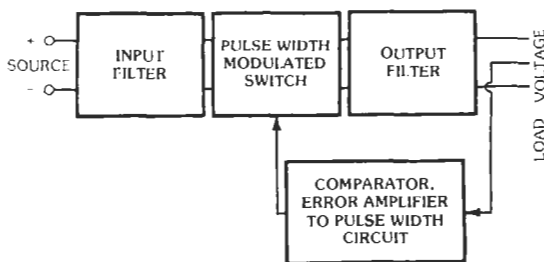


Fig. 28. Block diagram of switch-mode power supply.

(boost), the step-down/step-up (buck/boost), and two additional nonisolated converter topologies: the SEPIC and the Cuk converter. These five circuits do not provide ohmic isolation between the input and the output. Several other circuits have been developed that do provide transformer isolation between input and output. The typical examples are the flyback converter, the forward converter, the half-bridge, the full-bridge, the push-pull, and the Cuk converter with isolation transformer. Table 3 shows these circuits, their schematic diagrams, transfer functions, device rating equations, circuit waveforms, and advantages and disadvantages.

The vast majority of contemporary designs of both isolated and nonisolated converters is based upon the application of two basic concepts: the forward-mode and a flyback or boost-mode converter.

Forward-Mode Converter Basics

The *buck* or *step-down converter*, which can be seen in Fig. 29, is the basic forward-mode converter. Its operation period—one switching cycle—has two distinct parts, which occur when the series power switch is on and when it is off. When the power switch is on, the input voltage is connected to the input of the inductor. The output of the inductor is the output voltage, and the rectifier is back-biased. During this part of the switching period, since there is a constant voltage source connected across the inductor, the inductor current begins to linearly ramp upward, which is described by

$$i_{L(on)} = \frac{(V_{in} - V_{out}) \times t_{on}}{L}$$

where t_{on} is the on-time of the switching cycle when the power switch Q is on.

During the on-time, energy is being stored within the inductor core material in the form of magnetic flux. In steady-state operation, there is sufficient energy stored to carry the requirements of the load during the next off-time, when the input voltage source (V_{in}) is not connected with the load. The next off-time is the "off" period of the power switch. When the power switch turns off, the input voltage of the inductor flies below ground and is clamped at one diode drop below ground by the *catch diode*. The energy stored in the inductor core causes the current continue to flow in its previous direction through the inductor and the catch diode, thus maintaining the load current loop. This removes the stored energy from the inductor. The inductor current during this time is

$$i_{L(off)} = \frac{(V_{out} - V_D) \times t_{off}}{L}$$

where t_{off} is the off-time of the switching cycle when the power switch Q is off.

This period ends when the power switch is once again turned on. Regulation is accomplished by varying the on-to-off duty cycle of the power switch. The

TABLE 3. POWER TRANSISTOR AND DIODE REQUIREMENTS FOR SWITCH-MODE POWER SUPPLIES ²³

<p>Circuit Configuration</p>			
<p>Type of Converter</p>	<p>(A) Buck (Step-Down)</p>	<p>(B) Boost (Step-Up)</p>	<p>(C) Buck-Boost</p>
<p>Ideal Transfer Function</p>	$V_O/V_{IN} = \tau/T = D$	$V_O/V_{IN} = T/(T - \tau)$	$V_O/V_{IN} = -\tau/(T - \tau)$
<p>Collector Current (i_c) †</p>	$I_{C_{MAX}} = I_{RL} + \Delta I_{L1}/2$	$I_{C_{MAX}} = I_{RL}[T/(T - \tau)] + \Delta I_{L1}/2$	$I_{C_{MAX}} = I_{RL}[T/(T - \tau)] + \Delta I_{L1}/2$
<p>Collector Voltage Rating †</p>	$V_{CEO} = V_{IN}$	$V_{CEO} > V_O + 1$	$V_{CEO} > V_{IN} + V_O$
<p>Diode Currents †</p>	$I_{CR1} = I_{RL}(T - \tau)/T$	$I_{CR1} = I_{RL}$	$I_{CR1} = I_{RL}$
<p>Diode Voltages (V_{RM}) †</p>	$V_{CR1} = V_{IN}$	$V_{CR1} = V_O$	$V_{CR1} = V_O + V_{IN}$

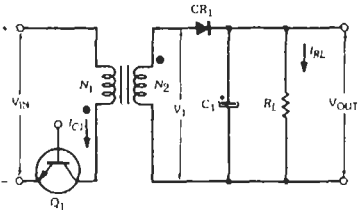
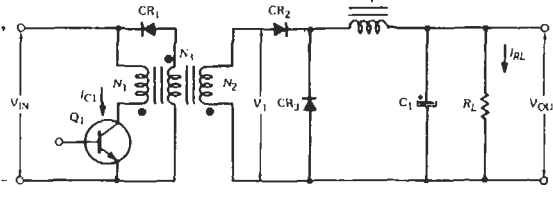
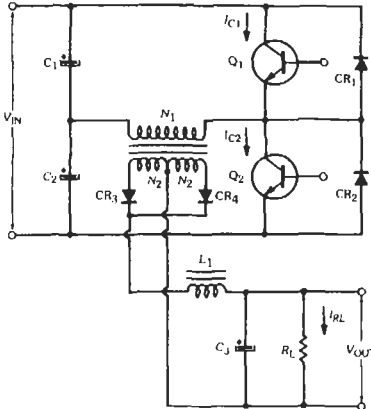
<p>Voltage and Current Waveforms</p>			
<p>Advantages</p>	<p>High efficiency. Simple. No transformer. High frequency operation. Easy to stabilize regulator loop.</p>	<p>High efficiency. Simple. No transformer. High frequency operation.</p>	<p>Voltage inversion without using a transformer. Simple. High frequency operation.</p>
<p>Disadvantages</p>	<p>No isolation between input and output. Requires a crowbar if Q_1 shorts. C_1 has high ripple current. Current limit difficult. Only one output is possible.</p>	<p>No isolation between input and output. High peak collector current. Only one output is possible. Poor transient response. Regulator loop hard to stabilize.</p>	<p>Q_1 must carry high peak current. No isolation between input and output. Only one output is possible. Poor transient response.</p>

* From Application Note 200.87. Auburn, NY: General Electric Co., Semiconductor Products Dept., 1979.

† For reliable operation, it is suggested and recommended that all voltage and current ratings be increased to 125% of the required maximum.

Continued on next page.

TABLE 3 (CONT). POWER TRANSISTOR AND DIODE REQUIREMENTS FOR SWITCH-MODE POWER SUPPLIES *

Circuit Configuration			
Type of Converter	(D) Flyback	(E) Forward	(F) Half Bridge
Ideal Transfer Function	$V_O/V_{IN} = (N_2/N_1)[\tau/(T-\tau)]$	$V_O/V_{IN} = (N_2/N_1)(\tau/T)$	$V_O/V_{IN} = (N_2/N_1)(\tau/T)$
Collector Current (i_c) [†]	$I_{C_{MAX}} = I_{RL}(N_2/N_1)[T/(T-\tau)] + \Delta I_{L1}/2$	$I_{C_{MAX}} = (N_2/N_1)(I_{RL} + \Delta I_{L1}/2) + \hat{I}_{MAG}$	$I_{C_{MAX}} = (N_2/N_1)(I_{RL} + \Delta I_{L1}/2) + \hat{I}_{MAG}$
Collector Voltage Rating [‡]	$V_{CEO} > V_{IN} + (N_1/N_2)V_{OUT}$	$V_{CEO} > V_{IN}(1 + N_1/N_2)$	$V_{CEO} = V_{IN}$
Diode Currents [‡]	$I_{CR1} = I_{RL}$	$I_{CR1} = (\hat{I}_{MAG}/2)(\tau/T)$ $I_{CR2} = I_{RL}(\tau/T)$ $I_{CR3} = I_{RL}(T-\tau)/T$	$I_{CR1} = I_{CR2} = (\hat{I}_{MAG}/2)(\tau/T)$ $I_{CR3} = I_{CR4} = I_{RL}/2$
Diode Voltages (V_{RM}) [‡]	$V_{CR1} = V_{IN}(N_2/N_1)$	$V_{CR1} = V_{IN}(1 + N_3/N_1)$ $V_{CR2} = V_{IN}(N_2/N_3)$ $V_{CR3} = V_{IN}(N_2/N_1)$	$V_{CR1} = V_{CR2} = V_{IN}$ $V_{CR3} = V_{CR4} = V_{IN}(N_2/N_1)$

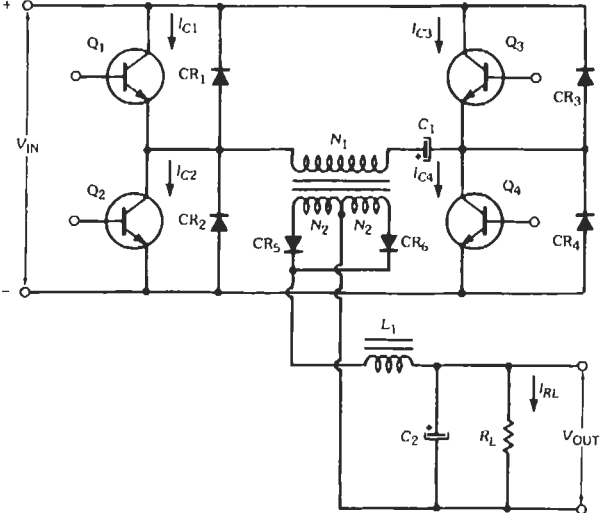
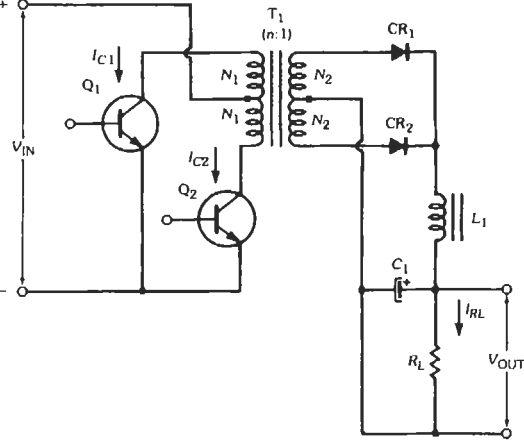
<p>Voltage and Current Waveforms</p>			
<p>Advantages</p>	<p>Simple. Multiple outputs are possible. Collector current reduced by turns ratio of transformer. Low parts count. Isolation.</p>	<p>Simple. Multiple outputs are possible. Collector current reduced by ratio of N_2/N_1. Low output ripple.</p>	<p>Simple. Good transformer utilization. Transistors rated at V_{IN}. Isolation. Multiple outputs. i_c reduced as a ratio of N_2/N_1. High power output.</p>
<p>Disadvantages</p>	<p>Poor transformer utilization. Transformer design critical. High output ripple.</p>	<p>Poor transformer utilization. Poor transient response. Parts count high. Transformer design is critical.</p>	<p>Poor transient response. High parts count. C_1 and C_2 have high ripple current. Limited dynamic range. Requires auxiliary power supplies for control circuits.</p>

* From Application Note 200.87. Auburn, NY: General Electric Co., Semiconductor Products Dept., 1979.

† For reliable operation, it is suggested and recommended that all voltage and current ratings be increased to 125% of the required maximum.

Continued on next page.

TABLE 3 (CONT). POWER TRANSISTOR AND DIODE REQUIREMENTS FOR SWITCH-MODE POWER SUPPLIES *

Circuit Configuration		
Type of Converter	(G) Full Bridge	(H) Push-Pull
Ideal Transfer Function	$V_O/V_{IN} = 2(N_2/N_1)(\tau/T)$	$V_O/V_{IN} = 2(N_2/N_1)(\tau/T)$
Collector Current (i_c)†	$I_{C_{MAX}} = (N_2/N_1)(I_{RL} + \Delta I_{LI}/2) + \hat{I}_{MAG}$	$I_{C_{MAX}} = (N_2/N_1)(I_{RL} + \Delta I_{LI}/2) + \hat{I}_{MAG}$
Collector Voltage Rating †	$V_{CEO} = V_{IN}$	$V_{CEO} = 2 V_{IN}$
Diode Currents †	$I_{CR1} = I_{CR2} = (\hat{I}_{MAG}/2)(\tau/T)$ $I_{CR3} = I_{CR4} = (\hat{I}_{MAG}/2)(\tau/T)$ $I_{CR5} = I_{CR6} = I_{RL}/2$	$I_{CR1} = I_{RL}/2$ $I_{CR2} = I_{RL}/2$
Diode Voltages (V_{RM})†	$V_{CR1} = V_{CR2} = V_{CR3} = V_{CR4} = V_{IN}$ $V_{CR5} = V_{CR6} = 2 V_{IN}(N_2/N_1)$	$V_{CR1} = 2 V_{IN}(N_2/N_1)$ $V_{CR2} = 2 V_{IN}(N_2/N_1)$

<p>Voltage and Current Waveforms</p>		
<p>Advantages</p>	<p>Simple. Good transformer utilization. Transistors rated at V_{IN}. Isolation. Multiple outputs. i_c reduced as a ratio of N_2/N_1. High power output. Preferred to circuit F where high power required.</p>	<p>Simple. Good transformer utilization. Collector current reduced as a function of N_2/N_1. Good at low values of V_{IN}.</p>
<p>Disadvantages</p>	<p>Poor transient response. High parts count. C_1 and C_2 have high ripple current. Limited dynamic range. Requires auxiliary power supplies for control circuit.</p>	<p>Cross conduction of Q_1, Q_2 possible. High parts count. Transformer design critical. Poor dynamic range. Poor transient response.</p>

* From Application Note 200.87. Auburn, NY: General Electric Co., Semiconductor Products Dept., 1979.

† For reliable operation, it is suggested and recommended that all voltage and current ratings be increased to 125% of the required maximum.

Continued on next page.

TABLE 3 (CONT). POWER TRANSISTOR AND DIODE REQUIREMENTS FOR SWITCH-MODE POWER SUPPLIES*

<p>Circuit Configuration</p>		
<p>Type of Converter</p>	<p>(I) Cuk (Boost-Buck Inverting)</p>	<p>(J) Cuk (With Transformer)</p>
<p>Ideal Transfer Function</p>	$V_O/V_{IN} = -\tau/(T-\tau)$ $D = \tau/T$	$V_O/V_{IN} = (N_2/N_1)[\tau/(T-\tau)]$ $D = \tau/T$
<p>Collector Current (i_c)†</p>	$I_c = 1.5 I_{RL} \text{ for } D = 0.33$ $I_c = 2 I_{RL} \text{ for } D = 0.50$ $I_c = 2.5 I_{RL} \text{ for } D = 0.60$	$I_c = 1.5(N_2/N_1)I_{RL} \text{ for } D = 0.33$ $I_c = 2(N_2/N_1)I_{RL} \text{ for } D = 0.50$ $I_c = 2.5(N_2/N_1)I_{RL} \text{ for } D = 0.60$
<p>Collector Voltage Rating †</p>	$V_{CEO} \geq 2 V_{IN}$	$V_{CEO} = 1.5 V_{IN} \text{ for } D = 0.33$ $V_{CEO} = 2 V_{IN} \text{ for } D = 0.50$ $V_{CEO} = 2.5 V_{IN} \text{ for } D = 0.60$
<p>Diode Currents †</p>	$I_{CR1} = 1.5 I_{RL} \text{ for } D = 0.33$ $I_{CR1} = 2 I_{RL} \text{ for } D = 0.50$ $I_{CR1} = 2.5 I_{RL} \text{ for } D = 0.60$	$I_{CR1} = 1.5 I_{RL} \text{ for } D = 0.33$ $I_{CR1} = 2 I_{RL} \text{ for } D = 0.50$ $I_{CR1} = 2.5 I_{RL} \text{ for } D = 0.60$
<p>Diode Voltages (V_{RM})†</p>	$V_{CR1} = 1.5 V_{IN} \text{ for } D = 0.33$ $V_{CR1} = 2 V_{IN} \text{ for } D = 0.50$ $V_{CR1} = 2.5 V_{IN} \text{ for } D = 0.60$	$V_{CR1} = 1.5(N_2/N_1)V_{IN} \text{ for } D = 0.33$ $V_{CR1} = 2(N_2/N_1)V_{IN} \text{ for } D = 0.50$ $V_{CR1} = 2.5(N_2/N_1)V_{IN} \text{ for } D = 0.60$

<p>Voltage and Current Waveforms</p>		
<p>Advantages</p>	<p>Continuous input and output current. Highest efficiency. Low ripple. Smallest number of switching components. Switching losses cut in half. Drive circuit referenced to ground. Highest operating frequency.</p>	<p>Continuous input and output current. Highest efficiency. Very low ripple. Smallest number of switching components. Switching losses low. Drive current referenced to ground. Highest operating frequency.</p>
<p>Disadvantages</p>	<p>High collector current. C_1 has high ripple current requirement. High voltage required for Q_1. Power output limited.</p>	<p>C_1 and C_2 have high ripple current requirements. Transformer design critical. Power output is limited.</p>

* From Application Note 200.87. Auburn, NY: General Electric Co., Semiconductor Products Dept., 1979.

† For reliable operation, it is suggested and recommended that all voltage and current ratings be increased to 125% of the required maximum.

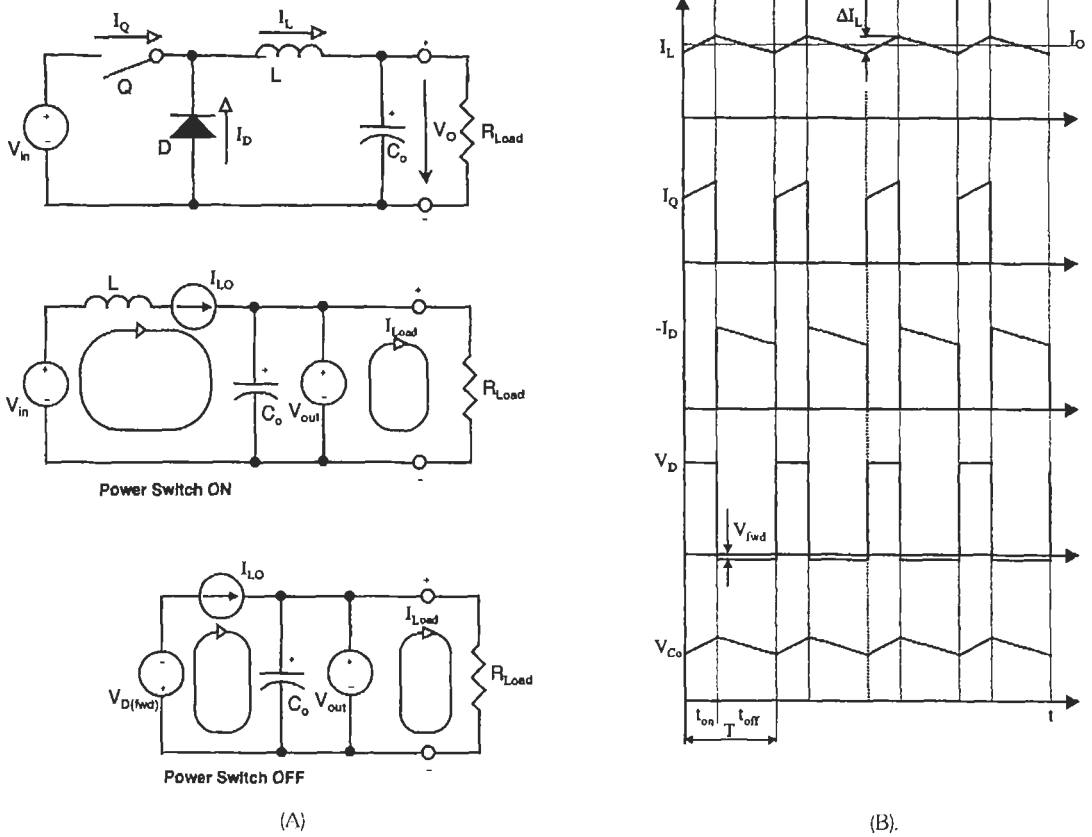


Fig. 29. Buck regulator principle of operation (A) and its waveforms (B).

relationship that approximately describes its operation is

$$V_{out} = d \times V_{in}$$

Here d is the duty cycle, $d = t_{on}/(t_{on} + t_{off}) = t_{on}/T$, where T is the period of one switching cycle, or $T = 1/f_{sw}$, and f_{sw} is the frequency of operation.

The buck converter is capable of kilowatts of output power. However, it suffers from one serious shortcoming: If the power switch were to fail short-circuited, the input power source would be connected directly to the load circuitry, which usually produces catastrophic results. The buck converter should only be used for board-level regulation. The most common buck derived transformer isolated converter topologies are: forward converter, two-transistor forward converter, push-pull converter, half-bridge and full-bridge converters.

Flyback or Boost-Mode Converters Basics

The most elementary flyback or boost-mode converter is the *boost* or *step-up converter*. Its schematic can be seen in Fig. 30. Its operation can also be broken into two distinct parts, one where the power switch is

on and another when it is off. When the power switch turns on, the input voltage source is placed directly across the inductor. This causes the current to begin linearly ramping upwards and is described by

$$i_{L(on)} = \frac{V_{in} \times t_{on}}{L}$$

where t_{on} is the on-time of the switching cycle when the power switch Q is on.

Once again, energy is being stored within the inductor core material. To maintain its steady-state operation the amount of energy stored during each switching cycle times the frequency of operation must be higher than the power demands of the load, or

$$P_{sto} = \frac{1}{2} LI_{pk}^2 \times f_{sw} > P_{out}$$

The power switch then turns off and the inductor voltage flies back above the input voltage and is clamped by the rectifier at the output voltage. The current then begins to linearly ramp downward until the new on-time begins. Its waveform, which is shown in Fig. 30, is determined by

$$i_{L(off)} = \frac{(V_{out} - V_{in}) \times t_{off}}{L}$$

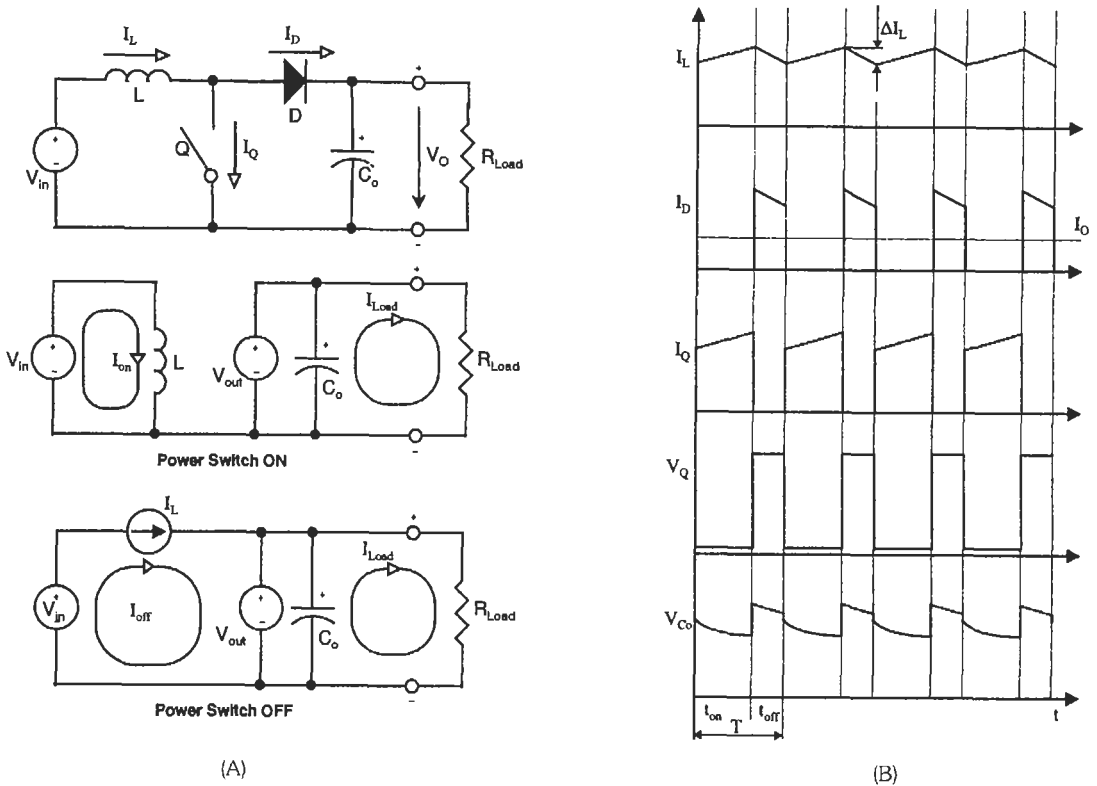


Fig. 30. Boost regulator principle of operation (A) and its waveforms (B).

where t_{off} is the on-time of the switching cycle when the power switch Q is off.

The operation of the boost converter can be approximately described by the following relation:

$$V_{in} V_{out} = (1-d)$$

Here d is the duty cycle, $d = t_{on} / (t_{on} + t_{off}) = t_{on} / T$, where T is the period of one switching cycle, or $T = 1 / f_{sw}$, and f_{sw} is the frequency of operation.

The boost converter should also be only used for board-level regulation.

Table 3 shows the most common switching converter topologies, their schematic diagrams, transfer functions, device rating equations, circuit waveforms, and advantages and disadvantages.

Power Semiconductors in Switch-Mode Power Supplies

The power semiconductors used in switch-mode power supplies are thoroughly described in the previous sections. Therefore the following paragraphs sum-

marize their properties in relation to the switch-mode power supplies.

Power MOSFETs

Power MOSFETs are very popular for use as power switches within switch-mode power supplies. MOSFETs have some advantages over the bipolar transistor, such as switching five to ten times faster than bipolar transistors and being easier to drive and use. The drive source, however, must be a well-bypassed low-impedance voltage source. This is because the gate of a MOSFET resembles a capacitor, which must be charged and discharged. Fig. 31 shows several power MOSFET gate drive circuits.

Bipolar Power Transistors

Bipolar power transistors offer several advantages over power MOSFETs: They have higher breakdown voltages, and they are somewhat less expensive for devices greater than 500 volts.

Bipolar power transistors are current-driven devices. That is, in order to have a current flowing from the col-

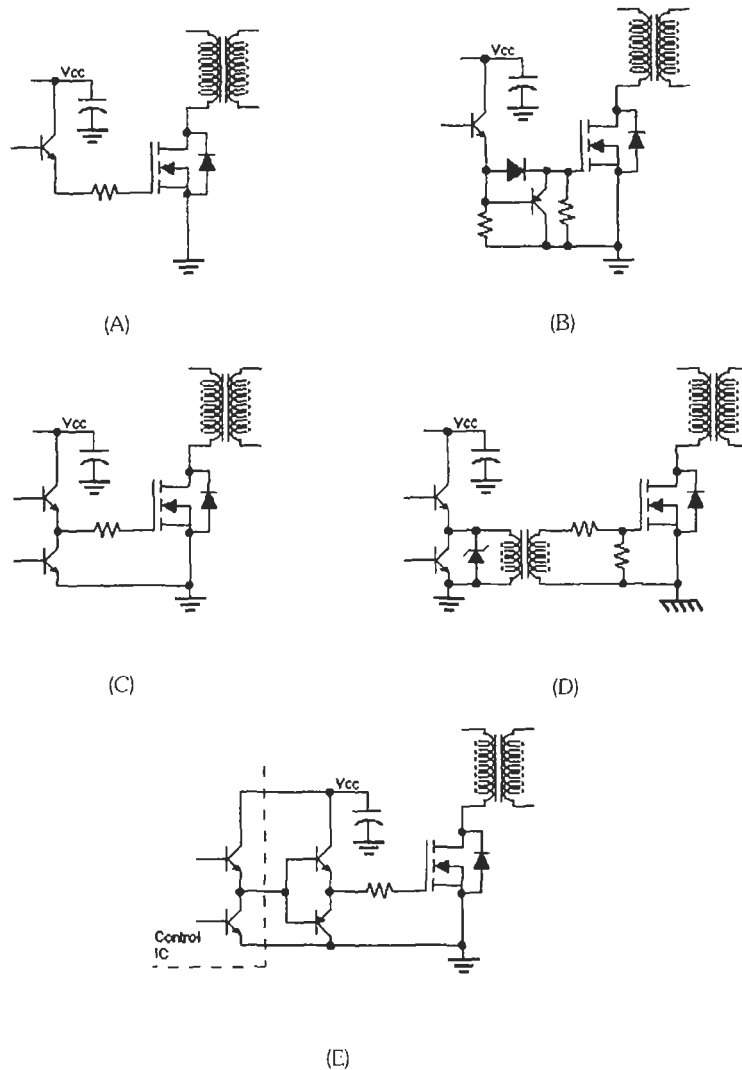


Fig. 31. Power MOSFET gate drive circuits: (A) passive base drive, (B) quasi-totem-pole, (C) totem-pole, (D) transformer-coupled totem-pole, (E) buffered totem-pole to drive MOSFETs with large gate-source capacitance.

lector to the emitter, one must drive current into the base. Within switch-mode power supplies, the objective is to have as small a collector-to-emitter voltage as possible when the transistor is "on." To do this the designer strives to have the transistor operate in, or close to, a saturated state. The saturation of bipolar transistors is defined as

$$I_B \geq \frac{I_{C(max)}}{h_{FE(min)}}$$

For the transistor to switch fast, the drive current should enter and exit the transistor fast. To accom-

plish this, special attention should be given to the drive circuitry.

Fig. 32. shows some common base drive circuits.

IGBTs

The insulated-gate bipolar transistor (IGBT) has gained popularity in the off-line, high-power switch-mode power supplies, because of its high voltage and current handling capabilities. It drives identically as a power MOSFET. It switches faster than the power bipolar transistor but still slower than a power MOSFET. The drive circuits are identical to those of power MOSFETs shown in Fig. 31.

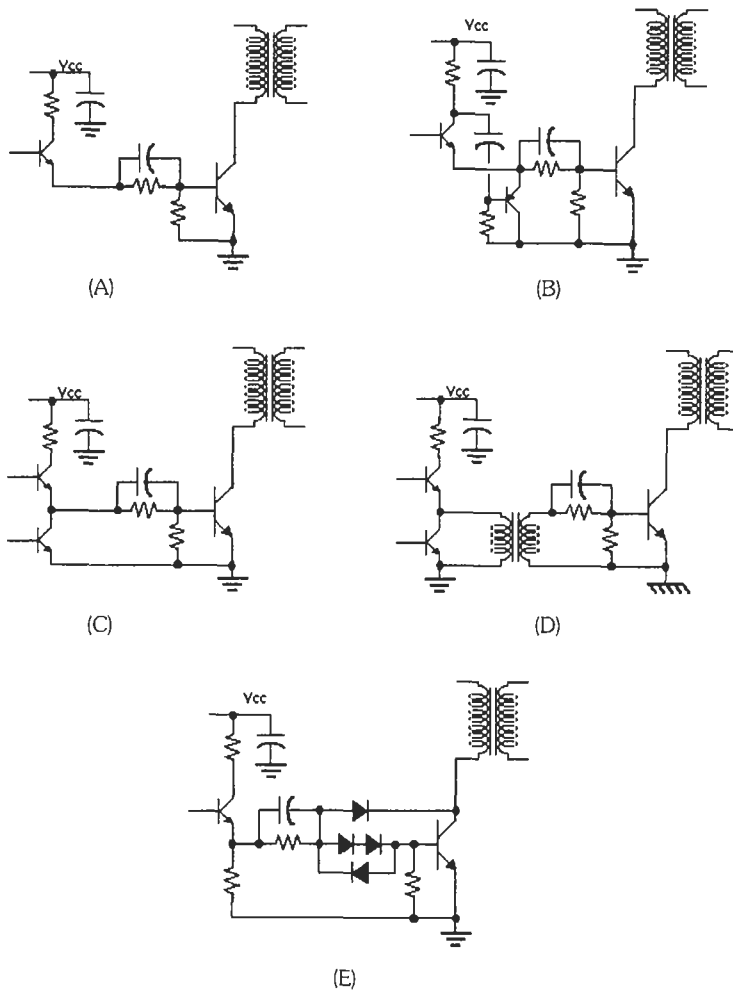


Fig. 32 Bipolar power transistor base drive circuits: (A) passive base drive, (B) quasi totem-pole, (C) totem-pole. (D) transformer-coupled totem-pole, (E) Baker clamp.

Rectifiers

Rectifiers are the largest source of loss within switch-mode power supplies. To choose the best rectifier, one must understand the parameters that affect their efficiency. The most important rectifier parameters are the *forward voltage drop* (V_f) and the *reverse recovery time* (t_{rr}). The reverse recovery loss is where the rectifier becomes reverse biased, and current appears to actually flow backward through the rectifier. It is actually the minority carriers being swept from the P-N junction. Nonetheless, it is a significant loss. This loss is minimized by selecting the rectifier with the shortest reverse recovery time (t_{rr}).

The summarized comparison of switching converter power semiconductors is given in Table 4.

Table 5 gives a comparison of power semiconductors in different switching converter applications.

Power Losses in Power Semiconductors

Conduction losses occur during the conducting period of the power semiconductor.

They can be roughly estimated as follows:

MOSFET:

$$P_{MOS(cond)} = I_{avg}^2 \times R_{DS(on)}$$

Since the $R_{DS(on)}$ increases with temperature, more conservative estimation can be made by doubling the $R_{DS(on)}$ value.

TABLE 4. COMPARISON OF VARIOUS TRANSISTORS AND RECTIFIERS

Transistors		Advantages	Disadvantages
Bipolars	Low $V_{CE(sat)}$ voltage		Low gain
	High voltage capability		Base drive is difficult
	Relatively inexpensive		Limited to approximately 75 kHz due to switching losses Nonsquare SOA*
MOSFETs	Easily driven		High voltage device
	Shares current well		Relatively high on-voltage
	Wide bandwidth		
	Low on-resistance for low voltage devices		
	Square SOA*		
IGBT	Easily driven		Existing devices only applicable up to 100 kHz. $T_f < 80$ ns
	High current capability		Future demands > 100 kHz and $T_f < 50$ ns
	Uses silicon more efficiently than MOSFETs		
	High voltage capability		
	Square SOA*		
Rectifiers		Advantages	Disadvantages
Schottky	Low V_f		200 V performance is marginal
	Very fast		High leakage
	No stored charge		High capacitance
Ultra-Fast	High voltage capability		Stored charge
	Low leakage		Very snappy reverse recovery

* SOA = safe operating area.

Bipolar and IGBT:

$$P_{BJT(cond)} = I_{avg} \times V_{CE(sat)}$$

Power rectifier:

$$P_{RECT} = I_{avg} \times V_f$$

Switching losses in power semiconductors occur during the transient turn-on and turn-off time. Switching losses are proportional to the switching frequency. Therefore to increase the switching frequency the turn-on and turn-off losses must be decreased.

Turn-on Losses

Fig. 33 shows the losses at turn-on. The rectifier reverse recovery effect puts stress upon the switching transistor. It is apparent that increasing di/dt in an attempt to reduce switching losses increases the severity of diode snap and current ringing; electromagnetic interference (EMI) is of great concern. The turn-on losses can be decreased by employing a fast power switch and a fast soft recovery rectifier.

Turn-off Losses

Turn-off losses are caused by the fact that the transistor current will not fall until an output rectifier is forward biased. It is preferable to turn-off as fast as the circuit parasitics allow, to reduce turn-off losses (see Fig. 34). Forward recovery of the output rectifier is very important.

Snubbers and Clamps

Snubber and clamp circuits are intended to reduce voltage spikes caused by parasitic elements of switching converters.

A *snubber* is used for reducing the level of a voltage spike and decreasing the rate of change of a voltage waveform. This has benefits in the Safe Operating Area (SOA) of semiconductors, and it lowers the spectral content of any conducted or radiated interference and thus less RF energy is radiated.

A *soft clamp* returns the spike energy to a medium-sized pulse capacitor. The capacitor is then discharged through the parallel resistor.

TABLE 5. COMPARISON OF POWER SEMICONDUCTORS IN DIFFERENT SWITCHING CONVERTER APPLICATIONS

Performance/ Topology	Bipolar < 80 kHz	MOSFET > 50 kHz	IGBT < 100 kHz	Schottky > 50 kHz	Ultra-Fast < 80 kHz
Buck or Step Down	N/A	Telecom Computer 3.3 V, 5 V out < 150 W	N/A	Telecom Computer 3.3 V, 5 V out < 150 W	N/A
Boost or Step Up	Seldom Used	PFC < 1000 W	PFC > 1000 W	Seldom Used	PFC
Flyback	TV/Monitors < 300 W	Computer < 150 W	TV/Monitors Computer < 300 W	Computer < 150 W	TV/Monitors Computer < 100 kHz < 150 W
Forward	Industrial < 500 W	Telecom Computer 5-48 V out < 500 W	Telecom Computer Industrial 5-48 V out < 500 W	Telecom Computer Industrial 5-48 V out < 500 W	Telecom Computer Industrial 48-400 V out < 500 W
Half & Full Bridge	Seldom Used	Telecom Computer Industrial 5-400 V out > 400 W	Telecom Computer Industrial 5-400 V out > 400 W	Telecom Computer Industrial 5-400 V out > 400 W	Telecom Computer Industrial 48-400 V out > 400 W

Note: N/A = not available; PFC = Power Factor Correction (see page 14-36).

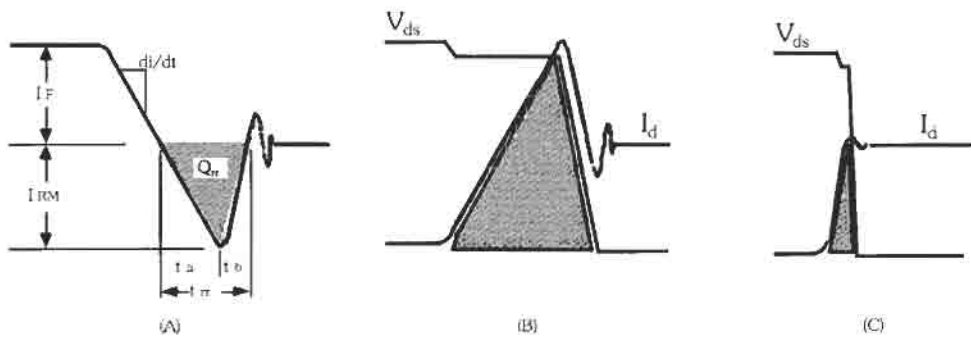


Fig. 33. The turn-on switching losses: (A) The rectifier reverse recovery effect. (B) The typical turn-on losses caused by slow di/dt and slow dv/dt . (C) Turn-on losses decreased by implementing a fast soft recovery rectifier.

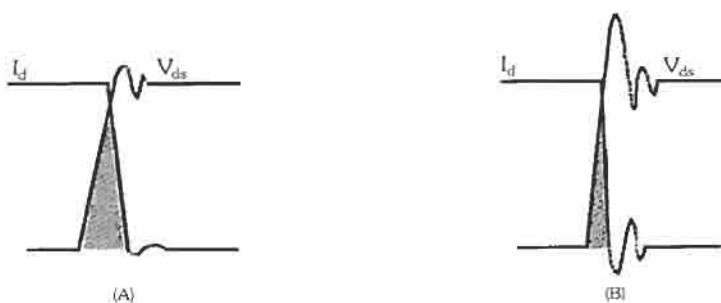


Fig. 34. (A) Losses during the turn-off. (B) Parasitic circuit inductances cause voltage overshoots at high di/dt .

A *hard clamp* is used only for reducing the level of a voltage spike. It has no effect on the dV/dt of the transition. Therefore it is not very useful for reducing EMI/RFI.* It is useful for preventing components such as semiconductors and capacitors from entering avalanche breakdown.

Typical snubber and clamp circuits are shown in Fig. 35. Their effects on a representative switching waveform are shown in Fig. 36.

* EMI—electromagnetic interference
RFI—radiated frequency interference

Switch-Mode Converter Control

Integrated circuits are available to control these various circuits properly. These signal-level integrated circuits contain reference voltages, voltage-error amplifiers, current-error amplifiers, linear timing ramps, error-voltage-to-pulse-width circuits, remote on-off, synchronization, and power-transistor base-drive signals. Only a few external components are required to set the basic operating frequency and voltage scaling to provide a complete control circuit. Isolation between the output voltage and the input voltage can be provided with optical couplers. Isolation of the

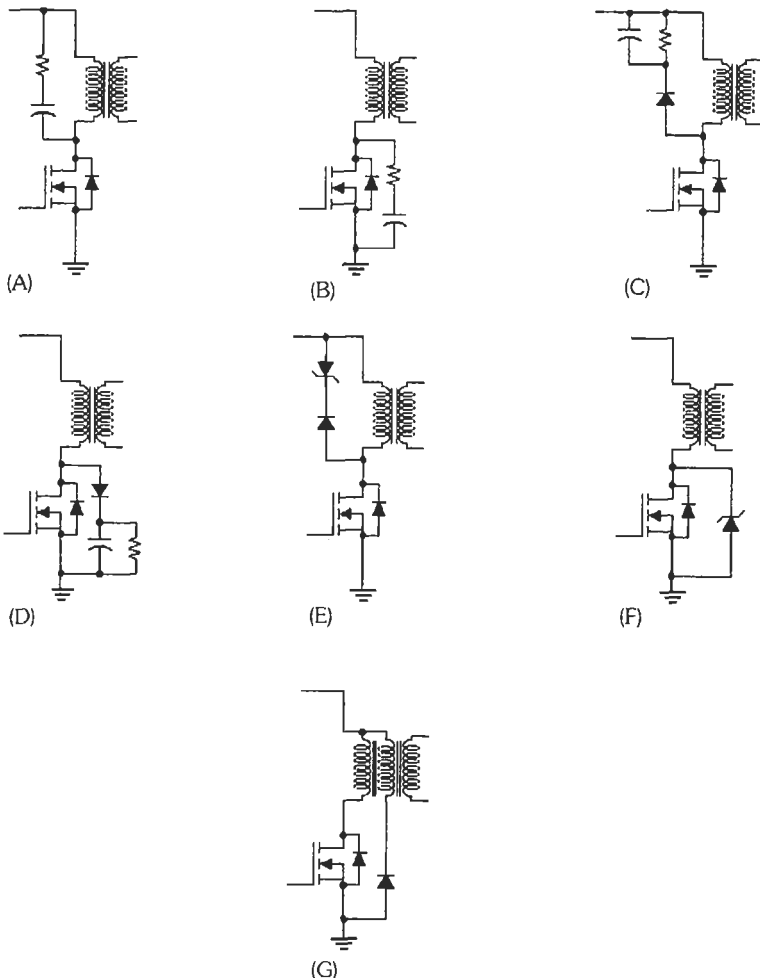


Fig. 35. Typical snubber and clamp circuits. (A) Snubber across transformer—parallel snubbing. (B) Snubber across switch—series snubbing. (C) Soft clamp across transformer. (D) Soft clamp across switch. (E) Zener clamp across transformer. (F) Zener clamp across switch. (G) Clamp winding.

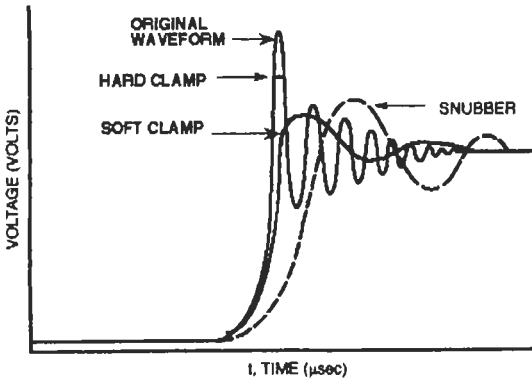


Fig. 36. The effect of snubbers and clamps upon a representative switching waveform.

power-transistor base-drive signals can be accomplished by pulse transformers with energy-storage capacitors to provide a reverse pulse of base current for fast turn-off of the power transistor or field-effect transistor. Refer to Figs. 31 and 32 for different drivers for bipolar switch-mode transistors and switch-mode power MOSFETs.

Different switching converters utilize different types of error voltage-to-switching pulse modulation. Refer to Fig. 37, which describes a basic switching converter control system.

In a Pulse-Width-Modulated (PWM) controller, the switching frequency is constant. The dc output voltage is determined by the ratio of pulse width to switching period.

In a Variable-Frequency-Modulated (VFM) controller, the pulse width is constant. The Off time will

change with load and input voltage changes. Therefore, the switching period is *not* constant. It is the sum of On time and Off time. The output voltage will be the ratio of pulse width to the sum of On time and Off time.

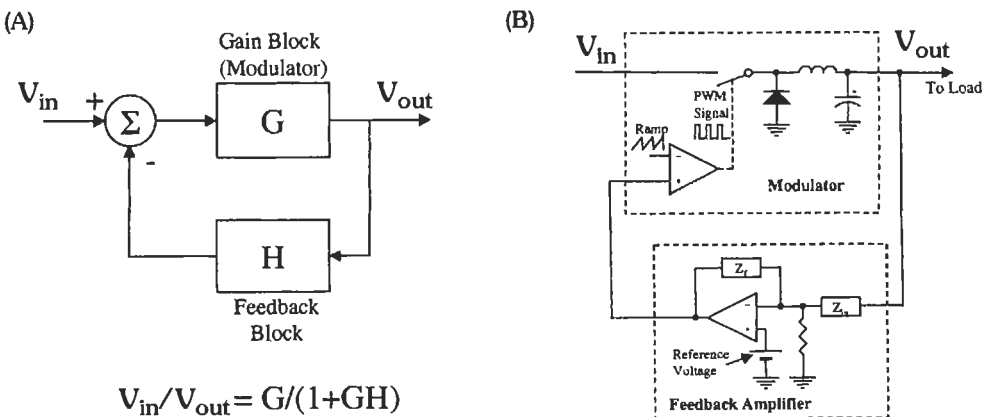
In a Pulse-Frequency-Modulated (PFM) converter, also called Pulse-Skipping or Cycle-Skipping Modulation or Ripple Mode converter, the switching frequency. On-time and Off-time are usually all variable with line and voltage changes and have relaxed stability criterion. At very light loads, its cycle-skipping characteristic consumes less power than the other methods while maintaining a regulated output voltage. However, the ripple on the output voltage is slightly higher than with PWM or VFM converters.

In noise-sensitive applications, the PWM controller is usually preferred because the switching frequency is constant and easier to filter out. PWM controllers are also preferred in systems using more than one switching converter in the equipment.

Voltage-Mode Control

Voltage-mode control is a single-loop control method. Only the output voltage of the converter is sensed and compared to the reference voltage in order to control the duty cycle of the power switch. Voltage-mode control provides fairly good control for the power supply but has a poor load transient response and line regulation. Another drawback is that it cannot protect against instantaneous overcurrent conditions in the power switch caused by transformer core saturation.

The voltage-mode control can be easily recognized by a comparator comparing the sawtooth ramp signal from the oscillator with the feedback signal from the error amplifier, thus creating a pulse-width modulated signal for the power switch driver. The typical voltage-mode controller is shown in Fig. 38.



$$V_{in}/V_{out} = G/(1+GH)$$

Fig. 37. Basic switching converter control system (A) and its principal block diagram (B).

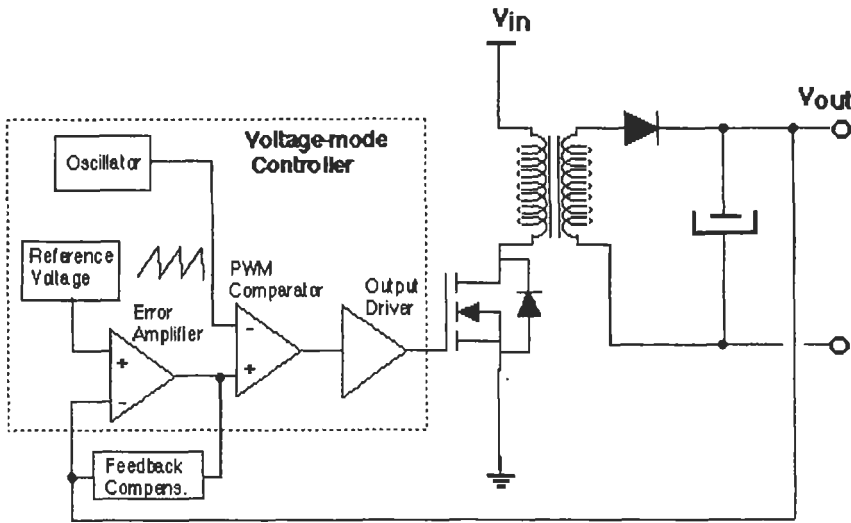


Fig. 38. Basic voltage-mode control.

Because of its simplicity this control method is still very popular among the power supply designers as well as among the manufacturers of switch-mode power supply (SMPS) control integrated circuits.

Current-Mode Control

Current-mode control adds a second control loop to the voltage feedback loop. Now, instead of the ramp from the oscillator, as in the voltage mode control, the current ramp created by magnetic elements is used for error-voltage-to-PWM conversion. The oscillator only serves to fix the frequency of operation. The basic current-mode controller diagram is shown in Fig. 39. This method of control offers significant improvements in terms of line regulation and load transient response as well as inherent cycle-by-cycle overcurrent protection.

Resonant Transistor Circuits

Unavoidable switching losses occur during each turn-on and turn-off switching event in a pulse-width-modulated power supply. In an effort to reduce the size and weight of power supplies, the frequency of operation has been increased up to the point that the switching losses are increasingly important. In an effort to reduce these switching losses, resonant circuits have been developed. Two general categories are: zero-voltage switching (ZVS), where the voltage across the switch is brought to zero before the switch is turned on, and zero-current switching (ZCS), where the power semiconductor does not switch unless the current through the switch is zero. With the switching losses greatly reduced, the frequency of operation can be increased up to, at present, 1 to 10 Mhz. The reso-

nant operation requires inductors and capacitors to set the resonant frequency. If the frequency is high enough, the capacitances of the semiconductor junctions and the leakage reactance of the transformers can be used as the resonant elements. The actual operating frequency can be higher or lower than the LC resonant frequency, depending on the characteristics of the power switches. Rather than use pulse-width modulation as a means for controlling the output voltage, many resonant converters vary the frequency of operation to control their output voltage. This area of the technology is the subject of ongoing research. Chapter 7 of Reference 2 provides additional discussion of resonant converters. Fig. 40 is a block diagram of a resonant converter configured as a dc-to-dc power supply.

Power Factor Correction Circuits

As stated in previous paragraphs, switch-mode power supplies, motor control frequency inverters, and lighting electronics pull current in narrow, high-amplitude pulses rather than sinusoidally from the ac mains. In switch-mode power supplies the problem lies in the input rectification and filter network. The typical input circuit and its associated waveforms are shown in Fig. 41. The input rectifiers can only conduct current when the ac line voltage exceeds the voltage on the bulk input filter capacitor. This typically occurs within 15 degrees of the crest of the ac voltage waveform. The resulting current pulses can be 5–10 times higher than the expected average current draw. This also can lead to distortion of the ac voltage waveform and an imbalance of the three-phase power lines feeding the circuits.

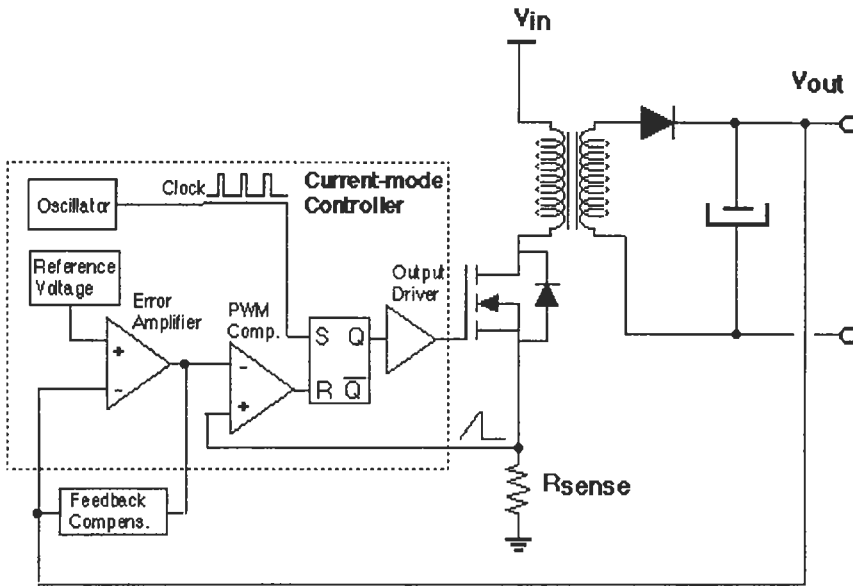


Fig. 39. Basic current-mode control.

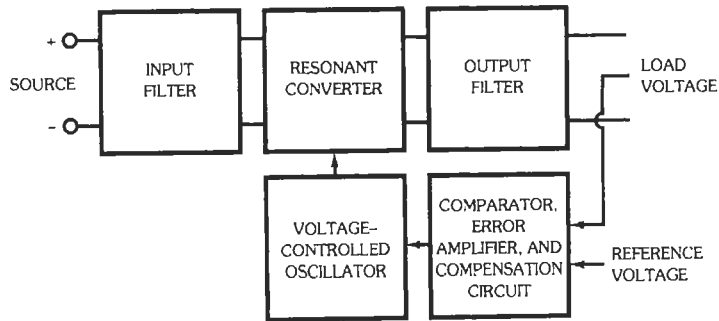


Fig. 40. Block diagram of resonant-mode power supply.

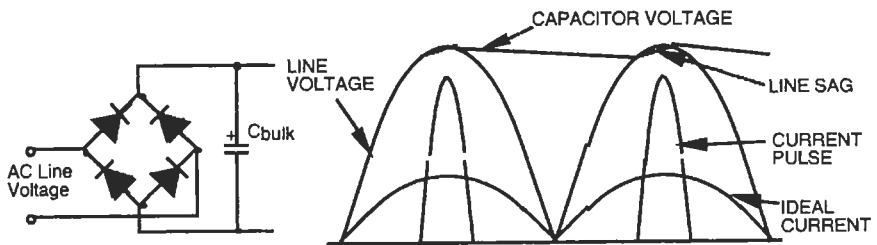


Fig. 41. Capacitive input filter and its waveforms.

Power Factor Definition

Ideally the ac current provided by local utilities should be purely sinusoidal and free of surges, transients, and electrical noise (see Fig. 42). The power factor (*PF*) is defined by the ratio of the real power P_R to the total apparent power P_A :

$$PF = \frac{P_R}{P_A} \tag{Eq. 2}$$

The real power P_R can be calculated as the average of the instantaneous product of voltage and current taken at each instance in time over the full period, T :

$$P_R = \frac{1}{T} \int_0^T e \cdot i dt \tag{Eq. 3}$$

where e = voltage and i = current. The apparent power, a sum of the real power and reactive power, is calculated as the product of the rms current measured over the full period and the rms voltage measured over the same full period:

$$P_A = V_{RMS} \cdot I_{RMS} \tag{Eq. 4}$$

In the case of an uncorrected switch-mode power supply, the current and voltage can be nearly in phase as the result of the capacitor charging near the peak of the voltage waveform (see Fig. 43(A)). Therefore using the definition $PF = \cos\phi$ would misleadingly yield a *PF* of 1.0 when the *PF* may be very low.

Probably the clearest definition of the power factor can be made by following way. For sinusoidal input voltage:

$$PF = \frac{P_R}{P_A} = \frac{P_{INPUT}}{V_{rms} \cdot I_{rms}} = \frac{V_1 I_1 \cos\phi}{V_1 \cdot I_{i,rms}} \tag{Eq. 5}$$

where

$$\frac{I_1}{I_{i,rms}} = DF \text{ is the distortion factor,}$$

I_1 is the first harmonic of the input current,
 $I_{i,rms}$ is the total harmonic current, and
 $\cos\phi$ is the displacement factor.

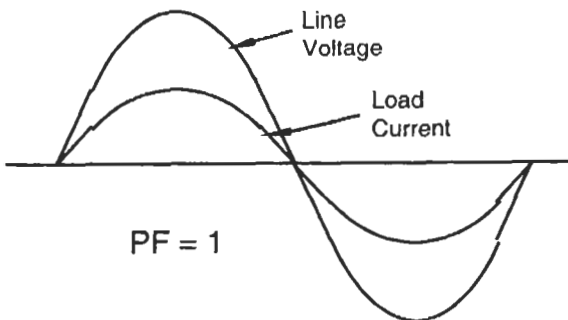


Fig. 42. The ac input current is sinusoidal and in phase with the voltage waveform. Power factor = 1.0

Then the distortion factor *DF* can be defined as follows:

$$DF = \frac{1}{\sqrt{1 + (THD)^2}} \tag{Eq. 6}$$

where the total harmonic distortion *THD* equals

$$THD = \frac{\sqrt{I_{i,rms}^2 - I_1^2}}{I_1} \tag{Eq. 7}$$

The situation is even more clearly seen in Figs. 43(A) and 43(B). In both cases the power factor does not equal 1.0.

International Regulatory Standards

The harmonic-current problems addresses the international community in the IEC 1000-3-2 document (the successor of the IEC 555-2), which establishes requirements that limit harmonic-current injection into the public supply system. The IEC 1000-3-2 addresses equipment with nominal line-to-neutral voltages of 220 V and higher. The aim of the standard is to prevent power-control devices from generating low-frequency harmonics.

IEC 1000-3-2 categorizes devices into four classes. Each class has unique limits. Class A covers balanced three-phase equipment and all others except those defined in specialized classes B, C, and D. Class B covers portable tools, Class C sets the limits for lighting equipment. Those devices that have a special "wave shape" (see Fig. 44) and are not motor driven or do not fit into Class A, B, or C are considered Class D.

Active Power Factor Correction

To improve power quality the mains power source should see a sinusoidal load current that appears resistive. Therefore an ideal power factor correction stage should appear as a resistive load to the ac line, pulling current sinusoidally and storing the power for the main converter so that the regulated voltage exists while the

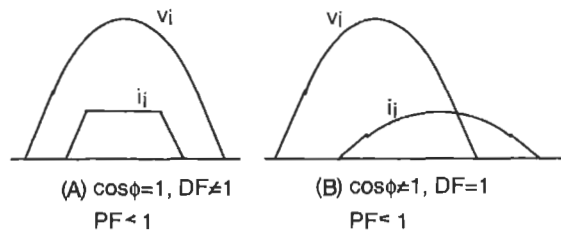


Fig. 43. Examples of non-unity power factor.

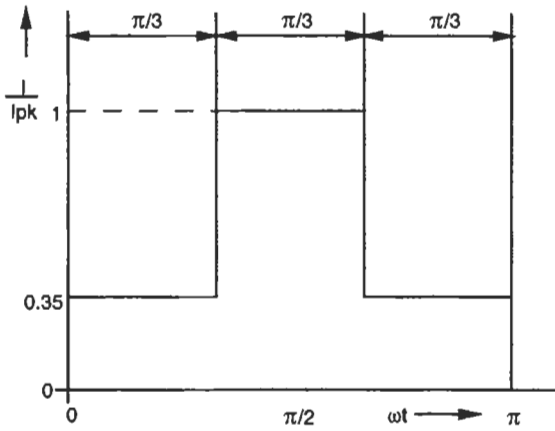


Fig. 44. Class D defines a "special wave shape" of the envelope of the input current.

line voltage sinusoidally varies. This can be approached actively or passively.

Active power factor correction circuits usually take the form of a modified boost converter (refer to Table 3, Column B, and Fig. 45). Three of these circuits can be arranged in a bridge circuit to shape the currents drawn from a multiphase ac supply.

There are several main implementations in the control of power factor correction (PFC) circuits: the fixed-frequency average current sensing method widely implemented by Unitorde under the Pioneer Magnetics patent (UC3854 control integrated circuit); the fixed-frequency, peak current sensing method by Micro Linear (ML4812); hysteric control used by Cherry Semiconductor (CS3810); and the borderline or critical conduction, peak current sensing method used by Siemens (TDA4818), Silicon General (SG3561), and Motorola (MC34262, MC33368). Fig. 46 shows current waveforms of two typical PFC control techniques.

FILTER CIRCUITS AND DESIGN

Rectifier filters may be classified into three types:

- (1) *Inductor Input* (Fig. 47): Have good voltage regulation, high transformer utilization factor, and low rectifier peak currents, but also give relatively low output voltage.
- (2) *Capacitor Input* (Fig. 48): Have high output voltage, but poor regulation, poor transformer utilization factor, and high peak currents. Used mostly in television and radio receivers.
- (3) *Resistor Input* (Fig. 49): Used for low-current and low-power applications.

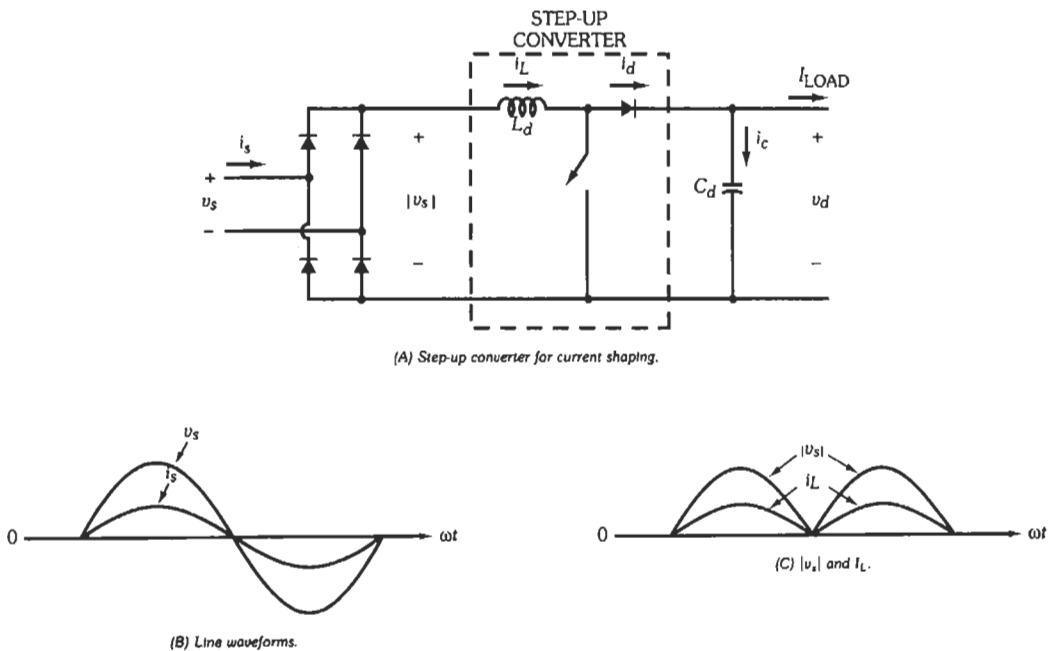


Fig. 45. Active power factor correction and harmonic filtering. (From Mohan, N., et al. *Power Electronics—Converters, Applications and Design*. New York: John Wiley & Sons, Inc., 1989.)

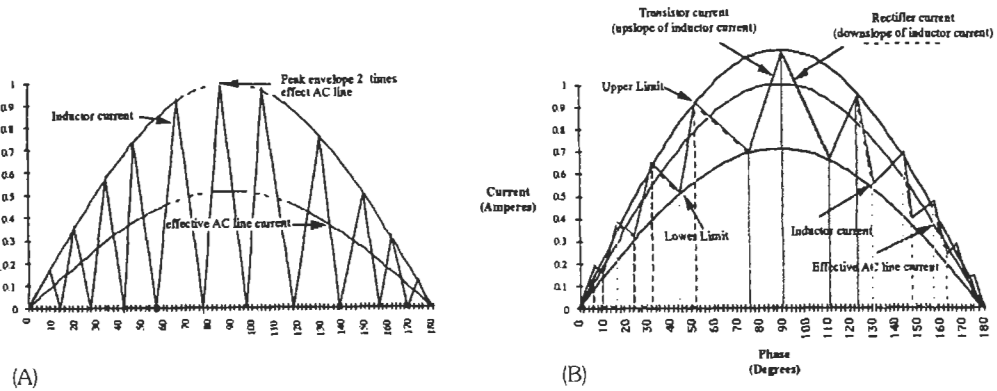


Fig. 46. (A) Critical conduction mode current waveform. (B) Hysteretic control current waveform.

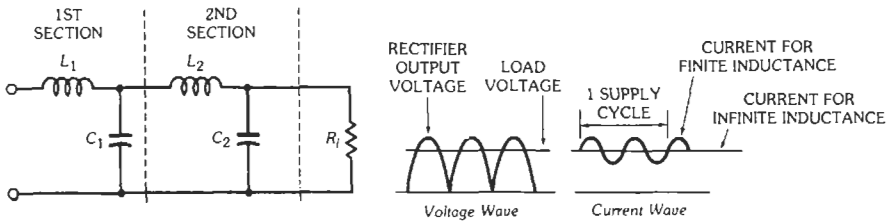


Fig. 47. Inductor-input filter.

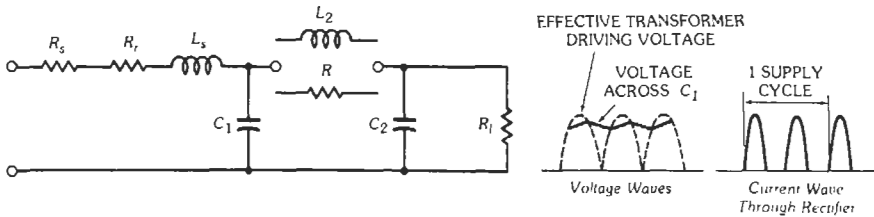


Fig. 48. Capacitor-input filter. C_1 is the input capacitor. $R_s = 1/2 \times$ (secondary-winding resistance). $L_s =$ leakage inductance viewed from $1/2$ secondary winding. $R_r =$ equivalent resistance of IR drop in rectifier element.

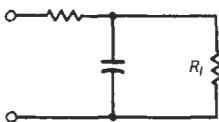


Fig. 49. Resistor-input filter.

Inductor-Input Filter Design

The constants of Fig. 47 are determined from the following considerations:

(A) There must be sufficient inductance to ensure continuous operation of rectifiers and good voltage regulation. When this critical value of inductance is

increased by a 25-percent safety factor, the minimum value becomes

$$L_{\min} = (K/f_s)R_l \text{ henry} \quad (\text{Eq. 8})$$

where

- $f_s =$ frequency of source in hertz.
- $R_l =$ maximum total load resistance in ohms,
- $K = 0.060$ for full-wave single-phase circuits
- $= 0.0057$ for full-wave two-phase circuits
- $= 0.0017$ for full-wave three-phase circuits.

At 60 hertz, single-phase full-wave

$$L_{\min} = R_l/1000 \text{ henry} \quad (\text{Eq. 8a})$$

(B) The LC product must exceed a certain minimum, to ensure a required ripple voltage factor

$$\begin{aligned} r &= E_r/E_{dc} \\ &= \left[\sqrt{2}/(p^2 - 1) \right] \left[10^6 / (2\pi f_s p)^2 L_1 C_1 \right] \\ &= K'/L_1 C_1 \end{aligned} \quad (\text{Eq. 9})$$

where, except for single-phase half-wave,

- p = effective number of phascs of rectifier,
- E_r = root-mean-square ripple voltage appearing across C_1 ,
- E_{dc} = direct-current voltage on C_1 ,
- L_1 is in henrys,
- C_1 is in microfarads.

For single-phase full-wave, $p = 2$ and

$$r = (0.83/L_1 C_1) (60/f_s)^2 \quad (\text{Eq. 9a})$$

For three-phase full-wave, $p = 6$ and

$$r = (0.0079/L_1 C_1) (60/f_s)^2 \quad (\text{Eq. 9b})$$

Equations (8) and (9) define the constants L_1 and C_1 of the filter, in terms of load resistor R_l and allowable ripple voltage factor r .

Swinging Chokes—Swinging chokes have inductances that vary with the load current. When the load resistance varies through a wide range, a swinging choke, with a bleeder resistor R_b (10 000 to 20 000 ohms) connected across the filter output, is used to guarantee efficient operation; i.e., $L_{min} = R_l'/1000$ for all loads, where $R_l' = (R_l R_b)/(R_l + R_b)$. Swinging chokes are economical because of their smaller relative size and result in adequate filtering in many cases.

Two-Section Filters—For further reduction of ripple voltage E_{r1} , a smoothing section (Fig. 46) may be added and will result in output ripple voltage E_{r2} .

$$E_{r2}/E_{r1} \approx 10^6 / (2\pi f_r)^2 L_2 C_2 \quad (\text{Eq. 10})$$

where f_r = ripple frequency.

Capacitor-Input Filter Design

The constants of the input capacitor (Fig. 47) are determined from the following:

(A) Degree of filtering required.

$$\begin{aligned} r &= E_r/E_{dc} \\ &= \sqrt{2}/2\pi f_s C_1 R_l \\ &= (0.00188/C_1 R_l) (120/f_s) \end{aligned} \quad (\text{Eq. 11})$$

where $C_1 R_l$ is in microfarads \times megohms or farads \times ohms.

(B) A maximum allowable C_1 (so as not to exceed the maximum allowable peak-current rating of the rectifier).

Unlike the inductor-input filter, the source impedance (transformer and rectifier) affects output dc and ripple voltages and the peak currents. The equivalent network is shown in Fig. 48.

Neglecting leakage inductance, the peak output ripple voltage E_{r1} (across the capacitor) and the peak rectifier current for varying effective load resistance are given in Fig. 50. If the load current is small, there may be no need to add the L-section consisting of an inductor and a second capacitor. Otherwise, with the completion of an $L_2 C_2$ or RC_2 section (Fig. 48), greater filtering is obtained, the peak output-ripple voltage E_{r2} being given by Eq. 10 or by

$$E_{r2}/E_{r1} = 1/2\pi f_r R C_2 \quad (\text{Eq. 12})$$

Harmonic Currents and Power Factor

In the case of the inductor-, resistor- and capacitor-input filters, the current drawn from the ac supply is nonsinusoidal. Since the definition of power factor includes a term associated with harmonic distortion, the power factor of all of the circuits is less than unity. The power factor given for the various rectifier circuits (Table 1) assumes that the inductance of the filter choke is infinite. For finite values of inductance, the power factor is decreased. With a three-phase supply, reducing the inductance of the filter reactor reduces the power factor and increases harmonic currents. The power factor as a function of filter inductance for a single-phase rectifier is shown in Fig. 51. The power factor as a function of filter inductance for a three-phase rectifier is shown in Fig. 52. The asymptote of 0.955 is the value given in Table 1, column 5. Standards are currently being prepared in Europe to limit the harmonic currents that can be drawn from the supply network. One approach to meeting these standards is to provide a pulse-width-modulated rectifier (Fig. 45).

Phase-Controlled-Thyristor Filter Design

The same general principles apply to filters for controlled rectifier circuits as to those for ordinary rectifier circuits. Capacitive-input filters, however, are rarely used, since a capacitive load restricts the range of conduction-angle control that it is possible to obtain with thyristors and results in large values of current at the instant of thyristor turn-on.

The two main differences to consider in filter design for controlled rectifiers, compared with that for ordinary rectifiers, are (A) the greatly increased values of critical inductance for the input choke and (B) the larger input ripple voltages.

Empirical equations for the critical inductance L_{min} were given that apply to various configurations of diode rectifiers. In controlled rectifiers, L_{min} rises as the conduction angle is decreased by gate control, i.e.,

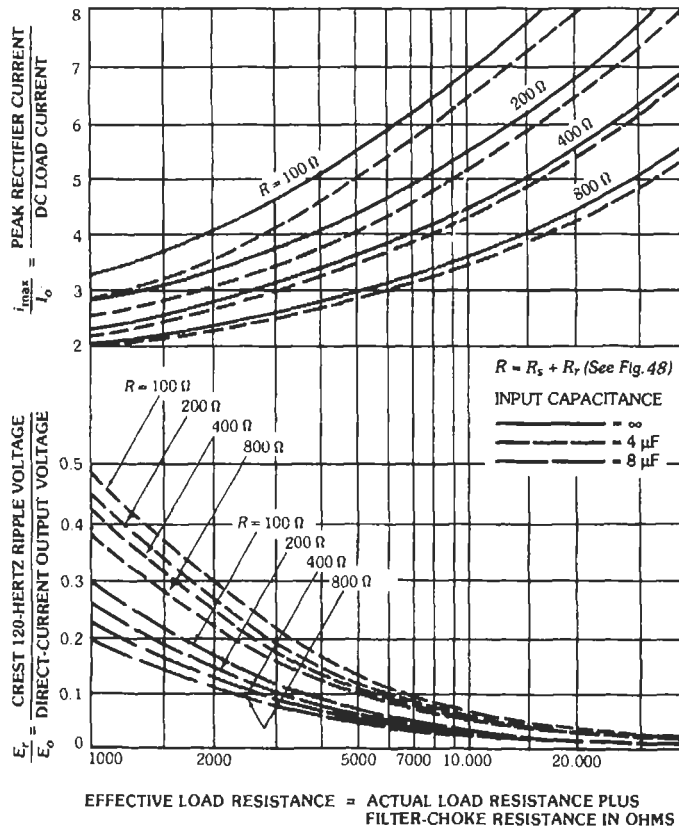


Fig. 50. Performance of capacitor-input filter for 60-hertz full-wave rectifier, assuming negligible leakage-inductance effect. (Adapted from *Radio Engineers Handbook*, by F. E. Terman, 1st ed., p. 603; 1943. By permission McGraw-Hill Book Co., New York.)

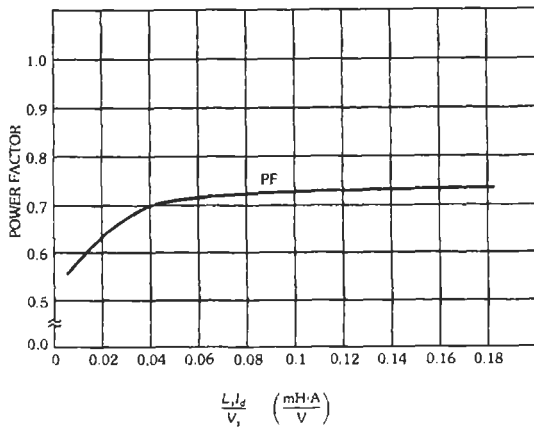


Fig. 51. Effect of L_s on power factor for a single-phase rectifier. L_s is in millihenries, I_d is the average current, and V_s is the rms voltage. (From Mohan, N., et al. *Power Electronics—Converters, Applications and Design*. ©1989. Reprinted by permission of John Wiley & Sons, Inc.)

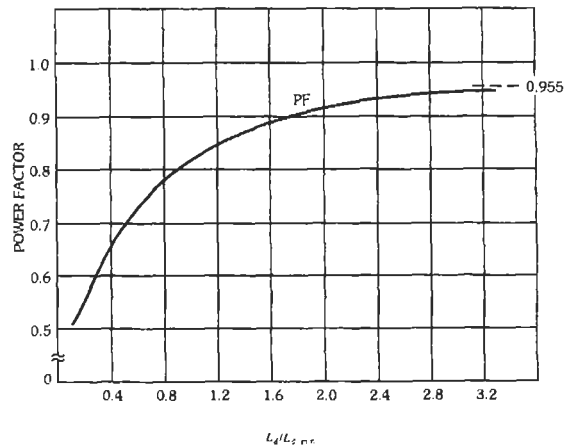


Fig. 52. Effect of L_d on power factor for three-phase rectifier. (From Mohan, N., et al. *Power Electronics—Converters, Applications and Design*. ©1989. Reprinted by permission of John Wiley & Sons, Inc.)

as the firing-angle delay, α , is increased. Fig. 53 shows the ratio L_{min}/R as a function of the percentage of maximum output voltage obtained from half- and full-controlled rectifier types discussed in the section on thyristors. Scales are included to show the corresponding values of α for various percentage outputs.

The ripple from a controlled rectifier is generally larger than that obtained from a conventional diode rectifier, and its value varies with the firing angle, α . Reference 3 provides graphs of the ripple voltage as a function of firing angle for the various single- and three-phase controlled rectifier circuits.

Resonant-Filter Design

Resonant filters are designed to reduce a specific frequency component of the waveform. A parallel-resonant section is placed in series with the load, and a series-resonant section is placed in parallel with the

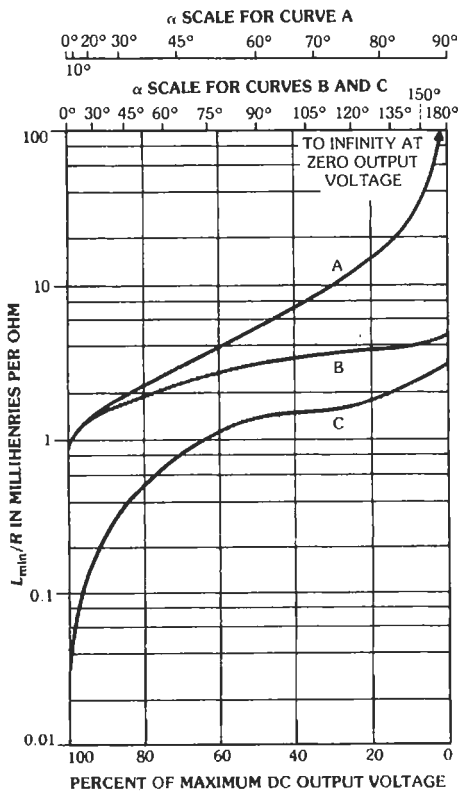


Fig. 53. L_{min}/R as a function of the percentage of maximum dc output voltage for thyristor rectifiers. The values of L_{min}/R given by the curves apply to a supply frequency of 60 hertz. Curve A is for full-controlled single-phase rectifiers such as those in Fig. 22 and in Fig. 24 with the bypass diode omitted. Curve B is for half-controlled single-phase rectifiers such as those in Fig. 23 and in Fig. 24 with the bypass diode included. Curve C is for a half-controlled three-phase rectifier such as that in Fig. 26.

load. Fig. 54 shows a resonant filter. Both L_1 and C_1 and L_2 and C_2 are selected to resonate at the unwanted frequency. These systems are applied to power-system filters where additional losses cannot be tolerated.

DC-OUTPUT POWER SUPPLIES

The previous sections described the various types of semiconductor devices, ac-to-dc converter circuits, and filter circuits. This section assembles these subsystems into complete dc power supplies. Block diagrams illustrate the various approaches to dc power supplies.

Unregulated Power Supplies

The most common type of unregulated power supply is the ac-to-dc rectifier with uncontrolled rectifiers followed by a suitable filter. Polyphase power produces less dc ripple voltage, less input harmonic current, and reduced dc voltage regulation. A power transformer allows for voltage matching between the incoming ac and the dc load voltage. If the load requires high-voltage dc, then a voltage-multiplier circuit can be chosen, since the design of the step-up transformer is more difficult with a large step-up ratio and severe insulation requirements. Load-voltage adjustment can be provided by transformer taps or an adjustable-voltage transformer.

Linear Regulated Power Supplies

A block diagram of a low-voltage linear regulated power supply is shown in Fig. 55. A utility-frequency transformer provides voltage matching between the incoming ac voltage and the desired load voltage. A diode rectifier converts the secondary ac voltage to dc voltage. A filter provides coarse filtering and fills in the gaps in the dc voltage caused by the rectification of single-phase power. The series transistor, operating in its linear region, absorbs the difference between its input voltage and the load voltage. The load voltage is sensed and fed back to the closed-loop voltage regulator. Load current limiting can also be provided to protect the series pass transistor from excessive power

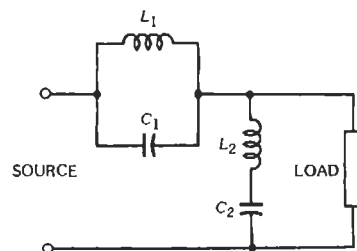


Fig. 54. Resonant filter circuit.

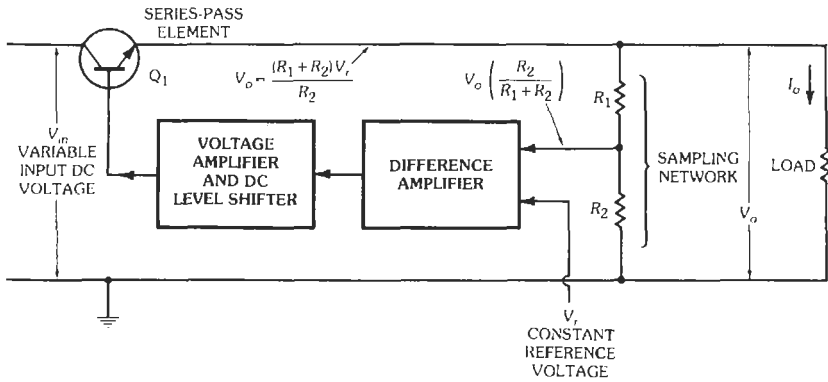


Fig. 55. Basic series-pass voltage regulator. Transistor Q_1 is an electronically controlled variable resistance in series with the load. (From Pressman, A. I. *Switching and Linear Power Supply, Power Converter Design*. Rochelle Park, NJ: Hayden Book Co., Inc., 1997.)

dissipation. Integrated and hybrid control circuits provide these and other functions in a single package. Instead of sensing output voltage, the output current can be sensed, and the system can provide a constant current to a variable load. A diode is sometimes provided to prevent load-voltage reversals caused by inductive energy in the load. Linear power supplies are used at low power (less than 50–100 watts); where multiple output voltages must be regulated; and where low noise, low ripple voltage, and fast response are essential.

Phase-Control Power Supplies

Two basic dc power-supply systems are configured with phase-controlled thyristors either on the primary or the secondary side of the isolation transformer. Figs. 56 and 57 show these two versions. The choice between the two systems is dependent on the cost of

thyristors as a function of their voltage and current ratings. The usual configuration is for the ac input to be high voltage and the dc load to be low voltage and high current. Placing the thyristors in the primary allows high-voltage, low-current thyristors to provide the control of the output voltage. If the thyristors are located on the low-voltage side of the transformer, they must be rated at high current and low voltage. The efficiencies of the two systems are comparable, as the double power conversion of one system balances the losses associated with a high-current, low-voltage phase-controlled rectifier. Control of the transformer primary voltage requires that the gating signals be symmetrical to avoid transformer saturation.

Systems without isolation transformers are used to supply adjustable dc voltage to dc motor armatures and field supplies. Isolation transformers are provided in general-purpose laboratory power supplies in the range of 1–50 kW.

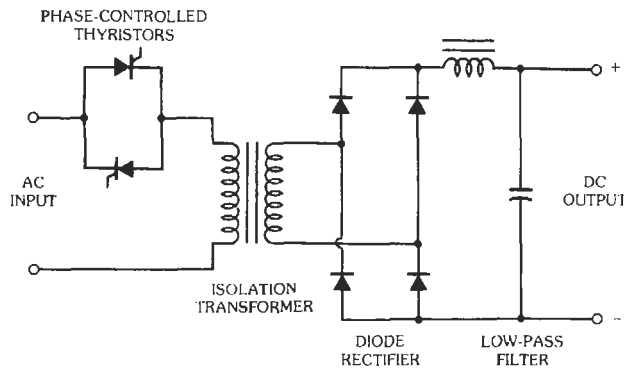


Fig. 56. Adjustable dc power supply with primary voltage control.

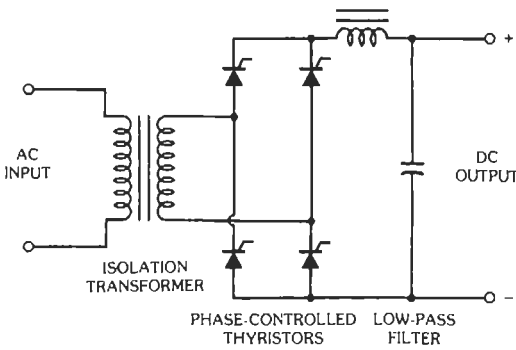


Fig. 57. Adjustable dc power supply with secondary voltage control.

Switching Power Supplies

Switching type power supplies are currently used for general-purpose dc power supplies in the 100–3000-watt range. It is expected that this range will increase in the future as the result of the technical features and cost reductions associated with this approach. Table 3 provides the design equations and device requirements for the most common types of dc-to-dc converters. To achieve small size and light weight, a 60-hertz isolation transformer is not provided. An uncontrolled diode rectifier and small capacitor filter provide an approximately 130-volt dc supply from 115-volt single-phase ac power. At light load, this voltage can increase to approximately 170–180 volts. A block diagram of this approach is shown in Fig. 58. A high-frequency pulse-width-modulated chopper or inverter is provided to regulate the voltage applied to the high-frequency transformer. The high-frequency transformer (20–500 kHz) provides voltage isolation and voltage matching at low weight, small size, and high efficiency. An uncontrolled ac-to-dc rectifier, at low voltages using Schotky diodes, converts from ac

to dc. A low-pass filter removes the high-frequency (40–1000 kHz) ripple from the output voltage. A closed-loop voltage regulator is provided by sensing the output voltage, comparing it to a reference, and controlling the on–off ratio of the primary transistors. Isolation of the feedback signal can be provided by optical couplers or voltage-to-frequency converters, pulse transformers, and frequency-to-voltage converter circuits.

A second approach is to separate the regulation and transformer functions into two separate circuits. A block diagram is shown in Fig. 59. Voltage control is provided by a series dc/dc chopper operating at high frequency followed by a low-pass filter. The dc-to-ac inverter operates with a fixed 180° waveform that allows the use of a lightweight non-air-gap transformer. After rectification, the output filter needs to provide only a small amount of high-frequency filtering. This system is useful when the output voltage is low and the current is high, so that a large filter inductor is indicated.

In addition to the pulse-width-modulated switching power supplies, resonant converters based on either zero-voltage switching or zero-current switching can be provided. These circuits can operate a very high frequency (1–5 MHz). These circuits are generally configured as dc–dc converters operating as part of a distributed power supply.

Distributed Power Supplies

Large electronic systems composed of many circuit cards containing digital logic chips require regulated dc power at 5 volts dc. Newer types of digital logic will require 3.3 volts dc. The connected load can require dc currents in the range of hundreds of amperes. It is difficult to supply this amount of current without resorting to large, heavy, and complex bus structures. The area of conductor must be sufficient to keep the load regulation less than the lowest voltage required by the digital logic. A way to reduce this

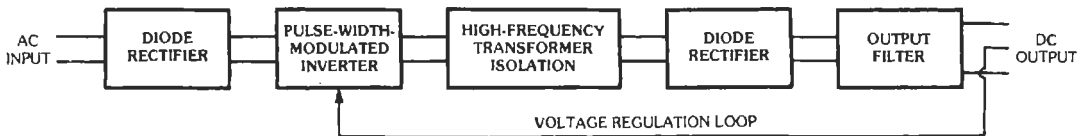


Fig. 58. Block diagram of dc power supply with pulse-width-modulated inverter for voltage control.

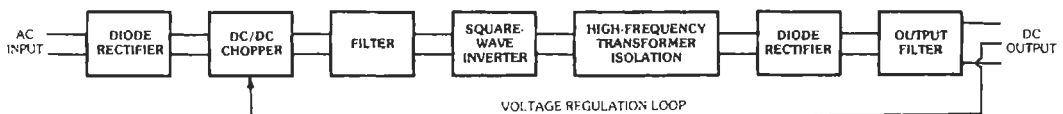


Fig. 59. Block diagram of dc power supply with dc/dc chopper for voltage control.

problem has been to use distributed power supplies. A block diagram is shown in Fig. 60. A bulk ac-to-dc converter, rated for full power, perhaps using a unity-power-factor ac-dc circuit, provides a regulated high-voltage dc output, in the range of 50 to 270 volts dc. For higher reliability, multiple converters can be connected to the high-voltage dc bus with diodes. This high-voltage dc power is supplied to the back plane of all of the electronic circuit cards and is brought on each printed circuit card. Each card has its own dc-dc converter to convert the 50 volts dc to 5 volts dc for that card. The power requirements for each card would be in the range of 10 to 100 watts, resulting in greatly reduced current demands. These "on-card" dc-dc converters operate at the highest possible frequency in order to reduce their size. Their aspect is such that their height is generally limited to less than 0.5 inch, so that the printed circuit cards can be closely stacked in their card rack.

Electromagnetic interference (EMI) can be controlled by shielding the dc-dc converter mounted on the printed circuit card. Conducted EMI can be controlled with series and parallel connected input filters. Heat rejection is accomplished by designing the converter to have the highest possible efficiency, providing fins for a larger heat surface, or conducting the heat from the printed circuit card to a cold plate.

AC-OUTPUT POWER SUPPLIES

The previous sections have discussed the aspects of converter and filter circuits and dc-output power supplies. This section describes some of the various types of power supplies that provide an alternating voltage and alternating current output. Some of the specific problem areas with an alternating output are associated with the power factor of ac loads, ac voltage magnitude control, polyphase output, adjustable-frequency output, and reverse power flow. Three basic ac-output power supplies are described, the ac voltage regulator, the dc-to-ac inverter, and the ac-to-ac cycloconverter. Greater detail and design equations for these generally complex systems are contained in References 2, 5, and 10-21.

AC Voltage Regulators

The ac voltage regulator is a static circuit that accepts a variable ac voltage input and produces a regulated ac voltage output. The output voltage is maintained for changes in input voltage, output load current, and load power factor. The input frequency is not changed, nor is the basic sinusoidal waveform modified. Since the input and output voltage is ac, thyristors with line commutation have been utilized. Fig. 61 shows a simple ac voltage-regulator circuit. This circuit is the static equivalent of the mechanical on-load tap changer. The thyristors are connected in inverse-parallel and provide a bidirectional connection from the tap on the autotransformer to the load. One tap is at a voltage lower than the input voltage; the other tap is at a higher voltage. In this manner, the input voltage can be either increased or decreased to produce a fixed voltage output. Output voltage control is provided by a closed-loop regulator that adjusts the firing angle of the two thyristors connected to the higher-voltage transformer tap. The gating on of these thyristors reverse biases the lower two thyristors, causing them to cease conducting the load current and transferring the load current to the higher-tap thyristors. More than two taps can be provided to develop a regulated output voltage with reduced harmonic content. The gating of the lower two thyristors is accomplished at the zero crossing of the load current. In this manner, the load current is initially in the lower tap irrespective of the load power factor. The transfer of load current at current zero allows the upper-tap thyristors to turn off and allows the lower-tap thyristors to conduct the load current. This transfer from upper to lower tap is at load-current zero crossing and is dependent on the load power factor. Therefore, the output voltage waveform and magnitude change with respect to the load power factor, putting an additional burden on the ac voltage regulator. Both three-phase and single-phase circuits with full isolation transformers have been fabricated. Where low harmonic distortion is required, an output filter can be provided.

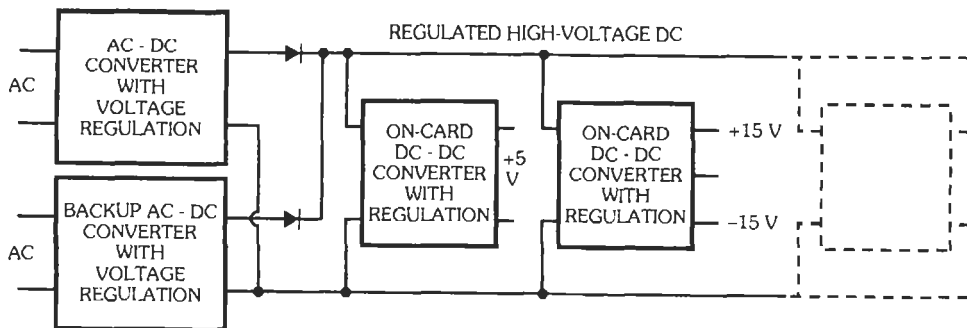


Fig. 60. Block diagram of distributed power supply.

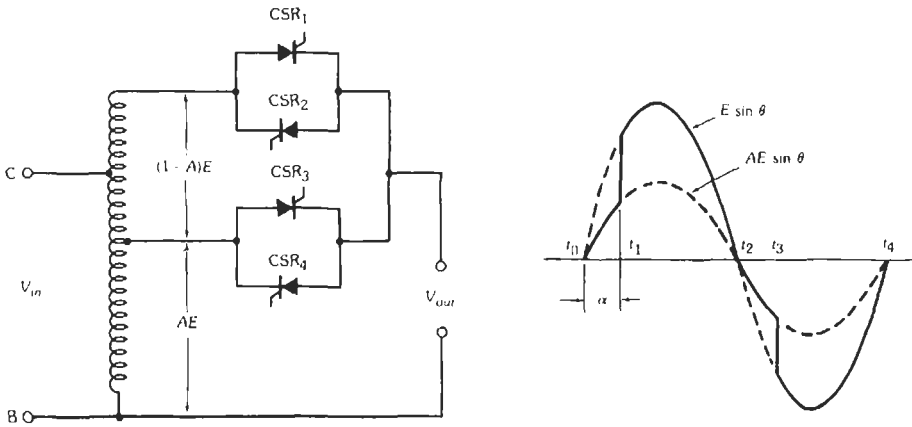


Fig. 61. Circuit diagram and waveforms for single-phase ac voltage regulator. (From Mazda, F. F. *Thyristor Control*. New York: John Wiley & Sons, Inc., 1973.)

DC-AC Inverters

This broad class of circuits is characterized by a dc input voltage and an ac output voltage. Since the input is dc, commutation of thyristors must be provided by techniques other than the reversal of the input voltage. Transistors are widely used in lower-power inverters because of their self-commutation capability. Four major classes of inverters are described in this section: first, a system called an “uninterruptible power supply (UPS),” which provides a fixed-voltage, fixed-frequency output; second, an adjustable-voltage, adjustable-frequency system for an adjustable-speed motor drive; third, an induction-heating supply with a high-frequency output; and fourth, an adjustable-ac-current, adjustable-frequency synchronous-motor drive operating at a leading power factor. These constitute typical applications for dc-to-ac inverters. Design equations and discussions of operating conditions in greater depth can be obtained in References 2, 5, and 10-14.

Uninterruptible Power Supply (UPS)—The purpose of this equipment is to provide a source of ac

power during outages of the normal source of utility supply. Uninterruptible power supplies are used in computer installations where power outages can mean loss of stored data (for example, in on-line reservations systems). Lower-power systems are provided to maintain continuous power to critical instrumentation (for example, a boiler-flame detector in a power plant). The source of power for these UPS installations is a battery that is kept charged from the utility. When the utility voltage is lost, the battery supplies power to the inverter and the connected load. These installations range in size from 1 kVA single-phase to 1000 kVA three-phase. The output frequency is generally 60 hertz, and the output voltage is the normal utilization voltage in order to use conventional ac-operated equipment and in some cases to operate in parallel with the ac utility.

Fig. 62 shows a block diagram of a large UPS installation. The separate battery charger, battery, inverter, utility bypass connection, and critical loads are shown.

High-power UPS systems use thyristors with forced commutation circuits as the power switches. Systems

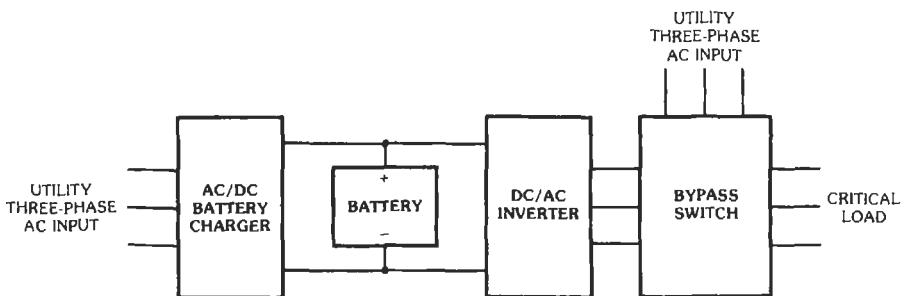


Fig. 62. Block diagram of UPS with utility bypass switch.

with ratings less than 200 kVA now use power transistors or insulated-gate bipolar transistors as the power switches. Fig. 63 shows a circuit diagram for a UPS system using a three-phase, pulse-width-modulated inverter supplied from a battery and feeding a transformer and filter. The transformer provides isolation and voltage matching. The filter reduces the harmonic content of the output voltage. A static bypass switch is shown. If the inverter must be taken off-line, the critical load can be supplied directly from the utility.

Drive for Adjustable-Speed AC Motor—A second application for an ac-output system is an adjustable-voltage, adjustable-frequency inverter operating from a dc source and providing power to an adjustable-speed ac motor. To operate at high efficiency over a wide range of speed, the ac motor must be provided with adjustable-voltage and adjustable-frequency power. The voltage and the frequency are not totally independent and are generally varied in such a way as to maintain a constant ratio of voltage to frequency. That is, higher ac voltage and higher frequency result in higher motor speed. Since the source is dc and induction motors operate at lagging power factor, forced commutation must be provided for the switching devices used in the inverter. If thyristors are used, they must be provided with auxiliary components to reduce the thyristor current to zero and provide a reverse voltage on the previously conducting device. Transistors with their self-commutation capability are becoming popular in inverters below 200 horsepower. The inverter operates from a relatively fixed source of dc voltage; therefore control of the output voltage must be provided in the inverter. The preferred technique is

to pulse-width modulate the output voltage by operating the transistors in an on-off mode many times during a given half-cycle of the motor fundamental frequency. Power transistors, with turn-on and turn-off times measured in microseconds, can be switched on and off many times during one half-cycle of a 60-hertz voltage, corresponding to a 3600-rpm synchronous speed with a two-pole induction motor. Various forms of pulse-width modulation (PWM) have been proposed, with a form of sine-wave shaping being a preferred approach. Rapid advancements in microprocessor-based control systems and stored-program memory devices have reduced their cost and complexity to the point where they are used to generate the pulse-width-modulation switching times.

A three-phase full-bridge inverter is shown in Fig. 64. Transistors T_1 - T_6 , together with feedback diodes D_1 - D_6 , operate in a pulse-width-modulated manner to generate a load voltage as shown in Fig. 65. The motor leakage reactance acts as a filter for the motor current. For industrial applications, the source of dc voltage shown in Fig. 64 is an uncontrolled diode rectifier from either a single-phase or polyphase ac supply.

Induction-Heating Supply—In this application, a source of high-frequency ac power is converted to localized heating for melting and heat treating metals and ferrous parts. The frequency is relatively fixed and is dependent on the application. Low frequency, 60-180 hertz, is used for large ferrous heating and melting applications, whereas 10-50 kHz is used for surface hardening of metallic parts and in ultrasonic cleaning tanks. This application requires that the ac power be coupled into the load with an induction coil or trans-

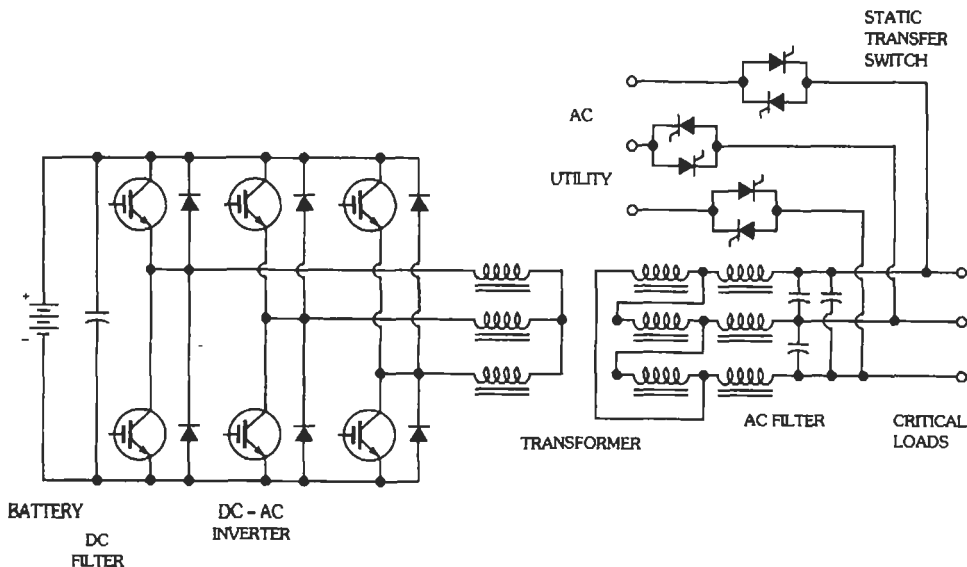


Fig. 63. Polyphase uninterruptible power supply with static transfer switch.

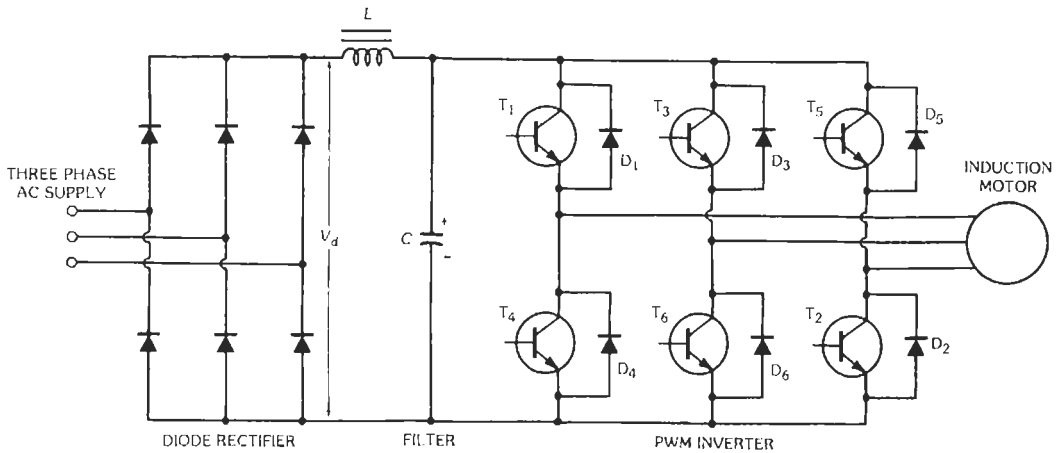


Fig. 64. Six-transistor PWM inverter/induction-motor drive. (From Bosc, Bimal K., ed. *Adjustable Speed AC Drive Systems*. New York: IEEE Press, 1981.)

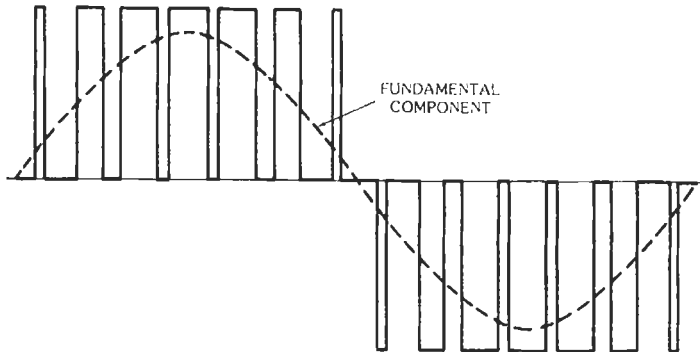


Fig. 65. Output-voltage waveform of a pulse-width-modulated inverter with sinusoidal modulation of the pulse width. (From Murphy, J. M. D. *Thyristor Control of A.C. Motors*. New York: Pergamon Press, 1973.)

ducer. Because of the lagging power factor of the work coil, a capacitor, either fixed or adjustable in steps, is provided to increase the power factor. High-power applications require the use of thyristors, and forced commutation is required because the operating frequency is higher than 60 hertz. The inductance of the work coil, an external capacitance, and the load resistor are selected to form an underdamped resonant circuit. Changes in the load resistance only change the damping ratio and barely affect the resonant frequency. The sinusoidal current waveform in the thyristors allows them to operate at higher frequencies with reduced losses. Fig. 66 shows an induction-heating inverter operating from a three-phase ac supply and supplying a single-phase load. High-power systems use this technique to balance the load on the three-phase supply.

Synchronous-Motor Drive—This system illustrates a fourth commutation technique for thyristors. In this case, the synchronous machine is operated at a leading fundamental-frequency power factor provided by overexciting the field of the machine. An adjustable-speed drive system of this type is shown in Fig. 67. The source of supply is an adjustable dc current supply. Thyristors T_1 – T_6 steer the dc current into the motor windings in sequence to provide the ac-motor line current. A suitable control system turns on the thyristors at the proper time with respect to the back electromotive force in the motor to provide a reverse voltage for commutation of the thyristors. Shaft-position sensors or electrical signal processing is used to determine the correct time to gate the thyristors. The source of adjustable dc current for the inverter is an ac-to-dc phase-controlled rectifier. Special commutating

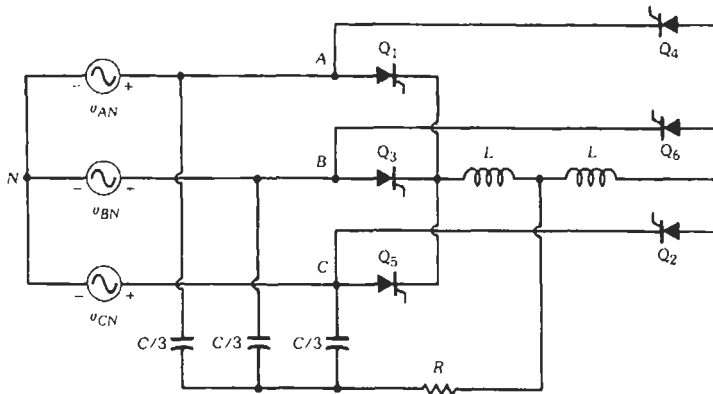


Fig. 66. Modified frequency multiplier with load commutation. (From Dewan, S. B., and Straughen, A. *Power Semiconductor Circuits*. New York: John Wiley & Sons, Inc., 1975.)

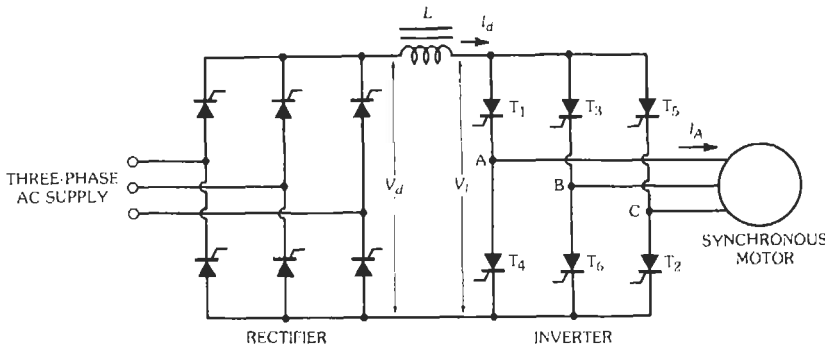


Fig. 67. Six-thyristor load-commutated inverter/synchronous-motor drive. (From Bose, Bimal K., ed. *Adjustable Speed AC Drive Systems*. New York: IEEE Press, 1981.)

circuits or turn-off intervals in the dc supply current are necessary to provide thyristor commutation during initial start-up and low-speed operation of the motor when the back emf is insufficient to commutate the thyristors. These drives have been fabricated to supply large synchronous machines up to 10 000 horsepower.

Cycloconverters

The preceding sections have described power-conversion systems that operate from dc power and produce an ac output. In industrial applications, the source is usually ac, thereby requiring a two-step conversion process from ac to dc to ac. The cycloconverter allows for a one-step conversion process from ac of a given voltage and frequency to adjustable-voltage ac (within limits) and adjustable-frequency ac (also within limits). This system uses the reversal of the line voltage to commutate the thyristors. The output voltage can be controlled by retarding the thyristor firing angle as in the phase-controlled rectifier. A second set of thyristors provides a path for the negative half-cycle of the

load current. A three-phase to single-phase cycloconverter is shown in Fig. 68. The firing angle of the thyristors is modulated in a sinusoidal manner to generate a cycle of the desired output frequency. The inverse-parallel connection of the thyristors provides a bidirectional path for current flow at any power factor. Since the commutation of the thyristors is provided by the input frequency, the maximum output frequency must be less than the input frequency. The maximum output voltage is also less than the input voltage because of the need to generate a sinusoidal output voltage waveform.

There are two major applications for cycloconverters. The first is very high horsepower, very low speed, multiple-pole induction and synchronous motors. The maximum frequency for these motors can be less than 5 hertz in the case of an 8000-hp motor enclosing a rotary cement ball mill. The second application is for variable-speed constant-frequency (VSCF) systems in aircraft. The speed of the aircraft generator, and hence its generated frequency, is a function of the engine speed. A cycloconverter is provided at the output of the

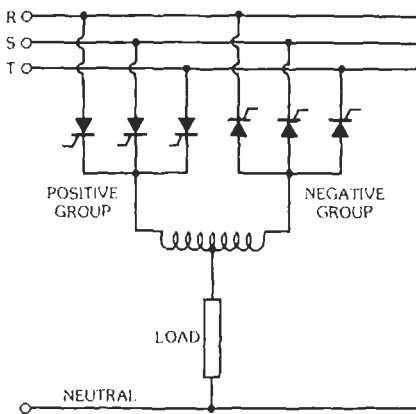


Fig. 68. Three-phase to single-phase cycloconverter power circuit. (From Murphy, J. M. D. *Thyristor Control of A.C. Motors*. New York: Pergamon Press, 1973.)

generator to convert the varying frequency into a fixed 400-hertz frequency for the aircraft power system. The adjustment of the thyristor gating times generates a sinusoidal output voltage waveform with minimum filtering. Output voltage regulation is provided by both the adjustment in thyristor gating times and adjustment of the field excitation on the synchronous generator.

REFERENCES

- Jahns, T. M., et al. "Circuit Utilization Characteristics of MOS-Controlled Thyristors." *IEEE Trans. Ind. Appl.*, Vol. 27, No. 3, May/June 1991, pp. 589-597.
- Mohan, N., et al. *Power Electronics—Converters, Applications and Design*. New York: John Wiley & Sons, Inc., 1989.
- "Harmonic Voltages in the Outputs of Controlled Rectifier Circuits." *Electronic Engineering*, Vol. 36, No. 442, December 1964.
- Baliga, B. J., and Chen, D. Y., Eds. *Power Transistors: Device Design and Applications*. New York: IEEE Press, 1984.
- Rashid, M. H. *Power Electronics—Circuits, Devices and Applications*. Englewood Cliffs, NJ: Prentice-Hall, Inc., 1988.
- Mazda, F. F. *Power Electronics Handbook*. London: Butterworths, 1990.
- Hnatek, E. R. *Design of Solid State Power Supplies*, 3rd ed. New York: Van Nostrand Reinhold Co., 1989.
- Billings, K. H. *Handbook of Switchmode Power Supplies*. New York: McGraw-Hill Book Co., 1989.
- Chryssis, G. C. *High-Frequency Switching Power Supplies*, 2nd ed. New York: McGraw-Hill Book Co., 1989.
- Murphy, J. M. D., and Turnbull, F. G. *Power Electronic Control of AC Motors*. Oxford: Pergamon Press, Inc., 1988.
- Bose, B. K. *Power Electronics and AC Drives, Inc.* Englewood Cliffs, NJ: Prentice-Hall, Inc., 1986.
- Kenjo, T. *Power Electronics for the Microprocessor Age*. Oxford: Oxford University Press, 1990.
- Bose, B. K., Ed. *Microcomputer Control of Power Electronics and Drives*. New York: IEEE Press, 1987.
- Dewan, S. B., and Straughen, A. *Power Semiconductor Circuits*. New York: John Wiley & Sons, Inc., 1975.
- Sen, P. C. *Thyristor DC Drives*. New York: John Wiley & Sons, Inc., 1981.
- Moltgen, G. *Line Commutated Thyristor Converters*. London: Pitman Publishing, 1972.
- Pelly, B. R. *Thyristor Phase-Controlled Converters and Cycloconverters*. New York: John Wiley & Sons, Inc., 1971.
- McMurray, W. *The Theory and Design of Cycloconverters*. Cambridge, MA: M.I.T. Press, 1972.
- Gyugyi, L., and Pelly, B. R. *Static Power Frequency Changers*. New York: John Wiley & Sons, Inc., 1976.
- Harnden, J. D., Jr., and Golden, F. B., Eds. *Power Semiconductor Applications*, Vol. 1 and Vol. 2. New York: IEEE Press, 1972.
- Bose, B. K., Ed. *Adjustable Speed AC Drive Systems*. New York: IEEE Press, 1981.
- Brown, M. *Practical Switching Power Supply Design*. San Diego: Academic Press, Inc., 1990.
- Brown, M. *Power Supply Cook Book*. Newton, MA: Butterworth-Heinemann, 1994.
- Mammano, R., Niedorf, R., *Improving Input Power Factor—A New Active Controller Simplifies the Task*. Unitrode Corp., Merrimack, NH.
- Tenti, P., Spiazzi, G., *Harmonic Limiting Standards and Power Factor Correction Techniques*. Sevilla, Spain: EPE Conference Tutorial, September 1995.

15 Feedback Control Systems

*Revised by
William R. Perkins**

Models of Control System Components	15-2
Error-Measuring Systems: Potentiometers, Synchros, Incremental Encoders	
DC Motors	
Two-Phase Servomotors	
Step Motors	
Rate Generators	
System Performance	15-4
Stability of Linear Time-Invariant Systems	
Tracking Accuracy	
Disturbance Inputs	
Sensitivity and Robustness	
Methods of Controller Design: Classical Design	
State Space Analysis and Design Techniques	15-12
Controller Design: Pole Placement	
Controller Design: LQ-Optimal Control	
Digital Control Systems	15-15
The z-Transform	
The z-Transfer Function	
The Inverse z-Transform	
State Variable Analysis of Digital Control Systems	
Stability of Linear Time-Invariant Digital Systems	
State Variable Design Methods: The Digital Case	
Phase-Locked Loop Servo Systems	15-21
Nonlinear Systems	15-22
Characteristics of Nonlinear Systems	
Describing-Function Technique	
Phase-Plane Method	
Computer-Aided Analysis and Design	15-33
Bibliography	15-33

* This revision draws heavily upon the previous version of this chapter prepared by B. C. Kuo.

Control is the use of feedback to achieve desired response of dynamic systems in the presence of uncertainties, disturbances, and constraints. Control systems are found widely in many processes and products, and include industrial and chemical process control; vehicle, aircraft, and spacecraft control; automation and manufacturing systems; and robotics, to name only a few. Even some issues in economic and social systems can be approached from the control system viewpoint. And feedback control systems exist naturally in many living systems, such as blood-pressure regulation and heart-rate control. Fig. 1 gives a simple block-diagram representation of a controlled process. For example, in a motor-speed control problem, the controlled process is the motor, the actuating signal is the voltage or current input to the motor, and the output in this case is the motor speed. The controlled process is called the "plant." Load disturbances may be present. The objective, then, is to produce a control input that causes the controlled output to respond in a desired manner, in the presence of disturbance inputs and uncertainties in the plant description.

A typical closed-loop (feedback) control system structure is shown in Fig. 2. In this system, the output is measured with sensors and is "fed back" to the controller. The controller generates an input to the actuator, which provides an actuating input to the plant. An actuator is needed in many systems because the controller signal cannot drive the plant directly. For example, consider an antenna position control system, in which the antenna is driven by a motor, which is the actuator for this case.

MODELS OF CONTROL SYSTEM COMPONENTS

For system analysis and design, the plant, sensors, actuators, and controllers must be modeled appropriately. Differential equations or transfer functions are frequently used as models. For digital (computer) control, the controller may be modeled by difference equations. As examples, we consider several components found frequently in electromechanical control systems.

Error-Measuring Systems: Potentiometers, Synchros, Incremental Encoders

Commonly used error-measuring systems, or comparators, are shown in Fig. 3. Rotary potentiometers are

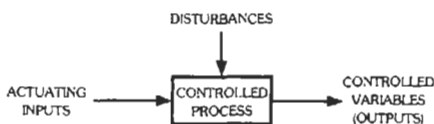


Fig. 1. Block diagram of a controlled process.

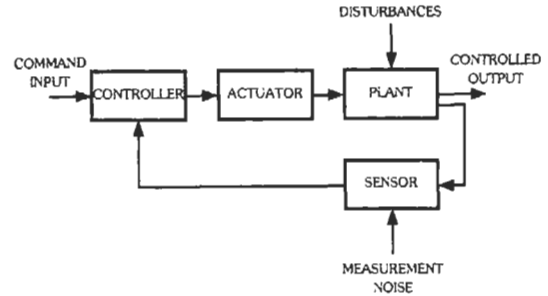


Fig. 2. Block diagram of a feedback control system.

available commercially in single-revolution or multi-revolution form. The potentiometers are commonly made with wirewound or conductive plastic resistance elements. Linear-motion potentiometers are also available for measuring linear motion properties.

The input-output relation of a potentiometer error detector is

$$e(t) = K_s[r(t) - c(t)]$$

where,

$e(t)$ is the output voltage,

$r(t)$ is the reference input,

$c(t)$ is the controlled variable,

K_s is the gain of the error detector,

Synchros are used in control systems as detectors and encoders due to their ruggedness of construction and high reliability. A synchro is basically a rotary device that operates on the same principle as a transformer and produces a correlation between an angular position and a voltage or set of voltages. See Fig. 3B.

Incremental encoders are available for rotary or linear motion. These devices typically convert linear or rotary displacement into digitally coded signals.

DC Motors

Types of dc motors include straight-series motors, split-series motors, shunt motors, compound motors (series-shunt fields), and permanent-magnet (PM) motors. Due to the advancements made in permanent-magnet materials, PM dc motors are one of the most widely used prime movers in industry today. The speed-torque characteristics of a PM dc motor are quite linear. However, the speed-torque characteristics of other types of dc motors are usually nonlinear.

The following variables and parameters are defined (SI units):

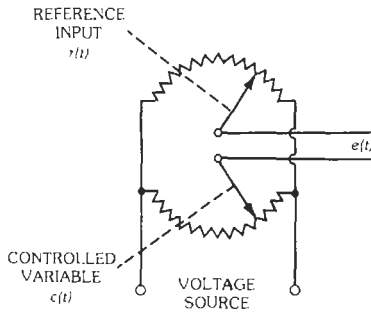
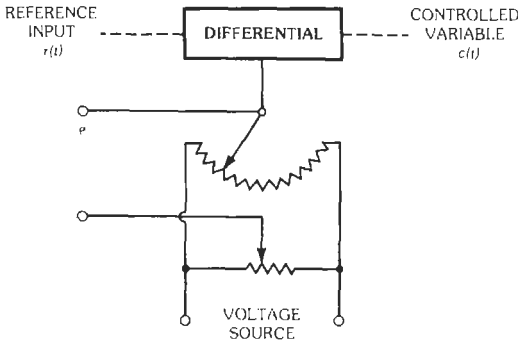
θ_m = motor angular position in radians

θ_i = load angular position in radians

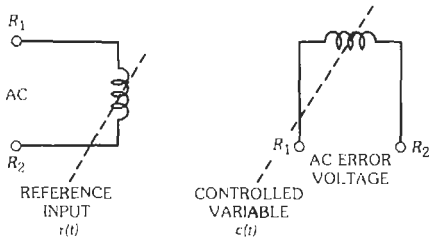
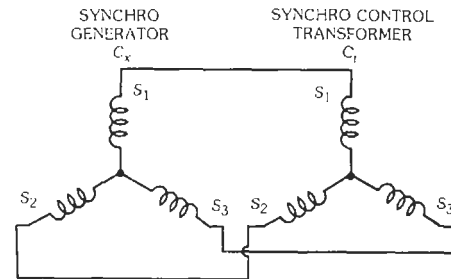
ω = angular velocity in radians/second = $d\theta/dt$

T_m = motor-developed torque in Newton-meters

J_m = motor moment of inertia in kilogram-meters²



(A) Potentiometer systems.



(B) Synchro system.

Fig. 3. Error-measuring systems.

- J_1 = load moment of inertia in kilogram-meters²
- E_m = applied voltage in volts
- K_m = motor torque constant in Newton-meters/ampere
- K_b = motor back emf constant in V/rad/s
- B_m = motor viscous-friction coefficient in Nm/rad/s
- B_1 = load viscous-friction coefficient in Nm/rad/s
- N = load-to-motor gear ratio = θ_l/θ_m

- B_{me} = viscous-friction coefficient reflected to motor shaft = $B_m + N^2 B_1$
- J_{me} = inertia reflected to motor shaft = $K_m + N^2 J_1$
- R_a = armature resistance of motor in ohms
- L_a = armature inductance of motor in henrys

The transfer function between the motor displacement and input voltage is

$$\theta_m(s)/E_m(s) = K_m / \{s[L_a J_{me} s^2 + (R_a J_{me} + B_{me} L_a)s + (K_b K_m + R_a B_{me})]\}$$

Two-Phase Servomotors

For low-power applications in control systems, ac motors are sometimes used because of their rugged construction. Most ac motors used in control systems are of the two-phase induction type. One of the two phases is excited from a constant-voltage source (the reference winding). The speed-torque curves shown in Fig. 4 are linearized for analytical purposes. Let

$$k = \frac{\text{Blocked-rotor torque at rated } E_m}{\text{Rated control voltage } E_1} = T_0/E_1$$

$$m = -\frac{\text{Blocked-rotor torque}}{\text{No-load speed}} = -T_0/\Omega_0$$

The transfer function between the control voltage and the motor displacement is

$$\theta_m(s)/E_m(s) = K_m/[s(1 + \tau_m s)]$$

where,

$$K_m = k/(B_{me} - m)$$

$$\tau_m = J_{me}/(B_{me} - m)$$

Step Motors

Step motors are electromechanical incremental actuators that convert digital pulse inputs to analog output shaft motion. The advantage of a step motor is that the

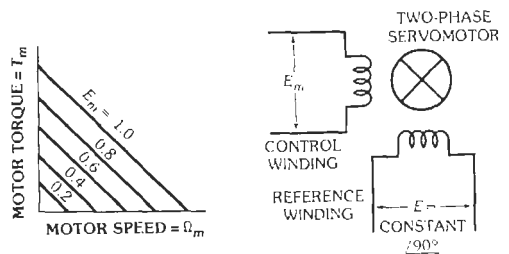


Fig. 4. Diagram of a two-phase servomotor and idealized torque-speed curves.

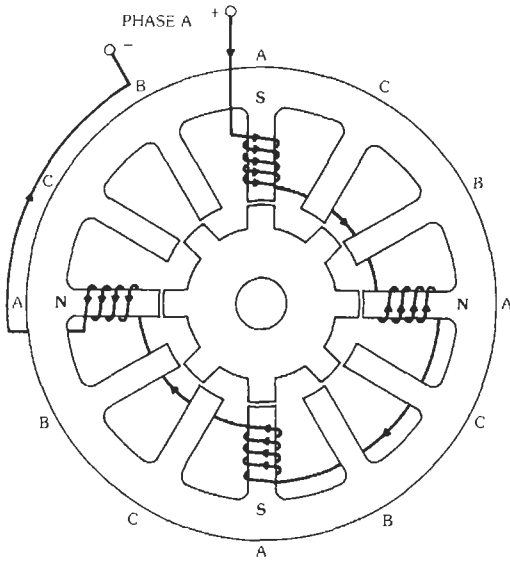


Fig. 5. Schematic diagram of a three-phase single-stack variable-reluctance step motor.

motor can be driven directly by digital inputs through the power drive, and no interface D/A conversion is necessary. In a rotary step motor, the output shaft of the motor rotates in equal increments in response to a train of input pulses.

The three most popular types of step motors are the variable-reluctance motor, the permanent-magnet motor, and the hybrid permanent-magnet motor. The latter two types have a permanent magnet in the rotor assembly.

Fig. 5 shows the schematic diagram of a single-stack three-phase variable-reluctance step motor. The stator in this case has 12 teeth, and the rotor has 8 teeth. There are four teeth per phase, and only the windings of phase A are shown. The rotor is shown to be at the detent position when phase A is energized with a dc current. If the dc excitation is shifted to the windings of

phase B, the rotor will rotate 15 degrees in the clockwise direction. If, instead, phase C is energized, the rotor will make a 15-degree step in the counter-clockwise direction. Therefore, the motor illustrated is a 24-step-per-revolution step motor.

Rate Generators

A rate generator (or tachometer) is a precision electromechanical component resembling a small motor and having an output voltage proportional to its shaft rotational speed. Rate generators have extensive applications both as computing instruments and as stabilizing components of feedback control systems. An example of the latter is illustrated in Fig. 6. The use of the rate generator produces an effective viscous damping and also tends to linearize the servomechanism by inserting damping of such magnitude that it is much larger than the nonlinear damping of the motor. Rate generators and servomotors are often constructed as integral units having a common shaft to eliminate the backlash. These units are available for dc or ac (either 400- or 60-hertz) operation.

SYSTEM PERFORMANCE

The performance of a closed-loop system (Fig. 2) is evaluated with respect to several desired objectives. These include:

- A. *Stability.* The system is bounded-input bounded-output (BIBO) stable if its output is bounded for any bounded input (zero initial conditions).
- B. *Steady-state tracking accuracy.* Frequently it is desired that the output $y(t)$ "follow," or track, a command or reference input $r(t)$, in the steady state ($t \rightarrow \infty$).
- C. *Transient response.* Usually "fast" transients are desirable, provided the control inputs remain realistic (satisfy amplitude, energy, or other constraints).
- D. *Disturbance rejection.* Ideally, the system should not respond to disturbance inputs. Such a goal is

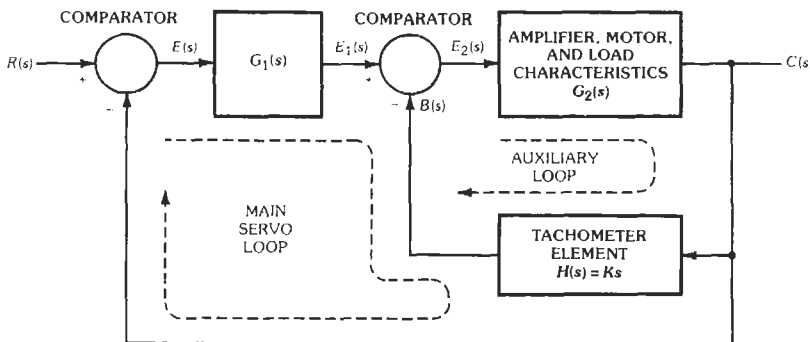


Fig. 6. Positioning servo.

unrealistic, and small (or zero) steady-state effects are sought, following rapid transients.

- E. *Sensitivity and robustness.* The system should perform well in the presence of uncertain plant models. Such modeling inaccuracies are always present, and a central purpose of feedback is to obtain satisfactory performance in the presence of such inaccuracies and uncertainties.

We consider each of these objectives in detail.

Stability of Linear Time-Invariant Systems

A linear time-invariant system is BIBO stable if and only if all roots of the characteristic equation lie strictly in the left-half s -plane. The location of these roots can be determined directly by finding the roots numerically (with computer assistance, if necessary), if all system parameter values are known. Several other methods are available for investigating stability. These methods also provide additional information concerning system performance that is useful in both analysis and design. We outline here the following methods:

1. Routh-Hurwitz criterion
2. Nyquist criterion
3. Root locus diagram
4. Bode diagram

Routh-Hurwitz Criterion—The stability of the common unity feedback signal-loop linear system modeled in Fig. 7 can be investigated by referring to the closed-loop transfer function

$$C(s)/R(s) = G(s)/[1 + G(s)]$$

The stability of the system depends on the location of the poles of $C(s)/R(s)$ or the zeros of $1 + G(s)$ in the complex s -plane. The zeros of $1 + G(s)$ are also known as the roots of the *characteristic equation*, which can be written

$$D(s) = \sum_{i=0}^n a_i s^i$$

where all the coefficients are real.

A *necessary* condition for $D(s)$ to have no roots on the imaginary axis or in the right half of the s -plane are that

all the coefficients of the equation must be of the same sign and that none of the coefficients is zero. This condition is not in general *sufficient*, however.

To check the necessary and sufficient conditions, we form the following Routh-Hurwitz tabulation:

a_n	a_{n-2}	a_{n-4}	a_{n-6}	·	·	·
a_{n-1}	a_{n-3}	a_{n-5}	a_{n-7}	·	·	·
b_1	b_2	b_3	b_4	·	·	·
c_1	c_2	c_3	c_4	·	·	·
d_1	d_2	d_3	·	·	·	·
e_1	e_2	·	·	·	·	·
f_1	·	·	·	·	·	·
·	·	·	·	·	·	·
·	·	·	·	·	·	·
·	·	·	·	·	·	·

where

$$\begin{aligned}
 b_1 &= (a_{n-1}a_{n-2} - a_{n-3}a_n)/a_{n-1} \\
 b_2 &= (a_{n-1}a_{n-4} - a_{n-5}a_n)/a_{n-1} \\
 b_3 &= (a_{n-1}a_{n-6} - a_{n-7}a_n)/a_{n-1} \\
 c_1 &= (b_1a_{n-3} - b_2a_{n-1})/b_1 \\
 c_2 &= (b_1a_{n-5} - b_3a_{n-1})/b_1 \\
 c_3 &= (b_1a_{n-7} - b_4a_{n-1})/b_1 \\
 d_1 &= (c_1b_2 - b_1c_2)/c_1 \\
 d_2 &= (c_1b_3 - b_1c_3)/c_1 \\
 d_3 &= (c_1b_4 - b_1c_4)/c_1 \\
 &\vdots
 \end{aligned}$$

The table will consist of n rows.

The system is stable, i.e., the characteristic polynomial has no zeros on the imaginary axis or in the right half of the s -plane, if every entry in the first column of the table has the same sign. The number of consecutive sign changes in the elements of the first column is equal to the number of zeros that are in the right half-plane.

A singular case develops when the first element in any one row is zero, or an entire row contains all zeros. Under such conditions, the table cannot be compiled by using the equations given above, and the tabulation must be modified.

Nyquist Stability Criterion—A second method for determining stability is known as the Nyquist stability criterion. This method involves obtaining the locus of the transfer function $G(s)$ in the complex G plane for values of $s = j\omega$ for ω from $-\infty$ to $+\infty$. Since the locus is always symmetrical about the real axis, it is sufficient to draw the locus for positive values of ω only. Fig. 8 shows loci for several simple systems. BIBO

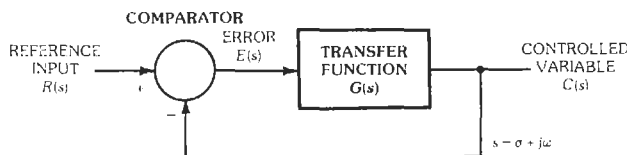


Fig. 7. Unity feedback single-loop system.

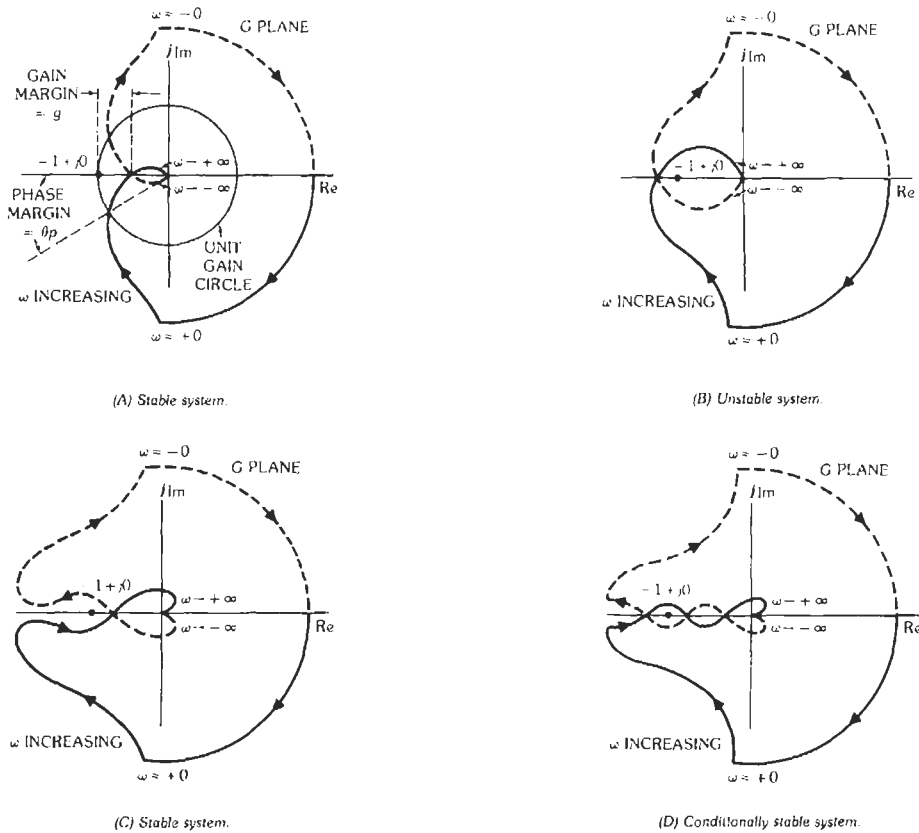


Fig. 8. Typical Nyquist loci.

stability can be determined from these plots by using the Nyquist stability criterion:

$$N = Z - P$$

where,

- N = number of counterclockwise encirclements of the point $-1 + j0$ in the G -plane,
- Z = number of zeros of $1 + G(s)$ in the right-half s -plane (same as number of poles of C/R in RHP),
- P = number of open-loop poles of $G(s)$ in RHP.

Clearly for BIBO stability, Z must be zero.

The curves in Figs. 8A and 8C represent stable systems (assuming $P = 0$) and are typical of the type-1 system;* the curve in Fig. 8B represents an unstable system. The system represented in Fig. 8D is conditionally stable; that is, for a particular range of values of gain K it is unstable. The system is stable for larger and smaller values of gain. *Note:* It is unstable as shown.

The *phase margin* θ_p and *gain margin* g are also

illustrated in Fig. 8A. The former is the angle between the negative real axis and $G(j\omega)$ at the point where the locus intersects the unit-gain circle. It is positive when measured as shown.

The gain margin g (in dB) is the negative dB value of $G(j\omega)$ corresponding to the frequency at which the phase angle is 180° (i.e., where $G(j\omega)$ intersects the negative real axis). Typical satisfactory values are -10 dB for g and an angle of 30° for θ_p . These values are selected on the basis of a good compromise between speed of response and reasonable overshoot. Note that for conditionally stable systems, the terms "gain margin" and "phase margin" may not have a simple interpretation.

Logarithmic Plots (Bode Diagrams)—The transfer function of a feedback control system can be described by separate plots of gain and phase versus frequency. Use of a logarithmic frequency scale permits simple straight-line (asymptotic) approximations for each curve. Fig. 9 illustrates the method for a transfer function with a single time constant. A comparison between approximation and actual values is included.

Transfer functions of the form $G = (1 + j\omega T)$ have similar approximations except that the attenuation curve slope is inverted upward ($+20$ dB/decade) and the values of phase are positive.

* System types are defined in the subsection "Tracking Accuracy."

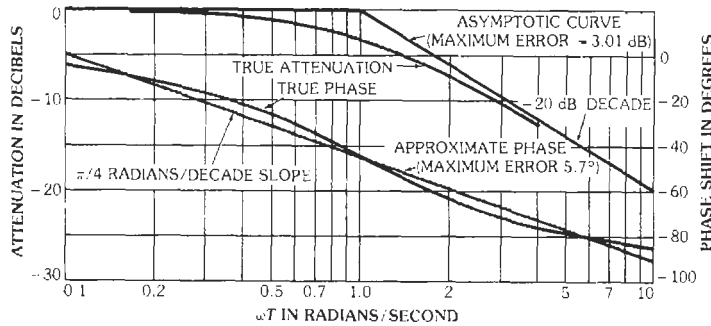


Fig. 9. Transfer-function plot. $G(j\omega) = 1/(1 + j\omega T)$.

The transfer function of feedback control systems can often be expressed as a fraction with the numerator and denominator each composed of linear factors of the form $(Ts + 1)$. Certain types of control systems, such as hydraulic motors where compressibility of the oil in the pipes is appreciable or some steering problems where the viscous damping is small, give rise to transfer functions in which quadratic factors occur in addition to the linear factors. The process of taking logarithms (as in making a dB plot) facilitates computation because only the addition of product terms is involved. The associated phase angles are directly additive. For example,

$$G(j\omega) = \frac{K(1 + j\omega T_2)}{[T^2(j\omega)^2 + 2\zeta T(j\omega) + 1](1 + j\omega T_1)(1 + j\omega T_3)}$$

where $s = j\omega$. The exact magnitude of G in decibels is

$$20 \log_{10} |G| = 20 \log_{10} K + 20 \log_{10} |1 + j\omega T_2|$$

$$- 20 \log_{10} |1 + j\omega T_1| - 20 \log_{10} |1 + j\omega T_3| - 20 \log_{10} |T^2(j\omega)^2 + 2\zeta T(j\omega) + 1|$$

Plots of attenuation and phase for quadratic factors as a function of the *relative damping ratio* ζ are given in Figs. 10 and 11. The low-frequency asymptote is 0 dB, but the high-frequency asymptote has a slope of ± 40 dB/decade (the positive slope applied to zero quadratic factors), twice the slope of the simple pole or zero case. The two asymptotes intersect at $\omega = 1/T$.

The difference between the asymptotic plot and the actual curves depends on the value of ζ with a variety of shapes realizable for the actual curve. Regardless of the value of ζ , the actual curve approaches the asymptotes at both low and high frequencies. In addition, the error between the asymptotic plot and the actual curve is geometrically symmetrical about the break frequency $\omega = 1/T$. As a result of this symmetry, the curves of Fig. 10 are plotted only for $\omega T \leq 1$. The error for $\omega = \alpha/T$ is identical with the error at $\omega = 1/\alpha T$.

From Bode plots, system stability can be determined by noting gain and phase margins. Nyquist diagrams

can be constructed for confusing cases, such as conditionally stable systems.

Nyquist diagrams and Bode plots can also be constructed with the use of computer-aided analysis software. (See the section "Computer-Aided Analysis and Design" at the end of this chapter.)

Root-Locus Method—Root locus is a method of design due to Evans. It is based on the relation between the poles and zeros of the closed-loop system function and those of the open-loop transfer function. The rapidity and ease with which the loci can be constructed form the basis for the success of root-locus design methods, in much the same way that the simplicity of the gain and phase plots (Bode diagrams) make design in the frequency domain so attractive. The root-locus plots can be used to adjust system gain, guide the design of compensation networks, or study the effects of changes in system parameters.

In linear time-invariant "lumped" systems, $G(s)$ is a rational algebraic function, the ratio of two polynomials in s :

$$G(s) = m(s)/n(s)$$

From Fig. 7

$$\begin{aligned} (C/R)(s) &= G(s)/[1 + G(s)] \\ &= \frac{m(s)/n(s)}{1 + [m(s)/n(s)]} \\ &= m(s)/[m(s) + n(s)] \end{aligned}$$

The zeros of the closed-loop system are identical with those of the open-loop system function. The closed-loop poles are the values of s at which $m(s)/n(s) = -1$. The root-locus method is a graphical technique for determination of the zeros of $m(s) + n(s)$ from the zeros of $m(s)$ and $n(s)$. Root loci are plots in the complex s plane of the variations of the poles of the closed-loop system function with changes in the open-loop gain. For the single-loop system of Fig. 7, the root loci constitute all s -plane points at which

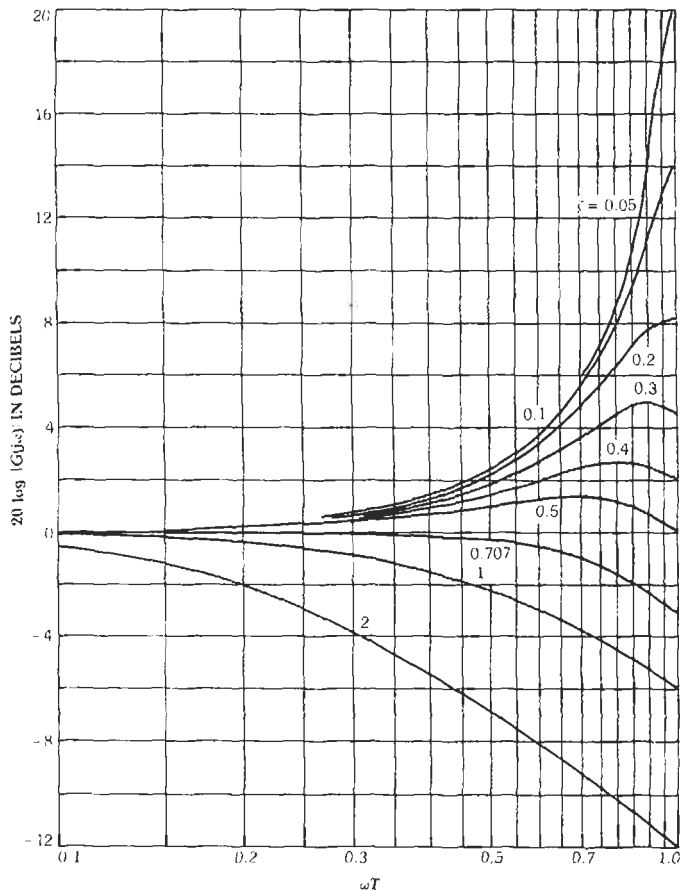


Fig. 10. Attenuation curve for quadratic factor. $G(j\omega) = 1/[T^2(j\omega^2) + 2\zeta T(j\omega) + 1]$. (By permission from Automatic Feedback Control System Synthesis, by J. G. Truxal. Copyright 1955, McGraw-Hill Book Co.)

$$\angle G(s) = 180^\circ + n360^\circ$$

where n is any integer including zero. A graphical interpretation for

$$G(s) = \frac{K(s + z_1)(s + z_2)}{s(s + p_1)(s + p_2)(q + p_3)}$$

is given in Fig. 12. Examples of root loci are given in Fig. 13 and 14.

For the example of Fig. 14, $K = K_1$ produces the case of critical damping. An increase in gain somewhat beyond this value causes a damped oscillation to appear. The latter increases in frequency (and decreases in damping) with further increase in gain. At gain K_3 , a sustained oscillation will result. Instability exists for gain greater than K_3 , as at K_4 . This corresponds to poles in the right half of the s plane for the closed-loop transfer function.

Various rules are available as aids in sketching root locus plots by hand. (Computer-aided packages can also be used; see section "Computer-Aided Analysis and Design.")

Intervals Along the Real Axis: The simplest portions of the plot to establish are the intervals along the negative real ($-\sigma$) axis, because then all angles are either 0° or 180° . Complex pairs of zeros or poles contribute no net angle for points along the real axis. Along the real axis, the locus will exist for intervals that have an *odd* number of zeros and poles to the right of the interval (Fig. 15).

Asymptotes: For very large values of s , $G(s) \sim K/s^{n-m}$. The locus will thus finally approach $(n - m)$ asymptotes at the angles (Fig. 16) given by the expression

$$(180^\circ + k360^\circ)/(n - m)$$

These asymptotes meet at a point s_1 (on the negative real axis) given by

$$s_1 = \frac{\Sigma(\text{poles}) - \Sigma(\text{zeros})}{(\text{finite poles}) - (\text{finite zeros})}$$

The other m branches of the locus will approach the zeros of $G(s)$, which are the zeros of $n(s)$.

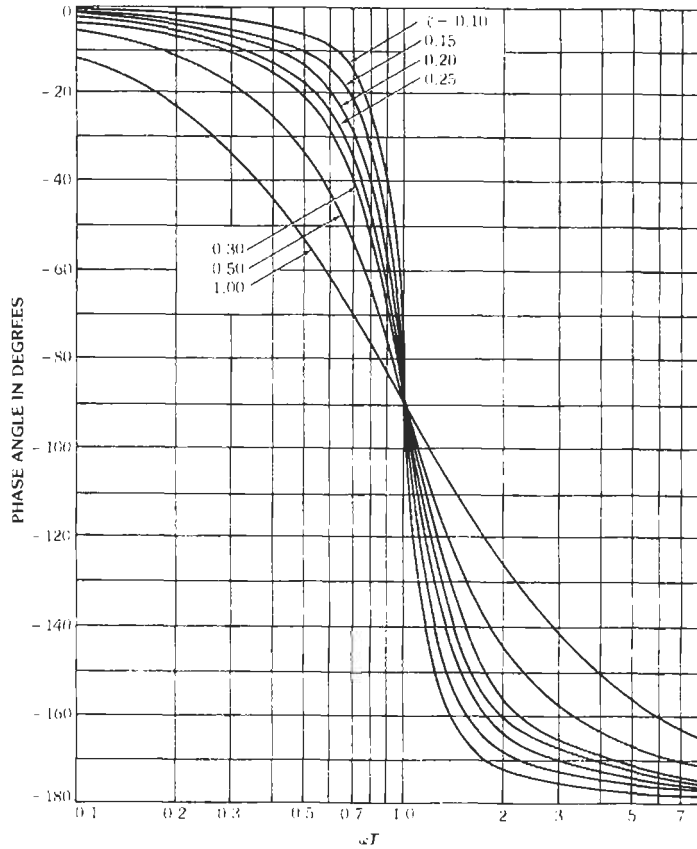


Fig. 11. Phase characteristics. (By permission from Theory of Servomechanisms, by H. M. James, N. B. Nichols, and R. S. Phillips. Copyright 1947, McGraw-Hill Book Co.)

Breakaway Points: Breakaway points from the real axis occur where the net change in angle caused by a small vertical displacement is zero. In Fig. 17, point p satisfies this condition at $1/x_0 = (1/x_1) + (1/x_2)$.

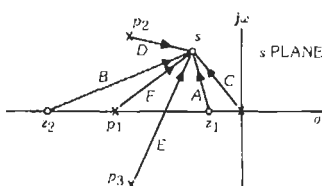
Intersections With $j\omega$ Axis: The Routh-Hurwitz test applied to the polynomial $m(s) + n(s)$ frequently permits rapid determination of the points at which the loci cross the $j\omega$ axis and the value of gain at these intersections.

Angles of Departure and Arrival: The angles at which the loci leave the poles and arrive at the zeros are readily evaluated from

$$\sum \angle \text{vectors from zeros to } s -$$

$$\sum \angle \text{vectors from poles to } s = 180^\circ + n360^\circ$$

For example, consider Fig. 18. The angle of departure of the locus from the pole at $(-1 + j1)$ is desired. If a test point is assumed only slightly displaced from the pole, the angles contributed by all critical frequencies (except the pole in question) are determined approximately by the vectors from these poles and zeros



$$G(s) = K(AB CDEF) / (A + \angle B - \angle C - \angle D - \angle E - \angle F)$$

Fig. 12. Graphical interpretation of $G(s)$.

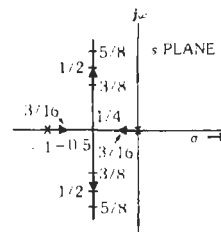


Fig. 13. Root locus for $G(s) = K/[s(s + 1)]$. Values of K as indicated by fractions.

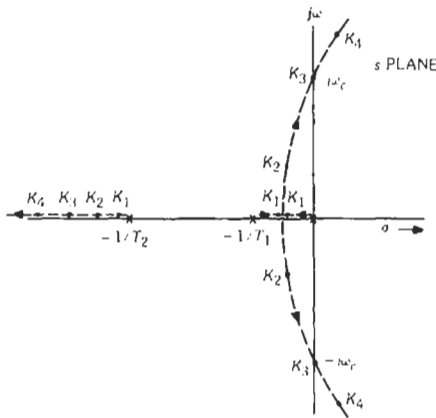


Fig. 14. Root locus for $G(s) = K/[s(T_1s + 1)(T_2s + 2)]$.

to $(-1 + j1)$. The angle contributed by the pole at $(-1 + j1)$ is then just sufficient to make the total angle 180° . In the example shown in the figure, the departure angle is found from the relation

$$+ 45^\circ - \left(\underbrace{135^\circ}_s + \underbrace{90^\circ}_{s+1+j1} + \underbrace{26.6^\circ}_{s+3} + \underbrace{\theta}_{s+1-j1} \right) = 180^\circ + n360^\circ$$

Hence, $\theta = -26.6^\circ$, the angle at which the locus leaves $(-1 + j1)$.

Tracking Accuracy

The tracking accuracy of various types of feedback control systems can be described most effectively in terms of the simple closed-loop unit-feedback system. Fig. 19 shows such a system; $R(s)$, $C(s)$, and $E(s)$ are the Laplace transforms of the reference input, controlled variable, and error signal, respectively.

For a typical linear system, $G(s)$ might appear as

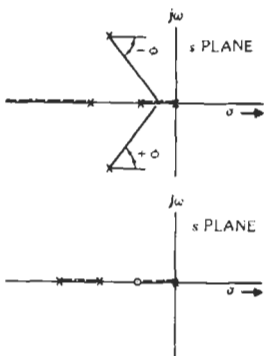


Fig. 15. Root-locus intervals along the real axis.

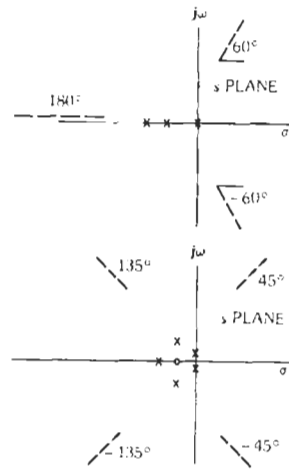


Fig. 16. Final asymptotes for root loci. Top, 60° asymptotes for system having three poles. Bottom, 45° asymptotes for system having an excess of four poles over zeros.

$$C(s)/E(s) = G(s) = [K(T_1s + 1)(T_3s + 1)]/[s^k(T_2s + 1)(T_4s + 1)]$$

The value of exponent k , an integer, designates the *type* of the system. This in turn reveals the nature of the steady-state performance of the system as follows:

For a type-0 system, a constant value of the controlled variable requires a constant error signal under steady-state conditions. A feedback control system of this type is generally referred to as a regulator system.

For a type-1 system, a constant rate of change of the controlled variable requires a constant error signal under steady-state conditions. A type-1 feedback control system is generally referred to as a servomechanism system. For reference inputs that change with time at a constant rate, a constant error is required to produce the same steady-state rate of the controlled variable. When applied to positions control, type-1 systems may also be referred to as "zero-displacement-error" systems. Under steady-state conditions, it is possible for the reference signal to have any desired constant position or displacement and the feedback signal or controlled variable to have the same displacement.

For a type-2 system, a constant acceleration of the controlled variable requires a constant error under steady-state conditions. Since these systems can main-

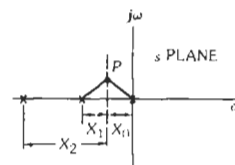


Fig. 17. Breakaway point.

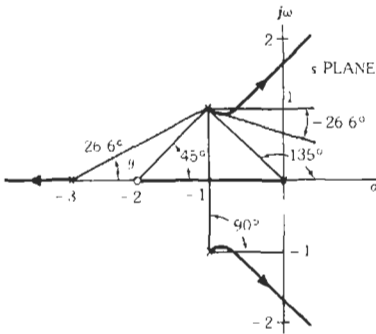


Fig. 18. Loci for $G(s) = K(s + 2)/[s(s + 3)(s^2 + 2s + 2)]$.

tain a constant value of controlled variable and a constant controlled variable speed with no actuating error, they are sometimes referred to as “zero-velocity-error” systems.

It is common to evaluate the steady-state tracking effectiveness using *error constants*. For unity feedback system (Fig. 19), these constants are defined as follows:

Positive Error Constant:

$$K_p = \lim_{s \rightarrow 0} [C(s)/E(s)] = \lim_{s \rightarrow 0} G(s)$$

Velocity Error Constant:

$$K_s = \lim_{s \rightarrow 0} [sC(s)/E(s)] = \lim_{s \rightarrow 0} sG(s)$$

Acceleration Error Constant:

$$K_a = \lim_{s \rightarrow 0} [s^2 C(s)/E(s)] = \lim_{s \rightarrow 0} s^2 G(s)$$

Disturbance Inputs

Frequently, systems are subjected to unwanted signals entering the system at points other than the input. Examples are load-torque disturbances, noise generated at a point within the system, etc. These may be represented as additional inputs to the system. Fig. 20 is a block diagram of such a situation. Transfer functions from the disturbance inputs to the output (with other inputs set to zero) can be found and used to analyze transient and steady-state disturbance responses, as has been done above for the command input transfer function.

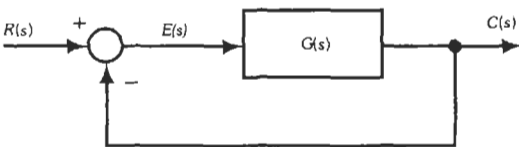


Fig. 19. Single-loop tracking system. $e = r - y$.

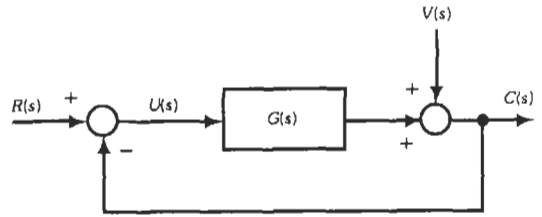


Fig. 20. Plant with “load disturbance” $v(t)$.

Sensitivity and Robustness

All system models used for analysis and design are approximations of the actual physical system. Ideally, the performance of the real physical system should be close to that predicted by our analysis and design methods, and it should be *insensitive* to the modeling approximations. A major reason for using feedback is the *potential* insensitivity of the performance of feedback structures to some variations in the plant. Bode suggested quantifying the study of plant parameter variations by defining the Bode Sensitivity Function,

$$S \triangleq (G/T)/(dT/dG)$$

where T is the system transfer function and G is the plant transfer function. Notice that S is the ratio of the percent (or logarithmic) change in T , with respect to the percent change in G . Ideally S should be “small”; that is, the (percent) change in the overall transfer function should be less than the (percent) change in the plant.

For the common “single-degree-of-freedom” unity feedback structure of Fig. 21,

$$S = 1/(1 + GG_c)$$

The quantity

$$GG_c \triangleq L$$

is called the *loop gain*, and $(1 + L)$ is the *return difference*. For S to be “small” (e.g., $|S(j\omega)| < 1$), the loop gain L should be “large.” Notice that

$$S + T = 1$$

This relationship provides a fundamental design constraint: S and T cannot be manipulated independently by

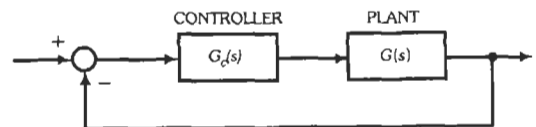


Fig. 21. Unity feedback “single-degree-of-freedom” structure.

choice of controller transfer function G_c for the system of Fig. 21.

Bode developed several other sensitivity results, of which two are especially important. One result reveals that S cannot be small for all frequencies. Suppose that $L = GG_c$ is of relative degree (number of poles - number of zeros) of at least 2, and that G and G_c are open-loop stable (no RHP poles). Then

$$\int_0^\infty \log_{10} |S(j\omega)| d\omega = 0$$

That is, the area under the sensitivity magnitude curve (in dB) is zero. Thus it is not possible for $|S(j\omega)| < 1$ (0 dB) for all frequencies. Typically $|S| < 1$ for low frequencies, but it becomes > 1 for higher frequencies.

The second important result relates the *slope* of the loop gain to stability. The exact statement of this result involves a complicated integral which is never evaluated directly. An accurate approximation, which is widely applied, is:

$$\angle L(j\omega) \cong n \times 90^\circ$$

where n is the slope of the loop amplitude curve in units of decade of amplitude per decade of frequency. Thus, if the phase margin is to be positive (stability), $\angle L(j\omega)$ must be greater than -180° at crossover, and so $n = 1$ (a slope of -20 dB/decade) at crossover. Greater slopes will produce unstable closed-loop systems.

Methods of Controller Design: Classical Design

The objective of controller design is to choose a controller transfer function to achieve desirable performance as described above. In the *classical design* approach, parameters are adjusted in several standard controller structures in an effort to meet performance specifications. Effects of parameter adjustments are commonly evaluated by using frequency response or root locus plots.

Methods of controller design for improving feedback-control-system response fall into the following basic categories:

- A. Series (cascade) compensation
- B. Feedback (parallel) compensation
- C. Load compensation

In many cases, any one of the above methods may be used to advantage, and it is largely a question of practical consideration as to which is selected. Fig. 22 illustrates the three methods.

Many systems can be controlled satisfactorily by using a proportional-integral-derivative (PID) controller:

$$G_c(s) = K[1 + (1/T_i)(1/s) + T_d s]$$

where K , T_i , and T_d are parameters to be adjusted. Such PID controllers are very common in process control and in many other control applications. The transfer function is easily implemented through analog or digital technology and can be obtained from many vendors in the form of an "off-the-shelf" controller.

Two other common classical controllers are the *lead* and the *lag* controllers. For the lead controller,

$$G_c(s) = K(T_s + 1)(\alpha T_s + 1)$$

where $\alpha < 1$. The lag controller has the same transfer function, but with $\alpha > 1$. Passive networks having these transfer functions are shown in Fig. 23. Op-amp-based active circuits for the same controllers are given in Fig. 24. Asymptotic attenuation and phase curves are shown in Fig. 25 and 26. The positive values of phase angle are to be associated with the phase-lead network, whereas the negative values are to be applied to the phase-lag network. Fig. 27 is a plot of the maximum phase shift for lag and lead networks as a function of the time-constant ratio.

Instead of direct feedback, the feedback path may contain frequency-sensitive elements. Typical of such frequency-sensitive elements are tachometers or other rate- or acceleration-sensitive devices that may be used for feedback directly or through suitable stabilizing circuitry.

The most common form of load stabilization involves the addition of an oscillation damper (tuned or untuned) to change the apparent characteristics of the load. Oscillation dampers can be used to obtain the equivalent of tachometric feedback. The primary advantages of load stabilization are the simplicity of instrumentation and the fact that the compensating action is independent of drift or the carrier frequency in ac systems.

STATE SPACE ANALYSIS AND DESIGN TECHNIQUES

State-variable methods are a modern approach to the analysis and design of control systems.

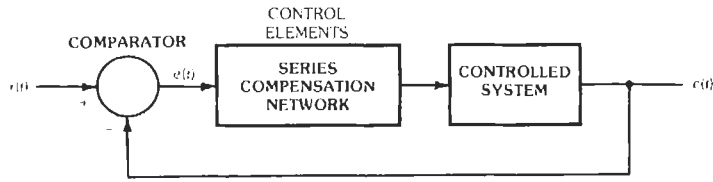
For linear time-invariant systems, the dynamic equations may be written in the following vector-matrix form:

$$\dot{x} = dx(t)/dt = Ax(t) + Bu(t)$$

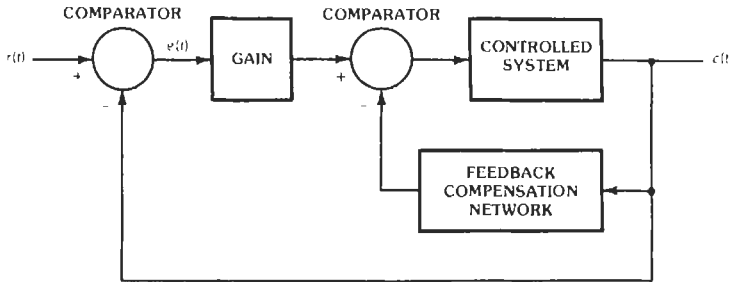
$$y(t) = Cx(t) + Du(t)$$

where

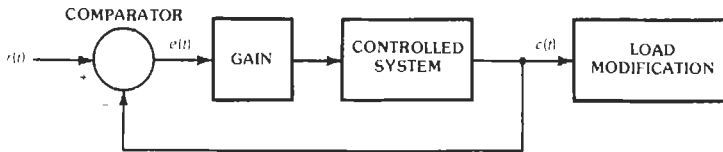
$$x(t) = \begin{bmatrix} x_1(t) \\ x_2(t) \\ \vdots \\ x_n(t) \end{bmatrix} = \text{state vector}$$



(A) Series compensation.



(B) Feedback compensation.



(C) Load compensation.

Fig. 22. Simple schemes for compensation.

$$u(t) = \begin{bmatrix} u_1(t) \\ u_2(t) \\ \vdots \\ u_m(t) \end{bmatrix} = \text{input vector}$$

$$y(t) = \begin{bmatrix} y_1(t) \\ y_2(t) \\ \vdots \\ y_p(t) \end{bmatrix} = \text{output vector}$$

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix}$$

= $n \times n$ coefficient matrix

$$B = \begin{bmatrix} b_{11} & b_{12} & \dots & b_{1p} \\ b_{21} & b_{22} & \dots & b_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ b_{n1} & b_{n2} & \dots & b_{nm} \end{bmatrix}$$

= $n \times m$ coefficient matrix.

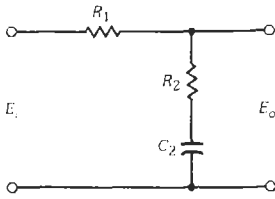
Similarly, C is a $p \times n$ coefficient matrix, and D is $q \times m$. The solution of the state equations is

$$x(t) = \phi(t - t_0)\xi(t_0) + \int_{t_0}^t \phi(t - \tau)Bu(\tau)d\tau$$

where $u(t)$ is the input applied at $t = t_0$, and $x(t_0)$ denotes the initial state vector. The $n \times n$ matrix $\phi(t)$ is the *state transition matrix*, and is defined as

$$\phi(t) = I + At + A^2t^2/2! + \dots + A^k t^k/k! + \dots$$

where I is the $n \times n$ identity matrix.



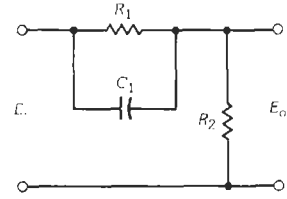
$$E_o/E = (T_2s + 1)/(T_1s + 1)$$

where,

$$T_2 = R_2C_2$$

$$T_1 = (R_1 + R_2)C_2$$

(A) Phase-lag network.



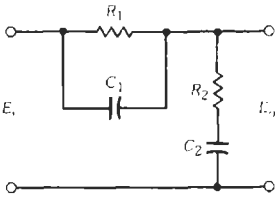
$$E_o/E = (T_2 + T_1)(T_1s + 1)/(T_2s + 1)$$

where,

$$T_1 = R_1C_1$$

$$T_2 = R_2R_1C_1/(R_1 + R_2)$$

(B) Phase-lead network.



$$E_o/E = \frac{(T_1s + 1)(T_2s + 1)}{T_1T_2s^2 + (T_1 + T_2 + T_1T_2)s + 1}$$

where,

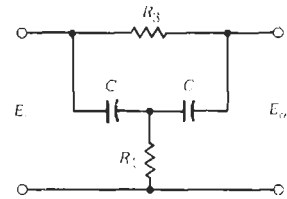
$$T_1 = R_1C_1$$

$$T_2 = R_2C_2$$

$$T_1T_2 = R_1C_2$$

$$G_1 = (T_1 + T_2)/(T_1 + T_2 + T_1T_2)$$

(C) Lead-lag network



$$E_o/E = \frac{T_1T_3s^2 + 2T_1s + 1}{T_1T_3s^2 + (2T_1 + T_3)s + 1}$$

where

$$T_1 = R_1C$$

$$T_3 = R_3C$$

(D) Bridged-T network.

Fig. 23. Passive lead, lag, and lead-lag networks.

The state transition matrix has the following properties:

$$\phi(0) = I$$

$$\phi^{-1}(t) = \text{matrix inverse of } \phi(t) = \phi(-t)$$

$$\phi(t_2 - t_1)\phi(t_1 - t_0) = \phi(t_2 - t_0) \quad \text{for any } t_0, t_1, t_2$$

$$\phi(t) = \mathcal{L}^{-1}[sI - A]^{-1}$$

where \mathcal{L}^{-1} denotes the "inverse Laplace transform of." The characteristic polynomial of A is defined as

$$\Delta = \text{determinant of } sI - A = |sI - A|$$

The roots of the characteristic equation $\Delta(s) = 0$ are the eigenvalues of A .

Controller Design: Pole Placement

One approach to control design is to feed back the state variables of the plant to adjust the closed-loop

system eigenvalues to obtain stability and desired transient response. State feedback is described by

$$u(t) = -Kx(t) + r(t)$$

where K is a $m \times n$ constant matrix of feedback controller gains, and $r(t)$ is an external command input. Then the closed-loop system is described by

$$\dot{x} = (A - BK)x + Br$$

The eigenvalues of the closed-loop system are given by

$$\det[sI - (A - BK)] = 0$$

and depend on K . These eigenvalues can be placed arbitrarily by choice of K if and only if the open-loop system is "controllable," that is, if

$$\text{rank } [B \ AB \ A^2B \ \dots \ A^{n-1}B] = n$$

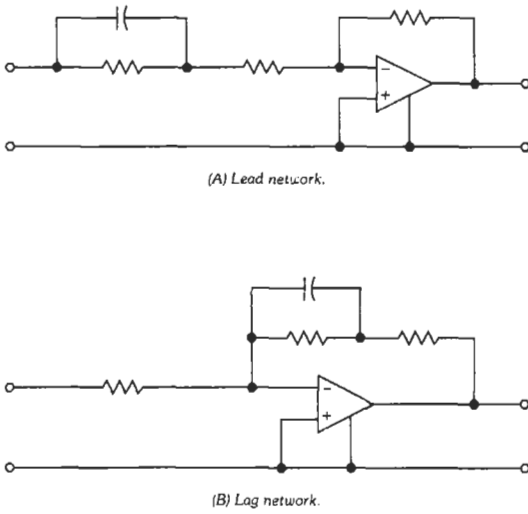


Fig. 24. Active lead and lag networks.

$$J = \frac{1}{2} \int_0^{\infty} [x^T(t)Qx(t) + u^T(t)Ru(t)]dt$$

This *integral* performance index *weights* quadratic forms in the state x and control u by symmetric constant matrices Q and R . The goal is to force the state from the initial state x_0 to the state $x = 0$ (a regulator problem) by *minimizing* J . That is, from among all possible controls u , the “best” or “optimum” control is that which results in the least value for J .

The weighting matrix Q is chosen to be positive-semidefinite, so that $x^T Q x \geq 0$. Thus a positive penalty is incurred in J for values of x (positive or negative) different from zero. The weighting matrix R , chosen to be positive definite, reflects a penalty for use of control. Typically, the designer would try various weightings Q and R until a desired state performance $x(t)$ was obtained, using an acceptable control $u(t)$.

The solution of this optimum control problem can be expressed in state feedback form,

$$u = -Kx$$

where $K = -R^{-1}B^T P$, the matrix P being the unique positive definite solution of the matrix algebraic Riccati equation (ARE):

$$A^T P + PA - PBR^{-1}B^T P + Q = 0$$

Controller Design: LQ-Optimal Control

In the optimum control approach to design, the system performance is manifested through a scalar *performance index* selected by the designer to reflect desired performance characteristics. In the LQ-optimal control approach, the linear plant is described in state form:

$$\dot{x} = Ax + Bu \quad x(0) = x_0$$

and the performance index is

DIGITAL CONTROL SYSTEMS

Digital control systems differ from continuous-data, or analog, systems in that the signals in one or more parts of the system are in the form of either pulses or

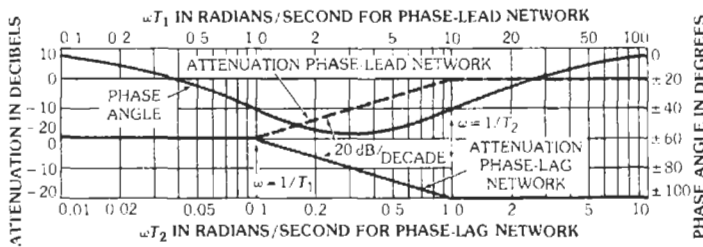


Fig. 25. Phase and magnitude for phase-lead and phase-lag networks. $T_1 = 10T_2$.

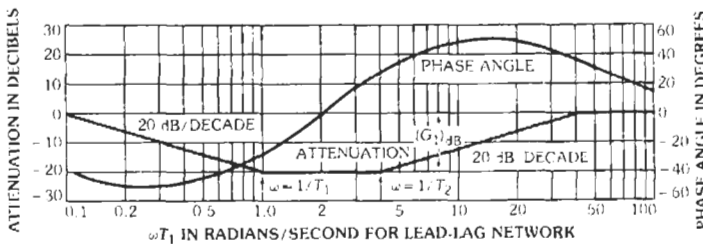


Fig. 26. Phase and magnitude for lead-lag network.

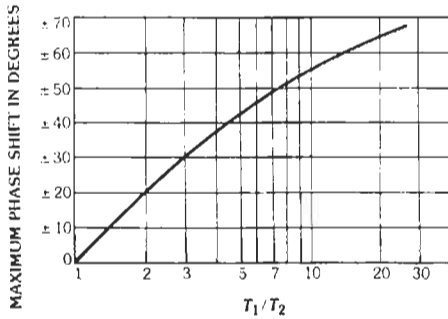


Fig. 27. Maximum phase shift for phase-lead (use positive angles) and phase-lag (negative angles) networks.

digitally coded signals. Due to the advances made in microcomputers, digital controllers have become very popular in control systems.

Fig. 28 shows the block diagram of a typical digital control system. Typically, the controlled process (plant) is an analog device, but the controller is digital. In this case, the reference input is shown to be an analog signal. The analog-to-digital converter (A/D) converts the analog signal into a digitally coded signal. The output of the digital controller is a digitally coded (such as binary-coded) signal. The digital-to-analog converter (D/A) converts the digital signal into an analog one for the controlled process.

From the analytical standpoint, the A/D operation can be represented by a *sampler* having T as the sampling period. The D/A operation can be represented by a sample-and-hold device. The block diagram representation of the sample-and-hold is shown in Fig. 29. The hold device simply holds the output of the sampler for one sampling period, T . Fig. 30 illustrates a typical set of waveforms of the inputs and the outputs of the sampler and the hold device.

The sampler defines the values of the function $f(t)$ at the sampling instants. In the time domain, the input-output relation of the sampler is written:

$$f^*(t) = \sum_{k=0}^{\infty} f(kT)\delta(t - kT)$$

where $\delta(t)$ is the unit impulse function.

The Laplace transform of the last equation is

$$F^*(s) = \sum_{k=0}^{\infty} f(kT)e^{-kTs}$$

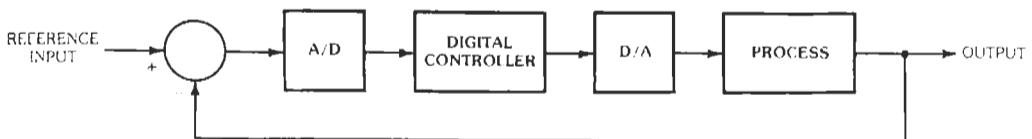


Fig. 28. Typical digital control system.

The transfer function of the hold device is

$$G_h(s) = (1 - e^{-Ts})/s$$

The z-Transform

Since the transfer functions of digital systems contain exponential terms of s , it is desirable to introduce the z-transform:

$$z = e^{Ts}$$

or

$$s = (1/T) \ln z$$

The z-transform of the sampled signal $f^*(t)$ is written:

$$\mathcal{Z}\{f^*(t)\} = F(z) = \sum_{k=0}^{\infty} f(kT)z^{-k}$$

The z-transforms of some common functions are given in Table 1. The role of the z-transform with respect to digital systems is similar to that of the Laplace transform with respect to continuous-data systems.

The z-Transfer Function

Fig. 31 shows the block diagram of a typical open-loop digital control system with a sample-and-hold. The transfer-function relation in the s -domain is written as:

$$C^* = [G_h(s)G_p(s)]^*F^*(s) \\ = \{(1 - e^{-Ts})/s\}G_p(s)^*F^*(s)$$

Thus, in terms of the z-transform,

$$C(z) = (1 - z^{-1})\mathcal{Z}\{G_p(s)/s\}F(z)$$

The open-loop transfer function is defined as

$$C(z)/F(z) = (1 - z^{-1})\mathcal{Z}\{G_p(s)/s\}$$

For the closed-loop digital control system shown in Fig. 32, the closed-loop transfer function is

$$C(z)/R(z) = G(z)/[1 + G(z)]$$

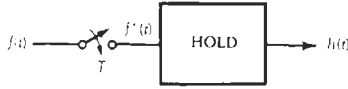


Fig. 29. Sample-and-hold.

where

$$G(z) = (1 - z^{-1})\mathfrak{Z}[G_p(s)/s]$$

The Inverse z-Transform

There are three methods of evaluating the inverse z-transform. These are the partial-fraction expansion method, the power-series method, and the inverse formula method.

Partial-Fraction Expansion Method—The function $F(z)/z$ is expanded into the form

$$F(z)/z = K_1/(z + a) + K_2/(z + b) + \dots$$

by partial-fraction expansion. Then

$$F(z)/z = K_1 z/(z + a) + K_2 z/(z + b) + \dots$$

The inverse z-transform of $F(z)$ is then taken term by term.

Power-Series Method—The function $F(z)$ is expanded into a power series in z^{-1} , i.e.,

$$F(z) = \sum_{k=0}^{\infty} f(kT)z^{-k}$$

Then, by the definition of the z-transform, the coefficient of z^{-k} , $k = 0, 1, 2, 3, \dots$ is $f(kT)$.

The Inversion Formula Method—The inverse z-transform of $F(z)$ can be expressed as the inversion formula

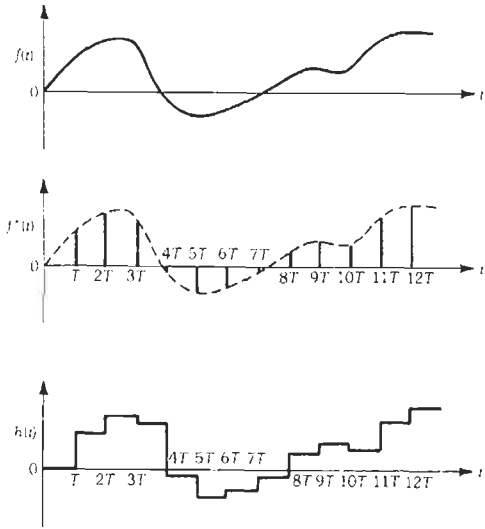


Fig. 30. Input-output relation of sample-and-hold.

$$f(kT) = (1/2\pi j) \oint F(z)z^{k-1} dz$$

where the integral is taken over a circle that encloses all the singularities of $F(z)z^{k-1}$. By use of the residue theorem of complex variable theory:

$$f(kT) = \sum \text{Residues of } F(z)z^{k-1} \text{ at the poles of } F(z)z^{k-1}$$

State Variable Analysis of Digital Control Systems

The state equations for continuous-data systems can be applied directly to digital systems with the input defined as

TABLE 1. z-TRANSFORMS

Time Function	Laplace Transform	z-Transform
$u(t)$ Unit-step function	$1/s$	$z/(z - 1)$
t	$1/s^2$	$Tz/(z - 1)^2$
$t^2/2$	$1/s^3$	$\frac{T^2 z(z + 1)}{2(z - 1)^3}$
e^{-at}	$1/(s + a)$	$z/(z - e^{-aT})$
te^{-at}	$1/(s + a)^2$	$(Tze^{-aT})/(z - e^{-aT})^2$
$\sin \omega t$	$\omega/(s^2 + \omega^2)$	$\frac{z \sin \omega T}{z^2 - 2z \cos \omega T + 1}$
$\cos \omega t$	$s/(s^2 + \omega^2)$	$\frac{z(z - \cos \omega T)}{z^2 - 2z \cos \omega T + 1}$

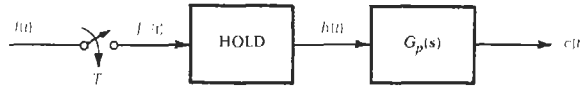


Fig. 31. Open-loop digital control system.

$$u(t) = u(kT) \quad kT \leq t < (k + 1)T$$

and

due to the action of the sample-and-hold. Thus, the state equations of an n^{th} -order digital system are written as

$$\phi(T) = e^{AT} = I + AT + \frac{A^2T^2}{2!} + \frac{A^3T^3}{3!} + \dots$$

$$dx(t)/dt = Ax(t) + Bu(kT) \quad kT \leq t < (k + 1)T$$

The solution of the state equation, the state transition equation, is

Stability of Linear Time-Invariant Digital Systems

The stability-analysis methods devised for linear continuous-data systems can all be extended to the stability study of digital systems. Since the z -transformation $z = e^{Ts}$ maps the imaginary axis in the s -plane onto the unit circle, $|z| = 1$, in the z -plane, the stability criterion of linear time-invariant digital systems is that all the roots of the characteristic equation must be found inside the unit circle in the z -plane.

$$x(t) = \phi(t - kT)x(kT) + \int_{kT}^t \phi(t - \tau)Bd\tau \cdot u(kT)$$

Let

The stability of the digital control system shown in Fig. 32 depends on the location of the poles of the closed-loop transfer function

$$\theta(t - kT) = \int_{kT}^t \phi(t - \tau)Bd\tau$$

Then,

$$x(t) = \psi(t - kT)x(kT) + \theta(t - kT)u(kT)$$

$$C(z)/R(z) = G(z)/[1 + G(z)]$$

To describe the state variables only at the sampling instants, let $t = (k + 1)T$. The last equation becomes

or of the zeros of $1 + G(z)$ in the complex z -plane.

The zeros of $1 + G(z)$ are also known as the roots of the characteristic equation, which can be written as

$$x[(k + 1)T] = \phi(T)x(kT) + \theta(T)u(kT)$$

which is the form of a vector-matrix difference equation with constant coefficients. The solution of this difference equation is

$$F(z) = a_{11}z^n + a_{n-1}z^{n-1} + \dots + a_2z^2 + a_0 = 0$$

where all the coefficients are real.

$$x(NT) =$$

The tabulation in Chart 1, known as the Jury Table, can be formed to check the necessary and sufficient conditions for the roots of the characteristic equations to be inside the unit circle.

The conclusion are:

$$[\phi(T)]^N x(0) + \sum_{k=0}^{N-1} \phi[(N - k - 1)T]\phi(T)r(kT)$$

$$= \phi(NT)x(0) + \sum_{k=0}^{N-1} \phi[N - k - 1)T]\phi(t)r(kT)$$

Number of *positive* calculated elements in the first column, $(b_0, c_0, \dots, p_0, q_0, r_0)$ = number of roots inside the unit circle.

Number of *negative* calculated elements in the first column, $(b_0, c_0, \dots, p_0, q_0, r_0)$ = number of roots outside the unit circle.

where,

$$\phi(NT) = \phi(T) \cdot \phi(T) \cdot \dots \cdot \phi(T) = [\phi(T)]^N$$

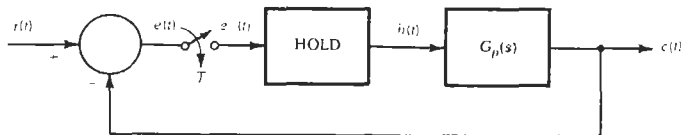


Fig. 32. Closed-loop digital control system.

CHART 1. CHART FOR CHECKING CONDITIONS FOR ROOTS OF CHARACTERISTIC EQUATIONS TO BE INSIDE UNIT CIRCLE

a_n $a_0 k_a$	a_{n-1} $a_1 k_a$	$a_{n-2} \dots a_2$ $a_2 k_a \dots a_{n-2} k_a$	a_1	a_0	$k_a = a_0/a_n$
b_0 $b_{n-1} k_b$	b_1 $b_{n-2} k_b$	$b_2 \dots b_{n-2}$ $b_{n-3} k_b \dots b_1 k_b$	b_{n-1}		$k_b = b_{n-1}/b_0$
c_0 $c_{n-2} k_c$	c_1 $c_{n-3} k_c$	$c_2 \dots c_{n-2}$ $c_{n-4} k_c \dots$			$k_c = c_{n-2}/c_0$
.....				
p_0 $p_2 k_p$	p_1 p_1/k_p	p_2			$k_p = p_2/p_0$
q_0 $q_1 k_q$	q_1				$k_q = q_1 q_0$
r_0					

The Nyquist stability criterion can be applied directly to determine the stability of digital control systems. The method involves the construction of the Nyquist locus of $G(z)$, with $z = e^{j\omega T}$, in the complex $G(z)$ plane for values of ω from $-\infty$ to $+\infty$. Actually, the locus of $G(z)$ repeats for every sampling frequency $\omega_s = 2\pi/T$. Thus, it is necessary to obtain only the locus of $G(z)$ for $0 \leq \omega \leq \omega_s$.

Fig. 33 shows the Nyquist loci of $G(z)$ of the system in Fig. 34. The open-loop transfer function is

$$G(z) = \frac{T^2 K_p (z + 1)}{2J_v z^2 + (2K_r T - 4J_v)z + 2J_v - 2K_r T}$$

with $J_v = 41\,822$, $K_r = 317\,000$, and $T = 0.1$ second. Or,

$$G(z) = \frac{1.2 \times 10^{-7} K_p (z + 1)}{(z - 1)(z - 0.242)}$$

Frequency response plots for the digital case present some difficulties not present in the continuous-time case. This is because the z -transform, with $z = e^{sT}$, is rational in z but not in s . So $s = j\omega$ leads to $z = e^{j\omega T}$, which means the frequency response transfer function is not rational in ω . Thus the Bode straight line approximations do not apply directly, making more difficult the sketching of the response curve by hand. Computer aids can be used to advantage. As an alternative, the Bode diagram of a digital control system can be obtained by use of the bilinear transformation

$$z = (1 + W)/(1 - W)$$

where W is a complex variable. For $z = e^{j\omega T}$,

$$W = j \tan(\omega T/2) = \sigma_W + j\omega_W$$

Thus,

$$\omega_W = \tan(\omega T/2)$$

Then the Bode plot of $G(z)$ can be made in the logarithmic coordinates using the definition of ω_W . For the system shown in Fig. 34,

$$G(j\omega_W) = \frac{1.583 \times 10^{-7} K_p (1 - j\omega_W)}{j\omega_W (1 + j1.636\omega_W)}$$

The Bode diagram of $G(j\omega_W)$ is shown in Fig. 35.

The rules on the construction of root loci in the s -plane for continuous-data systems can be applied

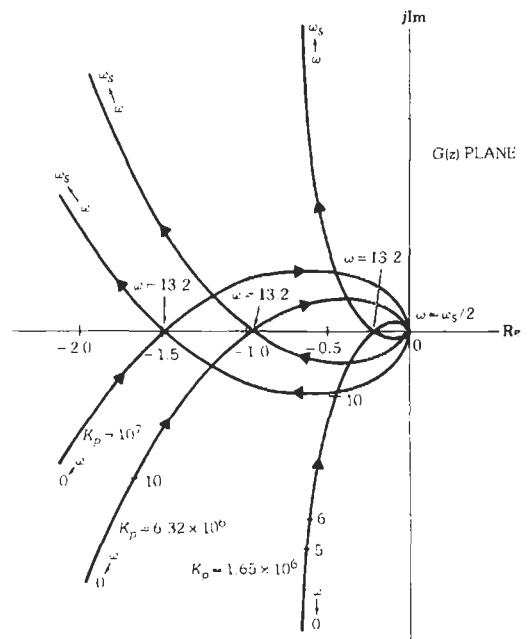


Fig. 33. Nyquist plot for system of Fig. 34.

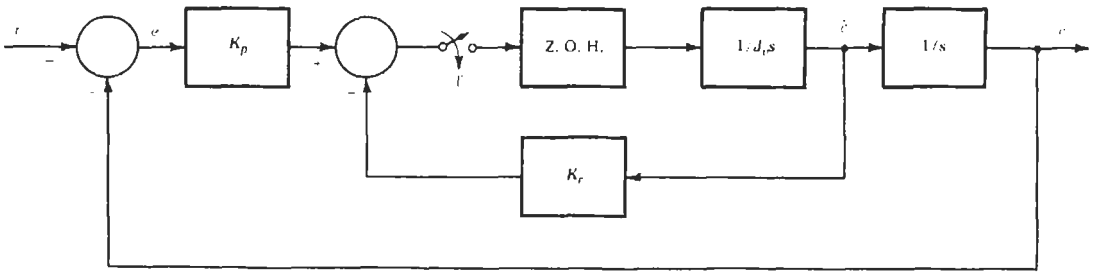


Fig. 34. Digital control of $1/s^2$ plant.

directly to the root loci in the z -plane. The only difference is that in the z -plane the stability of the digital control system must be investigated with respect to the unit circle $|z| = 1$.

For the digital control system in Fig. 34, the open-loop transfer function $G(z)$ has a zero at $z = -1$, and two poles at $z = 1$ and $z = 0.242$. The root loci of the characteristic equation for $0 \leq K_p < \infty$ are constructed as shown in Fig. 36.

State Variable Design Methods: The Digital Case

Both the pole placement and LQ-optimal control approaches can be applied to the digital control case. The results are similar to those for continuous-time. State feedback is given by

$$u(kT) = -Kx(kT) + r(kT)$$

and so, from

$$x[(k + 1)T] = \phi(T)x(kT) + \theta(T)u(kT)$$

The closed-loop system is described in state form by

$$x[(k + 1)T] = [\phi(T) - \theta(T)K]x(kT) + \theta(T)r(kT)$$

Thus the closed-loop eigenvalues can be placed arbitrarily by choice of K if and only if $[\phi(T), \theta(T)]$ satisfy the controllability condition

$$\text{rank}[\theta(T)\phi(T)\theta(T) \cdots \phi(T)^{n-1}\theta(T)] = n$$

For LQ design, the quadratic index to be minimized in the discrete (digital) case is

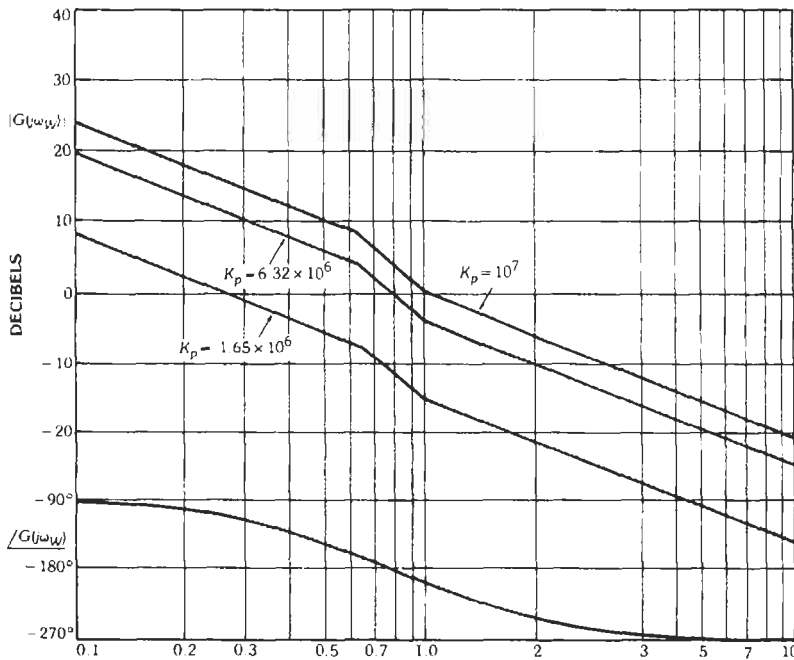


Fig. 35. Bode diagram of the digital control system in Fig. 34.

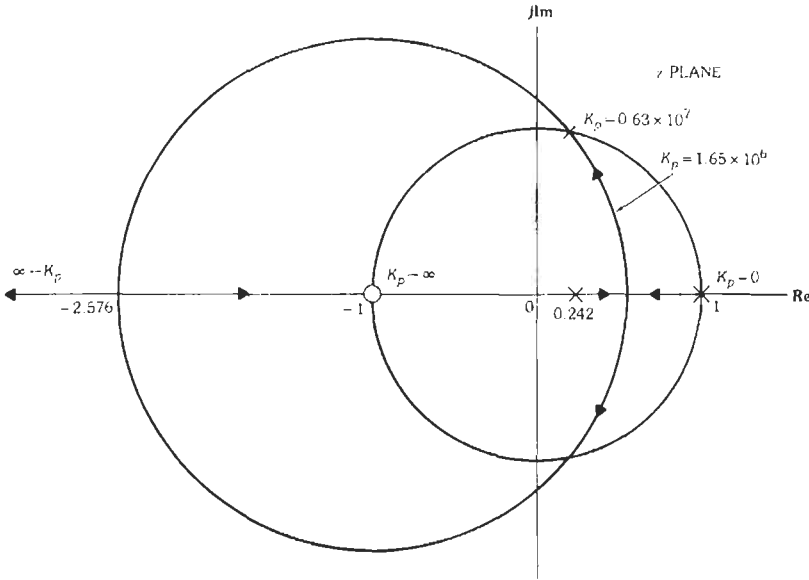


Fig. 36. Root locus for $G(z) = [(1.2 \times 10^{-7})(z + 1)K_p] / [(z - 1)(z - 0.242)]$.

$$J = \frac{1}{2} \sum_{k=0}^{\infty} [x^T(kT)Qx(kT) + u^T(kT)Ru(kT)],$$

where $Q \geq 0, R \geq 0$. (Note that, for the digital case, we require only positive *semidefiniteness* for R , in contrast to the positive definite requirement in the analog case). The optimum control is

$$u(kT) = -Kx(kT)$$

where

$$K = (\theta^T \theta + R)^{-1} \theta^T P \phi$$

with P being the unique positive definite solution of the discrete algebraic Riccati equation (DARE)

$$P = \phi^T [P - P \theta (\theta^T P \theta + R)^{-1} \theta^T P] \phi + Q$$

PHASE-LOCKED LOOP SERVO SYSTEMS

A phase-locked loop servo is a closed-loop control system that is used widely in communication systems for frequency demodulation and bit synchronization. Phase-locked loop servos are also used for velocity control, especially when a high degree of speed regulation and accuracy is desired.

A phase-locked loop in its basic form is represented by the block diagram of Fig. 37. The input signal and the feedback signal are sinusoidal, and, upon locking, both signals will have the same frequency and a constant phase difference. Any deviation from the desired phase difference is detected by the phase

detector and is transmitted to the voltage-controlled oscillator (VCO) to correct the error. Since the phase difference is constant, the frequencies of the input and the feedback signals are the same. This principle can be applied to servo systems for which the control objective is speed regulation.

In the case of a phase-locked loop servo for speed control, the VCO is replaced by an amplifier-motor-encoder combination. The output of the encoder is a sinusoidal signal with a frequency proportional to the angular velocity of the motor. Fig. 38 shows the block diagram of a phase-locked loop servo.

The main difference between a phase-locked servo and the conventional servo systems is that the reference input for a phase-locked servo is a periodic signal rather than a fixed voltage. The feedback signal from the encoder is a pulse train whose frequency is proportional to the speed of the motor. The phase detector compares the frequencies or phases of the reference and feedback signals, and generates an error voltage proportional to the difference. The error voltage is then filtered and sent to the motor; the motor speed changes to reduce the

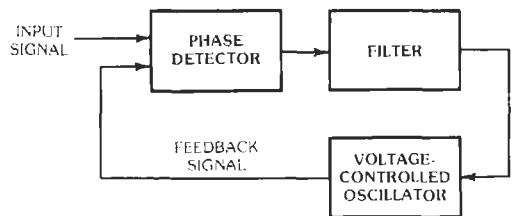


Fig. 37. Phase-locked loop.

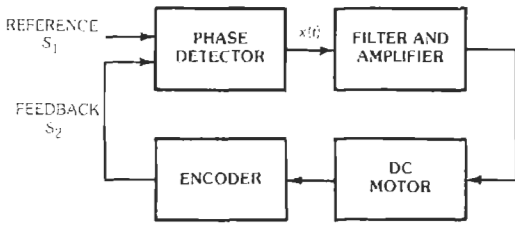


Fig. 38. Phase-locked servo system.

difference between the frequencies or phases of the two signals. Thus, the motor speed can be synchronized to the reference frequency.

A typical set of signals, $s_1(t)$ and $s_2(t)$, is shown in Fig. 39. The output of the phase detector, $x(t)$, switches between 0 and V_s as shown in Fig. 39.

The phase-locked loop servo can be modeled as a linearized digital control system as shown in Fig. 40.

The loop transfer function of the system in Fig. 40 is

$$G(z) = (V_s T n / 2\pi) \mathcal{Z}\{G_p(s) / s\}$$

Let

$$G_p(s) = K_m / (1 + \tau s)$$

Then,

$$\begin{aligned} \mathcal{Z}\{G_p(s) / s\} &= \mathcal{Z}\{K_m / [s(1 + \tau s)]\} \\ &= K_m z / (z - 1) - K_m z / (z - e^{-T\tau}) \end{aligned}$$

Let

$$K = (V_s T n / 2\pi) K_m$$

Then,

$$G(z) = \frac{K_z (1 - e^{-T\tau})}{(z - 1)(z - e^{-T\tau})}$$

The characteristic equation of the digital phase-locked loop servo is

$$z^2 + (-1 - e^{-T\tau} + K - K e^{-T\tau})z + e^{-T\tau} = 0$$

Applying the stability test to the last equation yields the condition of stability:

$$K < 2(1 - e^{-T\tau}) / (1 - e^{-2T\tau})$$

NONLINEAR SYSTEMS

All physical systems have nonlinearities to some degree. However, there are no general methods of analysis and synthesis for nonlinear systems. Therefore, they are often treated with linear approximations, and,

in many cases of small nonlinear effects, satisfactory results have been obtained. In other cases, however, linear methods become restrictive in their application and quite often unrealistic.

There are many approaches to nonlinear systems that may be applicable to a certain type of system but not to all. Among these methods, two techniques have proved especially useful. These are the describing function technique and the phase plane.

The describing-function technique was first applied to the analysis of nonlinear feedback control systems by Kochenburger. It is the object of the describing-function method of analysis to reduce the representation of the nonlinearity to an equivalent linear gain and phase angle. The representation of the nonlinearity is described in terms of the fundamental component of the distorted output waveform in response to a sinusoidal input. The result of the describing-function analysis is an approximate representation of the system in the frequency domain. The synthesis of nonlinear systems can be carried out with the describing-function technique in much the same way as is done with the linear systems.

The phase-plane method, on the other hand, is a time-domain representation of the behavior of the first- and second-order systems portrayed on the phase plane. The phase-plane diagrams are a useful tool in the study of transient response to any initial condition and in some cases to step and ramp inputs. This method can be extended conceptually to higher-order systems in the so-called phase space. However, an n th-order system in the n th-dimensional phase space compared with the phase plane is difficult to envision and interpret and may be less effective in its usefulness.

Characteristics of Nonlinear Systems

A comparison of linear and nonlinear systems yields the following representative characteristics.

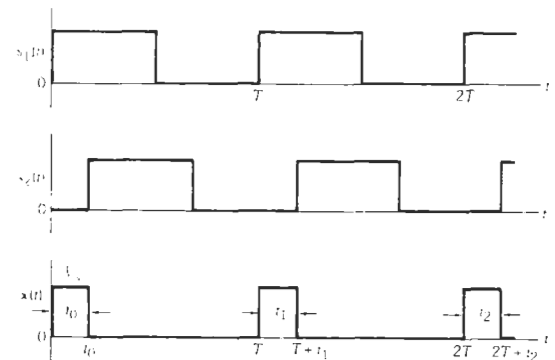


Fig. 39. Typical input and output signals of the phase detector.

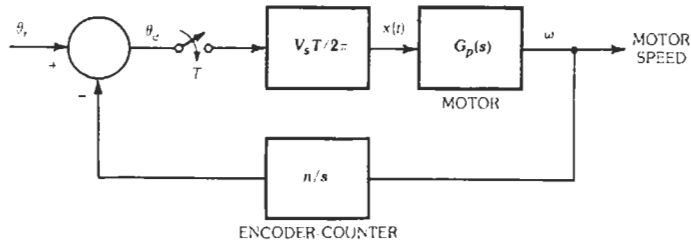


Fig. 40. Linearized model for a phase-locked loop servo.

Principle of Superposition Does not Apply—In nonlinear systems, the response to a combination of individual signals at the input will not be the same as the response to the sum of those same signals. The relationship of the impulse response (weighting function) to the frequency response in linear systems is therefore not applicable. The familiar Laplace transform operation ($s = \sigma + j\omega$) and the transfer-function concept used extensively in linear systems cannot be applied.

Nonlinear Response is Dependent on Input Signal—The response of nonlinear systems is dependent on the input-signal size and initial conditions as well as the system parameters. A stable response for one input signal may be unstable for another.

In nonlinear systems, phenomena exist that are not present in linear systems (generation of new frequencies, jump resonance, limit cycles, etc). The following are examples of such nonlinear phenomena.

Jump Resonance—The phenomenon called jump resonance is observed in certain closed-loop systems with saturation, where the input-output amplitude ratio and phase angle as a function of frequency exhibit sudden discontinuities. The typical closed-loop gain characteristics of a saturating system with jump resonance are shown in Fig. 41. Amplitude ratio $|\theta_o/\theta_i|$ is plotted as a function of frequency for fixed amplitude θ_i . As the frequency is increased from zero, the frequency response follows the curve along points A, B, and C. At point C a sudden discontinuous jump to D is observed with an incremental increase in frequency. Further increase in frequency leads to point E along the curve. If the frequency is reversed, the response retraces the path E, D and continues to point F, at which a sudden jump to point G occurs; it then continues on through B and A of the gain curve. The phase-angle response behaves similarly. The overall response curve exhibits a hysteresis-type property or jump resonance.

Limit Cycle—Limit cycle is an oscillation peculiar to nonlinear systems. The oscillatory behavior, unexplainable in terms of linear theory, is characterized by a constant amplitude and frequency determined by the nonlinear properties of the system. Limit cycles are distinguishable from linear oscillations in that their amplitude of oscillation is independent of initial conditions. For instance, if a system has a stable limit cycle,

the system will tend to fall into the limit cycle, with the output approaching the amplitude of that limit cycle regardless of the initial condition and forcing function. A limit cycle is easily recognized in the phase plane as an isolated closed path as shown in Fig. 42.

Generation of New Frequencies—In a nonlinear system, the output of the nonlinear device contains harmonic and subharmonic frequencies of the input signal. For instance, application of two sine waves of different frequencies f_1 and f_2 to the input will produce components corresponding to the input frequencies f_1 , f_2 , their sum and difference $f_1 \pm f_2$, their higher harmonics mf_1 , nf_2 , and various combinations of sums and differences $mf_1 \pm nf_2$. In a linear system, only the components of the input frequencies f_1 , f_2 will be reflected in the output. Clearly, the frequency-response concept of linear systems must be modified for nonlinear applications.

Hysteresis—Multivalued functions exist when two or more function values correspond to the value of the variable. Multivalued functions are intrinsically nonlinear. Examples of this are the hysteresis curves of magnetic materials and the backlash of a gear train.

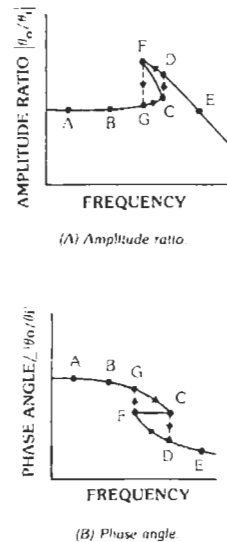


Fig. 41. Jump resonance.

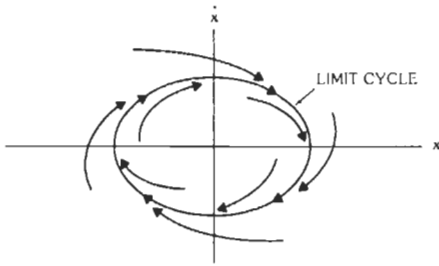


Fig. 42. Limit cycle.

Describing-Function Technique

The object of the describing-function technique is to represent the actual nonlinearity of the system in terms of an equivalent linear system by considering only the fundamental component of the output waveform of the nonlinear element subject to a sinusoidal input. This technique is valuable in the analysis and design of an important class of nonlinear feedback control systems, in which the output of the nonlinear element is acted upon as it travels around the closed loop by a linear low-pass filter.

The describing-function analysis is made on the following basic assumptions.

- A. The input to the nonlinear element n is a sinusoidal signal, and only the fundamental component of the output of n contributes to the input. The output response of a nonlinear element to a periodic signal consists of the fundamental frequency components of the input signal and its harmonics. Generally, the harmonic components are smaller in amplitude compared with the fundamental component. Further, in most control systems the system behaves as a low-pass filter, and the higher harmonics are attenuated. If the higher harmonics are sufficiently small, they can be neglected, and the equivalent linear approximation may be justified.
- B. There is only one nonlinear element in the system. All nonlinearities in the system are lumped into one single nonlinear element n . Fig. 43 shows a block diagram of a closed-loop system containing a nonlinear element n .

The describing function of a nonlinear element is defined as the ratio of the fundamental-frequency component of the output as a complex quantity (amplitude and phase angle) to the amplitude of the sinusoidal input signal. If the input signal as applied to the nonlinear element n is described by

$$e_m(t) = X \sin \omega t$$

the output response $e_o(t)$ may be expressed by the Fourier series

$$e_o(t) = (a_0/2) + a_1 X \sin \omega t + b_1 \cos \omega t + \sum_{n=2}^{\infty} a_n X \sin n\omega t + \sum_{n=2}^{\infty} b_n X \cos n\omega t + \dots$$

The $(a_0/2)$ term is the dc component; a_n and b_n are the harmonic components.

The fundamental-frequency component of the output may be expressed in terms of amplitude and phase angle as

$$e_{o1} = A(\omega, X) X \sin[\omega t + \phi(\omega, X)]$$

In this expression, $A(\omega, X)X$ is the amplitude, and $\phi(\omega, X)$ is the phase angle of the fundamental component. Both amplitude and phase angle are a function of the frequency and amplitude of the input signal. The describing function $N(\omega, X)$ by definition is

$$\begin{aligned} N(\omega, X) &= \{A(\omega, X) X \exp[j\phi(\omega, X)]/X\} \\ &= A(\omega, X) \exp[j\phi(\omega, X)] \\ &= A(\omega, X) \cos \phi(\omega, X) \\ &\quad + jA(\omega, X) \sin \phi(\omega, X) \quad (\text{Eq. 1}) \end{aligned}$$

The describing function $N(\omega, X)$ may be purely real or contain a phase angle depending on the type of nonlinearity. For single-valued nonlinear functions N is real, whereas for multivalued functions phase shift exists, generally lagging.

Calculation of a Describing Function—Calculation of the describing function involves performing a conventional Fourier analysis on the output waveform to obtain the fundamental component. The Fourier series expansion of the output waveform to an input sinusoidal $X \sin \omega t$ may be expressed as

$$e_o(t) = (a_0/2) + a_1 X \sin \omega t + b_1 X \cos \omega t + a_2 X \sin 2\omega t + b_2 X \cos 2\omega t + \dots$$

For the describing function, only the coefficients of the fundamental-frequency component are required. The coefficients may be obtained from the integrals

$$a_1 = (\pi X)^{-1} \int_0^{2\pi} f_0(t) \sin \omega t \cdot d(\omega t) \quad (\text{Eq. 2})$$

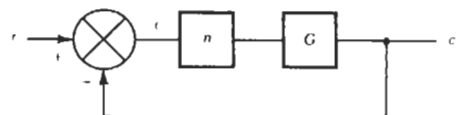


Fig. 43. Block diagram of nonlinear closed-loop system.

$$b_1 = (\pi X)^{-1} \int_0^{2\pi} f_0(t) \cos \omega t \cdot d(\omega t) \quad (\text{Eq. 3})$$

where $f_0(t)$ is the exact output of the nonlinear element expressed as a function of time. The describing function is then

$$|N(\omega, X)| = |a_1 + jb_1| = (a_1^2 + b_1^2)^{1/2} \angle N(\omega, X) = \tan^{-1}(b_1/a_1)$$

Where the exact output function $f_0(t)$ is known, the above method is applicable. If the function is not known, a graphical Fourier expansion can be performed on the output waveform. Two examples describing the procedure for calculation of the describing function using the graphical method are given.

Example 1: Saturation-Type Nonlinearity—A nonlinear element with saturation is shown in Fig. 44. Output y is held constant for input values greater than S . This region is called saturation or limiting. For input values less than S , the output behaves linearly with the input. The input-output relationship can be expressed by

- (A) $y = kx$ for $-S < x < S$
- (B) $y = kS$ for $x > S$
- (C) $y = -kS$ for $x < -S$

The output is an odd function, and thus only the sine term of the fundamental equation need be calculated. Furthermore, because of symmetry only the first quarter of the integration need be evaluated as follows:

$$a_1 = (4/\pi X) \int_0^{\pi/2} f_0(t) \sin \omega t \cdot d(\omega t)$$

If the input is $x = X \sin \omega t$, then the output $f_0(t) = y$ is expressed by

$$f_0(t) = kX \sin \omega t \quad X < S$$

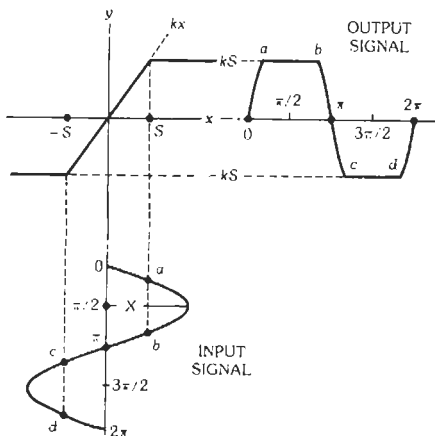


Fig. 44. Graphical representation of saturation.

and for X greater than S

$$f_0(t) = kX \sin \omega t \quad 0 < \omega t < \sin^{-1}(S/X)$$

$$= kS \quad \sin^{-1}(S/X) < \omega t < \pi/2$$

Therefore the coefficients a_1 are

$$a_1 = k \quad X < S$$

and for X greater than S

$$a_1 = (4/\pi X) \int_0^{\sin^{-1}(S/X)} kX \sin^2(\omega t) \cdot d(\omega t)$$

$$+ (4/\pi X) \int_{\sin^{-1}(S/X)}^{\pi/2} kS \sin \omega t \cdot d(\omega t)$$

$$= (2k/\pi)[\phi + (\sin 2\pi)/2]$$

where $\phi = \sin^{-1}(S/X)$.

The describing function N is given by

$$N = k \quad X < S$$

$$N = k(2/\pi)[\phi + (\sin 2\phi)/2] \quad X > S$$

The variation of amplitude of N with respect to X/S is plotted in Fig. 45. The phase angle is zero over the entire range.

Example 2: Backlash-Type Nonlinearity—For the second example, a simple backlash-type nonlinearity (Fig. 46) is evaluated. The backlash is a multivalued nonlinearity where the input-output relationship follows a different path dependent on the input-signal amplitude (curves 1, 2, and 3). After the steady state is established, the output $f_0(t)$ corresponding to different values of X of an input signal $x = X \sin \omega t$ are:

For $X < D$ (curve 1):

$$f_0(t) = 0$$

For $D < X < 2D$ (curve 2):

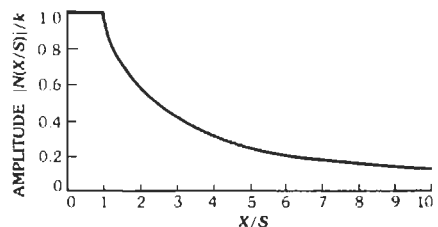


Fig. 45. Describing function for saturation (normalized amplitude).

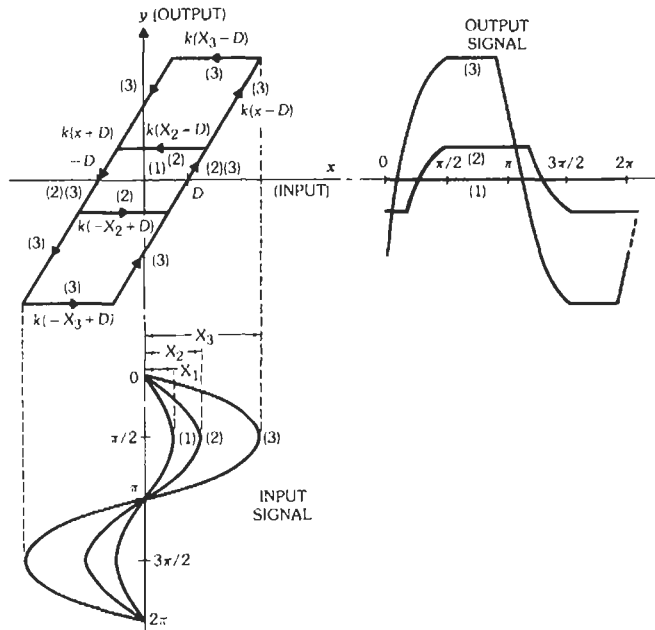


Fig. 46. Graphical representation of backlash nonlinearity. (Levinson, E. "Nonlinear Feedback Control Systems." Electro-Technology, September 1962; Fig. 36, p. 139.)

$$\begin{aligned}
 f_0(t) &= -k(X - D) & 0 < \omega t < \sin^{-1}[(2D/X) - 1] \\
 f_0(t) &= k(X \sin \omega t - D) & \sin^{-1}[(2D/X) - 1] < \omega t < \pi/2 \\
 f_0(t) &= k(X - D) & \pi/2 < \omega t < \pi + \sin^{-1}[(2D/X) - 1] \\
 f_0(t) &= k(X \sin \omega t + D) & \pi + \sin^{-1}[(2D/X) - 1] < \omega t < 3\pi/2 \\
 f_0(t) &= -k(X - D) & 3\pi/2 < \omega t < 2\pi
 \end{aligned}$$

For $X > 2D$ (curve 3):

$$\begin{aligned}
 f_0(t) &= k(X \sin \omega t - D) & 0 < \omega t < \pi/2 \\
 f_0(t) &= k(X - D) & \pi/2 < \omega t < \pi - \sin^{-1}[1 - (2D/X)] \\
 f_0(t) &= k(X \sin \omega t + D) & \pi - \sin^{-1}[1 - (2D/X)] < \omega t < 3\pi/2 \\
 f_0(t) &= -k(X - D) & 3\pi/2 < \omega t < 2\pi - \sin^{-1}[1 - (2D/X)] \\
 f_0(t) &= k(X \sin \omega t - D) & 2\pi - \sin^{-1}[1 - (2D/X)] < \omega t < 2\pi
 \end{aligned}$$

Solving for a_1 and b_1 of Eqs. 2 and 3, for the three conditions above, yields the describing-function terms

$$\begin{aligned}
 a_1 &= 0 & X < D \\
 a_1 &= (k/\pi)[\pi/2 + \theta + (\sin 2\theta)/2] & X > D \\
 b_1 &= 0 & X < D
 \end{aligned}$$

$$b_1 = (-k/\pi) \cos^2 \theta \quad X > D$$

where,

$$\theta = \sin^{-1}[1 - (2D/X)]$$

Or in terms of amplitude and phase angle

$$|N| = (k/\pi) \left\{ \left[\frac{1}{2}\pi + \theta + \frac{1}{2}(\sin 2\theta) \right]^2 + \cos^4 \theta \right\}^{1/2}$$

$$\angle N = \tan^{-1} \left[- \left(\frac{\cos^2 \theta}{\frac{1}{2}\pi + \theta + \frac{1}{2}(\sin 2\theta)} \right) \right]$$

The normalized amplitude and phase angle of the describing function N for different values of D/X are plotted in Fig. 47.

The describing function is calculated simply by determining the fundamental output component of the nonlinear element. Describing functions for some of the common nonlinear elements are given in Table 2.

Stability Analysis—The describing function N of the nonlinearity can be used to determine the stability of the system, providing the harmonics are sufficiently attenuated. In general, the describing function is a function of both frequency and amplitude of the input signal.

The closed-loop “transfer function” of the nonlinear feedback system in Fig. 48 is given by

$$(c/r)(\omega, X) = \frac{N(\omega, X)G(j\omega)}{1 + N(\omega, X)G(j\omega)}$$

The “characteristic equation” of the system is

$$1 + N(\omega, X)G(j\omega) = 0 \quad (\text{Eq. 4})$$

or

$$G(j\omega) = -[1/N(\omega, X)]$$

The condition of Eq. 4 must be satisfied for sustained oscillation of the output with zero input. Since $N(\omega, X)$ is a function of both frequency and amplitude, various combinations of ω and X can be found for oscillation. If there are no possible combinations satisfying the oscillation condition, the system is stable. In the case of sustained oscillation, the oscillatory mode may be either stable or unstable. If a slight disturbance in amplitude or frequency occurs and the oscillation returns to its original value, the oscillation is stable (stable limit cycle). If the oscillation amplitude increases or decreases from the original value, the oscillation is unstable (unstable limit cycle). The stability of the closed-loop system may be evaluated analytically by directly solving the characteristic equation by any one of the modified linear graphical methods.

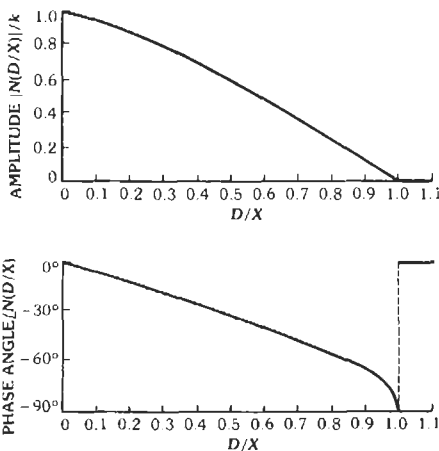


Fig. 47. Describing function for backlash (normalized amplitude and phase angle).

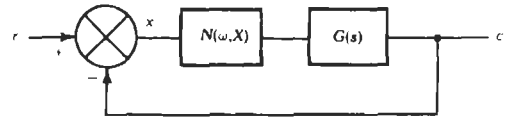


Fig. 48. Single-loop nonlinear system.

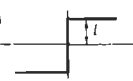


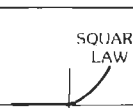
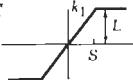
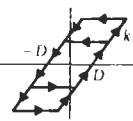
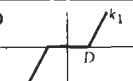
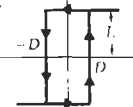
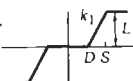
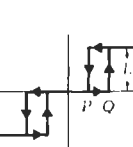
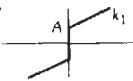
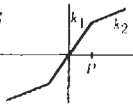
Polar Plot (Nyquist Diagram)—The conventional Nyquist diagram must be modified to apply the Nyquist stability criteria to the frequency-response plot. In a linear system, the crucial point on the Nyquist diagram is -1 . For nonlinear systems the $-[1/N(\omega, X)]$ locus corresponds to the critical point -1 . To evaluate the stability of the system, both $-[1/N(\omega, X)]$ and the $G(j\omega)$ function are plotted on the polar plane. The describing function $N(\omega, X)$ generally is a function of both ω and X . If N is only a function of X , there will be one locus $-[1/N(x)]$ plotted as a function of X . If N is also a function of ω , a family of constant-frequency loci are plotted for different values of ω (see Fig. 49).

The stability of the system is determined by the following relationship between the $-[1/N(\omega, X)]$ locus and the $G(j\omega)$ plot (Fig. 50). If the $-[1/N(\omega, X)]$ locus lies to the left of the $G(j\omega)$ plot or is not enclosed, the system is *stable*. Conversely, if the $-[1/N(\omega, X)]$ locus lies to the right of the $G(j\omega)$ plot or is enclosed, the system is *unstable*. If the $-[1/N(\omega, X)]$ locus intersects with the $G(j\omega)$ plot, the system may have a *sustained oscillation*. In the case where N is a function of ω , the condition for sustained oscillation is satisfied if the ω of the $G(j\omega)$ plot at the intersecting point is the same ω of the $-[1/N(\omega, X)]$ locus (see Fig. 51).

The oscillation may be either stable or unstable. If the $G(j\omega)$ intersects with the $-[1/N(\omega, X)]$ locus at one point only, the oscillation is stable (stable limit cycle). If more points of intersection exist, the limit cycle may be either stable or unstable. The stability of the limit cycle is determined by the direction of the two loci at the crossover point.

By establishing the $G(j\omega)$ locus pointing in the direction of increasing frequency as a reference, if the $-[1/N(X)]$ locus pointing in the direction of increasing amplitude X crosses the $G(j\omega)$ locus from right to left, the limit cycle is stable. If the crossover occurs from left to right, the limit cycle is unstable. A polar plot with both stable and unstable limit cycles is shown in Fig. 52.

TABLE 2. DESCRIBING FUNCTIONS FOR COMMON NONLINEAR ELEMENTS
 DESCRIBING FUNCTION: $N(X) = a_1 + jb_1$ or $|N(X)| = (a_1^2 + b_1^2)^{1/2}$, $\angle N(X) = \tan^{-1}(b_1/a_1)$

Characteristic	Describing-Function Coefficients		Characteristic	Describing-Function Coefficients	
A 	$a_1 = 4S/\pi X$ $b_1 = 0$		H 	$a_1 = k_1/2$ $b_1 = 0$	
B 	$a_1 = (4S/\pi X)\cos\theta$ $b_1 = 0$	$\theta = \sin^{-1}(D/X)$	I 	$a_1 = 4X/3\pi$ $b_1 = 0$	
C 	$a_1 = (2k_1/\pi)[\theta + (1/2)(\sin 2\theta)]$ $b_1 = 0$	$\theta = \sin^{-1}(S/X)$	J 	$a_1 = (k_1/\pi)[\pi/2 + \theta + (1/2)(\sin 2\theta)]$ $b_1 = -(k_1/\pi)\cos^2\theta$	$\theta = \sin^{-1}[1 - (2D/X)]$
D 	$a_1 = (2k_1/\pi)[\pi/2 - \theta - (1/2)(\sin 2\theta)]$ $b_1 = 0$	$\theta = \sin^{-1}(D/X)$	K 	$a_1 = (4L/\pi X)\cos\theta$ $b_1 = -(4L/\pi X)\sin\theta$	$\theta = \sin^{-1}(D/X)$
E 	$a_1 = (2k_1/\pi)[\psi - \theta + (1/2)(\sin 2\psi) - (1/2)(\sin 2\theta)]$ $b_1 = 0$	$\psi = \sin^{-1}(S/X)$ $\theta = \sin^{-1}(D/X)$	L 	$a_1 = (2L/\pi X)(\cos\theta + \cos\psi)$ $b_1 = (2L/\pi X)(\sin\psi - \sin\theta)$	$\psi = \sin^{-1}(P/X)$ $\theta = \sin^{-1}(Q/X)$
F 	$a_1 = k_1 + (4A/\pi X)$ $b_1 = 0$				
G 	$a_1 = k_2 - [(k_2 - k_1)/\pi](2\theta + \sin 2\theta)$ $b_1 = 0$	$\theta = \sin^{-1}(P/X)$			

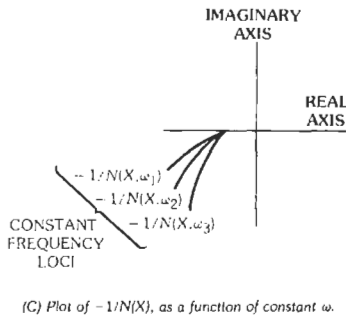
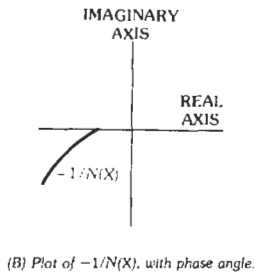
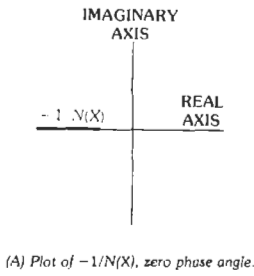


Fig. 49. Typical polar plots of various $N(\omega, X)$.

Gain-Phase Plot—The gain-phase plot is the direct transfer of the polar plot from polar coordinates to rectangular coordinates. The ordinate is the gain in decibels, and the abscissa is the phase angle in degrees.

The gain and phase angle of the two functions $G(j\omega)$ and $-[1/N(\omega, X)]$ are, for $G(j\omega)$

$$\text{Gain} \quad 20 \log |G(j\omega)|$$

$$\text{Phase angle} \quad \angle G(j\omega)$$

and for $N(X, \omega)$

$$\text{Gain} \quad -20 \log_{10} |N(\omega, X)|$$

$$\text{Phase angle} \quad -180^\circ - \angle N(\omega, X)$$

Typical gain-phase plots for various types of $N(\omega, X)$ are given in Fig. 53.

The system is stable if the $-(1/N)$ locus does not intersect with the $G(j\omega)$ plot. If the $-(1/N)$ locus intersects with the $G(j\omega)$ plot, the system has a sustained oscillation (Fig. 54).

In the case of sustained oscillation, there may be more than one point of intersection, as shown in Fig. 55. Points A and C are stable points (stable limit cycle), and point B is an unstable point (unstable limit cycle). The stability of the limit cycle is determined in a manner similar to that of the polar plot, except that if the $-(1/N(X))$ locus in the direction of increasing X crosses the $G(j\omega)$ locus point in the direction of increasing frequency from left to right, the limit cycle is stable; if it crosses from right to left, the limit cycle is unstable. Fig. 56 is a typical gain-phase plot of $G(j\omega)$ and $-[1/N(\omega, X)]$, where $N(\omega, X)$ is a function of ω . The family of $-(1/N)$ plots are the constant-frequency loci. Point A is the location for sustained oscillation.

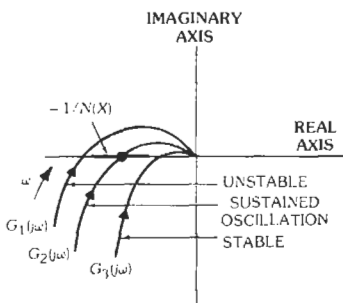


Fig. 50. Polar-plot stability criteria.

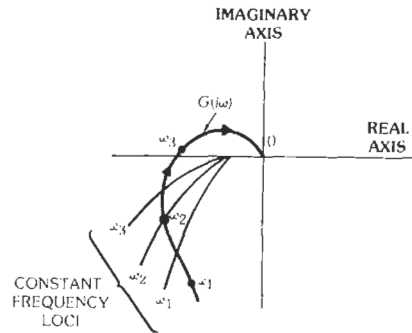


Fig. 51. Typical polar plot of $G(j\omega)$ and $-[1/N(\omega, X)]$ as a function of ω .

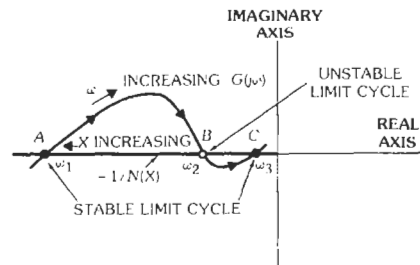
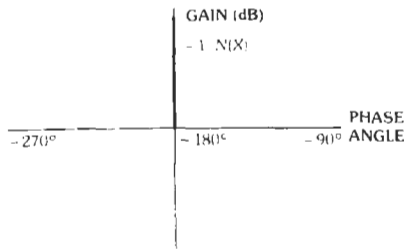
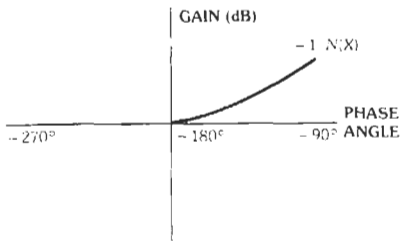


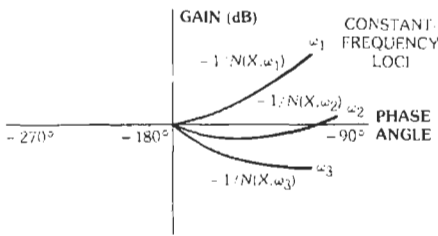
Fig. 52. Polar plot of stable and unstable limit cycles.



(A) Plot of $-1/N(X)$, zero phase angle.



(B) Plot of $-1/N(X)$, with phase angle.



(C) Plot of $-1/N(X)$ as a function of constant ω .

Fig. 53. Typical gain-phase plots of various $N(\omega, X)$.

Phase-Plane Method

The phase-plane method of analysis is used to study the transient behavior of second-order nonlinear systems. Differential equations for systems of higher order may be solved in the phase space; however, the results are complex and unwieldy, and not suited for graphical display. (Modern high-quality computer graphics can be used to advantage for higher-order problems).

Phase Plane—The differential equation describing a time-invariant second-order system may be expressed by

$$f_1(x, dx/dt)(d^2 x/dt^2) + f_2(x, dx/dt, t)(dx/dt) + f_3(x, dx/dt, t)x = g(t) \quad (\text{Eq. 5})$$

The equation in which t does not appear explicitly is called "autonomous." By defining

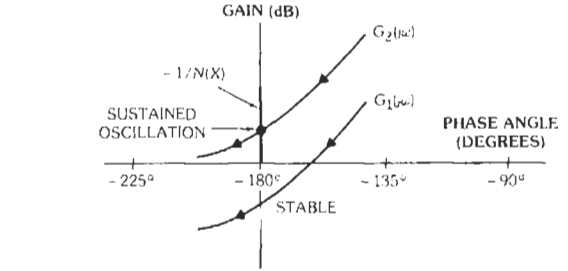


Fig. 54. Stability criteria for gain-phase plot.

$$\dot{x} = dx/dt$$

the equation may be rewritten as

$$f_1(x, \dot{x})(d\dot{x}/dt) + f_2(x, \dot{x})\dot{x} + f_3(x, \dot{x})x = 0$$

Or

$$dx/dt = P(x, \dot{x}) = \dot{x} \quad (\text{Eq. 6})$$

$$d\dot{x}/dt = Q(x, \dot{x}) = -[f_2(x, \dot{x})\dot{x} + f_3(x, \dot{x})x]/f_1(x, \dot{x}) \quad (\text{Eq. 7})$$

and, further:

$$\begin{aligned} (d\dot{x}/dt)/(dx/dt) &= d\dot{x}/dt = Q(x, \dot{x})/P(x, \dot{x}) \\ &= -[f_2(x, \dot{x})\dot{x} + f_3(x, \dot{x})x]/\dot{x}f_1(x, \dot{x}) \end{aligned} \quad (\text{Eq. 8})$$

The second-order differential equation with respect to time is converted to a first-order equation of x and \dot{x} .

The phase-plane diagram has the \dot{x} as its ordinate and x as its abscissa. The plot of \dot{x} as a function of x on the phase-plane diagram is termed *phase trajectory*. A family of phase trajectories is called the *phase portrait*.

The phase trajectory originates at a point corresponding to the initial condition (x_0, \dot{x}_0) and moves to a new location at each increment of time. Generally, the increments of time are not portrayed on the trajectory

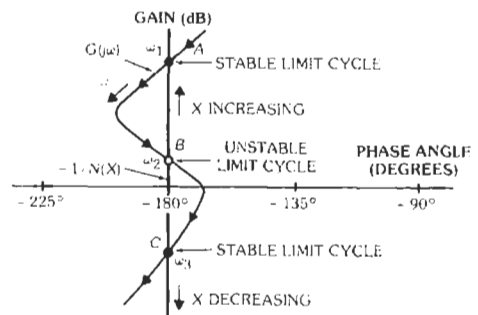


Fig. 55. Gain-phase plot of stable and unstable limit cycles.

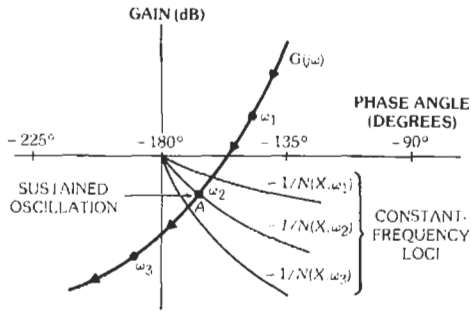


Fig. 56. Typical gain-phase plot of $G(j\omega)$ and $-1/N(\omega, X)$ as a function of ω .

and must be obtained by other means. If the value of time at each point on the trajectory is obtained, the time response of $\dot{x}(t)$ and $x(t)$ can be plotted. The phase trajectory has a definite direction associated with time. When \dot{x} is positive the trajectory moves from left to right, and for negative values of \dot{x} all paths move from right to left. If the trajectory approaches the origin or some finite point on the phase plane as time goes to ∞ , the system is stable. If the trajectory goes to ∞ with time, the system is unstable. If the trajectory approaches an enclosed path in the phase plane, the system has sustained oscillation. The enclosed path is called the limit cycle.

Construction of the Phase Portrait, Method of Isoclines—The slope of $d\dot{x}/dx$ of Eq. 8 is simply the slope of the trajectory in the phase plane. The locus of constants $d\dot{x}/dx$ is termed an *isocline* corresponding to the slope α ; that is

$$\alpha = d\dot{x}/dx = -[f_2(x, \dot{x})/f_1(x, \dot{x})] - [f_2(x, \dot{x})/f_1(x, \dot{x})]/(x/\dot{x}) = -g(x, \dot{x}) - h(x, \dot{x})(x/\dot{x})$$

The phase portrait is constructed by hand by plotting a large number of isoclines corresponding to the various slopes of the trajectory on the phase plane. All points located on the same isocline have the same slope α . Beginning at the location of the initial condition (x_0, \dot{x}_0) , the trajectory traverses in the clockwise direction, crossing each isocline at an angle corresponding to that slope α . Fig. 57 shows the isocline for a damped, linear, second-order system. Isoclines for first- and second-order linear differential equations are straight lines.

Singular Points: In a second-order system, the differential equation of the system may be described by two variables x and \dot{x} in the following form:

$$\begin{aligned} dx/dt &= P(x, \dot{x}) \\ d\dot{x}/dt &= Q(x, \dot{x}) \end{aligned}$$

The points where dx/dt and $d\dot{x}/dt$ vanish are called

singular points. At a singular point, the system is in a state of equilibrium.

The importance of a singular point in the phase plane is how the trajectories of the phase portrait behave in the vicinity of the singular point. When the trajectory converges to the singular point the system is stable, whereas if it diverges the system is unstable. Typical singular points are described below.

Types of Singular Points: Besides stable and unstable equilibria, the singular points may be classified into node, focus, center, and saddle points.

Consider, for example, a singular point at $x = a$ and $\dot{x} = b$ of Eqs. 6 and 7. At a singular point the derivatives dx/dt and $d\dot{x}/dt$ are both zero, and the location may be solved in the phase plane by setting equations 6 and 7 equal to zero. A singular point exists at $x = a$, $\dot{x} = b$, and the functions P and Q can be expressed in terms of the Taylor series about those points; then

$$\begin{aligned} dx/dt &= c_1(x - a) + c_2(\dot{x} - b) + c_3(x - a)^2 + \\ &\quad c_4(x - a)(\dot{x} - b) + c_5(\dot{x} - b)^2 + \dots \\ d\dot{x}/dt &= d_1(x - a) + d_2(\dot{x} - b) + d_3(x - a)^2 + \\ &\quad d_4(x - a)(\dot{x} - b) + d_5(\dot{x} - b)^2 + \dots \end{aligned}$$

If a sufficiently small region around the singular point is taken, the derivatives are dominated by the linear terms and hence quantities c_1 , c_2 , d_1 , and d_2 . By changing the variables, the singular point may be moved to the origin. Then the system equation may be rewritten as

$$\begin{aligned} dx/dt &= p_1x + p_2\dot{x} = 0 \\ d\dot{x}/dt &= q_1x + q_2\dot{x} = 0 \end{aligned}$$

and the characteristic equation is

$$\lambda^2 - (p_1 + q_2)\lambda + (p_1q_2 - p_2q_1) = 0$$

The roots of the characteristic equation determine the nature of the critical points. The roots are

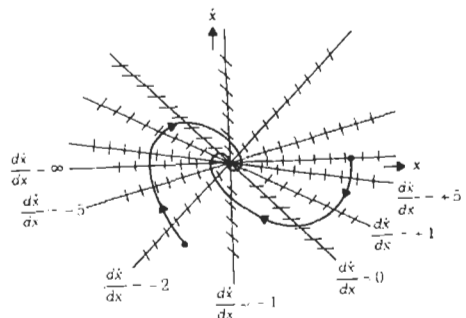


Fig. 57. Isocline method.

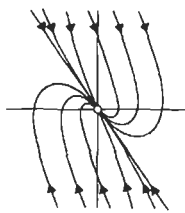
$$\lambda = \frac{1}{2}(p_1 + q_2) \pm [(p_1 + q_2)^2 - 4(p_1q_2 - p_2q_1)]^{1/2}$$

There are six possible cases for the six types of singular points.

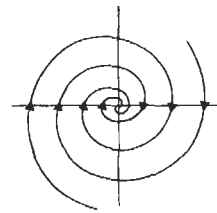
- A. The roots are real and are both negative. If the initial condition is $x_0 \neq 0, \dot{x}_0 = 0$, the trajectory approaches the singular point without an overshoot or oscillation and is called a *stable node* (Fig. 58A).
- B. The roots are complex conjugate with negative real parts. The trajectory displays a spiraling response as it converges to the singular point as shown in Fig. 58B. This is called the *stable-focus* type of singularity.
- C. The roots are conjugate and pure imaginary. The response exhibits a sustained oscillatory motion with the amplitude dependent on the initial condition. The trajectory displays a family of ellipses about the singular point, and is termed *center* (Fig. 58C).

- D. The roots are both positive real. The response in the time domain increases exponentially and is unstable. The portrait is the same as the stable node except the trajectory diverges from the singular point. This is termed *unstable node* (Fig. 58D).
- E. The roots are complex conjugate with positive real parts. The phase portrait is the same as the stable focus except the trajectory diverges from the singular point. This is termed the *unstable focus* (Fig. 58E).
- F. The roots are real with one negative and the other positive. The phase portrait consists of a family of curves of the hyperbolic type having $k_1 = (\lambda_1 - p_1)/p_2$ and $k_2 = (\lambda_2 - p_1)/p_2$ for its asymptotes. The direction of the paths is toward the singular point on the negative asymptote and away from the singular point on the positive asymptote. Singular points of this type are called *saddle points* and are unstable (Fig. 58F).

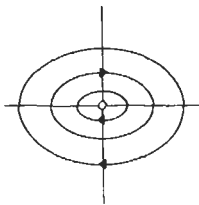
The six types of singularities correspond to the six regions of Fig. 59.



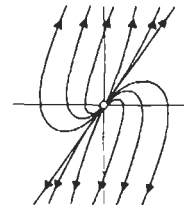
(A) Stable node.



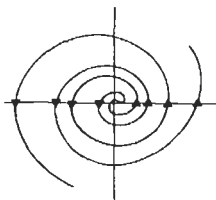
(B) Stable focus.



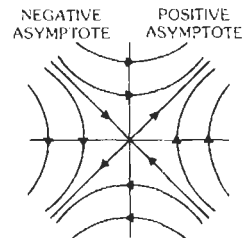
(C) Center



(D) Unstable node.



(E) Unstable focus



(F) Saddle point.

Fig. 58. Types of singular points.

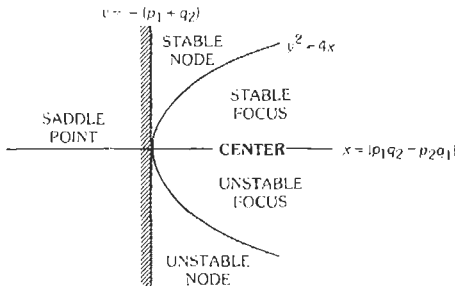


Fig. 59. Regions of various singular points.

COMPUTER-AIDED ANALYSIS AND DESIGN

The development of software for computer-aided control-system analysis and design has revolutionized control engineering. Several software packages are currently available commercially for carrying out the analysis and design methodologies described in this chapter. These packages provide high-quality graphics for displaying system time responses, as well as Bode plots and root loci. User-friendly interfaces provide the ease of interaction with these packages that is essential for an efficient design environment. Typically, these packages include both classical and state-space oriented methodologies. The underlying numerical algorithms are based on the latest research in numerical analysis. Currently available packages include Matlab, Matrix-X, Program CC, Simnon, Easy-5, Ctrl-C, and L-A-S, among others. See the end of the References section for further information on software packages.

BIBLIOGRAPHY

There are many good introductory texts in the field of control. Among these are:

Kuo, B. C. *Automatic Control Systems*, 6th ed. Englewood Cliffs, NJ: Prentice Hall, 1991.

Franklin, G. F., Powell, J. D., and Emami-Naeini, A. *Feedback Control of Dynamic Systems*. Reading, MA: Addison-Wesley, 1991.

Dorf, R. C. *Modern Control Systems*, 5th ed. Reading, MA: Addison-Wesley, 1989.

Phillips, C. L., and Harbor, R. D. *Feedback Control Systems*, 2nd ed. Englewood Cliffs, NJ: Prentice-Hall, 1991.

For control system components, sensors, and actuators, see:

deSilva, C. W. *Control Sensors and Actuators*. Englewood Cliffs, NJ: Prentice-Hall, 1989.

See also the journal *Sensor and Actuators* for up-to-date developments in this rapidly-moving field.

Modeling of a variety of physical systems and components for control purposes is discussed in:

Connon, R. H., Jr. *Dynamics of Physical Systems*. New York: McGraw-Hill, 1967.

Ogata, K. *System Dynamics*. Englewood Cliffs, NJ: Prentice-Hall, 1978.

The classical reference on sensitivity is:

Bode, H. W. *Network Analysis and Feedback Amplifier Design*. New York: Van Nostrand, 1945.

See also

Horowitz, I. *Synthesis of Feedback Systems*. New York: Academic Press, 1963.

An extension of the Bode sensitivity approach to multivariable systems (comparison sensitivity) is given in

Perkins, W. R. "Sensitivity Analysis," in J. B. Cruz, Jr., ed. *Feedback Theory*. Chapter 2, pp. 19-65. New York: McGraw-Hill Book Co., 1972.

Cruz, J. B., Jr. and Perkins, W. R. "A new approach to the sensitivity problem in multivariable feedback systems design." *IEEE Trans. Automat. Control*, July 1984, pp. 216-223.

Textbooks emphasizing the state-space approach include:

Friedland, B. *Control System Design: An Intro. to State-Space Methods*. New York: McGraw-Hill, 1986.

Anderson, B. D. O., and Moore, J. B. *Optimal Control: Linear Quadratic Methods*. Englewood Cliffs, NJ: Prentice-Hall, 1990.

Lewis, F. L. *Optimal Control*. New York: Wiley-Interscience, 1986.

Digital control is the subject for the following texts:

Kuo, B. C. *Digital Control Systems*, 2nd ed. Saunders, 1992.

Åström, K., and Wittenmark, B. *Computer-Controlled Systems: Theory and Design*, 2nd ed. Englewood Cliffs, NJ: Prentice-Hall, 1990.

Franklin, G. F., Powell, J. D., and Workman, M. L., *Digital Control of Dynamic Systems*, 2nd ed. Reading, MA: Addison-Wesley, 1990.

Auslander, D. M., and Tham, C. H. *Real-Time Software for Control: Program Examples in C*. Englewood Cliffs, NJ: Prentice-Hall, 1990.

Nonlinear systems are the subject of intense current interest. Useful texts include:

Vidyasagar, M. *Nonlinear Systems Analysis*. Englewood Cliffs, NJ: Prentice-Hall, 1978.

Khalil, H. K. *Nonlinear Systems*. New York: Macmillan, 1992.

Current research in control theory and applications is reported in journals such as *IEEE Trans. on Automatic Control*, *IEEE Trans. on Industrial Electronics, Automatica*, and the *International Journal on Control*, as a sampling. Very useful to designers are the "trade" magazine *Control Engineering* and the *IEEE Control Systems Magazine*. These magazines contain advertisements for the computer-aided design packages mentioned in the section "Computer-Aided Analysis and Design." Some control texts provide software on back-of-the-book disks.

16 Electron Tubes

Marvin Chodorow and Robert Symons

*Revised by David Abe, Richard Abrams, Bruce Danly, Henry Freund, Kevin Jensen,
Baruch Levush, Robert Myers, Robert Parker, Jonathan Shaw, Arnold Shih, Joan Yater*

- Electron Emission 16-3
 - Thermionic Emission
 - Secondary Emission
 - Field Emission
 - Macroscopic Current Density
 - Field Emitter Arrays
 - Current–Voltage Relationship, Transconductance, and Modulation
 - Protective Resistance
 - Field Emission Displays
 - Photocathodes

- Electrode Dissipation 16-16
 - Radiation Cooling
 - Water Cooling
 - Forced-Air Cooling
 - Evaporative Cooling
 - Conduction Cooling
 - Grid Temperature

- Noise in Tubes 16-18
 - Shot Effect
 - Partition Noise
 - Flicker Effect
 - Collision Ionization
 - Induced Noise
 - Miscellaneous Noise
 - Microwave Tubes

- Low-, Medium-, and High-Frequency Tubes 16-19
 Coefficients
- Materials and Structures 16-21
 Cathodes
 Grids
 Anodes
 Getters
 Tube Geometry
- Microwave Tubes 16-23
 Terminology
 Recent Trends
 Linear-Beam Tubes
 Crossed-Field Tubes
- Cyclotron Resonance Microwave Tubes 16-38
- Free-Electron Lasers 16-39
- Gas Tubes 16-41
 Characteristics of Gas Tubes
 Power Applications of Gas Tubes
 Microwave Applications of Gas Tubes
- Light-Sensing and Light-Emitting Tubes 16-45
 Radiometry and Photometry
 Flux Units
 Optical Imaging
 Typical Approximate Illumination Values at the Surface of the Earth
 Typical Approximate Brightness Values
- Light-Sensing Tubes 16-48
 Image Tubes and Image Intensifiers
 Gas Photodiodes
 Image Orthicons
 Vidicons
 Variations of the Vidicon
- Light-Emitting Tubes 16-55
 Cathode-Ray Tubes
 Storage Cathode-Ray Tubes
 Bistable Storage Tube

ELECTRON EMISSION

All electron tubes* depend for their operation on the flow of electrons within the tube, through either high vacuum or an ionized gas. The electrons are emitted from a cathode surface as a result of one of four processes that are distinguished on the basis of the mechanism by which the electrons are enabled to leave the surface. These processes are elevated temperature (thermionic or primary emission); bombardment by other particles, generally electrons (secondary emission); the action of a high electric field (field emission); or the incidence of photons (photoemission).

Thermionic Emission

Thermionic emission originates from the thermally excited electrons that have sufficient kinetic energy to overcome the vacuum barrier. The sum of all these thermal electrons gives rise to the Richardson-Dushman equation, which is

$$J = AT^2 \exp(-\phi/kT) \quad (\text{Eq. 1})$$

where

- J is emission density in amperes/cm²,
- A is 120 amperes/cm²,
- ϕ is work function in electron volts (eV),
- T is temperature in kelvins (K), and
- k , the Boltzmann constant, is 8.6164×10^{-5} eV K⁻¹.

Emission density is an exponential function of the work function. Between 1000 K and 2000 K, a decrease in the work function of 0.2 to 0.4 eV causes an order of magnitude increase in the emission density. The work function of commonly used thermionic cathodes varies from 4.6 eV (tungsten) to 1.5 eV (oxide-coated cathodes).

The first thermionic cathodes used in quantity were made from pure tungsten (W). Because of the high work function, a very high temperature is required to provide a reasonable emission density. For example, Eq. 1 indicates that an emission density of 0.3 A/cm² requires a temperature of 2150°C (or 2423 K).

To reduce the operating temperature, thoria (thorium oxide) was mixed with W. Thorium would diffuse to the surface of W and lower the work function. The cathode surface was sometimes carburized to reduce the thorium evaporation, extending the operating life. The resulting surface has a work function of ~ 3.2 eV.† For comparison with pure tungsten, the emission den-

sity of 0.3 A/cm² requires 1950°C (2223 K), a 200°C reduction in temperature.

A significant reduction in work function was realized in oxide-coated cathodes. The oxide-coated cathodes consist of an alkaline-earth-oxide coating on a metal substrate, which is usually nickel (Ni). A standard ASTM (American Society for Testing and Materials) mixture of the oxide coating is 49, 44, and 7 atomic percent of Ba, Sr, and Ca oxide, respectively. The substrate can be made of either active or passive Ni, depending on whether or not it contains reducing agents. A typical active Ni contains 4% W in the bulk. The work function of a fully activated oxide-coated cathode is 1.85 eV and 1.5 eV on passive and active Ni substrates, respectively.‡

The initial alkaline earths are in carbonate form and mixed by an organic binder. The carbonates and binder are either sprayed onto the substrate or they are cut and transferred from a self-supporting sheet, which is known as a "Sarong" coating. Activation of an oxide-coated cathode requires several hours of outgassing. A copious quantity of hydrocarbon gases is released initially from the decomposition of the organic binder, and then a large amount of CO₂ is released from the decomposition of the carbonates into the respective oxides. It is advisable to raise the cathode temperature gradually during the carbonate conversion to avoid CO₂ overpressure, which could cause melting of the oxide coating. On some oxide-coated cathodes, high work function patches were observed to form, which were identified as the formation of CaO slags on the surface caused by BaO/BaCO₃ eutectic melting during the conversion process.§

The activation and poisoning processes of an oxide-coated cathode are explained by a semiconductor model." The oxygen vacancies in the oxide lattice act as electron donors. During activation, the reducing agents, such as W in the Ni substrate, react and bind with the oxygen atoms in the oxides and generate oxygen vacancies. The creation of these donors shifts the Fermi level up from the intrinsic mid-band gap to near the bottom of the conduction band, and consequently, the work function is lowered. Thermal excitation lifts some of the electrons in the donors to the conduction band, and the mobile electrons are needed in the electron emission process. Cathode poisoning by oxidizing gases such as O₂, H₂O, and CO₂ occurs through filling the oxygen vacancies and shifting the Fermi level back toward the mid-band gap. Sulfur, which is present on the surface of many metals, can have a detrimental effect on the oxide-coated cathodes since S can also fill the oxygen vacancies. The well-known 10-volt slump occurs when S is transferred from the metal grids to the oxide coating through electron-stimulated

* J. W. Gewartowski and H. A. Watson, *Principles Of Electron Tubes* (Princeton, NJ: Van Nostrand, 1965). K. R. Spangenberg, *Vacuum Tubes* (New York: McGraw-Hill Book Co., 1948). A. H. W. Beck, *Thermionic Valves, Their Theory and Design* (London: Cambridge University Press, 1953). *Standards on Electron Tubes: Definitions of Terms* (New York: Institute of Radio Engineers, 1950).

† Reference 1.

‡ Reference 2.

§ Reference 3.

" Reference 4.

desorption. An oxide-coated cathode, once poisoned by S, cannot be reactivated.*

The major limitations of oxide-coated cathodes are low dc emission and susceptibility to damage. Because moderate electrical conductivity is inherent to semiconductive alkaline earth oxides, the dc emission of the oxide-coated cathodes is limited to less than 1 A/cm^2 . Pulsed emission densities are dependent on pulse length and can be as high as tens of A/cm^2 .† After the conversion from carbonates to the oxides, it is imperative to maintain a good vacuum since reactivation after exposure to oxidizing gases is often difficult.* Exposure to atmosphere is always detrimental since the coating expansion accompanying the hydroxide formation would cause the oxide coating to peel.

Dispenser cathodes are now the most commonly used high emission-density cathodes. Instead of a thick oxide coating, only a mono-layer or less of Ba and O covers the metal substrate of a dispenser cathode. Since the cathode body is made of W metal, the dc emission is no longer limited by the electrical conductivity. During the operation of a dispenser cathode, a fresh supply of active materials (Ba and O) is continuously dispensed onto the emitting surface, replenishing the loss due to evaporation. Consequently, the dispenser cathodes are quite robust in comparison to the oxide-coated cathodes. During the fabrication of a dispenser cathode, W powder with an average grain size of $4.5 \mu\text{m}$ is pressed to form the cathode body of a W-matrix with 20% porosity (i.e., 80% density of W). The pores of the W-matrix are then impregnated with a melt of mixed calcined $\text{BaO/CaO/Al}_2\text{O}_3$ drawn into the pores by capillary action in an atmosphere of dry hydrogen. Consequently, dispenser cathodes are also known as impregnated cathodes or matrix cathodes. Fig. 1 illustrates how the pores in the W-matrix form continuous paths, which allow the impregnation process to occur and which serve as a conduit in dispensing the active material. The most commonly used impregnant mixtures are B-type, S-type, and I-type.‡ The B-type impregnants consist of a mixture of 5 moles of BaO, 3 moles of CaO, and 2 moles of Al_2O_3 , with a designation of 5-3-2. The S-type impregnants consist of a 4-1-1 mixture. Less frequently used are the I-type impregnants, which consist of a 6-1-2 mixture. An important effect of the difference in the impregnant mixtures is the Ba evaporation rate. The higher the BaO-to-CaO ratio, the higher the Ba evaporation rate is. In applications where a fast turn-on after shelf storage is essential (e.g., expendable decoys), it is advantageous to choose impregnants with a higher Ba evaporation rate.§ However, excessive Ba evaporation is often undesirable. Ba deposition on the grids leads

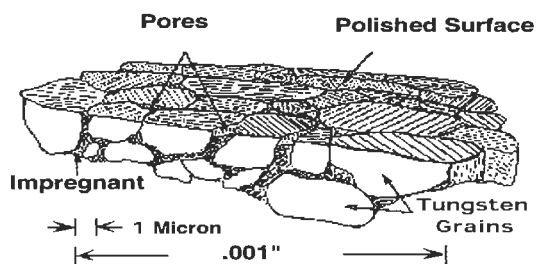


Fig. 1. Scale of dispenser-type cathode features near the emitting surface:

average W-grain diameter $\sim 4.5 \mu\text{m}$,

approximately 3×10^6 pores/ cm^2 ,

pore equivalent diameter $\sim 3 \mu\text{m}$,

average inter-pore distance $\sim 6 \mu\text{m}$.

(From M. C. Green, "Dispenser Cathode Physics," Final Technical Report RADC-TR-81-211, July 1981, 9, Fig. 1.)

to unwanted grid emission; a higher Ba ion concentration is also correlated with noisier tube performance.

There are two types of commercially available off-the-shelf dispenser cathodes, namely the standard-type and the *M*-type. The difference lies in the absence or presence of a metallic coating. As already discussed, the standard-type dispenser cathodes are classified as B-type, S-type, and I-type cathodes, depending on the impregnant mixture. *M*-type cathodes are standard-type cathodes coated with $\sim 5000 \text{ \AA}$ noble metal on top of the emitting surface. The noble metals can be Os, Os-Ru, Ir, or Re. The most common *M*-type cathode has a sputter-deposited Os/Ru coating with an atomic ratio of 80/20." At full activation, the effective work function of a standard-type dispenser cathode is about 2.1 eV and that of a *M*-type cathode is about 1.9 eV.# The low work functions are the effect of the surface layer, which consists of Ba, O, and substrate atoms. The charge transfer and the relative atomic positions result in a net dipole moment, which is directed outward from the surface and reduces the work function of the metal surface. The type of substrate atom changes the degree of charge transfer, the surface coverage, and the atomic arrangement, leading to the observed difference in the work functions.**

The emission from *M*-type cathodes can be further improved by the deposition of alloys of optimized composition, for example, 50% W and 50% Ir. The design of a structure that is stable at the optimum surface composition leads to MMM-type (Mixed Metal Matrix) or CD-type (Controlled Doping) cathodes.††

* Reference 4.

† Reference 1.

‡ Reference 5.

§ Reference 6.

" Reference 7.

Reference 7 & 8.

** Reference 9.

†† Reference 10.

The cathode body of a MMM-type cathode is a sintered plug of W and noble metal (usually Ir) particles. The surface coating of a CD-type cathode consists of deposited multiple layers of alloy material, each of different composition. However, these types of cathodes are no longer commercially available as off-the-shelf items.

Scandate cathodes have been under study for many years. In early years, their performance had not been stable or reproducible. However, recent developments at Philips Research Laboratories have demonstrated consistent emission for several thousand hours in diode tests from top-layer type scandate cathodes, which were fabricated by depositing a layer of Sc_2O_3 and W mixture on a S-type cathode. At 1030°C , emission densities of 100 A/cm^2 and 400 A/cm^2 have been demonstrated for top-layer scandate cathodes prepared by a mixed powder method* and laser ablation deposition,† respectively. These emission densities correspond to effective work functions of 1.63 eV and 1.48 eV, respectively. Even in applications that require an emission density lower than 100 A/cm^2 , the use of scandate cathodes allows operation at a lower temperature, which will result in lower Ba evaporation and longer life.

Fig. 2 shows the empirical relation between the emission density and life expectancy of B-type and M-type cathodes. The major life-limiting factor of a dispenser cathode is Ba depletion from the pores. A higher emission density requires a higher operating temperature, which results in a faster rate of Ba depletion, and consequently a shorter life. Since an M-type cathode has a lower work function than a B-type cathode, at a given emission density the operating temperature of the M-type cathode is lower, resulting in a longer life. At the high emission density end of the figure, the data of B-type and M-type cathodes converge.

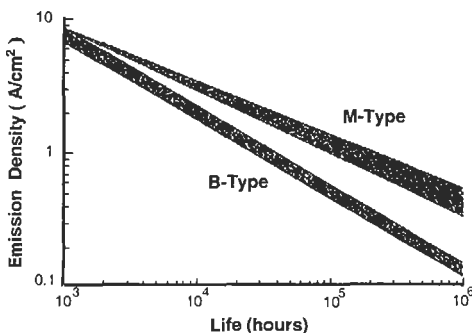


Fig. 2. Operating emission densities that can be expected as a function of desired cathode life. (From A. S. Gilmour, Jr., *Microwave Tubes*. Norwood, MA: Artech House, Inc., 1986, page 132, Fig. 5.27.)

* Reference 11.

† Reference 12.

At these operating temperatures, interdiffusion between the overlayer Os/Ru atoms and the substrate W atoms becomes rapid. The emission characteristics of M-type cathodes asymptotically approach that of B-type cathodes. At an emission density above 10 A/cm^2 , both B-type and M-type cathodes have too short a life expectancy to be of practical use. A dispenser cathode with a lower work function such as the scandate cathode is needed. However, further developments must be made before scandate cathodes find wide applications. One of the major obstacles is its susceptibility to ion damage. Ion bombardment could cause emission degradation through sputter removal of the surface atoms. With other types of dispenser cathodes, heating to a higher temperature completely restores the original emission characteristics. It was found that thermal treatments were unable to reactivate scandate cathodes after ion beam damage because the lost surface scandium atoms were not readily replenished. A method of replenishing scandium must be developed and incorporated in the cathode fabrication for scandate cathodes to have a reliable performance in a practical tube environment.

Secondary Emission

When the surface of a solid is bombarded by charged particles having appreciable velocity, electrons are emitted from the solid. This is the process of secondary emission.‡§# In the case where the bombarding particles are also electrons, which is the focus of this discussion, the incident and emitted electrons are differentiated by the terms primary and secondary, respectively. Secondary emission processes under electron bombardment play an essential role in vacuum electronic devices. The materials used in the devices may need to be judiciously selected in some cases to enhance the secondary emission or in other cases to suppress such emission. High secondary-emission materials are desirable for grids in electron multipliers and for cathodes in crossed-field devices, such as magnetrons and crossed-field amplifiers. In contrast, low secondary-emission materials are desirable for depressed collectors in microwave- and millimeter-wave power tubes to ensure a high efficiency in the energy conversion. Low-emission materials are also sought for coating the grids and tube walls to prevent RF vacuum breakdown.

‡ H. Bruining, *Physics and Applications of Secondary Electron Emission* (New York: McGraw-Hill Book Co., 1954).

§ O. Hackenberg and W. Brauer, *Secondary Electron Emission from Solids*, Advances in Electronics and Electron Physics, Vol. XI (New York: Academic Press, 1959).

A. J. Dekker, *Secondary Electron Emission*, Solid State Physics, Vol. 6 (New York: Academic Press, 1958).

R. O. Jenkins and W. G. Trodden, *Electron and Ion Emission* (New York: Dover, 1965) p. 54.

The secondary emission measured from a material includes *all* electrons emitted from the surface. These electrons comprise three groups: (a) true secondaries, (b) inelastically reflected primaries, and (c) elastically reflected primaries. True secondaries are considered to be those of the solid that have been excited above the energy level required for escape across the surface barrier. The three groups are separable to a degree on the basis of energy as indicated in the energy distribution curve of Fig. 3. True secondaries constitute the bulk of emitted electrons at moderate primary energies, and their distribution is almost independent of primary energy. Electrons in the relatively flat interval between *a* and *b* constitute a mixture of true secondaries and inelastically reflected primaries. It has become customary to arbitrarily designate those emitted electrons having energies less than 50 eV as true secondaries.

In practice the secondary emission characteristics of a material are described by the secondary yield coefficient δ (defined as the ratio of the secondary-electron current to the primary-electron current) and the emitted energy distribution. These emission characteristics are influenced by the properties of both the material and the primary electrons. To understand these influences, it is necessary to consider the three-step secondary-emission process: (1) generation of internal secondary electrons by kinetic impact of the primary electrons, (2) transport of the internal secondary electrons through the bulk to the surface, and (3) escape of the secondary electrons across the solid-vacuum interface. For example, the energy-dependent yield is a product of the competing generation and escape processes. First, primary electrons penetrate into the material and lose energy through collisions with electrons and ions. In the process, kinetic energy is transferred to internally generated secondary electrons. Each primary electron is able to generate more secondary electrons as its energy increases, and consequently the secondary yield exhibits an initial increase with increasing primary energy. However, as the primary energy increases, the primary electrons penetrate deeper into the material where most of the secondary electrons are generated. The secondary electrons cannot escape if they are generated too deeply in the material, and

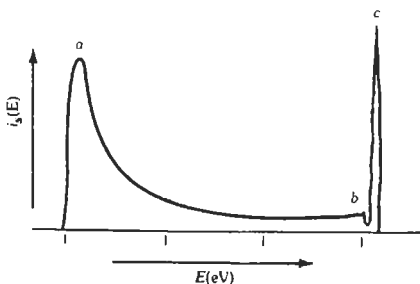


Fig. 3. Total energy distribution of secondary electrons.

hence the yield exhibits a decrease at very high primary energies. The resulting bell-shaped yield curve is seen in Fig. 4 and is observed for most materials. Significant points on the yield curve are the maximum yield, which occurs at a primary energy for which the penetration depth is approximately equal to the escape depth, and the first and second crossovers at which the yield becomes unity. The secondary yield also varies with incident angle, and this variation can be understood from the same process described above. If the primary electron is incident at an angle θ to the surface, the maximum penetration depth is reduced by a factor $\cos\theta$ relative to that at normal incidence ($\theta = 0$), assuming a forward path for the primary electron. Consequently, the yields are larger at more oblique angles, with this enhancement most pronounced at higher primary energies.

Energetic primary electrons generate many internal secondary electrons in all materials. However, high yields are not observed from most materials because the internal secondary electrons lose energy through collisions with electrons and ions and fall below the vacuum level. The dominant energy-loss mechanisms differ in metals and insulators, which results in different secondary emission characteristics for the two types of materials. In metals, secondary electrons lose energy primarily through strong inelastic scattering with conduction electrons. As a result, very few secondary electrons reach the surface with a sufficient kinetic energy to overcome the work function ($\sim 4-5$ eV). The secondary yields from metals are therefore small, with the maximum yield being on the order of unity and varying between 0.5 (for Li) and 1.8 (for Pt).^{*} Because the yield is relatively insensitive to changes in the metallic work function,[†] the secondary

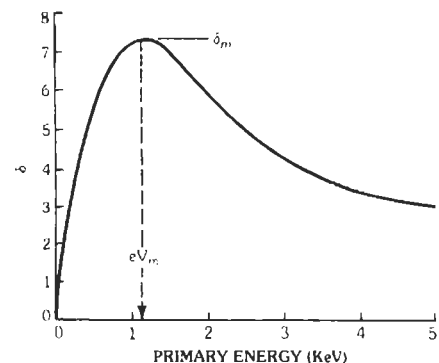


Fig. 4. Secondary-emission yield curve.

* A. J. Dekker, *Secondary Electron Emission*, Solid State Physics, Vol. 6 (New York: Academic Press, 1958).

† R. O. Jenkins and W. G. Trodden, *Electron and Ion Emission* (New York: Dover, 1965) p. 54.

emission from metals is relatively stable and long lived. Low yields ($\delta < 1$) are also obtained from semi-metallic graphite and amorphous carbon, and these materials are commonly used in low-yield applications.

Very few conduction electrons are present in insulators, and therefore secondary electrons lose energy primarily through the excitation of valence electrons into the conduction band. This electron-electron scattering interaction becomes forbidden if the secondary-electron energy is less than the energy gap, which is very large in insulators (~ 5 – 10 eV). As a result, secondary electrons lose energy more slowly as they move through insulators than they do through metals, and the escape depth becomes large. Upon reaching the surface, secondary electrons must have sufficient kinetic energy to overcome the electron affinity, which is the difference between the vacuum level and the conduction-band minimum. In insulators, the electron affinity is typically on the order of 1 eV. Because of the larger escape depth and lower surface barrier, the secondary yields from insulators are generally higher than from metals; for example, the maximum yield is 6.8 for NaCl and 25 for single-crystal MgO.*

The energy distribution of secondary electrons emitted from metals and insulators has generally been assumed to resemble the measurements shown in Fig. 3. From this curve, the secondary-electron energies appear to be relatively high (\sim tens of eV) and very broadly distributed. In the case of metals, this energy distribution is consistent with the broad internal distribution produced by inelastic scattering with conduction electrons. However, recent studies have revealed that secondary electrons generated in nonmetals (e.g., wide band gap materials) have a very narrow, low-energy distribution, and the emitted energy distribution depends strongly on the energy barrier at the surface. For example, in secondary emission measurements from diamond surfaces having *no* energy barrier (i.e., a negative electron affinity), the energy distribution is sharply peaked less than 1 eV above the conduction-band minimum with a FWHM of ~ 0.5 eV.† For primary energies above 1000 eV, the low-energy peak contains more than 90% of the total emitted electrons, and yields greater than 100 are measured. Once the electron affinity becomes greater than 1 eV, however, the very low energy electrons are blocked from emission and the yields decrease to values of about 1–3. In this case, only the high-energy tail of the low energy peak is able to escape and the emitted electron energy distribution resembles the “typical” distribution shown in Fig. 3.

It is clear that the secondary emission from insulators depends strongly on the nature of the surface bar-

rier, unlike the case of secondary emission from metals. At some insulator surfaces, the electron affinity can be made small or even negative, in which case the yields become extremely high. However, the secondary emission from such surfaces is very sensitive to surface contamination or damage and may not be stable or long lived. Secondary emission from insulators can also be affected by charging problems. Unlike metals, insulators do not have sufficient electrical conductivity to replenish the emitted electrons, and sample charging decreases the secondary yield to 1. To overcome charging problems, different approaches can be used to provide electrical conductivity in insulator materials. These include the fabrication of a thin insulator layer (e.g., < 50 Å metal-oxide film) on a metal substrate, the incorporation of metallic crystallites in the insulator bulk, and appropriate doping in single-crystal or polycrystalline insulator material. However, each of these approaches presents challenges in practical applications. For example, a thin oxide layer is eroded in the harsh environment of a crossed-field device, and hence a replenishment mechanism must be incorporated into the device. The size and distribution of crystallites in composite materials are difficult to control, and such materials suffer from erosion as well as localized charging problems. Insulator doping represents a promising approach, especially with recent advances in material growth capabilities. However, only a few selective material systems can be appropriately doped at this time, and the robustness and reliability of these materials have not yet been adequately established.

Secondary emitter materials typically used in microwave tubes include pure Pt metal and oxidized metals such as Mg, Be, and Al. For example, BeO is currently used as the secondary emission cathode in the crossed-field amplifiers of the AEGIS system. BeO produces high yields (~ 3 – 4), but the emitting layer must be replenished by use of an oxygen source. In other high-yield applications, Ag–Mg or Cu–Be alloys are processed to provide a high-yield, partly conductive surface film. Typical yields of 2.5 to 4 are obtained at 100 eV primary energy from the metal oxide (i.e., MgO or BeO) on the surface. In addition to these conventional materials, novel materials are being investigated for future device applications. In particular, doped semiconductor and insulator materials have produced extremely high yields (~ 100 – 1000) when the surface is prepared with a small or negative electron affinity. These include hydrogenated diamond and Cs-coated semiconductors (Si, GaAs, GaP). Although these surfaces are not sufficiently stable under heavy ion or electron bombardment, they may be suitable for operation in less harsh environments.

Field Emission

Field Emission From Metals—For current to flow from a metal into vacuum in the absence of heating (thermionic emission), an electric field must be sup-

* A. J. Dekker, *Secondary Electron Emission*, Solid State Physics, Vol. 6 (New York: Academic Press, 1958).

† J. E. Yater, A. Shih, and R. Abrams, *Phys. Rev. B* **56**, R4410 (1997).

plied of sufficient magnitude for electrons to quantum mechanically escape the surface barrier. When an electron leaves a metal surface, a redistribution of charge occurs to shield the electron's influence from the interior and maintain the surface at constant potential. The potential distribution the departing electron experiences is modified by an equal but opposite interior charge at the same distance from the surface and is called the *image charge potential*. It is given by

$$V(x) = \mu + \phi - Fx - Q/x$$

where

- x = the distance from the metal surface in nm
- μ = the chemical potential (approximately the Fermi energy E_F) in eV
- ϕ = the work function in eV
- F = the product of the electron charge and the electric field in eV/nm
- $Q = e^2/(16\pi\epsilon_0)$, 0.359991 eV-nm
- e = the electron charge, 1.60218×10^{-19} Coulombs
- ϵ_0 = the permittivity of free space, 8.85419×10^{-12} F/m

The chemical potential (Fermi energy) is related to the density of electrons in the metal (or "carrier concentration" in semiconductors) by the relationship*

$$n = N_c \frac{2}{\sqrt{\pi}} F_{1/2}(\beta\mu)$$

where

- n = the number density of electrons in $\#/nm^3$
- $N_c = 0.0250942 \#/nm^3$ for metals
- $\beta = 11604.5/T$ in eV^{-1} , where T is the temperature in Kelvin
- $F_{1/2}(x)$ = the Fermi-Dirac integral, equal to $x^{3/2} (2 + (\pi/x)^2)/3$ when $x \gg 1$

" F " is the force on an electron due to the local electric field, referred to here simply as "field." Barrier lowering by an amount $(4QF)^{1/2}$ occurs, and this increases the thermionic emission current, but the thermionic current is negligible by comparison to the tunneling current at room temperature for typical metals. The tunneling current is given by†

$$J(F) = \frac{em}{2\pi^2 \beta \hbar^3} \int_0^\infty \ln(1 + e^{\beta(\mu - E)}) \exp\left(-\frac{b_{fn}}{F} + c_{fn}(E - \mu)\right) dE$$

where

- J = the current density in A/cm^2
- E = the energy of the incident electron in eV
- m = electron rest mass, $510999 eV/c^2$
- c = the speed of light, 2.99792×10^{10} cm/s
- \hbar = Planck's constant over (2π) , 6.58212×10^{-16} eV-s

* Reference 13.
 † Reference 14.

The "ln" term arises from the Fermi-Dirac distribution for electron statistics, and the exponential term is from the WKB approximation to the transmission coefficient in which the log of the transmission coefficient is Taylor expanded to first order in energy about μ . The terms b_{fn} and c_{fn} are

$$b_{fn}(F) = \frac{4}{3\hbar} \sqrt{2m\phi^3} v(y); c_{fn}(F) = \frac{2}{\hbar F} \sqrt{2m\phi} t(y)$$

where

- $v(y)$ = function of elliptical integrals, approximated to within 0.1% for $0.40 \leq y \leq 0.92$ by $v(y) \approx 0.13119(1-y)(1.1321+y)(6.9464-y)$
- $t(y)$ = function of elliptical integrals, approximated to within 0.1% for $0.24 \leq y \leq 1.00$ by $t(y) \approx 0.023056(5.0842-y)(y^2 + 2.3317y + 8.4607)$
- $y = (4QF)^{1/2}/\phi$ and ranges between 0.4 to 0.85 for generic metal field emission

For metals, the integral lower limit in $J(F)$ can be extended to $(-\infty)$ with negligible error and analytically integrated to give the *Fowler Nordheim (FN) Equation*‡

$$J_{fn}(F) = a_{fn} F^2 \exp(-b_{fn}/F) \left[\frac{\pi c_{fn}/\beta}{\sin(\pi c_{fn}/\beta)} \right]$$

where

$$a_{fn}(F) = \frac{e}{16\pi^2 \hbar \phi t(y)^2}$$

(the present form is closer to that of Murphy and Good).§ The variation of the FN parameters a_{fn} , b_{fn} , and c_{fn} are shown in Fig. 5. The term in square brack-

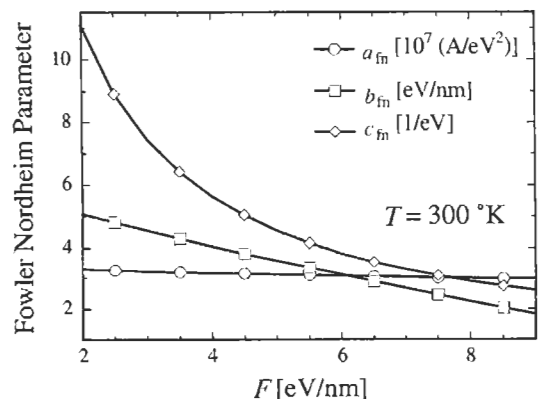


Fig 5. Variation of the Fowler-Nordheim current density parameters as a function of field F at the surface for $\phi = 4.41$ eV (photoelectric work function of Molybdenum) at room temperature. The weak dependence of a_{fn} on F motivates the common approximation that $t(y) \approx \text{constant}$.

‡ Reference 15.
 § Reference 16.

ets in $J(F)$ is within 6% of unity for metals under common field emission conditions and is often ignored, reflecting the relative independence of the electron energy distribution from temperature unless the temperature is sufficiently high, as in Fig. 6.

The plot of $\ln(J(F)/F^2)$ versus $(1/F)$ is approximately linear and is referred to as a "Fowler–Nordheim Plot." The slope is proportional to $\phi^{3/2}/F$ and is therefore sensitive to work function changes and geometrical field enhancement. Because of the exponential dependence of $J(F)$ on b_{in}/F , small changes in either ϕ or the local field result in substantial changes in the emitted current density, as in Fig. 7.

The function $v(y)$ is commonly taken to be linear in y^2 , for then $b_{in}(F)/F$ is linear in $(1/F)$ and $J_{in}(F)$ on a

Fowler–Nordheim plot is manifestly linear. The most widely used approximation is $t(y) \approx (1.1)^{-1/2}$ and $v(y) \approx 0.95 - y^2$.^{*} For emission from multidimensional structures (below), such approximations are similarly expedient, but because $v(y)$ versus y^2 exhibits a weak convexity, the tangent line at a particular $y_o = y(F)$ is a better approximation, for which

$$v(y)|_{\text{tangent}} = v_0 - v_1 y^2$$

$$v_0 = \frac{1}{4}(3t(y_o) + v(y_o))$$

$$v_1 = \frac{3}{4y_o^2}(t(y_o) - v(y_o))$$

where

- v_1 = the slope of the tangent line
- v_0 = the intercept of the tangent line, equal to 0.9368 when $v_1 = 1$
- y_o = the expansion point for a given F , equal to 0.5992 when $v_1 = 1$

Field Emission from Semiconductors—The carrier concentration in semiconductors is substantially less than the conduction-band electron density in metals, even under high field. Band bending occurs, so that electron affinity (energy difference between vacuum level and bottom of the conduction band) is a more useful parameter for semiconductors than work function (energy difference between vacuum level and Fermi level), as for semiconductors the latter is affected by applied field. Emission arises from either the valence or conduction band or both, as well as surface states and defects within the band gap. The nature of the current flow depends on the type of doping (*n*- or *p*-type) incorporated. Current resulting from hole transport can be significant. The effective mass in a semiconductor is generally less than the electron rest mass in vacuum, and it depends on the crystallographic plane. Surface layers, oxides, or adsorbates generally exist on the surface and contribute to emission fluctuations. If field emission is primarily from the valence band (e.g., diamond[†]), then the barrier height encountered by the valence electrons is augmented by the band gap. Because of these complications, a current density formula for semiconductors with the simplicity of the Fowler–Nordheim equation for metals is not available except for provisional approximations (more extensive treatments are available[‡]).

Electrons migrate to the surface to shield out an applied field, indicating that the electron density increases near the surface. The chemical potential is replaced by $\mu(x) = \mu_o + \phi(x)$, where μ_o is the bulk chemical potential (generally such that $\mu_o < 0$) and

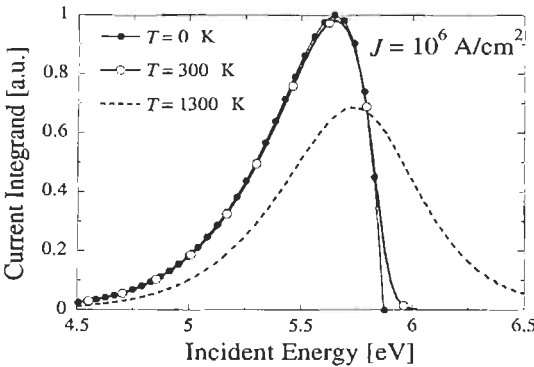


Fig. 6. Energy distribution of electrons incident on the surface barrier in the Fowler–Nordheim approximation. For a given current density, the effects of a raise in temperature are to increase the thermal "tail" near the Fermi level. At room temperature, the energy distribution for metals is well approximated by the 0 K limit.

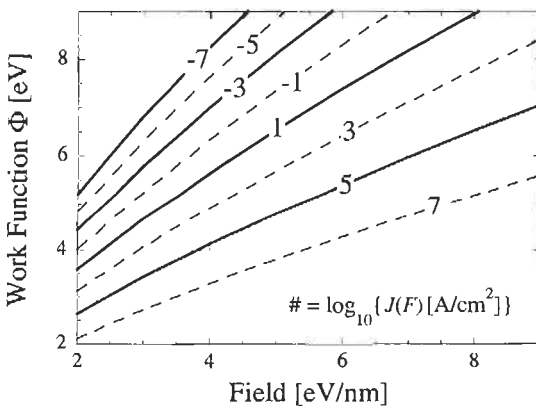


Fig. 7. Contours of constant current density, labeled by $\log_{10}(J(F)[A/cm^2])$, as a function of work function and field as evaluated using $J_{in}(F)$.

* Reference 17.

† Reference 18.

‡ Reference 14.

$\phi(x)$ is the solution to Poisson's equation. $\mu(x)$ is referred to as the "electrochemical potential." If emission is primarily from the conduction band and all other complications are ignored, then the equation for $J(F)$ for semiconductors is analogous to that for metals, but with the replacements:*

$$\phi \Rightarrow \chi - \mu$$

$$Q \Rightarrow \frac{K_s - 1}{K_s + 1} Q$$

$$\frac{c_{in}\pi/\beta}{\sin(c_{in}\pi/\beta)} \Rightarrow \frac{c_{in}\pi/\beta}{\sin(c_{in}\pi/\beta)} - (1 + c_{in}\mu)\exp(-c_{in}\mu)$$

where

- χ = the electron affinity in eV
- K_s = dielectric constant (e.g., 11.9 for silicon)

At the surface, $\mu(x)$ is related to F_{vac} by Poisson's equation for ϕ such that $K_s \partial_x \phi = F_{vac}$ at the surface, from which

$$F_{vac(\mu)} = \begin{cases} \left(\frac{2\pi^2 N_c K_s}{3\beta\epsilon_o} \left(\frac{\beta\mu}{\pi} \right)^{1/2} \left[\frac{8}{5} \left(\frac{\beta\mu}{\pi} \right)^2 + 1 \right] \right)^{1/2} & \beta\phi \text{ large} \\ \left(\frac{2N_c K_s}{\beta\epsilon_o} \exp(\beta\mu_o) [\exp(\beta\phi) - \beta\phi - 1] \right)^{1/2} & \beta\phi \text{ small} \end{cases}$$

where

- $N_c = 2 M_c (m/(2 \pi \beta \hbar^2))^{3/2}$
- M_c = number of equivalent minima in the conduction band (e.g., 6 for silicon)
- m = effective electron mass in the semiconductor (e.g., 0.3283 m_o for silicon)
- F_{vac} = applied field in vacuum in eV/nm

Macroscopic Current Density

Field emission sources come in a variety of geometries (e.g., conical [Spindt-type] and wedge [knife-edge], to planar thin films and mats of carbon nanotubes) and materials (e.g., metal and semiconductor with or without various coatings for robust performance or low work function). The current density is rarely uniform over extended emission areas because of geometrical and/or work function variation. The sharpest and lowest work function sites dominate, so that electron emission invariably occurs from localized emission sites of small size. For example, conical Molybdenum tips centered in gate holes of micron dimensions† can have an overall radius of curvature at the tip on the order of 25 nm, though emission appears to arise from sites of atomic dimensions.‡

* Reference 19.
 † Reference 20.
 ‡ References 21 and 22.

The total current from a single emission site is small, but that is offset by the magnitude of the number of sites. The packing density (alternately, the Emission Site Density [ESD]) is on the order of 10^7 – 10^9 tips/cm² for field emitter arrays and 10^4 – 10^6 sites/cm² for various carbon-based and wide band gap emitters.§ The combination of high current density in conjunction with the large number of sites counter the small emission site area and the tendency of only a subset of emitters or sites to be active. Current densities averaged over the cathode area vary between 0.01 and 0.1 A/cm² for carbon-based emitters and up to 10 to 2400 A/cm² for field emitter arrays. Often, macroscopic fields are cited, but such fields cannot be used in the Fowler–Nordheim equation unless modified by some estimate of the field at the emission site (the macroscopic field is related by a multiplicative factor dependent upon the geometry of the emission site).

Field Emitter Arrays

Field emitter arrays (FEAs) have been extensively investigated for vacuum electronics amplifiers,[¶] space-based applications,[#] and flat panel displays.** An example of a molybdenum FEA cathode is shown in Fig. 8. Concepts such as emission area and field enhancement treated here for FEAs have analogs for the other cold cathodes but are more ambiguous for the latter because a wider variety of mechanisms are believed to be responsible for electron emission.

The local field F is related to an applied (gate) voltage V by the "field enhancement factor" β_g such that $F = \beta_g V$. The field generated at the emission site by the gate significantly dominates other fields (such as that due to the collection anode). The field enhancement factor varies primarily as the inverse radius of the emission site, with other geometrical parameters contributing to a lesser extent. The emission is concentrated about the nanoprotusions and atomic-sized structures that exist on emitter surfaces.

In practice, β_g is not measured but is inferred from experimental data or simulated.†† β_g is analytically intractable for realistic configurations, but it is obtainable for idealized models. For a Spindt-type tip approximated by a gated hyperbolic cone, β_g is‡‡

$$\beta_g = \left\{ \frac{\pi}{\ln \left[\frac{a_g}{54a_s} \left(86 + \frac{a_g}{a_s} \right) \cot(\beta_o) \right]} - \tan^2(\beta_o) \right\} \frac{e}{a_s}$$

§ References 23 and 24.
 ¶ References 25 and 26
 # Reference 27.
 ** References 28 and 29.
 †† References 30 and 31.
 ‡‡ Reference 32.

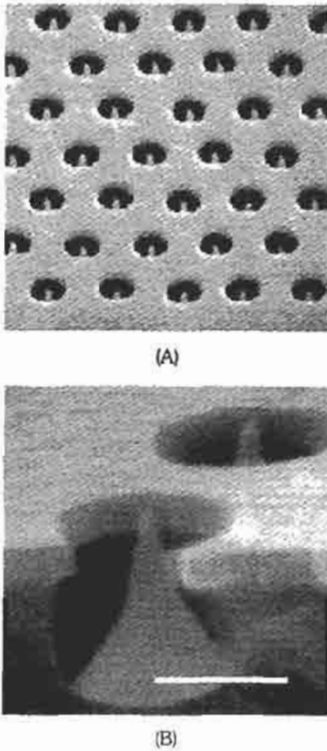


Fig. 8. (A) A $5 \mu\text{m}$ by $5 \mu\text{m}$ region of a Molybdenum field emitter array developed by SRI, corresponding to approximately 10^8 tips/ cm^2 . (Photo courtesy of C. A. Spindt.) (B) Close-up cross section of an individual SRI emitter. Bar corresponds to a length of $1 \mu\text{m}$. (Photo courtesy of C. A. Spindt.)

where (as shown in Fig. 9)

- β_g = the field enhancement factor in $\text{eV/Volt}\cdot\text{nm}$
- a_s = the radius of hyperbolic apex in nm
- a_g = the radius of the gate in nm
- β_o = the cone half-angle in radians

The emission is concentrated at the cone's apex. The characteristic emission area factor is defined by $b_{\text{area}} = I(V)/J(F)$ and is approximated by

$$b_{\text{area}}(F) = 2\pi a_s^2 \left(\frac{F \cos^2(\beta_o)}{b_{\text{in}}^o + F \sin^2(\beta_o)} \right)$$

where

b_{in}^o = the intercept of $b_{\text{in}}(y)$ as a linear function of y^2 .

Because the current falls off exponentially with field, the actual emission area is about a factor of 2.7 larger than b_{area} , and the average current density is smaller by the same factor. Generic values are given in Table 1 for a single molybdenum Spindt-type emitter operating at $1 \mu\text{A}$.

In a spherical coordinate system such that the symmetry axis of a single tip is parallel to the z axis, and

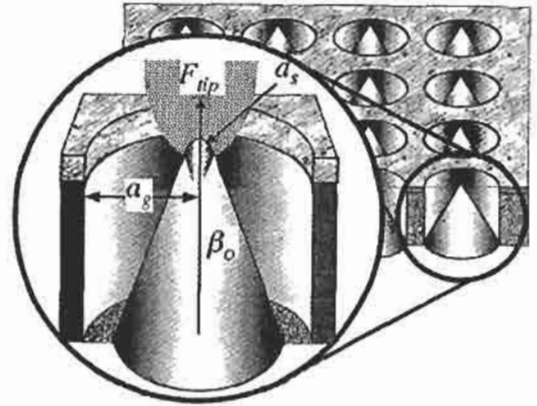


Fig. 9. Representation of an individual emitter as a hyperbolic cone, showing relevant parameters: gate radius (a_g), apex radius (a_s), and cone half-angle (β_o).

TABLE 1: GENERIC FEA MODEL PARAMETERS

Symbol	Description	Value
<i>Independent</i>		
T	Temperature	300 K
ϕ_o	Work Function (Mo)	4.41 eV
a_s	Emission Site Radius	4 nm
a_g	Gate Radius	$0.5 \mu\text{m}$
β_o	Cone Half-Angle	15°
V	Gate Voltage	72.2 V
<i>Dependent</i>		
β_g	Field Enhancement Factor	0.0867 e/nm
F	Field	6.26 eV/nm
y	Elliptical Integral Argument	0.681
b_{in}^o	$J(F)$ Exponential Term	58.2 eV/nm
b_{area}	Area Factor	10.0 nm^2
$b_{\text{area}} J_{\text{tip}}$	Single Tip Current	$1.00 \mu\text{A}$
χ	Emission Distribution Parameter	2.84
θ_{rms}	$\langle \theta^2 \rangle$ - rms Emission Angle	19.1°
E_t	Transverse Energy Component	7.32 eV

the polar angle is θ , then the distribution of electrons on a hemisphere of approximately the same diameter as the gate diameter (i.e., $r = a_g$) is given by

$$f(\theta) = \exp(-2\chi \sin^2 \theta)$$

$$\chi = \frac{1}{2} \left(1 + \frac{b_{\text{in}}^o}{2F} + \frac{1}{2} \sin^2 \beta_o \right)$$

where

F = field at apex in [eV/nm]

The root-mean-square (rms) value of θ and the transverse energy E_t for electrons emerging from the hemisphere is then, for a single tip,

$$\langle \theta^2 \rangle = \sqrt{\frac{1}{24} \left(15 - \frac{16I_1(\chi) - I_2(\chi)}{I_0(\chi)} \right)}$$

$$E_t = \frac{1}{2} \left(1 - \frac{I_1(\chi)}{I_0(\chi)} \right) V_g$$

where

$I_n(z)$ = hyperbolic Bessel function of order n

Asymptotically, $\theta_{rms} \approx (4\chi + 2)^{1/2}/(4\chi)$ and $E_t \approx \{(4\chi + 1)/(4\chi^2)\}V$. The presence of other emitters in close proximity, as in an array, will affect the electron distribution accordingly.

Current-Voltage Relationship, Transconductance, and Modulation

The proportion of active emission sites increases with applied voltage, that is, sites “turn on” as the voltage is raised. The proportion of operating tips scales approximately linearly with applied voltage* (and similarly for the emission sites of wide band gap emitters†). Further, as in b_{array} , the region over which emission occurs scales approximately linearly if the emission surface is curved. The commonly made assumption that total current $I(V)$ varies as $A J(F)$, where A is a constant area, is incorrect, as is the “derivation” of current-voltage relations by the substitution $J \Rightarrow I/A$ and $F \Rightarrow \beta_g V$ in $J(F)$. Nevertheless, experimental data plotted on axes defined by $\ln(I(V)/V^2)$ versus $1/V$ is more or less linear for FEAs, as well as for other field emission sources, so that

$$I(V) = A_{FN} V^2 \exp(-B_{FN}/V)$$

where

- I = the current measured in amperes
- V = the extraction potential in volts
- A_{FN} = the “intercept” parameter in A/V^2
- B_{FN} = the “slope” parameter in volts

The complex relationship between the parameters A_{FN} and B_{FN} and the parameters a_{fn} and b_{fn} is available only for idealized models, for which the majority of B_{FN}/V is due to b_{fn}/F . This implies that the slope on a Fowler-Nordheim plot of $I(V)$ data can give informa-

tion about the ratio $\phi^{3/2}/\beta_g$, (but not about ϕ or β_g , separately).

Because current is exponentially sensitive to changes in voltage, the large transconductance of FEAs and cold cathodes is of interest to vacuum electronics devices. Transconductance is the change in current with respect to voltage, or $g_m = \partial I/\partial V$. For Class A amplifiers, in which $I \propto V^n$ (n being of order unity) the transconductance does not vary greatly over a minimum to maximum voltage sweep (it is constant if $n = 1$). For space-charge limited thermionic emitters—where $n = 3/2$, or $I(V) = C V^{3/2}$, C being a constant of proportionality— g_m varies by only 41% when the maximum to minimum voltage ratio is 2. In contrast, for field emission, g_m is given by

$$g_m = \frac{(B_{FN} + 2V)}{V^2} I(V)$$

where

g_m = the transconductance in siemens

Transconductance for a field emission cathode therefore varies exponentially with inverse voltage. For a common Spindt-type 100-tip FEA‡ with $A_{FN} = 1.535 \text{ mA/V}^2$ and $B_{FN} = 544 \text{ V}$ operated between $V = 40 \text{ V}$ and 80 V , the max/min ratios for current and transconductance are 1000 and 3600, respectively, and g_m is $12 \text{ } \mu\text{A/V-tip}$ for $I = 110 \text{ } \mu\text{A-tip}$. Increasing g_m/C , where C is capacitance, is sought for amplifier designs that seek to modulate the electron beam.§

Three other parameters individually or in combination bear on the utility of a field emission cathode for vacuum electronic devices. They are beam turn-off conditions, average-to-peak current ratio, and pulse repetition frequency (PRF). For an emitter characterized by A_{FN} and B_{FN} , the current minimum to maximum ratio x identifies what minimum voltage is required to ensure that an electron beam is turned off. It is given by the large- n limit of x_n , defined by the recursion relation

$$x_n = \frac{\eta}{\eta - \ln\left(\frac{\zeta}{x_{n-1}}\right)}$$

where

- ζ = current minimum to maximum ratio (I_{min}/I_{max})
- $\eta = B_{FN}$ to voltage maximum V_{max} ratio B_{FN}/V_{max}
- x_n = voltage minimum to voltage maximum ratio V_{min}/V_{max} (n th iteration for $n \geq 1$)
- $= (\eta + 2)/(\eta + 2 - \ln(\zeta))$ ($n = 0$)

The $n = 0$ term is adequate to within 10% for general parameters. If “beam turn-off” is defined as $\zeta = 0.001$, then for $\eta = 6.8$, $V_{min} = 54.4\%$ of V_{max} (i.e., $x = 0.544$); compare to space-charge limited thermionic emitters,

* Reference 33.

† Reference 34.

‡ Reference 35.

§ Reference 36.

for which the minimum voltage is 1% of the maximum. Devices requiring beam turn-off (e.g., radar) may benefit.*

Emission gating of the electron beam is required for a class of rf devices known as *Inductive Output Amplifiers*.† If the extraction grid or gate potential is sinusoidally modulated, then the average-to-peak current ratio is given by

$$\frac{I_{ave}}{I_{pk}} = \frac{\frac{1}{2\pi} \int_0^{2\pi} I_{FN}(V(\omega t)) d(\omega t)}{I_{pk}}$$

$$= e^{-z} I_0(z) + \frac{(B_{FN} + V_{pk}) V_{pk}}{(B_{FN} + V_{pk})^2} \left\{ [(2z + 1) I_1(z) - 2z I_0(z)] z e^{-z} \right\}$$

where

- ω = the angular frequency of modulation in rad/s
- $V(\omega t) = V_{pk} - V_{rf}(1 - \cos(\omega t))$, with V_{pk} and V_{rf} constant voltages representing maximum and amplitude of oscillation, respectively
- $I_n(z)$ = hyperbolic Bessel function of order n
- $z = (V_{rf}/V_{pk})^2 [B_{FN} + 2V_{pk}]$ and is dimensionless

For generic parameters $77.5 V_{rf} = 6.8 V_{pk} = B_{FN}$, then $z = 0.772$, and $I_{ave}/I_{pk} \approx 0.5$.

Protective Resistance

To limit the damage caused by exponential rises in current (associated with arcing events, which can destroy an emitter) and to increase uniformity and homogenize the current density across the cathode,‡ protective resistance elements are often incorporated. Resistive protection suppresses current runaway and limits damage by providing a voltage drop in proportion to the emitted current. The $I(V)$ equation is modified according to

$$I(V) = A_{FN} (V - I(V)R)^2 \exp\left(-\frac{B_{FN}}{V - I(V)R}\right)$$

where

R = resistance in ohms

which may be solved graphically. Alternatively, an iterative solution is obtained from $I(V) = (V/R) z / (1 + z)$, where z is the ($n \rightarrow \infty$) limit of z_n defined by

$$z_n = \frac{RI_0(V) \exp\left(-\frac{B_{FN}}{V} z_{n-1}\right)}{V(1 + z_{n-1})}$$

* Reference 37.
 † Reference 38.
 ‡ Reference 39.

where

- $I_0(V)$ = current in the absence of resistance in amperes, or $A_{FN} V^2 \exp(-B_{FN}/V)$
- V = voltage drop from gate to ground in volts
- z_n = the ratio of $(I(V)R) / (V - I(V)R)$ (n th iteration for $n \geq 1$)
- $= (V I_0 R) / [B_{FN} I_0 R + V(V + I_0 R)]$ ($n = 0$)

When $I(V)R \ll V$, z_n rapidly converges, and so the $n = 0$ term is adequate. When $I(V)R \approx V$, convergence is slow.

Field Emission Displays

Flat panel cathodoluminescent displays utilizing pixel-size, addressed cathodes are available commercially. Most of these displays are called Field Emitter Displays (FEDs) because the cathodes are field emitter arrays. Monochrome grayscale displays with diagonal dimension of 5.2" are currently available,§ and larger color displays have been developed as prototypes. Similar displays using somewhat different kinds of cathodes are also under development. In addition to being flat, these displays can be very bright and energy efficient and appear similar to standard CRT displays.

FEDs are built by fabricating a "cathode plate" filled with field emitter arrays and placing it in close proximity to a printed phosphor screen similar to those used in cathode ray tube displays (see Fig. 10). The phosphor screen and the arrays are fabricated such that each phosphor pixel will be facing a separately gated array, which produces the electrons required for cathodoluminescence. The cathode and phosphor plates are manufactured separately, bonded together at the edges (typically using glass frit), and then evacuated. The emission from each cathode is modulated as needed to produce the desired visible light intensity at the adjacent pixel. The emission time or the emission intensity or both can be modulated to produce images with desired color and brightness. Because the field emitter arrays have a highly nonlinear current-voltage characteristic, individual arrays can be addressed when all the emitters in each row and all the gates in each column are connected together. Thus there is no need to fabricate transistors to address the pixels.

Each field emitter array contains many micron-scale field emitters and apertures. Ideally, all of the field emitters would have identical $I-V$ curves. However, state-of-the-art fabrication techniques produce emitters that vary both in shape and in surface electronic characteristics. Thus some means of homogenizing the emission is required. Many more emitters are placed at each pixel than are required, to provide a statistical distribution. In addition, the emitters are fabricated on top of a resistive film, such that excessive emission current will produce a voltage across the film. Thus the potential appearing between the gate and emission site

§ PixTech Corporation (www.pixtech.com).

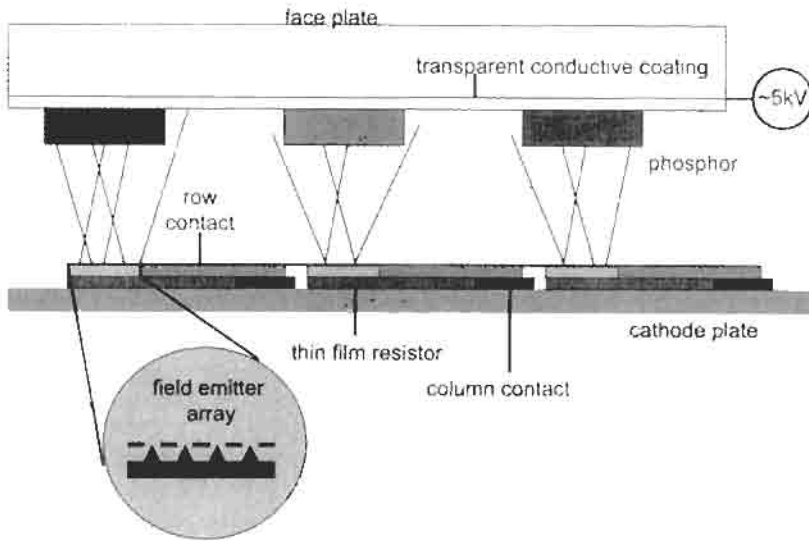


Fig. 10. Schematic diagram of a typical field emission display (FED). Individual pixels are addressed by applying positive and negative voltages to the corresponding row and column contacts. The resulting electrons strike the adjacent phosphor dots. Variations include use of a grid above the cathode plate (to partially focus the electrons), and various support structures between the cathode and face plate (not shown).

is less than the voltage applied between the row and column contacts. The emission characteristic becomes nearly linear and much more uniform from emitter to emitter as a result. Transistor structures can also be used in place of the thin-film resistors, producing saturated emission curves.

Since the thin glass plates are not mechanically strong enough to support atmospheric pressure over more than a few square inches, internal supporting structures are required. These supporting structures must not be visible to the viewer, and they must sustain the voltage applied between the cathode and phosphor plates. Slightly conductive materials can sustain high voltages more easily; however, larger currents dissipate precious power. Because the field emitter arrays emit electrons into a 20–30 degree angle, the distance between the cathode and phosphor plates must be kept small (~200 μm) to provide high resolution without crosstalk between adjacent pixels. The spacing can be increased by adding a grid to collimate the beams. Typical color phosphors used in CRTs work best with electron energies large enough to penetrate several hundred nanometers inside the phosphor particle. As in CRTs, a thin layer of aluminum is often applied to the backside of the phosphor, and this requires additional electron energy to penetrate. These issues translate to phosphor plate potentials of at least 5 kV. Thus the requirements on the support structures are severe.

Significant pressures of reactive gases are typically present inside the displays. These gases can be evolved when electrons strike the phosphor plate. These gases can react at the emission sites, reducing the emission current. Electrons traveling through the oxide tend to

break it down again, so that a kind of balance is struck between oxidation and electronic reduction. This process is poorly understood and depends on the type of phosphor as well as the emitter material.

Many of these issues could be resolved by using alternative emitter and phosphor technologies. For example, the emitter surfaces can be fabricated from less reactive materials. Some new phosphors can be excited at lower voltages and/or produce fewer reactive gases. Progress in these areas may produce improved FEDs in the future, provided that cost-effective production methods can be found.

Photocathodes

There is a renewed interest in photocathodes as electron sources for applications that require high current density and high brightness.* The primary focus of current research in photocathodes is the generation of high-brightness electron beams for either linear colliders or free electron lasers (FELs). Photocathodes are also of interest in high-power RF devices. Laser-driven photocathodes allow precise control over the electron peak current and the spatial and temporal profiles of the electron beam. Consequently, laser-driven photocathodes reduce the demands placed on beam structures and lead to a compact design.

Two types of materials have been used for photocathodes: metals and semiconductors. The primary advantages of metal photocathodes are fast response

* Reference 40.

time, tolerance to relatively poor vacuum (10^{-6} – 10^{-7} torr), and long lifetime (> 1 year). However, the main drawback is the low photoelectric yield, which is measured by the quantum efficiency (QE), i.e., the number of electrons produced per incident photon. The QE for metals is typically 10^{-5} – 10^{-7} *. Semiconductor photocathodes, in contrast, have higher QE, with the highest QE being around 10^{-1} . However, they are extremely sensitive to contamination and require an ultra-high vacuum (UHV) of 10^{-9} to 10^{-10} torr. Even then, the lifetimes are short, namely 1 to 100 hours.† Research efforts are ongoing to improve both types of photocathodes.

The characteristics of metal and semiconductor photocathodes can be understood by examining the photoemission process. Similar to the secondary-electron-emission process discussed earlier, the photoemission process can be visualized as a three-step process: photoelectron generation, transmission of the electrons through the material, and emission of the electrons over the vacuum barrier. Photoelectron generation is a quantum process, i.e., the energy of the absorbed photon is converted completely into the energy of the photoexcited electron. Of course, to be emitted the photoexcited electrons need to have sufficient energy to overcome the vacuum barrier. As a result, in a metal photocathode, the photon energy must be equal to or greater than the work function of the metal to achieve a high QE. Because the metal work function is typically greater than 4 eV, UV lasers are required to drive a metal photocathode (where the wavelength of a 4-eV photon is 310 nm). Before the photoexcited electrons can be emitted, however, they need to traverse the material to reach the surface, and in the transport process they suffer energy loss from collisions with conduction electrons. Because conduction electrons are abundant in metals, the collisions are frequent and the mean free path of the photoexcited electrons is short, which limits the QE. However, the photoexcited electrons have to be generated near the emitting surface of the metal to retain sufficient energy to overcome the high vacuum barrier. Consequently, metal photocathodes offer a short delay in the photon-to-electron conversion, resulting in a fast response time and narrowly bunched beam.

In semiconductor photocathodes, the photon energy must be at least as large as the energy gap (~ 1 – 2 eV in commonly used semiconductors) to generate the photoexcited electrons. A long mean free path for the photo-excited electrons and a low or negative electron affinity are then required to achieve a high QE. The conduction electron population is relatively low in semiconductor materials, and hence low-energy electrons (with insufficient energy to excite valence electrons over the band gap) experience infrequent electron–electron collisions. This results in a relatively

long mean free path for the low-energy electrons. Most of the high-QE semiconductor photocathodes are alkali-containing or cesiated semiconductors since the low or negative electron affinity of these materials allows the low-energy electrons that reach the surface to escape. However, adsorption of ambient gases often leads to an increase in the electron affinity, preventing many of the low-energy electrons from escaping. Because the contribution from low-energy electrons is essential to a high QE, the performance of semiconductor photocathodes is very sensitive to ambient gas adsorption.

The most promising metal photocathode materials have been Cu and Mg. They can be transported in air, are not sensitive to contaminants, and have lifetimes greater than a year. Various techniques have been used to improve the QE of metal photocathodes. High-power UV lasers are used to irradiate the sample to remove surface contamination, and the laser incident angle is near the Brewster's angle to confine the photons to the surface region. The electric field of the laser is oriented in the plane of incidence, and a high field is used to lower the work function by the Schottky effect. Recently, alkali ion implantations have also been used to reduce the work function. Among metal photocathodes, the highest QE obtained with a 266-nm (4.66 eV) laser beam has been 4×10^{-3} for a Mg photocathode, whereas the QE for Cu is 5×10^{-4} .‡

Alkali-based semiconductor cathodes, such as $\text{Na}_2\text{KSb:Cs}$ and K_2CsSb , have been used in photomultipliers. However, their sensitivity to contaminants prevents their wide use as photocathodes. Cs_3Sb , Cs_2Te , and GaAs:Cs are some of the materials that have been investigated as photocathodes. For potential application as high-brightness, high-current-density cathodes for induction LINAC FELs, various types of cathodes have been evaluated, including *M*-type dispenser and LaB_6 thermionic cathodes, graphite cathodes, velvet cathodes, and Cs_3Sb photocathodes. Although the highest brightness has been observed from some Cs_3Sb photocathodes,§ Cs_3Sb requires a vacuum of 10^{-10} torr and the lifetime is only tens of hours. Another material evaluated as a photocathode is Cs_2Te , which is typically fabricated in situ in a UHV system. Freshly prepared Cs_2Te samples grown on Cu and Mo substrates were shown to yield a QE of 6×10^{-2} – 1.2×10^{-1} and 1.6×10^{-1} using photons of 4.66 eV (266 nm) and 4.94 eV (251 nm), respectively. However, the QE decreased roughly exponentially with a decay time of about 100 hours." A very high QE ($\sim 1.4 \times 10^{-1}$) has been achieved from GaAs:Cs photocathodes with 2.55 eV (486 nm) photons. Typically, these cathodes have been used to generate polarized electrons, but recently there has been interest in using them as photocathodes in RF cavities. With the advent of GaAs:Cs photocath-

* References 40 and 41.

† Reference 40.

‡ Reference 40.

§ Reference 42.

" Reference 40.

odes, it may be possible to generate electron bunches with visible photons instead of UV photons in the foreseeable future.

ELECTRODE DISSIPATION

After the electron stream has given up the useful component of its energy, the remainder is dissipated as heat in some suitable part of the tube. Five processes are commonly used to remove this heat. The amount that can be removed depends on the area available, the temperature differential, and, in the cases of forced cooling, the coolant flow.

In computing cooling-medium flow, a minimum velocity sufficient to assure turbulent flow at the dissipating surface must be maintained. The figures for specific dissipation (Table 2) apply to clean cooling surfaces and may be reduced to a small fraction of the values shown by heat-insulating coatings such as scale or dust.

Radiation Cooling

In a radiation-cooled system, that portion of the tube on which the heat is dissipated is allowed to reach a temperature such that the heat is radiated to the surroundings. The amount of heat that can be removed in this manner is given by the equation

$$P = \epsilon_r \sigma (T^4 - T_o^4)$$

where

P = radiated power in W/cm²

ϵ_r = total thermal emissivity of the surface

σ = Stefan-Boltzmann constant = 5.67×10^{-12} W cm⁻² K⁻⁴

T = temperature of radiating surface in kelvins

T_o = temperature of surroundings in kelvins

Total thermal emissivity varies with the degree of roughness of the surface of the material and the temperature. Values for typical surfaces are listed in Table 3.

TABLE 3. TOTAL THERMAL EMISSIVITY ϵ_r OF ELECTRON-TUBE MATERIALS

Material	Temperature (K)	Thermal Emissivity
Aluminum	450	0.1
Anode graphite	1000	0.9
Copper	300	0.07
Molybdenum	1300	0.13
Molybdenum, quartz-blasted	1300	0.5
Nickel	600	0.09
Tantalum	1400	0.18
Tungsten	2600	0.30

Water Cooling

For water cooling, the water is circulated through a suitably designed structure. The amount of heat that can be removed by this process is given by

$$P = 264Q_w(T_2 - T_1)$$

where

P = power in watts

Q_w = flow in gallons per minute

T_2, T_1 = outlet and inlet water temperatures, respectively, in kelvins

This same relationship is given in the nomogram of Fig. 11 with the temperature rise in degrees Fahrenheit or Celsius and the power in kilowatts.

Forced-Air Cooling

With forced-air cooling, a stream of air is forced past a suitable radiator. The heat that can be removed by this process is given by

$$P = 169Q_A [(T_2/T_1) - 1]$$

TABLE 2. TYPICAL OPERATING DATA FOR COMMON TYPES OF COOLING

Type	Average Cooling Surface Temperature (°C)	Specific Dissipation of Cooling Surface (W/cm ²)	Cooling-Medium Supply
Radiation	400-1000	4-10	
Water	30-150	30-110	0.25-0.5 gallon/minute/kilowatt
Forced air	150-200	0.5-1	50-150 ft ³ /minute/kilowatt
Evaporative	100-120	80-125	Water-, air-, or convection-cooled condenser. A water-cooled condenser would require 0.07-0.1 gallon/minute/kilowatt.
Conduction	100-250	5-30	Heat sink operating at 50-100°C

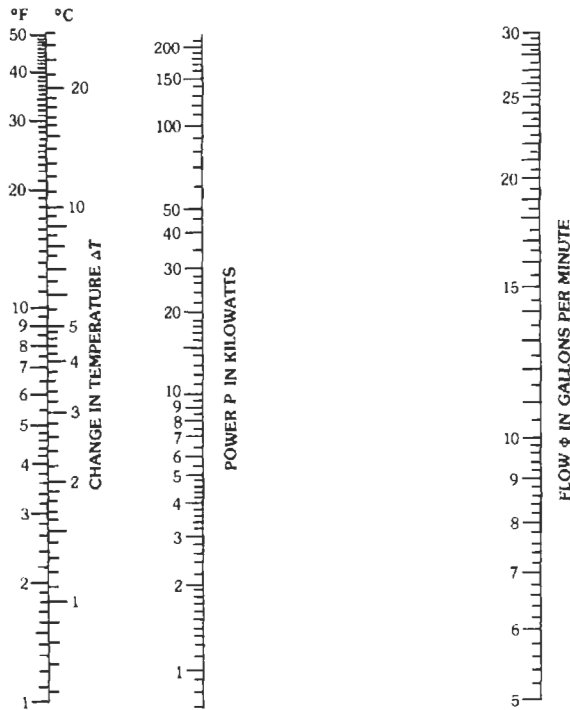


Fig. 11. Heat transfer in cooling water. $P = 0.1466\phi\Delta T$, for T in Fahrenheit. (Courtesy Clyde G. Haehnle.)

where

$$Q_A = \text{air flow in ft}^3/\text{minute}$$

Evaporative Cooling

A typical evaporative-cooled system consists of a tube with a specially designed anode immersed in a boiler containing distilled water. When power is dissipated on the anode, the water boils and the steam is conducted upward through an insulating pipe to a condenser. The condensate is then gravity fed back to the boiler, thus eliminating the pump required in a circulating water system.

For some transmitter applications, the steam is directed downward to leave the space above the tube available for other components. Such a system requires a pump to return the condensate to the boiler, but even then the pump has to handle only about 0.05 of the amount of water required for a water-cooled system because of the exploitation of the latent heat of steam.

The size of the heat-exchanger equipment for an evaporative-cooled system is less than one-third of that required for a water-cooled system because of the greater mean temperature differential between the cooled liquid and the secondary coolant. Typical temperature differentials for the two systems are 75°C and 30°C, respectively.

The anode dissipation should not exceed 135 watts per square centimeter of external anode surface because at this point, often referred to as the “Leidenfrost” or “calefaction” point, the surface becomes completely covered with a sheath of vapor, and the thermal conductivity between the anode and the cooling liquid drops to 30 watts per square centimeter, with resultant overheating of the anode. Special designs of the external anode surface (such as the “pineapple”) allow up to 500 watts to be dissipated per square centimeter of internal anode surface.

Conduction Cooling

When an external heat sink is available, heat may be removed from the tube by conduction. Since the electrode where the heat appears is usually at an elevated potential, it is often necessary to conduct the heat through an electrical insulator.

Because of its relatively high thermal conductivity, beryllia ceramic can be used as a common insulator and thermal conductor between the anode of a tube and a heat sink. Beryllia has the following properties:

- Breakdown strength = 10 kV/mm
- Dielectric constant = 6–8
- Thermal conductivity = 2.62 W/cm/°C at 20°C, 1.75 W/cm/°C at 200°C
- Dielectric loss factor = 4×10^{-5}

Tensile strength = 18 000 lbs/square inch
 Compressive strength = 150 000 lbs/square inch

The temperature drop in degrees Celsius across the beryllia ceramic is given by

$$t_1 - t_2 = dW_a / KA \text{ (for a parallel configuration)}$$

where

- t_1 = temperature of tube anode (typical maximum 250°C)
- t_2 = temperature of heat sink (typically 100°C)
- d = thickness of beryllia in cm
- A = cross-sectional area of beryllia perpendicular to direction of heat flow
- K = thermal conductivity of beryllia in W/cm/°C
- W_a = power dissipated on anode in watts

To the temperature drop across the beryllia ceramic must be added the temperature drop across the interfaces between the ceramic and the anode and heat sink, typically 20°C for clamped surfaces at a loading of 25 W/cm².

Because of its toxic nature, care must be taken in handling and disposal of beryllia ceramic.

Grid Temperature

Operation of grids at excessive temperatures will result in one or more harmful effects: liberation of gas, high primary (thermal) emission, contamination of the other electrodes by deposition of grid material, and melting of the grid. Grid-current ratings should not be exceeded, even for short periods.

NOISE IN TUBES

There are several sources of noise in electron tubes,* some of which are associated with the nature of electron emission and some of which are caused by other effects in the tube.

Shot Effect

The electric current emitted from a cathode consists of a large number of electrons and consequently exhibits fluctuations that produce tube noise and set a limit to the minimum signal that can be amplified. The root-

mean-square value of the fluctuating (noise) component of the plate current I_n , is given in amperes by

$$I_n^2 = 2eI\Gamma^2\Delta f$$

where

- I = plate direct current in amperes
- e = electron charge = 1.6×10^{-19} coulomb
- Δf = bandwidth in hertz
- Γ^2 = space-charge reduction or smoothing factor

For temperature-limited cases, $\Gamma^2 = 1$. For space-charge-controlled regions

$$\Gamma^2 = 2kT_c g \theta / \sigma e I$$

where

- k = Boltzmann's constant = 1.380×10^{-23} joule/kelvin
- T_c = cathode temperature in kelvins
- g = conductance or transconductance in mhos, which relates the output signal current to the input signal voltage†
- θ = a factor that in most practical cases is nearly equal to its asymptotic value of $3[1 - (\pi/4)] = 0.644$
- σ = a tube parameter, related to the amplification factor and electrode spacings, that has a value of unity for diodes and varies between 0.5 and 1.0 for negative-grid tubes

Partition Noise

Excess noise appears in multicollector tubes because of fluctuations in the division of the current between the different electrodes. In a grid-controlled tube, these fluctuations in current division reduce the effectiveness of the space-charge smoothing of the shot noise in the plate current. For a screen-grid tube, the root-mean-square noise currents in the cathode lead, the screen-grid lead, and the plate lead (I_{nk} , I_{nc2} , and I_n , respectively) are given by

$$I_{nk}^2 = 2eI_k\Gamma^2\Delta f$$

$$I_{nc2}^2 = 2eI_{c2}[(\Gamma^2 I_{c2} + I)/I_k]\Delta f$$

$$I_n^2 = 2eI[(\Gamma^2 I + I_{c2})/I_k]\Delta f$$

where I_k and I_{c2} are the cathode and screen-grid currents, respectively.

Flicker Effect

The mechanism involved in the flicker effect is not completely understood but appears to depend on the field distribution in the surface layer of the cathode

* S. Goldman, *Frequency Analysis, Modulation and Noise* (New York: Dover Publications, 1967). A. van der Ziel, *Noise* (Englewood Cliffs, NJ: Prentice-Hall, 1954). L. D. Smullin and H. A. Haus, eds., *Noise in Electron Devices* (The Technology Press of Massachusetts Institute of Technology and John Wiley & Sons, New York, 1959). D. H. Bell, *Electron Noise* (London: D. Van Nostrand Co., 1960). W. R. Bennett, *Electrical Noise* (New York: McGraw-Hill Book Co., 1960). D. K. C. MacDonald, *Noise and Fluctuations: An Introduction* (New York: John Wiley & Sons, 1962).

† For diodes, g is the conductance; for triode and pentode amplifiers, g is the transconductance g_m ; and for triode or pentode mixers and converters, g is the conversion conductance g_c .

resulting from its porous structure. Because this same field distribution also will influence the cathode activity and temperature, flicker noise will depend on cathode activity and temperature in a complicated manner.

The flicker noise spectrum is usually of the form f^{-a} with a close to unity and thus is important only at low frequencies. The sensitivity of audio, subaudio, and direct-current amplifiers is limited by the flicker noise generated in the first tube.

Collision Ionization

Free gas ions can be generated by collisions with the electron stream. The electrons thus liberated and collected by the anode will appear as noise in the anode circuit. The ions that travel to the cathode will travel slowly through the potential minimum and reduce the space charge, which in turn will reduce the space-charge smoothing effect. This also will increase the noise in the anode circuit.

Induced Noise

At high frequencies it is not necessary for electrons to reach an electrode for induced current to flow in the electrode leads. This noise is an important consideration in miniature tubes above 15 megahertz and becomes the principal limiting factor in low-noise amplifier design above about 100 megahertz. For microwave tubes, this is the dominant method by which beam noise is coupled to the output circuit.

Miscellaneous Noise

Other noise may be present as the result of microphonics, hum, leakage, charges on insulators, poor contacts, and secondary emission.

Microwave Tubes

The noise appearing in the output circuit of a microwave tube is due in part to induced noise from the beam. Also, some of the electrons may be intercepted by the radio-frequency structure (microwave cavity, slow-wave circuit, etc.) giving rise to partition noise. In well-designed low-noise tubes, however, the latter effect is kept negligibly small.

For lossless linear beam tubes (traveling-wave amplifiers, klystron amplifiers, backward-wave amplifiers), the minimum obtainable noise figure F_{\min} for one-dimensional single-velocity small-signal theory and high gain has been found to be given by

$$F_{\min} = 1 + (2\pi/kT_0)(S - \pi)$$

where $S - \pi$ is the basic noise parameter and is established in the region of the potential minimum of the beam. If certain assumptions concerning the potential minimum are made, such as full shot noise and uncor-

related current and velocity fluctuations,* then values for S and π can be obtained. They are given as

$$\pi = 0$$

$$S = [1 - (\pi/4)]^{1/2} (kT_c/\pi)$$

Therefore

$$F_{\min} = 1 + (4 - \pi)^{1/2} (T_c/T_0)$$

For $T_c/T_0 = 4$, $F_{\min} = 4$. The assumptions made are not entirely valid, as shown by the fact that noise figures of less than 4 have been obtained experimentally. At the present time values of S and π/S are obtained by measurement.

LOW-, MEDIUM-, AND HIGH-FREQUENCY TUBES

This section applies particularly to triodes and multigrid tubes operated at frequencies where electron-inertia effects are negligible. Traditionally, the vacuum envelope of such tubes has been of glass with metal, usually copper, for the anode in larger sizes. In recent years, the trend has been toward using ceramic in place of glass for the external insulating portions of such tubes. Fig. 12 shows a typical construction of a medium-power transmitting tube.

Ceramic-envelope tubes have the following advantage over glass tubes.

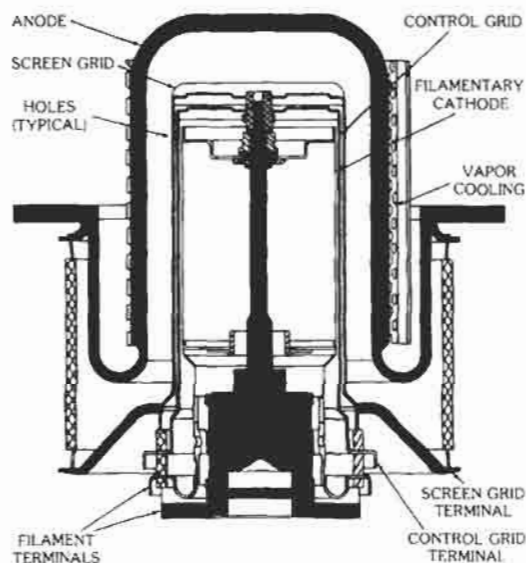


Fig. 12. Typical medium-power transmitting tube.

* J. R. Pierce, "A theorem concerning noise in electron streams," *Journal of Applied Physics*, vol. 25, p. 931, 1954.

(A) The radio-frequency loss P_{rf} in the seals of a tube is given by

$$P_{rf} = Kf^{5/2}R^{1/2}/\mu^{1/2}$$

where

- K = constant
- f = frequency
- R = resistivity of the conducting material
- μ = permeability of the conducting material.

In glass-to-metal seals, the metal is normally of a magnetic material such as Kovar. As Kovar has high resistivity and permeability, the radio-frequency losses at the seals are therefore high, and at high frequencies cracking and/or glass suck-in near the seals can result. With ceramic-to-metal seals, this problem is minimized because the radio-frequency circulating currents at the seals flow through the metallizing and plating on the ceramic. The resistivity is low, and the permeability is unity.

- (B) Ceramics have a lower dielectric loss than glass. Furthermore, the loss factor of glass rapidly rises with temperature. This leads to a "runaway" condition, glass suck-in, and hence severe limitation of maximum frequency of operation of glass tubes.
- (C) The safe operating temperature of a ceramic-to-metal seal may be between 220 and 250°C as against 180°C for Kovar-glass seals.
- (D) The high bakeout temperature of ceramic-envelope tubes during evacuation increases reliability and life.
- (E) Ceramic tubes withstand higher thermal and mechanical shocks than those with glass envelopes. They can also be manufactured to closer dimensional tolerances.

Coefficients

Amplification factor μ : Ratio of incremental plate voltage to control-electrode voltage change at a fixed plate current with constant voltage on other electrodes.

$$\mu = \left[\frac{\delta e_b}{\delta e_{c1}} \right]_{i_b, E_{c2}, \dots, E_{cn} \text{ constant}}$$

Transconductance s_m : Ratio of incremental plate current to control-electrode voltage change at constant voltage on other electrodes.

$$s_m = \left[\frac{\delta i_b}{\delta e_{c1}} \right]_{E_b, E_{c2}, \dots, E_{cn} \text{ constant}}$$

When electrodes are plate and control grid, the ratio is the mutual conductance g_m .

$$g_m = \mu/r_p$$

Variational (ac) plate resistance r_p : Ratio of incremental plate voltage to current change at constant voltage on other electrodes.

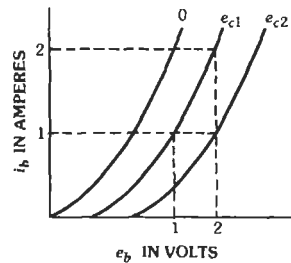
$$r_p = \left[\frac{\delta e_b}{\delta i_b} \right]_{E_{c1}, \dots, E_{cn} \text{ constant}}$$

Total (dc) plate resistance R_p : Ratio of total plate voltage to current for constant voltage on other electrodes.

$$R_p = \left[\frac{E_b}{I_b} \right]_{E_{c1}, \dots, E_{cn} \text{ constant}}$$

A useful approximation of these coefficients may be obtained from a family of anode characteristics, Fig. 13. Relationships between the actual geometry of a tube and its coefficients are given roughly in Table 4.

When the operating frequency is increased, the operation of triodes and multigrid tubes is affected by electron-inertia effects. The design features that distinguish the high-frequency tube shown in Fig. 15 from the lower-frequency tube (Fig. 12) are: reduced cathode-to-grid and grid-to-anode spacings, high emission density, high power density, small active and inactive capacitances, heavy terminals, short support leads, and adaptability to a cavity circuit.



- Amplification factor $\mu = (e_{b2} - e_{b1}) / (e_{c2} - e_{c1})$
- Mutual conductance $g_m = (i_{b2} - i_{b1}) / (e_{c2} - e_{c1})$
- Total plate resistance $R_p = e_{b2} / i_{b2}$
- Variational plate resistance $r_p = (e_{b2} - e_{b1}) / (i_{b2} - i_{b1})$

Fig. 13. Graphic method of determining coefficients.

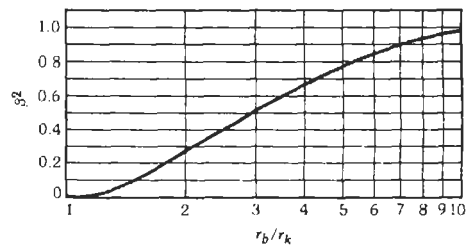


Fig. 14. Values of β^2 for values of $r_b/r_k < 10$.

TABLE 4. TUBE CHARACTERISTICS FOR UNIPOTENTIAL CATHODE AND NEGLIGIBLE SATURATION OF CATHODE

Function	Parallel-Plane Cathode and Anode	Cylindrical Cathode and Anode
Diode anode current (amperes)	$G_1 e_b^{3/2}$	$G_1 e_b^{3/2}$
Triode anode current (amperes)	$G_2 [(e_b + \mu e_c)/(1 + \mu)]^{3/2}$	$G_2 [(e_b + \mu e_c)/(1 + \mu)]^{3/2}$
Diode perveance G_1	$2.3 \times 10^{-6} (A_p/d_b^2)$	$2.3 \times 10^{-6} (A_p/\beta^2 r_b^2)$
Triode perveance G_2	$2.3 \times 10^{-6} (A_p/d_b d_c)$	$2.3 \times 10^{-6} (A_p/\beta^2 r_i r_c)$
Amplification factor μ	$2.7 d_c [(d_b/d_c) - 1] / \rho \log(\rho/2\pi r_g)$	$(2\pi d_c/\rho) [\log(d_b/d_c) / \log(\rho/2\pi r_g)]$
Mutual conductance g_m	$1.5 G_2 \mu / (\mu + 1) (E'_g)^{1/2}$ $E'_g = (E_b + \mu E_c)/(1 + \mu)$	$1.5 G_2 [\mu / (\mu + 1) (E'_g)^{1/2}]$ $E'_g = (E_b + \mu E_c)/(1 + \mu)$

A_b = effective anode area in square centimeters; d_b = anode-cathode distance in centimeters; d_c = grid-cathode distance in centimeters; β = geometric constant (a function of the ratio of anode-to-cathode radius), $\beta^2 = 1$ for $r_b/r_c > 10$ (Fig. 14); ρ = pitch of grid wires in centimeters; r_g = grid-wire radius in centimeters; r_i = anode radius in centimeters; r_c = cathode radius in centimeters; r_g = grid radius in centimeters.

Note:

These equations are based on theoretical considerations and do not provide accurate results for practical structures; however, they give a fair idea of the relationship between the tube geometry and the constants of the tube.

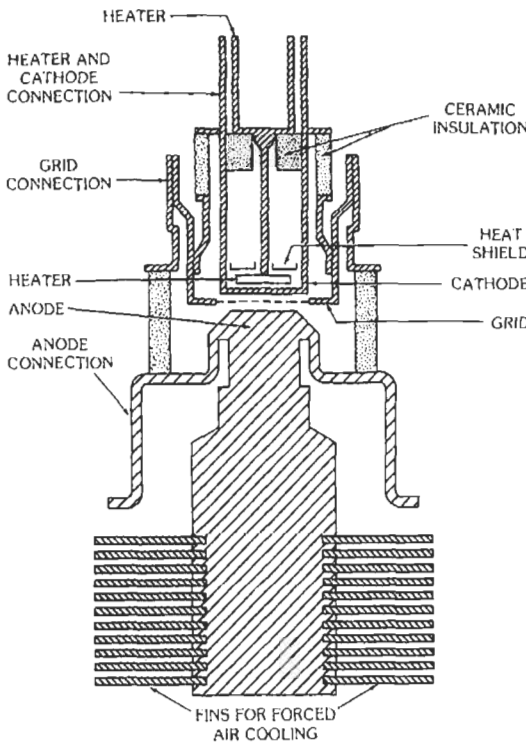


Fig. 15. Electrode arrangement of a small high-frequency external-anode triode.

MATERIALS AND STRUCTURES

Cathodes

Early cathodes were all directly heated, that is, the filament, which was heated by passing a current

through it, was also the cathode. The filaments were made of mostly tungsten (W) or thoriated W (i.e., W mixed with a small percentage of thoria). Thoriated W cathodes are still commonly used in high-power switch tubes, transmitter tubes, and microwave magnetrons at operating temperatures of 1600–1900°C.

In a directly heated cathode, the heating current causes a voltage drop along the current path and across the surface. Most applications require indirectly heated cathodes so that the surface of the cathode is at a constant potential. Oxide-coated cathodes and dispenser-type cathodes use mostly indirect heaters, with the heater filament made of W. A potted heater, shown in Fig. 16, uses alumina (Al_2O_3) potting material to fill the space in the cavity that contains the heater, whereas a nonpotted heater has alumina coating on the surface of the heater filament only. A potted heater maintains good thermal contact with the cathode body and provides uniform heating across the cathode. A nonpotted heater has a low heat capacity and is suitable for applications that require fast warm-up. In addition, nonpotted heaters cost much less than potted heaters.

The emitting surfaces of oxide-coated cathodes are initially made of Ba, Ca, and Sr carbonates, which are converted thermally to respective oxides before electron emission can occur. The conversion occurs around 800°C and the operating temperatures are 750–800°C. The emitting surfaces of dispenser-type cathodes are made of pressed tungsten powder impregnated with BaO, CaO, and Al_2O_3 . The surfaces of M-type dispenser cathodes are coated with ~5000 Å of noble metal, most commonly an Os/Ru mixture. The operating temperatures of dispenser-type cathodes are typically 900–1250°C.

Grids

Grids are conductors, semitransparent to electrons, usually in the form of meshes or parallel bars. Their

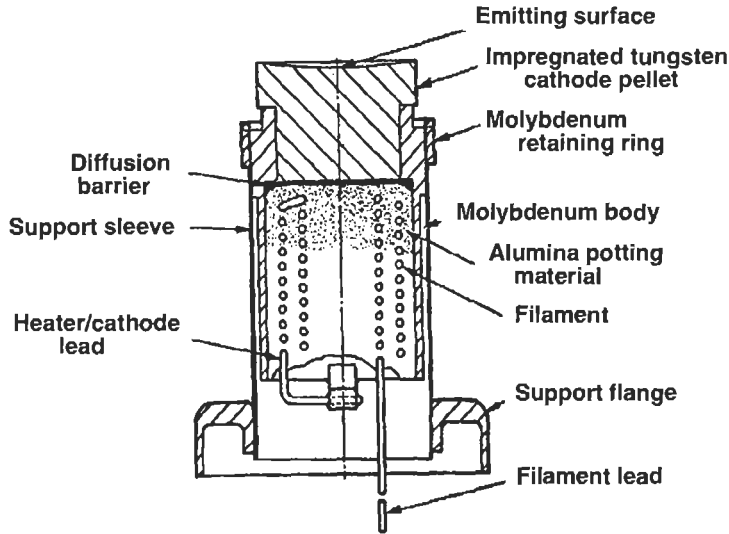


Fig. 16. Cross-sectional view of a typical cathode and heater assembly. The materials used are indicated by the labels in the figure. (From A. S. Gilmour, Jr., *Microwave Tubes*, Norwood, MA: Artech House, Inc., 1986, p. 126, fig. 5.20.)

function is to control or accelerate electrons by establishing appropriate electric fields. An ideal grid would intercept no electrons; in practice, interception is minimized by careful design but still occurs, causing heating that usually sets the upper limit to power output of the tube.

In tubes with oxide cathodes, grids may be tungsten or molybdenum wire structures, coated with gold to reduce primary emission caused by deposition of material evaporated from the cathode during operation. In tubes with thoriated tungsten cathodes, grids may be made of tungsten or molybdenum coated with proprietary compounds to reduce primary emission. Relatively recently, pyrolytic graphite has been found to be an excellent grid material because of its combination of good thermal and mechanical properties and its low primary and secondary emission.* Grids made of this material behave elastically up to very high temperatures and as a result are more dimensionally stable than refractory metal wire grids. This material is anisotropic, is formed by chemical vapor deposition at high temperature (2000°C or higher), and should not be confused with ordinary graphite, which is isotropic.

Anodes

Anodes are usually made of copper, cooled by air or water. Some glass-envelope tubes have molybdenum internal anodes. Carbon has been used. Insulating structures inside the vacuum envelope are usually high-purity aluminum oxide or beryllium oxide. Con-

ductors are molybdenum, tungsten, nickel, Kovar, cupro-nickel, Hastalloy, or copper.

Getters

Getters are used to maintain a high vacuum during the life of the tube. They are so called because they "get" or trap and hold gases that may evolve. Typical materials are zirconium, cerium, barium, and titanium.

Tube Geometry

Gridded tubes at all power levels for frequencies up to about 1 GHz are invariably cylindrical in form. The tetrode in Fig. 12 is a typical example. The anode is vapor cooled by means of a wire mesh, or "wick," held in contact with the anode. This is a substitute for the knobs or fins usually found on vapor-cooled anodes and is satisfactory at medium power density. At higher frequencies, the "planar" triode structure is almost universal. An example is shown in Fig. 15.

The number of electrodes is usually three (triode) or four (tetrode) and rarely five (pentode). That is, the number of grids may be one, two, or three. An unusual hybrid tube invented in 1939 by Haeff†‡ and recently rediscovered and improved has been named Klystrode TM. Its geometry is part triode, part klystron. Its performance is described in Chapter 17, which covers the

† A. V. Haeff, "An UHF power amplifier of novel design," *Electronics*, Feb. 1939, pp. 30–32.

‡ D. Priest and M. B. Shrader, "The klystrode—An unusual transmitting tube with potential for UHF-TV," *Proc. IEEE*, Nov. 1982.

* W. H. Smith and D. H. Leeds, *Pyrolytic Graphite*, *Modern Materials*, Vol. 7 (New York: Academic Press, 1970), pp. 139–221.

performance of tubes and some typical associated circuits.

MICROWAVE TUBES*

The reduced performance of space-charge control tubes in the microwave region has fostered the development of other types of tubes for use as oscillators and amplifiers at microwave frequencies. Such tubes generally function on the basis of the modulation of the velocity of an electron stream rather than of its density. They may be roughly divided simply into linear-beam devices and crossed-field devices. In the former, the electron stream flows essentially linearly, often with a collimating magnetic field to counteract space-charge spreading; in the latter, the electron stream follows a curved path under the action of orthogonal electric and magnetic fields. The linear-beam devices are often referred to as *O*-type, while the crossed-field devices are referred to as *M*-type.

Terminology

Bunching: Any process that introduces a radio-frequency conduction-current component into a velocity-modulated electron stream as a direct result of the variation in electron transit time that the velocity modulation produces.

Cavity resonator: Any region bounded by conducting walls within which resonant electromagnetic fields may be excited.

Circuit efficiency: The ratio of (A) the power of the desired frequency delivered to the output terminals of the circuit of an oscillator or amplifier to (B) the power of the desired frequency delivered by the electron stream to the circuit.

Coherent-pulse operation: Method of pulse operation in which the phase of the radio-frequency wave is maintained through successive pulses.

Drift space: In an electron tube, a region substantially free of externally applied alternating fields

in which a relative repositioning of the electrons is determined by their velocity distributions and the space-charge forces.

Duty cycle: The product of the pulse duration and the pulse repetition rate. It is also the ratio of the average power output to the peak power output.

External Q : The reciprocal of the difference between the reciprocals of the loaded Q and the unloaded Q .

Frequency pulling of an oscillator: The change in the generated frequency caused by a change of the load impedance.

Frequency pushing of an oscillator: The change in frequency caused by a change in anode current (or in anode voltage).

Loaded Q of a specific mode of resonance of a system: The Q when there is external coupling to that mode. Note: When the system is connected to the load by means of a transmission line, the loaded Q is customarily determined when the line is terminated in its characteristic impedance.

Mode: One of the components of a general configuration of a vibrating system. A mode is characterized by a particular geometric pattern of the electromagnetic field and a resonant frequency (or propagation constant).

Noise Figure: The ratio in decibels of the total available output noise from an amplifier to the available noise that would be present at the output if the amplifier itself were noiseless, assuming a source temperature of 290 K.

Perveance: The ratio of electron beam current to the $3/2$ power of the beam voltage ($IV^{-3/2}$), which is an invariant for a particular electron-gun design.

Pulling figure of an oscillator: The difference in megahertz between the maximum and minimum frequencies of oscillation obtained when the phase angle of the load-impedance reflection coefficient varies through 360 degrees, while the absolute value of this coefficient is constant and is normally equal to 0.20.

Pulse: Momentary flow of energy of such short time duration that it may be considered as an isolated phenomenon.

Pushing figure of an oscillator: The rate of frequency pushing in megahertz per ampere or megahertz per volt.

Q : The Q of a specific mode of resonance of a system is 2π times the ratio of the stored electromagnetic energy to the energy dissipated per cycle when the system is excited in this mode.

Reflector: Electrode whose primary function is to reverse the direction of an electron stream. It is also called a *repeller*.

Reflex bunching: Type of bunching that occurs when the velocity-modulated electron stream is made to reverse its direction by means of an opposing direct-current field.

Slow-wave structure: A microwave circuit, as used in beam-type microwave tubes, capable of propa-

* A. H. W. Beck, *Space-Charge Waves and Slow Electromagnetic Waves* (New York: Pergamon Press, 1958); J. W. Gewartowski and H. A. Wason, *Principles of Electron Tubes* (Princeton, NJ: D. Van Nostrand Co., 1965); A. S. Gilmour, Jr., *Microwave Tubes* (Boston: Artech House, 1986); R. G. E. Hutter, *Beam and Wave Electronics in Microwave Tubes* (Princeton, NJ: D. Van Nostrand Co., 1960); W. J. Kleen, *Electronics of Microwave Tubes* (New York: Academic Press, 1958); L. Sivan, *Microwave Tube Transmitters* (London: Chapman and Hall, 1994); J. C. Slater, *Microwave Electronics* (Princeton, NJ: D. Van Nostrand Co., 1950); J. F. Hull, "Microwave tubes of the mid-sixties," 1965 *IEEE International Convention Record*, IEEE, New York; "Special section on high-power microwave tubes," *Proceedings of the IEEE*, vol. 61, no. 3, March, 1973. "Special issue: New vistas for vacuum electronics," *Proceedings of the IEEE*, vol. 87, no. 5, May, 1999.

gating radio-frequency waves with phase velocities appreciably less than the velocity of light.

Recent Trends

Dramatic progress in device performance and capability has been achieved in the 1980s and 1990s. These advances can be credited to a combination of device innovation, enhanced understanding gained through improved modeling and design, introduction of superior materials and subassembly components, and development of advanced vacuum processing and construction techniques. The most familiar types of microwave power tube, invented in the decade spanning the mid-1930s and 1940s, may be broadly classified as "slow-wave" devices. This description as being slow-wave emphasizes the fundamental requirement in these devices that the phase velocity of the electromagnetic wave on the interaction circuit must be approximately equal to the velocity of the electron stream if energy transfer between the beam and the wave is to occur. The electron velocity is less than the speed of light and hence the wave phase velocity must also be slower than the speed of light in free space. From the richness of device concepts investigated through the 1960s, the helix and coupled cavity traveling-wave tubes (TWTs), the klystron, magnetron, and crossed field amplifier (CFA) have emerged as the primary products of today's technology. The continuing vitality of this relatively mature segment of the technology is demonstrated by the recent development of a hybrid amplifier architecture, the microwave power module (MPM), that combines the best features of solid-state and vacuum technologies. This hybrid approach yields a cost-competitive amplifier that offers the power, efficiency, and bandwidth of a TWT with the reduced noise and functionality of a solid-state amplifier with either discrete or monolithic components. Most significantly, this performance is obtained in a miniaturized package that can be a small fraction of the size of a comparable TWT amplifier (TWT-A) or an all solid-state power amplifier (SSPA).

The most dramatic of the modern trends, however, is the extension of device operation into the spectral range well beyond the microwave. Beginning in the mid-1960s, two new classes of powerful "fast-wave" device, the gyro-oscillator/amplifier and the free electron laser (FEL), have experienced ongoing development. In fast-wave devices, the phase velocity of the electromagnetic wave exceeds the speed of light in the interaction region. Although the physical interaction mechanism is not the same, gyro-devices are technologically very similar to other microwave power tubes. They differ primarily by their intrinsic need for magnetic field strengths typically available only from superconducting magnets. In recent years, gyro-oscillators have produced output powers approaching a megawatt in the frequency range of 100–200 GHz with the pulse duration nearing a second. In the guise of the ubitron, free electron masers

are much akin to their slow-wave cousins. By using electrons at relativistic energies, the operational domain of the FEL is extended to much shorter wavelength. For this extension, a high-energy electron accelerator replaces the electron gun and optical cavities replace the interaction circuits of the microwave power tube. Free electron lasers are now generating usable levels of tunable power at submillimeter wavelengths and in the infrared and promise to provide high-power, continuously tunable coherent sources over a spectral range extending into the ultraviolet and eventually to x rays.

For very general comparison of device limits, the average power density, $P_{av}f^2$, in units of (MW)(GHz)², has been used as a figure of merit providing an insightful basis for comparing device types and portraying growth trends in performance. The physical significance of $P_{av}f^2$ as a figure of merit derives from the fact that the maximum beam or charge-carrier power that can be transported through a device is proportional to the circuit cross-sectional area, which is inversely proportional to the operating frequency. Upper limits in "power density" were therefore considered as intrinsic to a device design. With this in mind, the progression of power density for major device types is plotted in Fig. 17. In general, the power density profile for each device type exhibits the characteristic S-shaped curve associated with technological maturity. The slow progress typical of a period of innovation and preliminary study is followed by the rapid growth associated with intense development. As the fundamental limits of a particular device are approached, progress becomes more difficult and the curve flattens. Escape from the limits imposed by a particular device concept requires finding a better approach. Following this argument, we conclude that the vitality of a technology is measured by its ability to provide research opportunities leading to the next step.

As seen by the sustained growth in achieved power density, RF vacuum electronics continue to provide such opportunity. The history of RF vacuum electronics can be viewed as successive waves of scientific innovation driven by technical opportunity coupled with market demand. First, the gridded tube is overtaken by the magnetron under the impetus of wartime demands. Then, the high-power linear-beam tubes, the klystron and the coupled cavity TWT, move to the fore in response to accelerator and radar interests. Finally, the gyro-oscillator and the FEL surpass slow-wave devices in the mid-1970s. With these transitions, the envelope of vacuum device performance, as defined by the straight line in Fig. 17, continues to grow exponentially at an approximate rate of one-and-a-half orders of magnitude per decade, a rate sustained for the past 60 years. In contrast to the fast-wave devices, the figures of merit for the slow-wave devices have grown very little since the mid-1960s.

One significant exception to limited growth with slow-wave devices is the periodic-permanent-magnet (PPM-) focused helix TWT, whose pace of growth

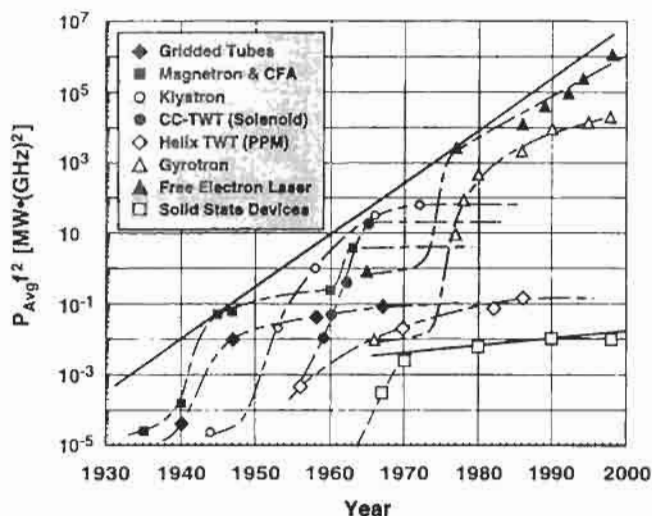


Fig. 17. Progression of device power density, $P_{avg} f^2$, for major device types.

continues but at a somewhat diminished rate. This continued growth of the PPM-focused TWT is not surprising given its dominance in the marketplace. As measured in terms of the basic parameters of information theory, the capabilities of the conventional TWT are extraordinary. The TWT is a key building block in the rapidly expanding telecommunication systems and retains a major role in satellite-based transmitters. With their increasing demand for bandwidth, military radar, electronic countermeasures, and communication data links depend heavily on TWTs. In a package weighing less than a pound, the modern TWT provides average powers from tens to hundreds of watts combined with a gain of ~ 60 dB and an instantaneous bandwidth of well over an octave. The basic distortion parameters can satisfy the most exacting requirements and the overall efficiency can be as high as 70%. Moreover, the demonstrated mean lifetime of TWTs in the most demanding environment of space exceeds 15 years ($> 130\,000$ hrs).

In general, the character of the research and developmental effort associated with fast- and slow-wave devices is commensurate with the relative state of maturity (as defined within the context of power density limits) and the market served. Fast-wave efforts retain a strong research flavor while having become progressively more developmental since the mid-1980s. The areas of application (thermonuclear plasma heating; high-power, millimeter-wave radar; RF drivers for high-energy particle accelerators; and biological, chemical, and medical research in the submillimeter and IR) served by fast-wave devices still require higher power performance at high frequency. Hence, market pull and technological opportunity coalesce to sustain the pursuit of power at high frequency with fast-wave devices.

In contrast, the development of conventional microwave power tubes is strongly driven by the market.

Most current military, commercial, and scientific applications do not require simply power, but rather some combination of more bandwidth, efficiency, and linearity or lower noise and intermodulation. As new systems become more complex and expensive, the components must become smaller and lighter as well as more reliable. These demands have increased the importance of technology in determining the opportunity for advance. Revolutionary change has given way to iterative improvement. Although these steps may be evolutionary, the impact can be dramatic and important advances continue to be achieved within this broader view of performance. The current capabilities of the major device types is given in Fig. 18 in terms of average power with frequency.

Linear-Beam Tubes

The principal types of linear-beam tubes are the klystron, the traveling-wave amplifier, and the backward-wave oscillator.

Klystrons—A klystron* is an electron tube in which the following processes may be distinguished:

- (A) Periodic variations of the longitudinal velocities of the electrons forming the beam in a region confining a radio-frequency field

* D. R. Hamilton, J. K. Knipp, and J. B. H. Kuper, *Klystrons and Microwave Triodes* (New York: McGraw-Hill Book Co., 1948); A. H. W. Beck, *Velocity-Modulated Thermionic Valves* (London: Cambridge University Press, 1948); A. H. W. Beck, *Thermionic Valves, Their Theory and Design* (London: Cambridge University Press, 1953); A. H. W. Beck, *Space-Charge Waves and Slow Electromagnetic Waves* (New York: Pergamon Press, 1958).

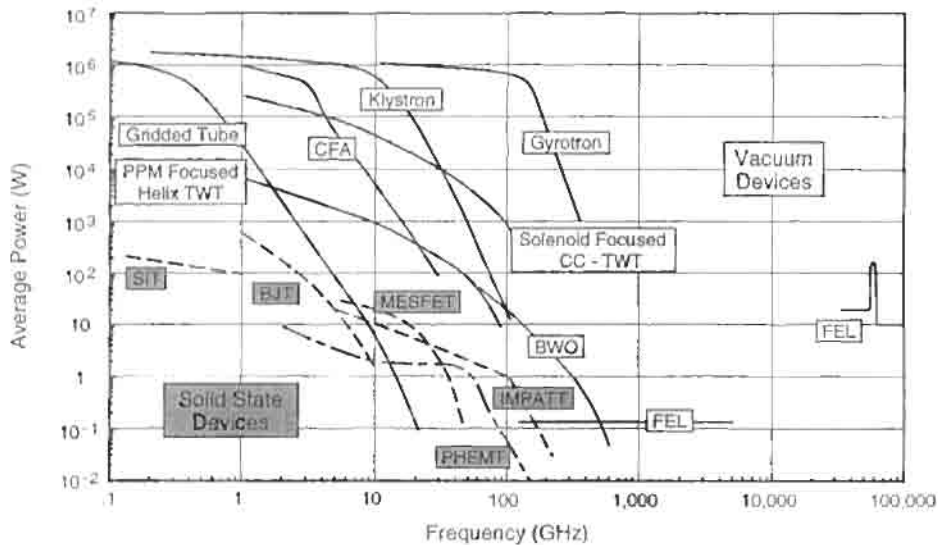


Fig. 18. RF Vacuum electronic and solid-state devices: current state of technology for single device performance.

- (B) Conversion of the velocity variation into conduction-current modulation by motion in a region free from radio-frequency fields
 (C) Extraction of the radio-frequency energy from the beam in another confined radio-frequency field

The transit angles in the confined fields are made short ($\delta \cong \pi/2$) so that there is no appreciable conduction-current variation while the electrons traverse them.

Several variations of the basic klystron exist. These include the two-cavity amplifier or oscillator, the reflex klystron, and the multicavity high-power amplifier.

Two-Cavity Klystron Amplifiers: An electron beam is formed in an electron gun and passed through the gaps associated with the two cavities (Fig. 19). After emerging from the second gap, the electrons pass to a collector designed to dissipate the remaining beam power without the production of secondary electrons. In the first gap, the electron beam is alternately accelerated and decelerated in succeeding half-periods of the radio-frequency cycle, the magnitude of the change in speed depending on the magnitude of the alternating voltage impressed on the cavity. The electrons then move in a drift space where there are no radio-frequency fields. Here, the electrons that were accelerated in the input gap during one half-cycle catch up with those that were decelerated in the preceding half-cycle, and a local increase of current density occurs in the beam. Analysis shows that the maximum of the current-density wave occurs at the position, in time and space, of those electrons that passed the center of the input gap as the field changed from negative to posi-

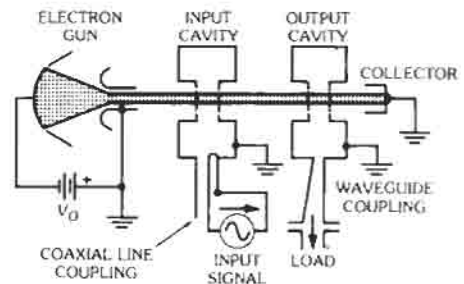


Fig. 19. Two-cavity klystron amplifier. (From J. W. Geurtowski and H. A. Watson, *Principles of Electron Tubes*, 1965; p. 296. Courtesy D. Van Nostrand Co., Inc.)

tive. There is therefore a phase difference of $\pi/2$ between the current wave and the voltage wave that produced it. Thus at the end of the drift space, the initially uniform electron beam has been altered into a beam showing periodic density variations. This beam now traverses the output gap, and the variations in density induce an amplified voltage wave in the output circuit, phased so that the negative maximum corresponds with the phase of the bunch center. The increased radio-frequency energy has been gained by conversion from the direct-current, beam energy.

The two-cavity amplifier can be made to oscillate by providing a feedback loop from the output to the input cavity. A much simpler but less efficient structure results if the electron-beam direction is reversed by a negative electrode, termed the reflector.

*Reflex Klystrons:** A schematic diagram of a reflex klystron is shown in Fig. 20. The velocity-modulation process takes place as before, but analysis shows that in the retarding field used to reverse the direction of electron motion, the phase of the current wave is exactly opposite to that in the two-cavity klystron. When the bunched beam returns to the cavity gap, a positive field extracts maximum energy from the beam, since the direction of electron motion has now been reversed. Consideration of the phase conditions shows that for a fixed cavity potential, the reflex klystron will oscillate only near certain discrete values of reflector voltage for which the transit time measured from the gap center to the reflection point and back is given by

$$w\tau = 2\pi(N + 3/4)$$

where N is an integer called the mode number.

By varying the reflector voltage around the value corresponding to the mode center, it is possible to vary the oscillation frequency by a small percentage. This fact is made use of in providing automatic frequency control or in frequency-modulation transmission.

Reflex-Klystron Performance Data: The performance data for a reflex klystron are usually given in the form of a reflector-characteristic chart. This chart displays power output and frequency deviations as a function of reflector voltage. Several modes are often displayed on the same chart. A typical chart is shown in Fig. 21.

Table 5 shows typical reflex-klystron performance.

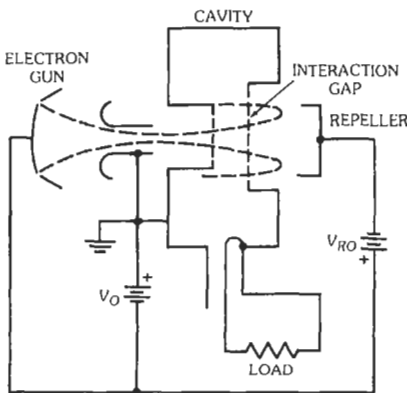


Fig. 20. Schematic diagram of reflex klystron with power supply. (From J. W. Gevartowski and H. A. Watson, *Principles of Electron Tubes*, 1965; p. 311. Courtesy D. Van Nostrand Co., Inc.)

* J. R. Pierce and W. G. Shepherd, "Reflex oscillators," *Bell System Technical Journal*, vol. 26, pp. 460-681, July, 1947.

Multicavity Klystrons: Multicavity klystrons† have been perfected for use in two rather different fields of application: 1. those requiring extremely high pulse powers and 2. continuous-wave (CW) systems in which a few kilowatts to hundreds of kilowatts are required. Examples of the first application are power sources for nuclear-particle acceleration and radar; ultra-high-frequency television, troposcatter, satellite-groundstation, and space-communication transmitters are examples of the latter.

A multicavity klystron amplifier is shown schematically in Fig. 22. The example shown has three cavities, all coupled to the same beam, although as many as seven cavities are sometimes used. The radio-frequency input modulates the beam as before. The bunched beam induces an amplified voltage across the second cavity, which is tuned to the operating frequency. This amplified voltage remodulates the beam with a certain phase shift, and the now more strongly bunched beam excites a highly amplified wave in the output circuit. It is found that the optimum power output is obtained when the second cavity is slightly detuned. When four or more cavities are used, optimum efficiency is obtained with both the penultimate and the antepenultimate cavities tuned above the operating frequency. Moreover, when increased bandwidth is required, the earlier cavities may be loaded. Modern multicavity klystrons use magnetically focused, high-perveance beams, and under these conditions high

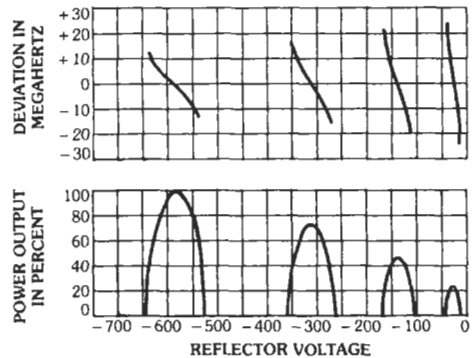


Fig. 21. Klystron reflector-characteristic chart. (Courtesy Sperry Gyroscope Co.)

† M. Chodorow, E. L. Ginzton, I. R. Neilson, and S. Sonkin, "Design and performance of a high-power pulsed klystron," *Proceedings of the IRE*, vol. 41, pp. 1584-1602, November, 1953; D. H. Priest, C. E. Murdock, and J. J. Woerner, "High-power klystrons at UHF," *Proceedings of the IRE*, vol. 41, pp. 20-25, January, 1953; A. Staprans, E. McCune, and J. Ruetz, "High-power linear beam tubes," *Proceedings of the IEEE*, vol. 61, pp. 299-329, March, 1973.

TABLE 5. CLASSES OF REFLEX KLYSTRONS

Frequency (MHz)	Power Output (mW)	Useful Mode Width Δf_{3dB} (MHz)	Operating Voltage
Local Oscillators			
3 000	150	40	300
9 000	40	40	350
24 000	35	120	750
35 000	> 15	50	2 000
50 000	10-20	60-140	600
Master Pumps			
35 000	500-1 500	70	2 000
45 000	500-1 000	80	2 000
Frequency-Modulation Transmitters			
4 000	10 000	40	1 100
7 000	10 000	37	750
9 000	6 000	60	500

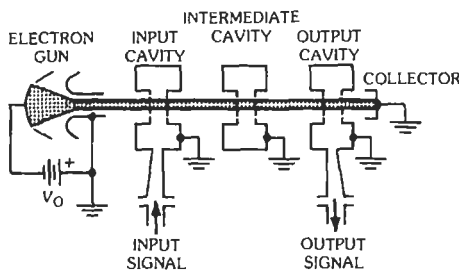


Fig. 22. Three-cavity klystron amplifier. (From J. W. Gewartowski and H. A. Watson, *Principles of Electron Tubes*, 1965; p. 340. Courtesy D. Van Nostrand Co., Inc.)

gains, large power output, and reasonable values of efficiency are readily obtained.

Fig. 23 shows the power output and tuning range versus frequency for a number of typical CW klystrons. Efficiencies from 30% to 75% are typical. The lower efficiencies apply for high-perveance ($3 \times 10^{-6} A/V^{3/2}$) stagger-tuned, broad-band klystrons, and the higher efficiencies apply for narrow-band, lower-perveance ($0.5 \times 10^{-6} A/V^{3/2}$) klystrons. Average output power scales as about λ^3 ($\lambda =$ wavelength) at constant perveance if cavity losses limit output power. At the lowest frequencies, the limit on power is imposed by single-mode-waveguide, output window capability, and power scales as A . At the highest frequencies, cathode current and cavity losses impose the limit, and the power scales as $\lambda^{9/2}$.

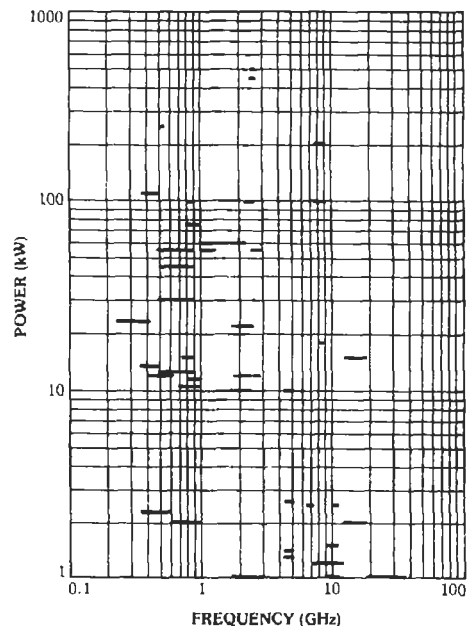


Fig. 23. Output power vs frequency for a representative group of CW klystrons.

Fig. 24 shows the power output and tuning range versus frequency for pulsed klystrons. Because most pulsed klystrons use higher-perveance beams (1.5 to $3 \times 10^{-6} A/V^{3/2}$) to keep the beam voltage low, efficiencies range from 30% to 50%. As much as 11.8% band-

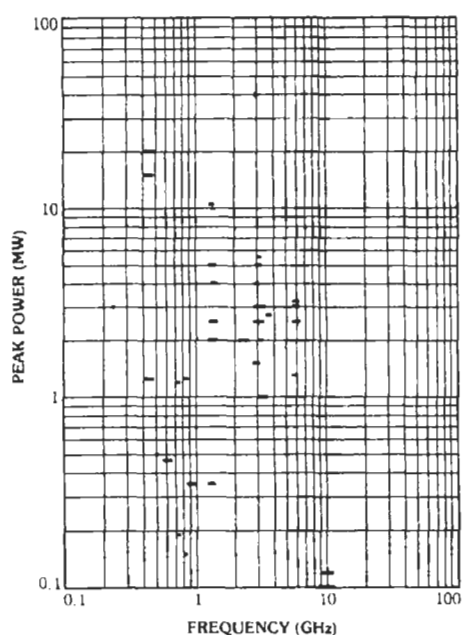


Fig. 24. Output power vs frequency for a representative group of pulsed klystrons.

width has been achieved from a 10-MW klystron with a 1000-ohm "beam impedance." Even wider bandwidth, at lower power, has been achieved recently by replacing the individual intermediate cavities along the beam of a multicavity klystron with groups or clusters of two or three closely spaced artificially loaded low- Q cavities. In this way, two or three times the bandwidth can be achieved while preserving the overall dimensions of the tube as well as the gain.*

Bandwidth scales inversely with the ratio of dc beam voltage to dc beam current, or "beam impedance," in both CW and pulse klystrons. Peak output power from pulsed klystrons is usually limited by the beam voltage that can be applied to the electron gun without giving rise to arcing problems and by the current that can be put through the beam tunnel while keeping the cathode current density reasonable. For this reason, maximum peak power from klystrons is relatively constant up to about 3000 MHz and then falls as λ^2 .

Multiple-Beam Klystrons: Multiple-beam klystrons (MBKs) are useful in systems that require high peak/average power with relatively high instantaneous bandwidths (up to approximately 15%) and/or systems

where relatively low operating voltages are desired. They are used in a broad range of applications, ranging from compact, low-power C-, X-, and Ku-band 10–500 watt amplifiers for communications, airborne radar, and missile seekers, to high power L-, S-, and X-band 100–1000 kilowatt amplifiers for radar and broadcasting.

The concept of the multiple-beam klystron appeared in the early 1940s,† soon after the invention of the conventional (single-beam) klystron by the Varian brothers. The development of the multiple-beam klystron was driven by the desire of microwave source developers to produce a given amount of microwave power at as low an operating voltage as possible, since, in general, lower voltages make the design of power supplies simpler, lighter, cheaper, and more reliable. However, to keep the beam power constant, any reduction in the beam voltage must be accompanied by a corresponding increase in beam current. In conventional klystrons (as well as in many other microwave tubes), an increase in electron current density increases the role of space-charge forces. When these forces are negligibly small, electrons will ballistically bunch in the drift channel of the klystron because of the initial modulation of electron velocities at the input cavity gap. At high current densities, however, space-charge forces are not negligible and the repulsive forces will interfere with the ballistic bunching, leading to a degradation in efficiency.

It occurred to klystron designers that higher current beams could be achieved by separating the single electron beam into multiple beamlets, each with a sufficiently low current density to ensure efficient electron bunching. Each beamlet is transported down its own drift channel, parallel to, but isolated from, the other beamlets. The beamlets are allowed to interact only over small axial distances corresponding to cavity gap locations. Klystrons with such an electron beam generation and transport design are termed "multiple-beam klystrons." Over the years, many multiple-beam klystron configurations have been considered.‡

An important advantage of a multiple-beam klystron is the potential for enhanced instantaneous bandwidths relative to single-beam klystrons. Since the total current and perveance of the beam is higher in multiple-beam klystrons than in conventional klystrons, the quality factors of the cavities in the multiple-beam klystrons will be lower. The lower quality factors increase the width of the resonance curves of the modes in the cavities, thus increasing the operating bandwidth.§ Another significant advantage of the multiple-beam configuration is a reduction in operating voltage relative to conventional klystrons. This reduc-

* R. S. Synons, B. Arfin, R. E. Boesenberg, P. E. Ferguson, M. Kirshner, and J. R. M. Vaughan, "An experimental clusiced-cavityTm klystron," Paper 7.4, *Technical Digest, International Electron Devices Meeting*, Washington, DC, Dec. 6–9, 1987, IEEE.

† J. Bernier, Patent No. 992-853, September 15, 1944 (France).

‡ M. R. Boyd, R. A. Dehn, J. S. Hickey, and T. G. Mhiran, "The multiple-beam klystron." *IRE-ED* vol. 9, pp. 247–252 (1962).

tion in voltage can be illustrated by examining the cases of a multiple-beam klystron and a single-beam klystron (SBK) generating the same level of microwave power with the following assumptions: (i) the microperveance of a single beamlet in the multiple-beam klystron is the same as the microperveance of a whole beam in the single-beam klystron; (ii) the power and efficiency of the multiple-beam and single-beam klystrons are identical; and (iii) the electron current density in both cases is limited by space-charge effects. Under these assumptions, the voltage for the multiple-beam klystron will scale as $V_{MBK} = V_{SBK}/N^{2/5}$, where N is number of beamlets. In addition, the total length of the interaction region will scale as $V_{MBK}^{1/2}$, and so the reduced voltage also leads to a more compact device.

Traveling-Wave Tubes—The traveling-wave tube* differs from the klystron in that the radio-frequency field is not confined to a limited region but is distributed along a wave-propagating structure. A longitudinal electron beam interacts continuously with the field of a wave traveling along this wave-propagating structure. In its most common form, it is an amplifier, although there are related types of tubes that are basically oscillators.

The principle of operation may be understood by reference to Fig. 25. An electron stream is produced by an electron gun, travels along the axis of the tube, and is finally collected by a suitable electrode. Spaced closely around the beam is a circuit, in this case a helix, capable of propagating a slow wave. The circuit is proportioned so that the phase velocity of the wave is small with respect to the velocity of light. In typical low-power tubes, a value of the order of one-tenth of

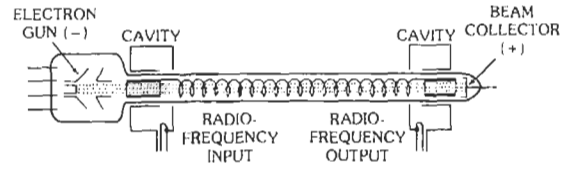


Fig. 25. Basic helix traveling-wave tube. The magnetic beam-focusing system between input and output cavities is not shown.

the velocity of light is used; for higher-power tubes the phase velocity may be two or three times higher. Suitable means are provided to couple an external radio-frequency circuit to the slow-wave structure at the input and output. The velocity of the electron stream is adjusted to be approximately the same as the axial phase velocity of the wave on the circuit.

When a wave is launched on the circuit, the longitudinal component of its field interacts with the electrons traveling in approximate synchronism with it. Some electrons will be accelerated and some will be decelerated, resulting in a progressive rearrangement in phase of the electrons with respect to the wave. The electron stream, thus modulated, in turn induces additional waves on the helix. This process of mutual interaction continues along the length of the tube with the net result that direct-current energy is given up by the electron stream to the circuit as radio-frequency energy, and the wave is thus amplified.

By virtue of the continuous interaction between a wave traveling on a broad-band circuit and an electron stream, traveling-wave tubes do not suffer the gain-bandwidth limitation of ordinary types of electron tubes. By proper circuit design, such tubes can have instantaneous bandwidths approaching three octaves.

The helix is an extremely useful form of slow-wave circuit because the impedance that it presents to the wave is relatively high and because, when properly proportioned, its phase velocity is almost independent of frequency over a wide range.

An essential feature of this type of tube is the approximate synchronism between the electron stream and the wave. For this reason, the traveling-wave tube will operate correctly over only a limited range in voltage. Practical considerations require that the operating voltages be kept as low as is consistent with obtaining the necessary beam input power; the voltage, in turn, dictates the phase velocity of the circuit. The electron velocity v in centimeters/second is determined by the accelerating voltage V in volts in accordance with

$$v = 5.93 \times 10^7 V^{1/2}$$

Fig. 26 shows a typical relationship between gain and beam voltage. The small-signal gain $G_{s,s,r}$ in decibels, of a traveling-wave tube is given approximately by

$$G_{s,s,r} = A + BCN$$

* E. A. Gelvich, L. Borisov, Y. Zhara, A. Zakurdayev, A. Pebedonostsev, and V. Poogin, "The new generation of high-power multiple-beam klystrons," *IEEE Transactions on Microwave Theory and Techniques*, vol. 41, no. 1, pp. 15–19, January, 1993.

* J. R. Pierce, *Traveling-Wave Tubes* (New York: D. Van Nostrand Co., 1950). R. Kompfner, *Reports on Progress in Physics*, Vol. 15 (London: The Physical Society, 1952; pp. 275–327); a bibliography is given in a survey paper by J. R. Pierce, "Some recent advances in microwave tubes," *Proceedings of the IRE*, vol. 42, pp. 1735–1747, December, 1954; B. N. Basu, *Electromagnetic Theory and Applications in Beam-Wave Electronics* (Singapore: World Scientific Publishing Co., 1996); A. H. W. Beck, *Space Charge Waves and Slow Electromagnetic Waves* (New York: Pergamon Press, 1958). A. S. Gilmour, Jr., *Principles of Traveling-Wave Tubes* (Boston: Artech House, 1994); R. G. E. Hutter, *Traveling-Wave Tubes*, *Advances in Electronics and Electron Physics*, Vol. 6 (New York: Academic Press, 1954); S. Sensiper, "Electromagnetic wave propagation on helical structures," *Proceedings of the IRE*, vol. 43, pp. 149–161. L. Sivan, *Microwave Tube Transmitters* (London: Chapman and Hall, 1994); D. A. Watkins, *Topics in Electromagnetic Theory* (New York: John Wiley & Sons, 1958).

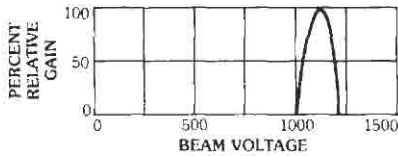


Fig. 26. Traveling-wave tube gain vs accelerating voltage.

where

- A is the initial loss due to the establishment of the modes on the helix and lies in the range from -6 to -9 decibels,
- B is a gain coefficient that accounts for the effect of circuit attenuation and space charge,
- C is a gain parameter that depends on the impedances of the circuit and the electron stream, and
- N is the number of active wavelengths in the tube.

$$C = \left\{ \left[E^2 / (\omega/v)^2 P \right] \times (I_0 / 8V_0) \right\}^{1/3}$$

and

$$N = (l/\lambda_0)(c/v)$$

where

- I_0 is the beam current,
- V_0 is the beam voltage,
- l is the axial length of the helix,
- λ_0 is the free-space wavelength,
- v is the phase velocity of the wave along the tube, and
- c is the velocity of light.

The term $E^2/(\omega/v)^2 P$ is a normalized wave impedance that may be defined in a number of ways.

In practice, the attenuation of the circuit will vary along the tube, and consequently the gain per unit length will not be constant. The total gain will be a summation of the gains of various sections of the tube.

Commonly, C is of the order of 0.02 to 0.2 in helix traveling-wave tubes; typical gains vary from 20 to 70 decibels, depending on the application. The gain in a tube designed to produce appreciable power will vary somewhat with signal level when the beam voltage is adjusted for optimum operation. Fig. 27 shows a typical characteristic.

To restrain the physical size of the electron stream as it travels along the tube, it is necessary to provide a focusing field, either magnetic or electrostatic, of a strength appropriate to overcome the space-charge forces that would otherwise cause the beam to spread. For tubes requiring a strong confining field (in the range of 0.5 to 2 tesla) and/or a high degree of spatial field uniformity, a solenoid electromagnet is typically used to generate the longitudinal magnetic field. Many traveling-wave tube designs (e.g., for electronic warfare, communications, and airborne/space-based appli-

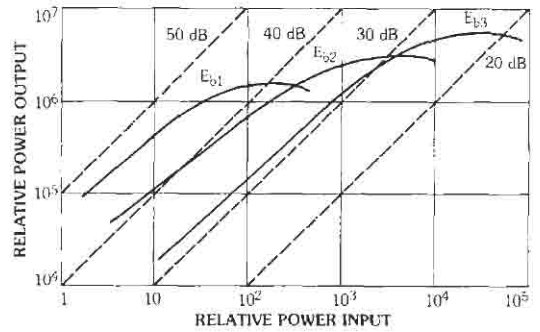


Fig. 27. Gain of traveling-wave tube as a function of input level and beam voltage. $E_{b1} < E_{b2} < E_{b3}$.

cations) require more modest bore fluxes in the range of 0.2 to 0.5 tesla, allowing the use of periodic-permanent-magnet (PPM) focusing. A PPM stack consists of rings of permanent magnets with alternating axial polarities separated by iron pole pieces and/or nonconducting shims. Typically, the permanent magnets are fabricated from a rare-earth compound of samarium-cobalt or neodymium-iron-boron. The high-energy products ($[BH]_{max} > 200$ kilojoules/cubic meter) and coercivities ($H_c > 100$ amperes/meter) of these materials enable amplifier designs with substantial weight and size reductions compared to solenoid-focused systems along with improved efficiency and reliability owing to the elimination of external electromagnet power supplies.

After interacting with the fields in the slow-wave circuit of the traveling-wave tube, the spent beam electrons are collected in an electron collector. If the collector is maintained at the same potential as the slow-wave circuit, the kinetic energy of the electrons will be dissipated in the walls of the collector as heat. However, if the collector is operated at a reduced (depressed) potential relative to the slow-wave circuit, much of the remaining energy in the beam can be recovered by decelerating the electrons prior to collection, greatly enhancing the efficiency of the device.

In practice, the beam-wave interaction process introduces a spread of velocities in the beam electrons; for efficient collection, the collector must have multiple electrodes (stages) maintained at a range of depressed potentials to efficiently sort electrons traveling at different velocities. With an optimized two- or three-stage depressed collector, the efficiency of the traveling-wave tube can be increased by as much as a factor of two or three. For space-based applications, where efficiency is critical, a traveling-wave tube may use up to a five-stage depressed collector to achieve narrow-band electronic efficiencies approaching 80%.

Helix traveling-wave circuits of the type shown in Fig. 25 are widely used in low- to medium-power applications (≤ 200 watts), where gain, bandwidth, and efficiency are important performance criteria. For

applications requiring high peak and/or high average power, the helix can be replaced by periodic slow-wave structures such as ring-bar, ring-loop, and coupled-cavity circuits. These structures have the advantage of being mechanically robust with high intrinsic thermal conductivities and are readily adaptable to liquid cooling. For example, ladder circuit coupled-cavity traveling-wave tubes have been built that operate at peak output powers ranging from 50 to 50,000 W at frequencies from 30 GHz to 95 GHz and bandwidths ranging from 14% to 20%. In general, the nonhelix circuits are capable of higher power operation at the expense of bandwidth.

Traveling-Wave-Tube Performance Data—Traveling-wave tubes are designed to emphasize particular inherent characteristics (e.g., gain, bandwidth, peak/average power, linearity, efficiency) for specific applications. Principle application areas include radar, electronic countermeasures, telecommunications, and satellite communications.

Radar: For this application, the traveling-wave tube may be an intermediate or final output stage of a transmitter; power output may be pulsed, continuous-wave, or dual-mode (user selectable low-power continuous-wave/high-power pulsed operation). Gains of 30 to 65 decibels are typical. Designs using coupled-cavity circuits have achieved peak output powers ranging from > 1000 kilowatts in S-band to one kilowatt in W-band; helix circuits operate at lower peak powers (ranging from tens to thousands of watts) but with broader bandwidth.

Electronic Countermeasures: In this application, the principle performance requirements are high bandwidth and low harmonic distortion at moderate output power levels; power output may be pulsed or continuous-wave. Bandwidths approaching three octaves have been achieved; gains range from 25 to 60 decibels. For applications in L-band through Ku-band, peak powers

are in the several kilowatt range for pulsed tubes and in the several hundreds of watts range for continuous-wave devices.

Telecommunications: Terrestrial telecommunications applications include transmitters for satellite ground stations, point-to-point microwave links, and broadband, high-speed data transmission. Power output is typically continuous-wave and falls into one of three classes: low power (≤ 200 watts), medium power (300 to 1,200 watts), or high power ($\geq 1,200$ watts). All commercial/defense/scientific bands are covered, including C-, X-, Ku-, DBS, K-, Ka-, Q-, and V-bands. Gains range from 30 to 65 decibels.

Satellite Communications: Because of their reliability, lifetime, and high efficiency, traveling-wave tubes are the amplifier of choice for many space-based applications. Typical operating lifetimes for space tubes are in excess of 15 years. Power output is typically continuous-wave and can range from tens of watts to hundreds of watts depending on the application. Gains are generally high, in the range of 45 to 65 decibels, with efficiencies from 55% to approaching 80%.

Microwave Power Modules*—The microwave power module (MPM) and millimeter-wave power module (MMPM) represent a class of highly compact, fully integrated, self-contained amplifiers that combine both solid-state and vacuum electronics technologies. As shown in Fig. 28, the power module consists of a solid-state driver amplifier (constructed from either discrete or monolithic microwave integrated circuit components) combined with a highly efficient, low-gain traveling-wave tube; power conditioning is inte-

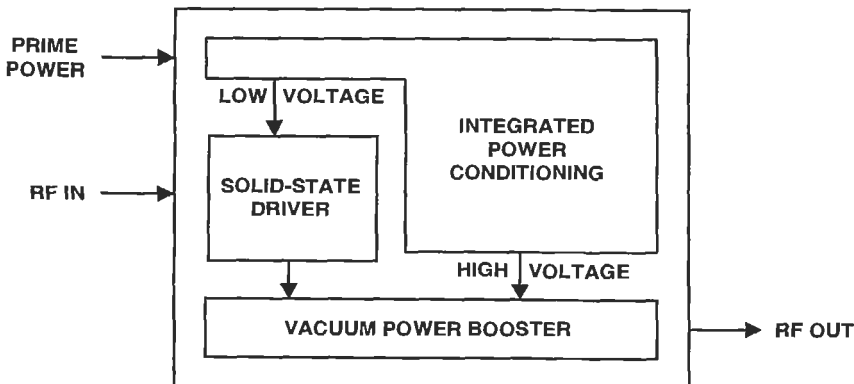


Fig. 28. Block diagram of a microwave power module.

* C. R. Smith, C. M. Armstrong, and J. Duthie, "The microwave power module: A versatile RF building block for high-power transmitters," *Proceedings of the IEEE*, vol. 87, no. 5, May, 1999.

grated with the amplifier chain. The overall gain of the power module is partitioned approximately equally between the solid-state driver and the traveling-wave tube, taking maximum advantage of the low-noise, signal conditioning capabilities of solid-state technology and the high power, high efficiency, and broad bandwidth capabilities of vacuum electronics technology. To emphasize its complementary function, the traveling-wave tube is often referred to as a "vacuum power booster" in the context of the power module.

The optimal partitioning of gain between the solid-state driver and vacuum power booster stages results in a number of performance benefits, including (i) reduced amplifier size and weight, owing to the reduction in gain, and hence length, of the vacuum power booster; (ii) increased broad-band power handling capability, again owing to the reduction in gain of the vacuum power booster; (iii) extremely high DC-to-RF transmitter efficiency, owing to the use of a multi-stage depressed collector in the vacuum power booster; (iv) improved noise performance (noise figure < 10 decibels), owing to the use of the solid-state driver; (v) the ability to incorporate gain and temperature compensation, as well as phase shifting, in the driver; (vi) the opportunity to incorporate linearization circuitry in the solid-state driver that can be matched to the characteristics of the vacuum power booster and the needs of the system; and (vii) an anticipated increase in reliability owing to the reduction of the high voltage required for the reduced-gain vacuum power booster. Depending on the application, power modules have been developed that can deliver 50 to 200 watts of continuous-wave radio frequency output power at power gains of > 50 decibels over bandwidths approaching three octaves. For optimized narrow-band designs, DC-to-RF conversion efficiencies of 50% have been achieved (using a multi-stage depressed collector); for multi-octave frequency designs, efficiencies of > 35% have been achieved.

Modeling of the Helix Traveling-Wave Tube—

The need for experimental prototyping and testing cycles in traveling-wave tube development can be reduced or eliminated through the use of modern computer-aided design tools,* reducing the time and cost of development. The design tools are based on physical models of the interaction between the electron beam and the electromagnetic fields in the traveling-wave tube structure. Implemented into design codes using modern numerical techniques and advanced algorithms, these sophisticated models can simulate actual hardware geometries and incorporate realistic material properties. Advances in computationally efficient codes and processor speeds enable intricate three-dimensional structures with advanced, inhomogeneous,

anisotropic materials to be simulated on modern desktop computers, where complex parameter optimization trade-off analyses can be made to develop tube designs that meet specific goals for frequency, bandwidth, power, efficiency,† gain, noise, and linearity.

Large-signal codes of different levels of sophistication are employed to model helix traveling-wave tube performance.‡ Electron gun design codes accurately model electron beam formation and transport to the interaction region, whereas the collector codes are used in the design of efficient multi-stage depressed collectors (examples of 2-D electron gun codes include EGUN§ and UGUN"; 3-D gun/collector codes include MICHELLE#). Sophisticated finite-element codes are employed to design magnetic circuits (e.g., MAXWELL***) and are used to address thermo-mechanical design issues (e.g., ANSYS††). The effective exercise of computational design tools can also provide insight into the dependence of device performance on the variation of critical parameters, such as dimensional variations or variations in materials properties. This information can significantly impact device design, as well as fabrication and assembly processes.

Backward-Wave Oscillators—A member of the traveling-wave-tube family, the *O*-type backward-wave oscillator‡‡ makes use of the interaction of the electron stream with a radio-frequency-circuit wave whose phase and group velocities are 180° apart. The group velocity, and thus the direction of energy flow, is

† D. R. Whaley, C. M. Armstrong, B. Gannon, G. Groshart, E. Hurt, J. Hutchins, M. Roscoe, T. M. Antonsen, Jr., and B. Levush "Sixty-percent-efficient miniature C-band vacuum power booster for the microwave power module," *IEEE Transactions on Plasma Science*, vol. 26, no. 3, pp. 913–921, 1998.

‡ T. M. Antonsen, Jr. and B. Levush "CHRISTINE: A multi-frequency parametric simulation code for traveling wave tube amplifiers", NTIS no. ADA327112, U.S. Naval Research Laboratory report NRL/FR/6840-97-9845; V. Srivastava, R. G. Carter, B. Ravinder, A. K. Sinha, and S. N. Joshi "Design of helix slow-wave structures for high efficiency TWTs," *IEEE Transactions on Electron Devices*, vol. 47, no. 12, pp. 2438–2443, 2000; H. P. Freund, "Three-dimensional nonlinear theory of helix traveling wave tubes," *IEEE Transactions on Plasma Science*, vol. 28, pp. 748–759, 2000; D. Chernin, T. M. Antonsen, Jr., B. Levush, and D. R. Whaley "A three-dimensional multifrequency large signal model for helix traveling wave tubes," *IEEE Transactions on Electron Devices*, vol. 48, no. 1, pp. 3–11, 2001.

§ W. B. Herrmannsfeldt, "Electron trajectory program." Stanford University, Stanford Linear Accelerator Center Report No. 166, September, 1973.

" N. Dionne, "UGUN gun/collector design tool."

J. Petillo et al. "The new 3-D electron gun and collector modeling tool: MICHELLE," *International Vacuum Electronics Conference 2000*, May 2–4, 2000. Monterey, CA.

* T. M. Antonsen, Jr., A. A. Mondelli, B. Levush, J. P. Verboncoeur, and C. K. Birdsall "Advances in modeling and simulation of vacuum electronic devices," *Proceeding of the IEEE*, vol. 87, no. 5, pp. 804–839, 1999.

directly opposed to the direction of electron motion. Fig. 29 shows schematically a backward-wave tube with connection to both ends of the slow-wave structure, so that operation as either oscillator or amplifier could be achieved. An electron beam is produced by the electron gun, traverses the slow-wave structure, and is dissipated in the collector structure. During its transit, the beam is confined by a longitudinal magnetic field. With a beam current of sufficient magnitude, the beam-structure interaction will produce oscillations, and microwave power will be delivered from the end of the structure adjacent to the electron gun. At beam-current levels below the "start-oscillation" value, a radio-frequency signal may be introduced at the collector end of the device, and the tube will operate as an amplifier.

To improve interaction efficiency, electron beams with hollow cross sections are usually used. This places all the electrons as close as possible to the slow-wave structure in the region of maximum radio-frequency field. This is beneficial because the strength of the -1 space harmonic field goes to zero on the axis. To produce this hollow-cross-section beam, it is necessary to use magnetically confined electron flow from the cathode, and thus the electron gun is entirely immersed in the magnetic field.

An O -type backward-wave device is voltage tunable, with the frequency being proportional to the $1/2$ power of the cathode-helix voltage as well as dependent on the dimensions of the structure. Typically, tuning over a full octave range is possible, and in special cases a range of two or more octaves can be achieved.

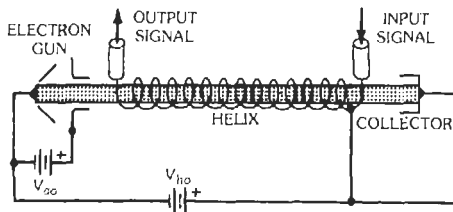


Fig. 29. A traveling-wave tube in operation as a backward-wave amplifier. A separate power supply connected to the anode permits beam-current control independent of the helix voltage. (From J. W. Gewartowski and H. A. Watson, *Principles of Electron Tubes*, 1965; p. 398. Courtesy D. Van Nostrand Co., Inc.)

** Ansoft Corporation, Maxwell 3D structure electromagnetic simulator.

†† ANSYS Inc.

‡‡ H. Heffner, "Analysis of the backward-wave traveling-wave tube," *Proceedings of the IRE*, vol. 42, pp. 930-937, June, 1954; A. H. W. Beck, *Space-Charge Wave and Slow Electromagnetic Waves* (New York: Pergamon Press, 1958; pp. 241-255); R. Kompfner and N. T. Williams, "Backward-wave tubes," *Proceedings of the IRE*, vol. 41, pp. 1602-1611, November, 1953.

However, where confined limits are desired on power variation or other special characteristics, more restricted frequency ranges may be necessary. Where full octave tuning is used, power output variation of 6 to 10 decibels across the range is usual. In most cases, a separate control element in the gun permits adjustment of beam-current amplitude and thus provides control of power output. Oscillators of this type have very low pulling figures, but the pushing figure is often substantial. Frequency stability is generally excellent, with the achievable value normally depending on power-supply capabilities rather than on inherent tube limitations.

Generally, O -type backward-wave oscillators are low-power devices, with 10 to 50 milliwatts being typical. However, in the range from 1 to 4 gigahertz, up to several hundred milliwatts is feasible, while in the range from 50 to 100 gigahertz, 5 to 10 milliwatts is relatively difficult to achieve reliably. Typical performance data for low-power helix-type permanent-magnet-focused backward-wave oscillators are listed in Table 6.

Crossed-Field Tubes

Microwave crossed-field tubes are often called M -type devices to distinguish them from linear-beam tubes that are often called O -type devices. The earliest type of crossed-field tube* was the magnetron oscillator. Crossed-field amplifiers (CIA) and voltage-tunable crossed-field oscillators have been developed more recently. Crossed-field tubes generally operate with higher electronic conversion efficiency for dc power to microwave power than linear-beam tubes, so they are especially attractive for high-power applications.

Magnetrons—A magnetron is a high-vacuum tube containing a cathode and an anode, the latter usually divided into two or more segments. A constant magnetic field modifies the space-charge distribution and the current-voltage relations. In modern usage, the

TABLE 6. PERFORMANCE OF TYPICAL LOW-POWER BACKWARD-WAVE OSCILLATORS

Frequency Range (GHz)	Tuning Voltage (V)	Cathode Current (mA)	Minimum Power Output (MW)
1.0-2.0	250-1150	15	100
2.0-4.0	300-1800	10	100
4.0-8.0	250-2400	12	25
5.3-11.0	245-2400	10	25
8.0-12.4	550-2400	10	25

* E. Okress, ed., *Crossed-Field Microwave Devices* (New York: Academic Press, 1961).

term "magnetron" refers to the magnetron oscillator in which the interaction of the electronic space charge with the resonant system converts direct-current power into alternating-current power, usually at microwave frequencies.

Many forms of magnetrons have been made in the past, and several kinds of operation have been employed. The type of tube that is now almost universally employed is the multicavity magnetron generating traveling-wave oscillations. It possesses the advantages of good efficiency at high frequencies, capability of high outputs either in pulse or continuous-wave operation, moderate magnetic-field requirements, and good stability of operation. A section through the basic anode structure of a typical fixed-frequency, conventional magnetron is shown in Fig. 30.

In magnetrons, the operating frequency is determined by the resonant frequency of the separate cavities arranged around the central cylindrical cathode and parallel to it. A high direct-current potential is placed between the cathode and the cavities, and radio-frequency output in the conventional type of magnetron is brought out through a suitable transmission line or waveguide, usually coupled to one of the resonator cavities. Under the action of the radio-frequency voltages across these resonators and the axial magnetic field, the electrons from the cathode form a bunched space-charge cloud that rotates around the tube axis, exciting the cavities and maintaining their radio-frequency voltages.

Pulsed magnetrons have been developed to operate with peak power output from a few kilowatts to several megawatts and at frequencies from a few hundred megahertz to tens of gigahertz. The primary application for these tubes is the RF power generator for radar transmitters. Continuous-wave magnetrons have also been developed with power output of several hundred watts at frequencies from a few hundred megahertz to about 10–12 GHz. The principal use of many of these CW tubes is in countermeasure equipment used for electronic jamming. In addition, many fixed-frequency CW magnetrons at 2450 MHz have been used in microwave ovens.

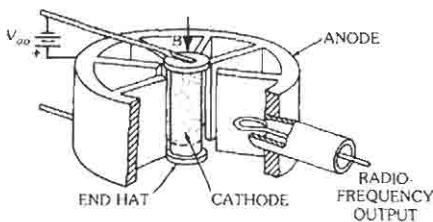


Fig. 30. Magnetron oscillator. (From J. W. Gewartowski and H. A. Watson, *Principles of Electron Tubes*, 1963; p. 428. Courtesy D. Van Nostrand Co., Inc.)

The output frequency of conventional magnetrons can be altered in an undesired fashion by variations in the microwave load impedance (frequency pulling) or by changes in cathode current (frequency pushing) caused by fluctuations in voltage applied to the magnetron. Significant improvement in operational stability for the magnetron is obtained by coupling to a high- Q stabilizing cavity. A convenient method for doing this is shown in Fig. 31. The stabilizing cavity is made an integral part of the magnetron by surrounding the magnetron anode with a resonant coaxial cavity (hence, the name "coaxial magnetron"). Alternate resonators of the anode circuit are slotted to provide coupling to excite the TE_{011} circular electric mode in the coaxial cavity. Power output from a coaxial magnetron is obtained through a slot in the outer wall of the cavity.

The circulating space charge in both conventional and coaxial magnetrons contains wide-band noise-frequency components that can couple to the output. In conventional magnetrons, this spurious noise can couple directly to the output waveguide. Spurious noise power measured in a 1-MHz bandwidth is typically greater than 30 to 40 dB below the carrier. The coaxial cavity in the coaxial magnetron provides some isolation between the spurious noise coupled to the vanes and the output waveguide. The spurious noise power from coaxial magnetrons is typically 10 to 20 dB lower than the noise power from conventional magnetrons of comparable peak power level.

Magnetron frequencies can be tuned by various means. One method used for conventional magnetrons employs metallic probes that are inserted uniformly into the multiple resonant cavities of the anode. This changes the resonant frequency of the cavities that control the magnetron frequency. Coaxial magnetrons are usually tuned by changing the resonant frequency of the coaxial stabilizing cavity by moving an end plate in the cavity. Other methods are also used in each magnetron type.

Many conventional and coaxial magnetrons are designed for very rapid tuning. Fast frequency operation is employed in radar transmitters to improve radar signal detection and/or as an electronic countermeasure to avoid jamming. Fast tuning is performed by several methods. Some tubes have tuners that are driven mechanically by an electromagnet in a manner similar to the voice coil of a radio speaker. Tuning rates, tuning range, and oscillator frequency-versus-time profiles are determined by the waveform of the voltage applied to the electromagnet. Other tubes have a mechanical linkage between the tuner and a rotating motor and produce approximately sinusoidal variation in frequency-time profile. Tuning rates and range are limited by acceleration forces in the reciprocating tuner. For example, X-band, 200-kW tubes with narrow-band, rapid tuning of 30–50 MHz are tuned through this range at 200 Hz. Wider-band tuner excursions of 250–500 MHz are cycled at 25–40 Hz. Other tuning mechanisms in conventional and coaxial magnetrons utilize continuously rotating tuners. More

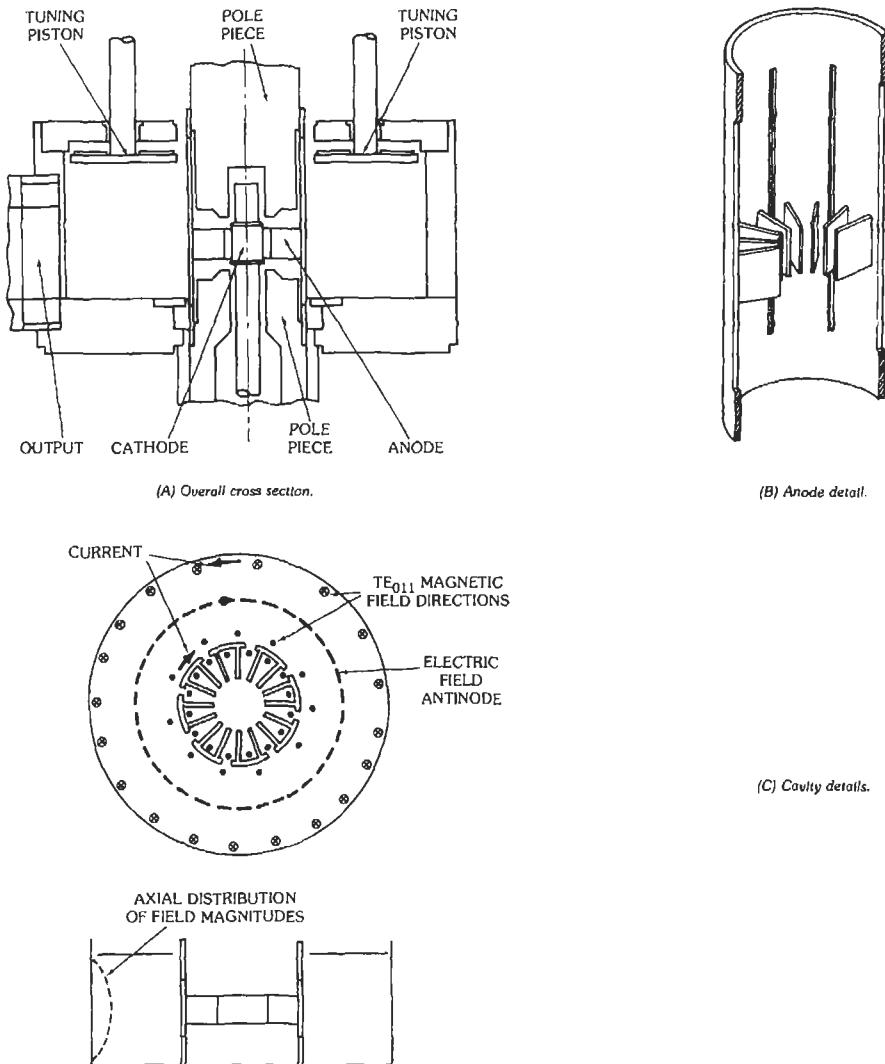


Fig. 31. Geometry of the composite anode. (From E. Okress, *Crossed-Field Microwave Devices*, Vol. 2, New York: Academic Press, Inc., 1961; p. 125.)

rapid frequency excursions through a fixed tuning range can be obtained with rotating tuners than from reciprocating tuners, but the rotational inertia of the tuners prohibits rapid changes in tuning profiles.

Voltage-Tuned Crossed-Field Oscillators—There are two types of voltage-tunable, crossed-field oscillators. These are the carcinotron and the voltage-tuned magnetrons (VTM).

The carcinotron is an *M*-type backward-wave oscillator in which the electron stream traverses the tube and interacts with the fields on the slow-wave structure under conditions where the electric and magnetic fields are perpendicular to each other. Fig. 32 shows schematically a linear version of the carcinotron. In the electron gun, current is drawn from the cathode when

the accelerator voltage is applied. Because of the presence of the magnetic field, directed as shown, the electron paths are curved approximately 90° so that they enter the interaction region between the slow-wave structure and the sole. If the voltages and the magnetic field strength are proper, the electrons will travel along a path approximately parallel to the structure until they reach the collector.

Although Fig. 32 shows a linear arrangement, carcinotrons are conventionally designed in a circular arrangement to conserve magnet size and weight. In this arrangement, the sole approximates the appearance of the cathode of a magnetron, and the slow-wave structure is in the position of the magnetron anode, but neither the sole nor the structure is reentrant.

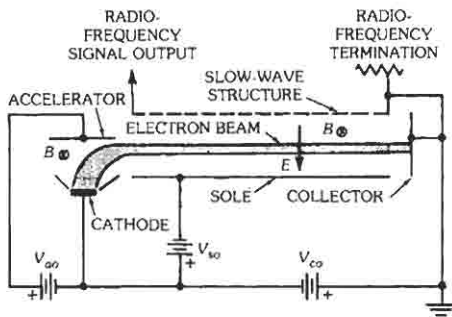


Fig. 32. Linear version of an *M*-type carcinotron oscillator. (From J. W. Gewartowski and H. A. Watson, *Principles of Electron Tubes*, 1965; p. 459. Courtesy D. Van Nostrand Co., Inc.)

The carcinotron gives performance similar to that of the *O*-type backward-wave oscillator, but it offers several of the advantages of crossed-field devices. High-efficiency operation is possible, with values of 20% to 30% being readily obtained. This efficiency capability makes the carcinotron useful as a high-power device with continuous-wave capabilities of hundreds of watts through the X-band. Its construction permits direct scaling to very high frequencies, with several milliwatts of power having been achieved at frequencies beyond 300 gigahertz.

The carcinotron, like the *O*-type backward-wave oscillator, is voltage tunable with the oscillation frequency being approximately directly proportional to the voltage between the cathode and the slow-wave structure. This linear relationship simplifies the associated electronic tuning circuit considerably. Frequency pushing is considerably lower than in *O*-type backward-wave oscillators. The *M*-type carcinotron has the disadvantage, however, that it is relatively noisy, with spurious power output often not more than 10 to 15 decibels below the main signal output.

Voltage-tuned magnetrons (VTMs) are distinguished from conventional magnetrons in two critical ways:

1. The output is coupled very heavily to the circuit, and so the loaded Q is very low, typically between 1 and 10.
2. The normal cathode is replaced by a nonemitting post, and a cylindrical beam is injected into the interaction space from an emitting filament at one end; the injected current is controlled by the voltage on a control electrode surrounding the filament.

Under these conditions, the oscillation frequency is no longer controlled by the resonator, but by the ratio of electric to magnetic fields in the interaction space. The magnetic field is held constant by a permanent magnet so that the frequency becomes directly proportional to the applied voltage, which is normally between 1 and 4 kV dc. The control-electrode voltage is normally about 20% of the cathode-anode voltage.

Power output of VTMs ranges from a few watts to a few hundred watts at frequencies from about 1 to 5 GHz. Tuning range is more than one octave at low powers (under 10 watts) but reduces to about 30% at 100 watts.

Efficiency is remarkably high, exceeding 60% in the higher-power VTMs. The tubes are more compact and more efficient than carcinotrons, but they do not reach as high an absolute power level or as high a frequency.

Because the heavy coupling would make the VTM very susceptible to load changes, a ferrite isolator is normally made an integral part of the VTM package.

The application of high-power, voltage-tuned, crossed-field oscillators is exclusively in the field of electronic countermeasures. Low-power VTMs are also used as test oscillators and are capable of extremely high modulation rates.

Crossed-Field Amplifiers—Crossed-field amplifiers (CFAs) can be divided into two general classes, injected-beam tubes and distributed-emission (or sometimes emitting-sole) tubes. The former are similar in appearance to the carcinotron in that they employ separate electron guns, interaction regions, and beam-collector elements. For amplifiers, the slow-wave circuit is provided with both input and output terminals. The RF input is close to the electron gun, and the RF wave is amplified as it travels along the slow-wave circuit to the output while maintaining synchronism with the electron stream. Most injected-beam CFAs are fabricated in a linear format and are often compared to beam-type traveling-wave tubes. The principal application of injected-beam CFAs is in electronic-countermeasures equipment.

Distributed-emission CFAs are similar to magnetrons in that electron current for interaction is obtained from the cathode (or sole electrode) throughout the interaction space. Most distributed-emission tubes are fabricated in a circular format like magnetrons. As may be seen from Fig. 33, the major difference is that the slow-wave structure is not reentrant, whereas in the magnetron both the beam and the circuit are reentrant. In Fig. 33, voltage and magnetic field are applied as for the magnetron. A radio-frequency signal is applied to the structure and progresses in a clockwise direction toward the output terminal. Current spokes, produced in the cathode-circuit region by the radio-frequency electric fields, also progress in a clockwise direction synchronously with the circuit wave. The interaction between the beam and circuit wave results in a growing circuit wave and, hence, gain. If desired, interaction with a backward mode may also be accomplished with this device.

Because the beam is reentrant, the crossed-field amplifier will oscillate if the circuit gain becomes high. Gain is usually limited to 10 to 15 decibels. If only a portion of the circumference is used for the slow-wave structure and a drift area is left between the two ends of the structure, the feedback mechanism is

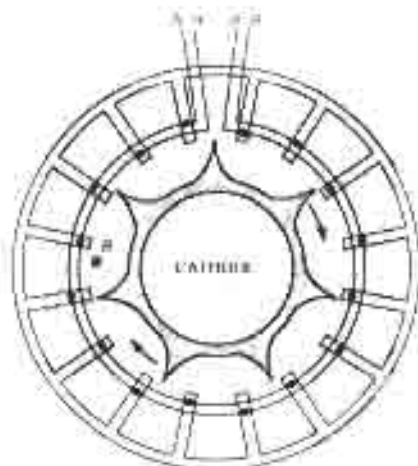


Fig. 33 Schematic drawing of a crossed-field amplifier. (From J. W. Geary and H. A. Wilson, *Principles of Electron Tubes*, 1963, p. 449. Courtesy: D. Van Nostrand Co., Inc.)

disrupted, and gains of 15 to 20 decibels may be realized.

The power output of the crossed-field amplifier is essentially independent of the radio-frequency drive signal, and the device thus operates as a saturated amplifier. This characteristic makes it unsuitable for amplifying amplitude-modulated signals.

Distributed-emission crossed-field amplifiers offer the advantage of relatively high efficiency (40% to 60% or even higher), and they may be designed to provide very high peak output powers. Their disadvantages are their low gain, limited bandwidth, high noise, and saturated-amplifier characteristics.

CYCLOTRON RESONANCE MICROWAVE TUBES* †

One of the highest power device technologies for the generation and amplification of power in the millimeter wave bands (30–300 GHz) are the cyclotron resonance devices.

In cyclotron resonance devices (gyromons, gyroklystrons, gyro traveling-wave tubes, and gyro backward-wave oscillators), the electrons of the electron beam have substantial motion perpendicular to the axis of the beam and the solenoidal focusing magnetic field. They interact with radio-frequency electric fields

perpendicular to the magnetic focusing field. The RF fields correspond to the fields of TE modes in cavities that, to a first approximation, are cylindrical, and the operating modes need not be the fundamental mode of the cavity. Operation of higher-order modes permits use of larger diameter cavities and the attainment of higher output power. As the electrons rotate and the fields alternate in synchronism, there is a cumulative interaction. Some electrons gain energy, and other electrons lose energy. The electrons that gain energy undergo a relativistic mass increase, and the ones that lose energy undergo a relativistic-mass decrease. Because the cyclotron frequency is equal to $\Omega_c = eB/m\gamma$, in which m is the rest mass of the electrons and γ is the relativistic factor ($\gamma = [1 - v^2/c^2]^{-1/2}$, where v is the electron velocity), electrons that gain energy have larger gyroradii and thus smaller angular rotation rates and electrons that lose energy have larger angular rotation rates. Hence the electrons gather into rod-like bunches parallel to the axis of the helical electron trajectories. As a result of this bunching, there is a net exchange of energy between the electron beam and the resonant electromagnetic field. This resonance condition is given by

$$\omega = n\Omega_c$$

where

$$\omega = 2\pi\nu$$

ν is the frequency, and

n is the harmonic number ($n = 1, 2, 3$).

In most instances gyro devices operate on the fundamental resonance ($n = 1$), but many successful devices have also been made at higher harmonics ($n > 1$). For a solenoidal magnetic field of 1 T and $\gamma = 1$, a gyro device operating on the fundamental cyclotron resonance ($n = 1$) will have a frequency of $\nu = \Omega_c/2\pi = 28$ GHz. The magnetic field of the solenoid serves two purposes for the cyclotron resonance maser: 1, to propagate the electron beam and 2, more importantly, to permit this resonance with a particular cavity mode. Because the resonant nature of the interaction permits the designer to select an interaction with a higher-order cavity mode without necessarily having interference from lower-order modes, cyclotron resonance devices will produce higher powers at higher frequencies than will other microwave tubes. Although the bunching in this device is based on a relativistic effect, the electron energies need not be that relativistic; electron beam kinetic energies typical in gyro devices are 20–80 keV. Some cyclotron resonance devices designed for very high peak powers do employ much higher energy beams.

Cyclotron resonance devices that have been built so far fall into four major categories: (A) gyromon oscillators, † which are single-cavity devices that oscillate because of the negative interaction impedance of the electron beam with the cavity; (B) gyroklystrons, in which azimuthal bunching of the electrons is achieved

* R. S. Symons and H. R. Job, "Cyclotron resonance devices," in *Advances in Electronics*, vol. 58, pp. 1–75 (New York: Academic Press, Inc., 1980).

† K. L. Felch, H. O. Daily, H. R. Job, K. E. Kreschler, W. Lawson, H. Levush, and R. J. Temlett, "Characteristics and applications of fast wave gyrodevices," *Proceedings of the IEEE*, vol. 87, no. 5, pp. 752–781, 1989.

by first passing them through an input cavity wherein the microwave signal to be amplified is injected, then through successive beam-excited bunching cavities, and then through a final output cavity in which the energy extraction is achieved; (C) gyro-traveling-wave tubes (gyro-TWTs), in which the electrons travel through a uniform or tapered waveguide supporting a transverse electric mode; and (D) gyro backward-wave oscillators, which are similar in structure to gyro-TWTs but in which the group velocity of the wave in the uniform waveguide is opposed to electron flow and power comes out of the circuit at the end where the electron beam enters.

Gyrotron oscillators are by far the best developed of these four categories. The principle applications are scientific and include applications to the cyclotron resonance heating of fusion plasmas and material processing. Representative of present state of the art for gyrotron oscillators is a 110-GHz gyrotron capable of operation at 1 MW for 0.6-s pulses, 400 kW for 6.5-s pulses, and 100 kW CW.* To achieve these unprecedented powers without thermal damage to the interaction circuit, this device was operated in a very high order cavity mode, the $TE_{22,6,1}$ mode. Consequently, gyrotrons often now have internal mode converters that convert the TE mode used in the interaction to a more usable output mode, such as a TEM_{00} Gaussian beam. Gyrotrons have been developed or are under development at frequencies ranging from 8 to 170 GHz. Gyrotrons at even higher frequencies (up to 800 GHz) have been developed for scientific applications.

Though they have received less development than gyrotron oscillators, gyro-amplifiers such as the gyro-klystron have been developed in several frequency bands with radar and accelerator applications as the primary motivation. As an example, a 90-kW peak power, 10-kW average power gyro-klystron amplifier has been developed at 94 GHz for radar applications. The device has an instantaneous bandwidth of 600 MHz, an efficiency of 33%, and operates with a 65-kV, 6-A electron beam.† A second device delivered the same average power at 700 MHz bandwidth, and it offers a 1-GHz bandwidth at 4-kW average power. This represents more than an order of magnitude higher average power than is available from any other

amplifier technology at this frequency. Gyro-klystron amplifiers behave as coherent amplifiers of the input signal, and they have been demonstrated to have low phase noise comparable to that obtained from klystrons, with -147 dBc/Hz SSB phase noise at 10 MHz offset from the carrier in one experimental device when operated with a carrier at 200-kW output power at 35 GHz.‡

Experimental gyro-klystrons at high peak power levels in the microwave bands have also been developed for powering RF linear accelerators. Experimental devices designed for this application have achieved 27-MW peak power at 9.87 GHz at the fundamental ($n = 1$) cyclotron resonance device, and 32-MW peak power has been obtained at 19.76 GHz in a second harmonic ($n = 2$) device.§

Future efforts in the development of the cyclotron resonance or gyrotron-class of devices are directed toward higher power oscillators for fusion applications, more compact and efficient oscillators for industrial processing applications, and gyro-amplifiers such as the gyro-klystron and gyro-TWT for high-power millimeter-wave radar and accelerator applications.

FREE-ELECTRON LASERS

Free-electron lasers (which are also called ubitrons when applied to microwaves millimeter waves) consist of an electron beam propagating through a periodic transverse magnetic field (called a *wiggler* or *undulator*) that induces motion perpendicular to the axis of symmetry." A schematic illustration of a free-electron laser is shown in Fig. 34.

The wiggler and the radiation combine to produce a beat wave, which is essentially an interference pattern, that exerts a slowly varying ponderomotive force on the electrons. Because the electrons experience a near-constant field the interaction can be extremely strong. Stimulated emission occurs when electrons form coherent bunches over a wavelength. Bunching occurs because the ponderomotive wave has the same frequency (ω) as the radiation, but the wavenumber is the sum of the wavenumbers of the electromagnetic (k) and wiggler (k_w) fields. Hence, the beat wave travels more slowly than the light wave and can be in synchronism with the electrons when the phase velocity of the ponderomotive wave equals the electron beam velocity (v_b), i.e., $\omega/(k + k_w) = v_b$. The amplitude of the transverse velocity (v_{\perp}) is proportional to the product of the wiggler amplitude and period and inversely proportional to the relativistic dilation factor. Specifically,

‡ *Gyrotron Oscillators: Their Principles and Practice*, ed. C. J. Edcombe (London: Taylor & Francis Ltd, 1993).

* K. L. Felch, B. G. Danly, H. R. Jory, K. E. Kreischer, W. Lawson, B. Levush, and R. J. Temkin, "Characteristics and applications of fast-wave gyrodevices," *Proceedings of the IEEE*, vol. 87, no. 5, pp. 752-781, 1999.

† B. G. Danly, M. Blank, J. P. Calame, B. Levush, K. T. Nguyen, D. E. Pershing, R. K. Parker, K. L. Felch, B. G. James, P. Borchard, P. Cahalan, T. S. Chu, H. R. Jory, T. A. Hargreaves, R. B. Truc, W. G. Lawson, and T. M. Antonsen, Jr., "Development and testing of a high-average power 94-GHz gyro-klystron", *IEEE Trans. Plasma Science*, vol. 28, no. 3, pp. 713-726, 2000.

‡ J. P. Calame, B. G. Danly, M. Garven, and B. Levush, "Studies of electronic noise in gyro-klystrons," *Phys. Plasmas*, vol. 7, no. 5, pp. 2180-2185, 2000.

§ V. L. Granatstein and W. Lawson, "Gyro-amplifiers as candidate RF drivers for TeV linear colliders," *IEEE Trans. Plasma Science*, vol. 24, pp. 648-665, 1996.

" Reference 44.

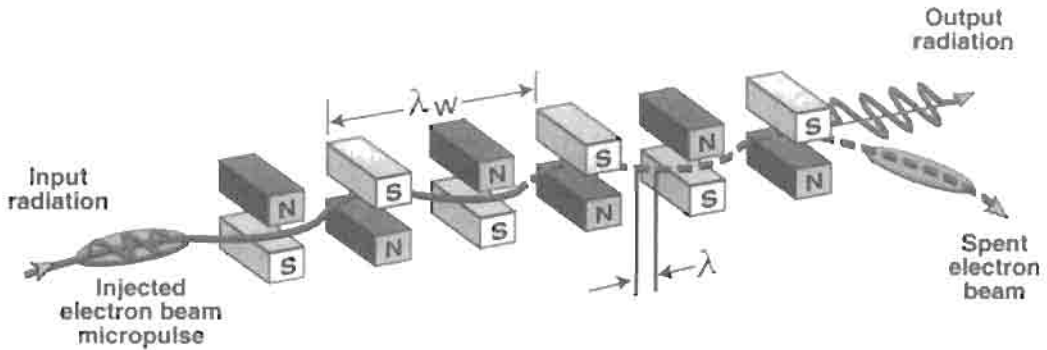


Fig. 34. A schematic of a typical free-electron laser.

$v_w = 0.934B_w \lambda_w / \gamma_w$, where B_w denotes the wiggler amplitude in kG, λ_w is the wiggler period in cm, and $\gamma_w = 1 + E_w/m_e c^2$ (where E_w is the electron beam energy and $m_e c^2$ is the electron rest mass energy). As a result, for frequencies well above the waveguide cutoffs, the resonant wavelength is given by $\lambda = (1 + 0.872B_w^2 \lambda_w^2) \lambda_w / 2\gamma_w^2$.

Because of the wavelength scaling, the free-electron laser is capable of producing very high frequency radiation when either high-energy beams or short-period wigglers are used. In practical applications, however, high gain requires high wiggler field amplitudes, which, in turn, precludes very short wiggler periods;

hence, high-frequency operation can entail relatively high beam energies. This often precludes the use of standard power supplies and modulators, requiring instead relatively large and complex electron accelerators. As a result, free-electron lasers are principally useful when infrared or shorter wavelengths are needed.

At the present time, 10 free-electron laser user facilities are in operation (see Table 7) for a wide range of applications in materials and biomedical research. Most existing free-electron lasers operate at infrared wavelengths where the high-intensity picosecond pulses they produce are used to study principle excita-

TABLE 7. SUMMARY OF CURRENTLY AVAILABLE FEL USER FACILITIES

Country	Institution	Device	λ (μm)	τp (ps)	$E_b I_b$ (MeV/A)	P_{peak} (MW)	P_{ave} (W)
USA	SU	FIREFLY	19-65	1-5	20/14	.3	.4
		SCA/FEL	3-10	0.7	37/10	10	1.2
	VU	FEL1	2-10	2	43/50	10	10
		DU	Mark III	3-10	3	44/20	2
		OK-4	0.3	10	1000/350	1000	0.1
	JL	IR Demo	3-6.6	1-2	42/50	10	1000
UCSB	PU	CIRFEL	10-43	4	14/150	4	0.05
Japan	FELI	FELI 1	5-22	1.7	33/42	5	2
		FELI 2	1-6	1.7	75/50	5	0.5
		FELI 3	0.23-1.2	1.7	165/60	5	0.5
		FELI 4	20-80	1.7	30/40	5	1
		FELI 5	40-100	3-5	20/40	(2)	(1)
France	LURE	SuperACO	0.3	20	800/10	12	0.8
		CLIO	1.8-17.5	1.5-6	70/80	10	9
NL	FOM	FELIX-2	5-35	0.5-10	45/70	2	1
		FELIX-1	25-110	1-10	25/70	2	0.5
China	Beijing Univ.	BFEL	10-16	4	30/14	20	2

tions in condensed matter systems where it is possible to access the principal excitations such as plasmons, phonons, magnons, and inter-sub-band transitions. Direct linear probing of defect modes and buried interfaces with bond specificity is possible. Free-electron lasers have been used for pump probe observation of coherent transient grating effects of narrow-gap semiconductors and third-order nonlinearity coefficients and electron relaxation times in GaAs/AlGaAs quantum wells. Studies in kinetics include vibrational energy transfers in molecules and resonant excitations of molecular vibrations. Applications to neurosurgery take advantage of particular absorption bands to produce an exceptionally fine cut in bone and tissue. At shorter wavelengths, free-electron lasers have been used to make measurements of time-resolved fluorescence decays of biological molecules, to study time-resolved photoemission, and photoionization.

There are two principal areas of future development for free-electron laser (FEL) technology: high average powers and ultra-short wavelengths. The high average power goal is several tens of kilowatts at infrared or shorter wavelengths, and the most likely configuration is an oscillator driven by a radio-frequency linear accelerator (rf linac). Rates of progress in the development of vacuum generators of coherent radiation can be gauged by using the power density (the product of the average power and the square of the frequency) as a figure of merit. The evolution of power densities in magnetrons, klystrons, gyrotrons, and FELs is shown in Fig. 35.

Significantly, although all the other vacuum tube types listed have reached their technological limits, the ultimate limit for free-electron lasers has not yet been reached. The record FEL average power of 1.7 kilo-

watts at a 3-micron wavelength was achieved at the Thomas Jefferson Laboratory,* and an upgrade is currently under construction to reach 10 kW. Short-wavelength free-electron lasers are primarily directed toward research applications of X-rays at wavelengths down to 1 Å. Currently, the principal X-ray sources are synchrotrons, but FELs hold promise for substantial increases in the average brilliance. Since neither coherent sources nor mirrors exist at X-ray wavelengths, these free-electron lasers cannot operate either as coherent amplifiers or oscillators. Instead, these X-ray free-electron lasers are intended to operate in a mode referred to as self-amplified spontaneous emission (SASE) in which shot noise on the beam is amplified to high power levels during a single pass through a long undulator. Current designs† call for beam energies of the order of 14 GeV, peak beam currents of about 3400 A, and wigglers of the order of 100 m in length. Obviously, therefore, construction of these sources are clearly major financial undertakings.

GAS TUBES

A gas tube‡ is an electron tube in which the pressure of the contained gas is such as to affect substantially the electrical characteristics of the tube. Such effects are caused by collisions between moving electrons and

* Reference 45.

† References 46 and 47.

‡ B. E. Cherrington, *Gaseous Electronics and Gas Lasers* (Elmsford, NY: Pergamon Press, 1979); B. N. Chapman, *Glow Discharge Processes* (New York: John Wiley & Sons, Inc., 1980).

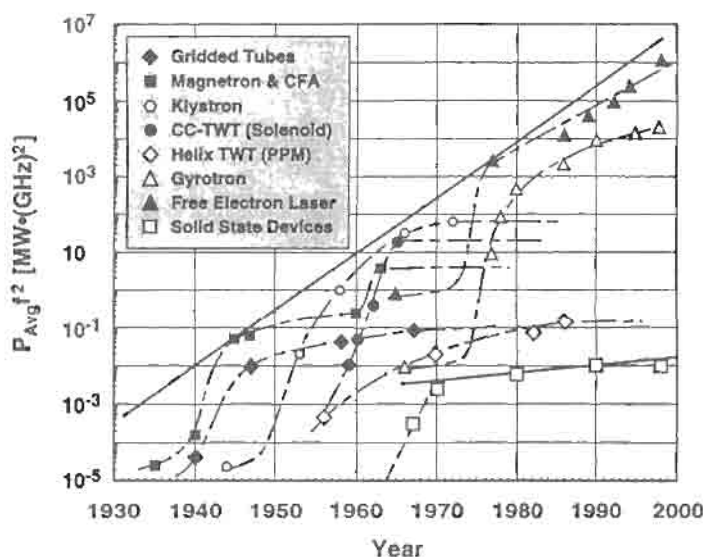


Fig. 35. Progression of device power density, $P_{avg} f^2$, for major device types.

gas atoms. These collisions, if of sufficient energy, may dislodge an electron from the atom, thereby leaving the atom as a positive ion. The electron space charge is effectively neutralized by these positive ions, and comparatively high free-electron densities are easily created.

Table 8 gives the energy in electron volts necessary to produce ionization. The column headed P_c contains the kinetic-theory collision probability per centimeter of path length for an electron in a gas at 15°C at a pressure of 1 mm of mercury. The collision frequency is given by

$$v_c = P_c p_o v$$

where

- v_c = collisions per second
- P_c = collision probability in collisions per centimeter per torr of pressure
- p_o = reduced gas pressure in torr
- v = electron velocity

Characteristics of Gas Tubes

Gas tubes may be generally divided into two classes, depending on whether the cathode is hot or cold and thus on the mechanism by which electrons are supplied.

Hot-Cathode Gas Tubes—The electrons in the hot-cathode gas tube are produced thermionically. The voltage drop across such tubes is that required to produce ionization of the gas and is generally a few tens of volts. The current conducted by the tube depends primarily on the emission capability of the cathode.

Fig. 37 shows the effect of the ionized gas on the voltage distribution in a hot-cathode tube.

Cold-Cathode Gas Tubes—The electrons in a cold-cathode tube are produced by bombardment of the cathode by ions and/or by the action of a localized high electric field. The voltage drop across such a tube is higher than in the hot-cathode tube because of this mechanism of electron generation, and the current is limited. Fig. 38 shows the effect of tube geometry and gas pressure on the voltage required to initiate the discharge.

Fig. 39 shows a typical volt-ampere characteristic of a cold-cathode discharge. Cold-cathode gas tubes may be divided into two categories, depending on the region of this characteristic in which they operate. *Glow discharge tubes* require a drop of several hundred volts across the tube and operate in region II. The current is of the order of tens of milliamperes. *Arc discharge tubes* operate in region III. They are not, strictly speaking, cold-cathode tubes since the current is drawn from a localized spot on the cathode, which is consequently heated and provides a large thermionic current. The voltage drop is thus lowered. Such a tube is capable of conducting currents of thousands of amperes at voltage drops of tens of volts. Mercury-pool cathodes are used in one common form of arc discharge tube, supplying the electron current from an arc spot on the mercury-pool surface. The mercury vapor evaporated from the surface provides the gas atmosphere that is ionized.

Power Applications of Gas Tubes

Power rectifier and control tubes include mercury-vapor rectifiers, thyratrons, and ignitrons. These tubes employ the very high current-carrying capacity of gas discharge tubes with low power losses for rectification and control in high-power equipment. The operation of tubes depends on temperature insofar as mercury-vapor tube voltage drop and peak inverse voltages are concerned. (See Fig. 40.)

Hydrogen thyratrons are hot-cathode hydrogen-filled triodes designed for use as electronic switching devices where short anode delay time is important. In pulsing service, they are capable of switching tens of megawatts at voltages of tens of kilovolts. Anode delay time and time jitter are in the nanosecond range, and the tubes do not depend on ambient temperature for proper operation. Hydrogen thyratrons are also used in crowbar applications to protect other circuit components against fault voltages or currents and are capable of handling peak currents of several thousand amperes.

Triggered spark gaps are cold-cathode gas tubes operating in the arc discharge region (region III). The gaps contain two high-power electrodes and a trigger electrode, which is generally fired through a step-up pulse transformer by a simple low-energy pulse. The gaps are used as electronic switching devices for peak currents of tens of thousands of amperes and voltages of tens of kilovolts. They can discharge stored energies

TABLE 8. IONIZATION PROPERTIES OF GASES

Gas	Ionization Energy (eV)	Collision Probability P_c
Helium	24.5	12.7
Neon	21.5	17.5
Nitrogen	16.7	37.0
Hydrogen (H ₂)	15.9	20.0
Argon	15.7	34.5
Carbon monoxide	14.2	23.8
Oxygen	13.5	34.5
Krypton	13.3	45.4
Water vapor	13.2	55.2
Xenon	11.5	62.5
Mercury	10.4	67.0

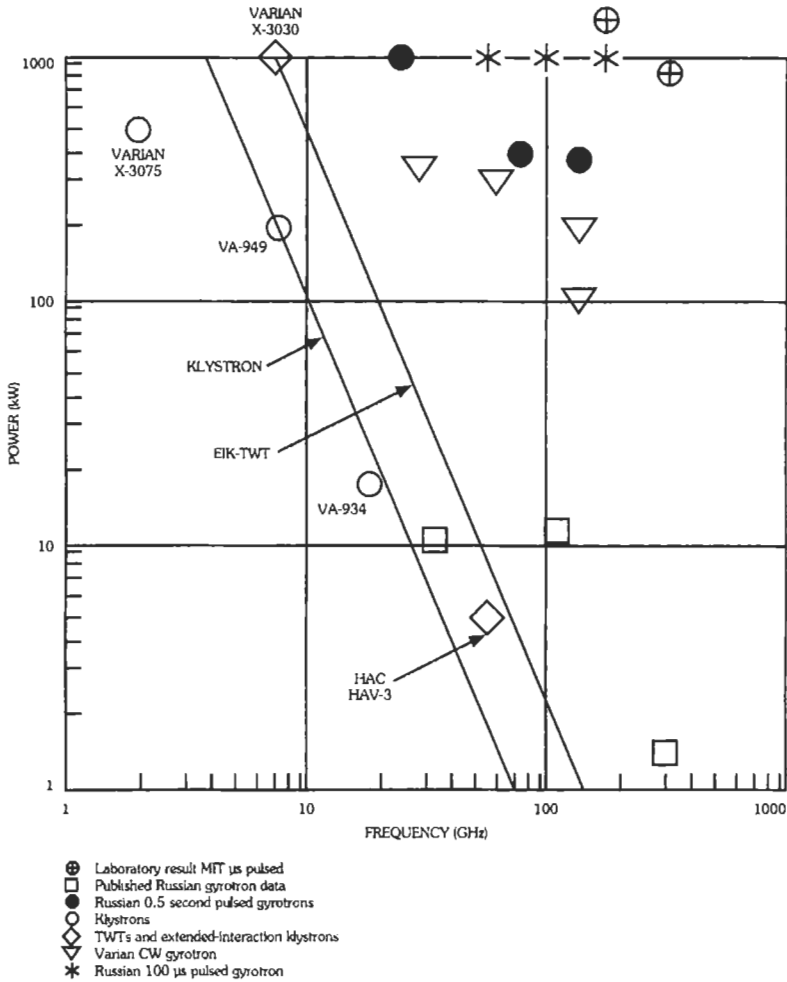


Fig. 36. Selected power output of high-frequency microwave tubes circa 1980.

(For an overview of the current state-of-the-art in high power millimeter-wave devices, see "Characteristics and applications of fast-wave gyrodevices," K. L. Felch, B. G. Danly, H. R. Jory, K. E. Kreischer, W. Lawson, B. Levush, and R. J. Temkin, *Proc. of the IEEE*, vol. 87, no. 5, pp. 752-781, May, 1999.)

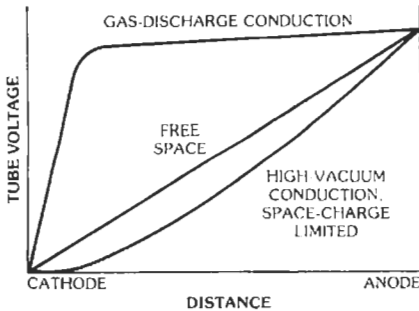


Fig. 37. Voltage distribution between plane parallel electrodes showing effect of space-charge neutralization in a hot-cathode gas tube.

of several thousand joules and are used for energy transfer in exploding-bridge-wire circuits, gas plasma discharges, spark chambers, and Keff cells. They are also used in crowbar applications for fast-acting protection of other circuit components against fault voltages and currents. Before conduction, the gap presents a low capacitance and a very high impedance to the circuit. After triggering, when the gap is conducting, the impedance drops to a few ohms or less.

Voltage regulators of the glow-discharge type take advantage of the volt-ampere characteristic in region II, where the voltage is nearly independent of the current. They operate at milliamperes and up to a few hundred volts.

Voltage regulators of the corona-discharge type operate at currents of less than a milliampere and at voltages up to several thousand volts.

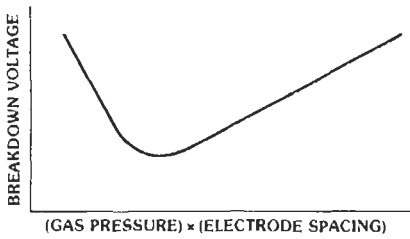


Fig. 38. Effect of gas pressure and tube geometry on gap voltage required for breakdown to occur in a cold-cathode gas tube.

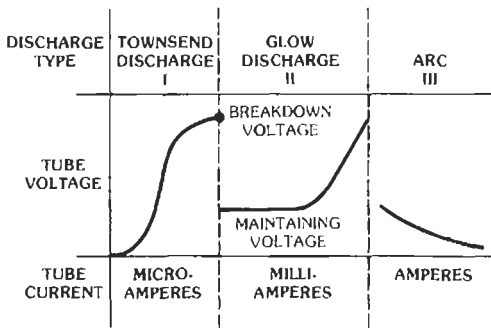


Fig. 39. Typical volt-ampere characteristic of cold-cathode gas discharge.

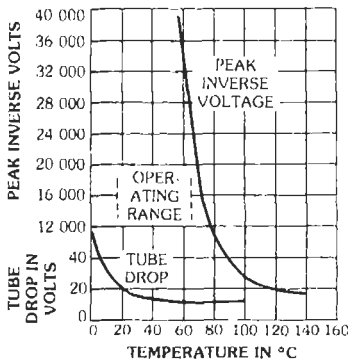


Fig. 40. Tube drop and arc-back voltages as a function of the condensed mercury temperature in a hot-cathode mercury-vapor tube. (Courtesy McGraw-Hill Book Co.)

Microwave Applications of Gas Tubes

Noise Sources—Gas-discharge devices possess a highly stable and repeatable effective noise temperature when in the fired condition. This feature provides a convenient and accurate means for determining the

noise figure. The microwave energy radiated from a gas-discharge plasma is coupled into a radio-frequency transmission line with which it is used. The amount of radio-frequency power available from a gas-discharge tube depends mainly on the nature of the gas fill, the geometric characteristic of the discharge tube, and the electron temperature of the positive column or plasma. The design parameter that most strongly determines the noise temperature is the type of gas employed. Any of the noble gases may be used in a noise source. In practice, however, only two or three are normally used (see Table 9).

In referring to a noise source or generator, the ratio of its noise power output to thermal noise power is called the Excess Noise Ratio (ENR) and is defined by

$$F = ENR = \frac{[(T_2/T_0) - 1] - Y[(T_1/T_0) - 1]}{Y - 1}$$

where

Y = ratio of the noise output power of the receiver with the noise generator on to that with the noise generator off

T_0 = 290 K

T_1 = temperature (in kelvins) of the termination

T_2 = effective noise temperature (in kelvins) of the noise generator in the fired condition

The expression $[(T_2/T_0) - 1]$ is termed the excess noise power of the noise source. When $T_1 = T_0 = 290$ K

$$ENR = [(T_2/T_0) - 1] / (Y - 1)$$

$$ENR(\text{dB}) = 10 \log_{10} ENR$$

The effective temperature of the noise source is equal to the temperature of the discharge only if the coupling of the transmission line to the discharge is complete. Otherwise, there is a reduction in the noise power output that can best be determined by measuring the fired and unfired insertion loss of the unit at the frequency of interest. The relation between these factors is given by

$$[(T_e/T_0) - 1] / [(T_2/T_0) - 1] = 1 - (L_u/L_f)$$

where

$[(T_e/T_0) - 1]$ is the effective excess noise power of the generator,

$[(T_2/T_0) - 1]$ is the excess noise power, and

L_u and L_f are the insertion losses in the unfired and fired conditions, respectively.

TABLE 9. NOISE FIGURES FOR SOME NOBLE GASES

Gas	$F = ENR(\text{dB})$
Helium	21.0
Neon	18.5
Argon	15.3

This correction should be subtracted from the apparent measured noise figure. The noise figure is always measured with reference to a standard temperature of 290 K (T_0). If the ambient temperature (T_1) of the noise-generator termination differs from the standard temperature, the noise figure calculated must be corrected. To find the correction factor, substitute the ambient temperature of the noise-generator termination for T_1 , in the following equation, and add the temperature factor (F_T) to the noise figure calculated.

$$F_T = [Y/(Y - 1)] [(T_1/T_0) - 1]$$

TR Tubes—Transmit-receive (TR) tubes are gas-discharge devices designed to isolate the receiver section of radar equipment from the transmitter during the period of high power output. A typical TR tube and its circuit are illustrated in Fig. 41. The cones in the waveguide form a transmission cavity tuned to the transmitter frequency, and the tube conducts received low-power-level signals from the antenna to the receiver. When the transmitter is operated, however, the high-power signal causes gas ionization between the cone tips, which detunes the structure and reflects all the transmitter power to the antenna. The receiver is protected from the destructively high level of power, and all of the available transmitter power is useful output.

Microwave Gas-Discharge Circuit Elements—Because of the high free-electron density, the plasmas

of gas discharges are capable of strong interaction with electromagnetic waves in the microwave region. In general, microwave phase shift and/or absorption result. If used in conjunction with a magnetic field, these effects can be increased and made nonreciprocal. Phase shift is a result of the change in dielectric constant caused by the plasma according to

$$\epsilon_p/\epsilon_0 = 1 - (0.8 \times 10^{-4} N_0/f_s^2)$$

where

- ϵ_p = dielectric constant in the plasma
- ϵ_0 = dielectric constant in free space
- N_0 = electron density in electrons/cm³ = signal frequency in megahertz
- f_s = signal frequency in megahertz

Absorption of microwave energy results when electrons, having gained energy from the electric field of the signal, lose this energy in collisions with the tube envelope or neutral gas molecules. This absorption is a maximum when the frequency of collisions is equal to the signal frequency and the absolute magnitude is proportional to the free-electron density.

LIGHT-SENSING AND LIGHT-EMITTING TUBES

Radiometry and Photometry*

Radiometric and photometric systems are generally based on the concept of radiated flux, where flux is defined as the total amount of radiation passing through a unit area per unit time. If a flux is measured in terms of its thermal heating ability, the most common unit is the watt, and the resultant measurement system is referred to as radiometry. If a flux is measured in terms of its ability to stimulate the standard photopic human eye, the resultant unit is the lumen, and the resultant measurement system is called photometry. A third choice for the measurement of flux is the number of photons per unit time. These three choices, in conjunction with the mks system of units, lead to the three mutually compatible systems of units shown in Table 10. Table 11 gives equivalents among units in different photometric measurement systems.

Flux Units

The number of lumens dL_λ and the number of photons per second dN_λ associated with a monochromatic flux dW_λ in watts are given by

$$dL_\lambda = 680 E_\lambda dW_\lambda \text{ and } dN_\lambda = (\lambda/hc) dW_\lambda$$

where

680 = number of lumens per watt of radiation at the peak photopic eye response

* See also Chapter 21.

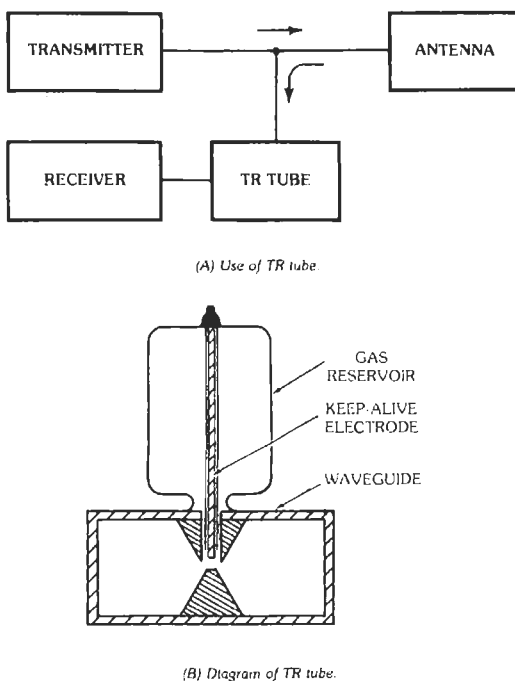


Fig. 41. Transmit-receive tube.

TABLE 10. COMPATIBLE SYSTEMS OF RADIATION UNITS

Parameter	Radiometric System	Photometric System	Photon System
Flux	watt	lumen	photon s ⁻¹
Source intensity	watt sr ⁻¹	lumen sr ⁻¹	photon s ⁻¹ sr ⁻¹
Incidence	watt m ⁻² (irradiance)	lumen m ⁻² (illuminance)	photon s ⁻¹ m ⁻²
Excitance	watt m ⁻² (emittance)	lumen m ⁻² (emittance)	photon s ⁻¹ m ⁻²
Scrance	watt sr ⁻¹ m ⁻² (radiance)	lumen sr ⁻¹ m ⁻² (luminance)	photon s ⁻¹ sr ⁻¹ m ⁻²
Energy	watt second	lumen second	photon

Note: The terms in parentheses are often used to characterize a measurement as either radiometric or photometric.

TABLE 11. PHOTOMETRIC EQUIVALENTS

Photometric Unit	Equivalent Unit Based on the Lumen (lm) as the Unit of Flux	Equivalent Lumen-MKS Unit
Source Intensity, <i>C</i>		
1 candela	1 lm sr ⁻¹	1 lm sr ⁻¹
1 Hefner candle	0.92 lm sr ⁻¹	0.92 lm sr ⁻¹
Surface Luminance, <i>B</i>		
1 candle cm ⁻²	1 lm sr ⁻¹ cm ⁻²	10 ⁴ lm sr ⁻¹ m ⁻²
1 candle m ⁻²	1 lm sr ⁻¹ m ⁻²	1 lm sr ⁻¹ m ⁻²
1 candle in ⁻²	1 lm sr ⁻¹ in ⁻²	1.55 × 10 ³ lm sr ⁻¹ m ⁻²
1 candle ft ⁻²	1 lm sr ⁻¹ ft ⁻²	10.8 lm sr ⁻¹ m ⁻²
1 nit	10 ⁻⁴ lm sr ⁻¹ cm ⁻²	1 lm sr ⁻¹ m ⁻²
1 stilb	1 lm sr ⁻¹ cm ⁻²	10 ⁴ lm sr ⁻¹ m ⁻²
1 apostilb	π ⁻¹ lm sr ⁻¹ m ⁻²	π ⁻¹ lm sr ⁻¹ m ⁻²
1 lambert	π ⁻¹ lm sr ⁻¹ cm ⁻²	10 ⁴ π ⁻¹ lm sr ⁻¹ m ⁻²
1 millilambert	10 ⁻³ π ⁻¹ lm sr ⁻¹ cm ⁻²	10 π ⁻¹ lm sr ⁻¹ m ⁻²
1 footlambert	π ⁻¹ lm sr ⁻¹ ft ⁻²	10.8 π ⁻¹ lm sr ⁻¹ m ⁻²
Illuminance of a Surface, <i>I_L</i>		
1 lux	1 lm m ⁻²	1 lm m ⁻²
1 phot	1 lm cm ⁻²	10 ⁴ lm m ⁻²
1 milliphot	10 ⁻³ lm cm ⁻²	10 lm m ⁻²
1 footcandle	1 lm ft ⁻²	10.8 lm m ⁻²
Energy, <i>U</i>		
1 talbot	1 lm s	1 lm s

E_λ = normalized (to unity maximum) photopic human eye response (Fig. 42)
 λ = wavelength of the monochromatic radiation (m)
 h = Planck's constant $\approx 6.6 \times 10^{-34}$ (J · s)
 c = velocity of light $\approx 3.0 \times 10^8$ (m/s)

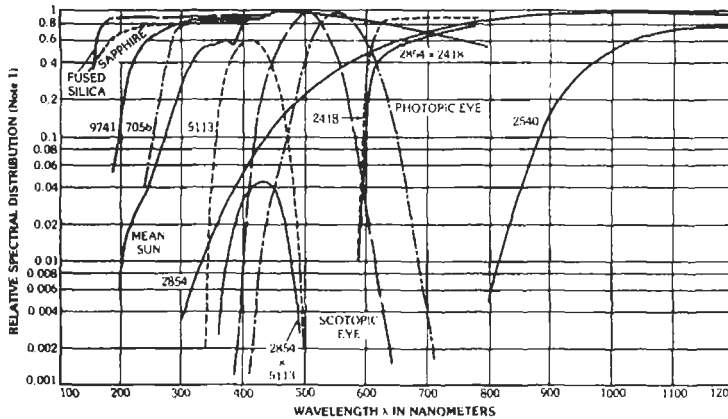
having a wattage W between the wavelength limits λ_1 and λ_2 are given by

$$L/W = 680 \int_0^\infty E_\lambda w_\lambda d\lambda / \int_{\lambda_1}^{\lambda_2} w_\lambda d\lambda$$

and

$$Nhc/W = \int_{\lambda_1}^{\lambda_2} \lambda w_\lambda d\lambda / \int_{\lambda_1}^{\lambda_2} w_\lambda d\lambda$$

The number of lumens L and the number of photons per second N between the wavelength units of λ_1 to λ_2 associated with a distributed spectral radiation source



Fused silica: transmission through polished "Suprasil" (1 millimeter thick)
 Sapphire: transmission through polished Sapphire (1 millimeter thick)
 9741: transmission through polished Corning 9741 glass (1 millimeter thick)
 7056: transmission through polished Corning 7056 glass (1 millimeter thick)
 Mean sun: mean solar distribution at Earth's surface
 5113: transmission through polished Corning 5113 filter (CS-5-58) (half-stock thickness)

2854 × 5113: product of 5113 curve and 2854 curve
 2854: spectral density distribution of 2854 K color-temperature tungsten lamp
 Scotopic eye: relative response of dark-adapted eye
 Photopic eye: standard eye response
 2418: transmission through polished Corning 2418 filter (CS-2-62) (stock thickness)
 2854 × 2418: product of 2854 and 2418 curves
 2540: transmission through polished Corning 2540 filter (CS-7-56) (stock thickness)

Note 1: "Relative spectral distribution" designates the relative radiant-energy density distribution w_λ for sources, the relative visual stimulation for equienergy inputs for the eye response, and the spectral transmission t_λ for windows and filters. The transmission characteristics of individual filter and window samples can be expected to depart appreciably from these typical values.

Fig. 42. Useful spectral distributions.

where

$$W = w_{\lambda_{max}} \int_{\lambda_1}^{\lambda_2} w_\lambda d\lambda$$

and where

$w_{\lambda_{max}}$ = maximum spectral density in watts per unit wavelength in the spectral band between λ_1 and λ_2 ,
 w_λ = relative spectral distribution of the radiation source on thermal-energy basis, normalized to a maximum value of unity.

Typical w_λ spectral distributions are in Fig. 42.

Optical Imaging

In an optical lens system of flux-gathering diameter D_f in meters, focal length f in meters, and optical transmittance T , the ratio $f/D_f = n_f$ is called the f -number of the lens. If the surface of an object of radiance or luminance B in flux units per steradian per meter² is imaged by this system with a linear magnification m , and if we assume Lambertian emittance characteristics over the solid angle subtended by the optical system,

the image will be subjected to an irradiance or illuminance I_L in flux units per meter² given by

$$I_L = \pi B T / [4n_f^2 (m+1)^2 + m^2]$$

For objects at infinity, $m = 0$, and

$$I_L (\text{object at infinity}) = \pi B T / 4n_f^2$$

If the irradiance (or illuminance) I_L in flux units per meter² is allowed to fall on a nonabsorbing Lambertian diffusing surface, the resultant image radiance (or luminance) B_i in flux units per steradian per meter² is given by

$$\pi B_i = I_L$$

Any desired method of measuring flux units, such as watts, lumens, or photons/second (Table 11), can be selected for expressing the object radiance (or luminance) B in flux units steradian⁻¹ meter⁻² and the irradiance (or illuminance) I_L in flux units meter⁻² in these relationships. Thus, a radiance B in watts steradian⁻¹ meter⁻² would be paired with an irradiance I_L in watts meter⁻², a luminance B in lumens steradian⁻¹ meter⁻² with an illuminance I_L , in lumens meter⁻², and a radi-

ance B in photons second⁻¹ steradian⁻¹ meter⁻² with an irradiance I_L in photons second⁻¹ meter⁻².

Any spectral distribution modifications, if present, would be included in the numerical magnitude of the lens transmission T , defined as the ratio of the total output flux from the optical system to the corresponding input flux.

Selection of appropriate alternative pairs of luminance and illuminance units when the flux units are not explicitly stated (first column of Table 11) must be made with care. Thus candle centimeter⁻² (or stilb) would be paired with phot, candle meter⁻² (or nit) with lux, and candle foot⁻² with footcandle. Even greater difficulty arises when the factor Π in the preceding relationships is absorbed or included in the units of luminance. Thus the product ΠB in apostilbs would be paired with I_L in lux, the product ΠB in lamberts with I_L in phots, the products ΠB in millilamberts with I_L in milliphots, and the product ΠB in footlamberts with I_L in footcandles. These difficulties are avoided by the use of the compatible systems of radiation units shown in Table 10.

Typical Approximate Illumination Values at the Surface of the Earth

- Sun at zenith $\approx 10^4$ footcandles
 $\approx 10^5$ lumens meter⁻²
- Full moon $\approx 3 \times 10^{-2}$ footcandles
 $\approx 3 \times 10^{-1}$ lumens meter⁻²

Typical Approximate Brightness Values

	footlamberts	lm sr ⁻¹ m ⁻²
Highlights, 35-millimeter movie	≈ 4	≈ 100
Page brightness for reading fine print	≈ 10	$\approx 3 \times 10^2$
November football field	≈ 50	$\approx 1.5 \times 10^3$
Surface of moon seen from Earth	$\approx 1.5 \times 10^3$	$\approx 5 \times 10^4$
Summer baseball field	$\approx 3 \times 10^3$	$\approx 10^5$
Surface of 40-watt frosted lamp bulb	$\approx 8 \times 10^3$	$\approx 2.5 \times 10^5$
Crater of carbon arc	$\approx 4.5 \times 10^7$	$\approx 10^9$
Sun seen from Earth	$\approx 5.2 \times 10^8$	$\approx 1.5 \times 10^{10}$

LIGHT-SENSING TUBES

Image Tubes and Image Intensifiers

An image tube* is an optical-image-in to optical-image-out electron-tube device, combining an input photocathode and an output phosphor screen such that

photoelectrons emitted from each point on the photocathode subsequently excite a corresponding individual image "point" on the phosphor screen. Various focusing means, including magnetic and electrostatic electron lenses, may be used to assure maximum point-to-point correlation between the input and output images. The principal operating requirements are a lens to form the input image and a high-voltage supply, typically 5–25 kilovolts, to provide sufficient electron-beam energy to excite the output phosphor screen.

If means are provided within the image tube to amplify the photoelectrons before they strike the output phosphor screen, or if the tube without such means produces a much brighter output image than the input image would produce on a diffusing screen, the tube is commonly called an image-intensifier tube.

Image-intensifier tubes are used to amplify the brightness of a faint input image for better visual or photographic viewing, whereas image tubes without amplification are used to convert radiation from one spectral region to another (image conversion) or to perform such control operations as optical shuttering by programming the applied high voltage.

The total output flux dW_o , in watts exiting (through 2Π steradians) from the phosphor-screen faceplate of an image-intensifier tube for an input monochromatic flux dW_λ in watts at a wavelength λ is given by

$$dW_o = s_\lambda G \mathcal{E}_w (V - V_k) dW_\lambda = G_\lambda dW_\lambda$$

where

- G_λ = monochromatic wattage gain of the image-intensifier tube at a wavelength λ -ratio of the total output flux dW_o in watts to the input monochromatic flux dW_λ in watts
- s_λ = radiant sensitivity of the input photocathode in amperes per watt (see Fig. 13, Chapter 21)
- G = internal current gain ratio of the image-intensifier tube = ratio of the current bombarding the output phosphor screen to the corresponding photocurrent leaving the input photocathode
- \mathcal{E}_w = absolute phosphor efficiency = ratio of the total radiated flux in watts to the exciting electron-beam power in watts dissipated in the particles of the output phosphor screen
- V = energy of the electron beam in volts bombarding the output phosphor screen
- V_k = extrapolated knee voltage of the output phosphor screen.

If the phosphor screen radiates flux according to Lambert's law (usually only approximately valid), the

* *Photo-Electrical Image Devices*, Advances in Electronics and Electron Physics, Vols. 12, 16, 22A, and 22B (New York and London: Academic Press, 1960, 1962, and 1966), H. V. Soule, *Electro-Optical Photography at Low Illumination Levels* (New York: John Wiley & Sons, 1968).

corresponding output image radiance R_o , in watts steradian⁻¹ meter⁻² is given by

$$R_o = G_{\lambda_1} I_{\lambda_1} / \Pi m^2$$

where

- I_{λ_1} = input image irradiance on the photocathode expressed in watts meter⁻² at the wavelength λ_1 .
- m = differential magnification ratio of the image tube = output incremental image size divided by the corresponding input incremental image size.

For a spectrally distributed input flux having a known relative spectral distribution w_λ and a known total radiated power $W_{\lambda_1\lambda_2}$ in watts between the wavelength limits λ_1 and λ_2 , the resulting total output flux W_o , in watts exiting from the image tube, is given by

$$W_o = s_{\lambda_{max}} \left(\int_0^{\infty} \sigma_\lambda w_\lambda d\lambda / \int_{\lambda_1}^{\lambda_2} w_\lambda d\lambda \right) \times G \mathcal{E}_w (V - V_k) W_{\lambda_1\lambda_2} = G_{\lambda_1\lambda_2} W_{\lambda_1\lambda_2}$$

where

- $s_{\lambda_{max}}$ = peak radiant sensitivity of the input photocathode in amperes per watt
- σ_λ = relative radiant sensitivity of the input photocathode as a function of wavelength λ normalized to unity maximum
- w_λ = relative spectral distribution of the power density spectrum of the input flux normalized to unity maximum
- $G_{\lambda_1\lambda_2}$ = wattage gain of the image tube for the relative spectral distribution w_λ and the wavelength limits λ_1 and λ_2 .

Typical values for the magnitude of the dimensionless spectral-matching-factor ratio are given by

$$\int_0^{\infty} \sigma_\lambda w_\lambda d\lambda / \int_{\lambda_1}^{\lambda_2} w_\lambda d\lambda$$

The total output flux L_o , in lumens exiting from an image tube, corresponding to the total output flux W_o in watts, can be computed from the flux conversion relationships given in the section on radiometry and photometry, or from the following relationship

$$L_o = s_{\lambda_{max}} \left(\int_0^{\infty} \sigma_\lambda w_\lambda d\lambda / \int_0^{\infty} E_\lambda w_\lambda d\lambda \right) \times G \mathcal{E}_w \left(\int_0^{\infty} E_\lambda w_{o\lambda} d\lambda / \int_0^{\infty} w_{o\lambda} d\lambda \right) (V - V_k) L_i = G_L I_i$$

where

- G_L = luminous gain of the image intensifier tube = ratio of the output flux in lumens to the cor-

responding input flux in lumens for the spectral input distribution w_λ

- E_λ = standard tabulated average relative photopic eye response
- $w_{o\lambda}$ = relative spectral density distribution of the output flux
- L_i = input flux in lumens

The typical values of the dimensionless spectral matching factors can be used to determine the magnitude of the dimensionless integral ratios appearing in these relationships. For the special case where the input flux L_i (2854) in lumens is generated by a 2854 K color-temperature tungsten-filament lamp, the output flux L_o (2854) in lumens is given by

$$L_o(2854) = S(2854) G \mathcal{E}_w \left(\int_0^{\infty} w_{o\lambda} E_\lambda d\lambda / \int_0^{\infty} w_{o\lambda} d\lambda \right) \times (V - V_k) L_i(2854) = G_L(2854) L_i(2854)$$

where

- $G_L(2854)$ = luminous gain of the image intensifier for 2854 K tungsten-lamp radiation
- $S(2854)$ = luminous sensitivity of the input photocathode for 2854 K tungsten-lamp radiation

The magnitude of the luminous gain $G_L(2854)$ is commonly used to characterize the image intensification properties of an image-intensifier tube.

If the output phosphor screen radiates flux according to Lambert's law (usually only approximately valid), the output image luminance (or brightness) B_o in lumens steradian⁻¹ meter⁻² is given by

$$B_o = G_L I_i / \pi m^2$$

where

- G_L = luminous gain of the image intensifier for input spectral distribution w_λ
- I_i = input illuminance (or illumination) on the photocathode in lumens meter⁻² for the spectral distribution w_λ
- m = differential magnification ratio of the image tube = output incremental image size divided by the corresponding input incremental image size

Internal current gain G of an image-intensifier tube can be obtained by the use of an internal sandwich electrode, in which an auxiliary or sandwich photocathode is mounted in close proximity to and following an auxiliary or sandwich phosphor. Photoelectrons from the input photocathode of the tube are then imaged onto this sandwich phosphor screen, and the flux from this screen is coupled to the sandwich photocathode, generating an enhanced photocurrent. The current gain ratio G of this sandwich phosphor-photocathode combination, defined

as the ratio of output photocurrent to input photocurrent, is given by

$$G = S_{\lambda_{max, sand}} \left(\int_0^{\infty} w_{\lambda(sand)} \sigma_{\lambda(sand)} d\lambda / \int_0^{\infty} w_{\lambda(sand)} d\lambda \right) \mathcal{E}_{w, sand} (V_{sand} - V_{k, sand}) \gamma$$

where

- $S_{\lambda_{max, sand}}$ = peak monochromatic responsivity of the sandwich photocathode in amperes watt⁻¹
- $w_{\lambda(sand)}$ = relative spectral distribution of the flux emitted by the sandwich phosphor screen
- $\sigma_{\lambda(sand)}$ = the relative spectral distribution of the sandwich photocathode
- $\mathcal{E}_{w, sand}$ = absolute efficiency of the sandwich phosphor screen in watts watt⁻¹
- V_{sand} = electron beam energy in volts bombarding the sandwich phosphor screen
- $V_{k, sand}$ = extrapolated knee voltage in volts for the sandwich phosphor screen
- γ = optical coupling efficiency of the sandwich electrode = ratio of the flux falling onto the sandwich photocathode to the corresponding flux emitted by the sandwich phosphor screen

The combination of a phosphor screen and a photocathode to produce current G can also be achieved by optically coupling the output flux from one image tube to the input of a second tube.

Resolution in image tubes and image-intensifier tubes is a subjective parameter describing the number of pairs of equally spaced illuminated and unilluminated bars per unit distance at the photocathode imaged onto the input photocathode surface that can just be distinguished visually by a trained observer under stated test conditions.

Distortion is a parameter describing any change in the geometric shape of the output image compared with the input image. Radially increasing magnification leads to "pincushion" distortion; radially decreasing magnification leads to "barrel" distortion; and radially changing image rotation leads to "S" distortion.

Gas Photodiodes

In diode phototubes not containing a high vacuum, ionization by collision of electrons with neutral molecules may occur so that more than one electron reaches the anode for each originally emitted photoelectron. This "gas amplification factor" has a value of between 3 and 5; a higher factor causes instabilities. Gas-tube operation is restricted to frequencies below about 10 000 hertz.

Image Orthicons

The image orthicon* is a camera tube that is widely used for commercial television. This fact derives from its high sensitivity, its close spectral-sensitivity match to the human eye, and its relatively fast response. Good quality commercial television pictures can be generated by an image orthicon viewing a 5-to-20-footlambert ($\approx 15\text{--}50$ lumen steradian⁻¹ meter²) scene through an $f/5.6$ lens. The image orthicon is generally available with either S-10 or S-20 spectral response (see Fig. 13, Chapter 21) and is capable of 500 picture elements per raster height (9.9 line pairs/millimeter) at 30% video-amplitude response.

Principle of Operation—Fig. 43 is a diagram of an image orthicon. In the image section, a light image incident on the translucent photocathode liberates photoelectrons into the adjacent vacuum region in proportion to the light intensity (gamma is unity) on each

* A. Rose, P. K. Weimer, and H. B. Law, "The image orthicon—A sensitive television pickup tube," *Proceedings of the IRE*, vol. 34, no. 7, pp. 424-432, July, 1946.

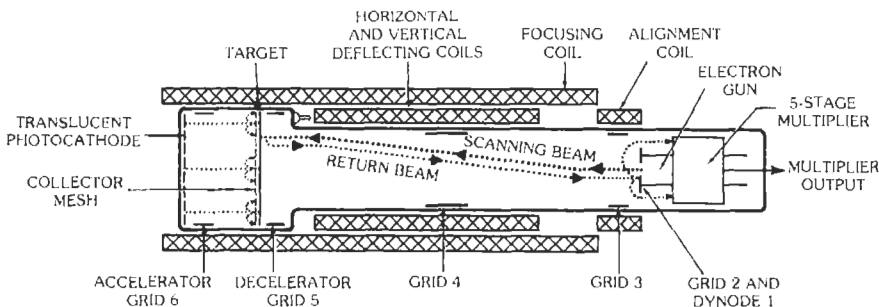


Fig. 43. Image orthicon. (By permission of RCA, copyright proprietor.)

element of the cathode. These photoelectrons are accelerated toward and magnetically focused onto the surface of a thin semiconducting target. Electrons strike this target with sufficient energy to liberate a larger number of secondary electrons (typically five) for each incident primary. The secondary electrons are collected by a mesh closely spaced from the target membrane. Hence, by depletion of electrons from the thin membrane, incremental areas become positive in proportion to the number of photoelectrons striking each element. In cases of high-light-level operation, parts of the target may become charged to target (collector) mesh potential, and saturation charge results. This phenomenon accounts for the so-called "knee" in the signal-versus-illumination transfer curve (Fig. 43).

Because the target membrane is very thin, of the order of microns, a charge distribution pattern formed on the image-section surface appears nearly simultaneously and identically on the scanning-section surface.

In the scanning section, an electron gun generates a highly apertured electron beam from a fraction to tens of microamperes in intensity. A solenoidal magnetic-focus coil and saddle-type deflection coils surrounding the scan section focus this beam on the insulator target and move it across the target. Scan-beam electrons impinge on the target at very low velocity, giving rise to relatively few secondary electrons. The target acts somewhat as a retarding-field electrode and reflects a large number of the beam electrons that have less than average axial velocity. These two phenomena—small but finite secondary emission and reflection of slow beam electrons—limit scan-beam modulation to a maximum of about 30% at high light levels, and to 2 orders less at threshold. As will be shown later, the large unmodulated return beam current is the primary source of noise in the image orthicon.

Another problem created by the retarding-field aspect of low-velocity target scanning appears when the deflected beam does not strike the target normally. Since the entire beam-velocity component normal to the surface is now reduced by the cosine of the angle of incidence, the effective beam impedance is greatly increased. To overcome this problem, the decelerating field between grids 4 and 5 is shaped such that the electron beam always approaches normal to the plane of the target at a low velocity. If the elemental area on the target is positive, then electrons from the scanning beam deposit until the charge is neutralized. If the elemental area is at cathode potential (corresponding to a dark picture area), no electrons are deposited. In both cases, the excess beam electrons are turned back and focused into a five-stage electron multiplier. The charges existing on either side of the semiconductive target membrane will, by conductivity, neutralize each other in less than one frame time. Electrons turned back at the target form a return beam that has been amplitude-modulated in accordance with the charge pattern of the target.

The return beam is redirected by the deflection and focus fields toward the electron gun where it origi-

nated. Atop the electron gun, and forming the final aperture for that gun, is a flat secondary-emitting surface comprising the first dynode of the electron multiplier. The return beam strikes this surface, generating secondary electrons in a ratio of approximately 4:1.

Grid 3 facilitates a more complete collection by dynode 2 of the secondary electrons emitted from dynode 1. The gain of the multiplier is high enough that in operation the limiting noise is the shot noise of the returned electron beam rather than the input noise of the video amplifier.

Signal and Noise—Typical signal output currents for tube types 5820 and 5826 are shown in Fig. 44. The tubes should be operated so that the highlights on the photocathode bring the signal output slightly over the knee of the signal-output curve.

The spectral response of type 5820 and 5826 image orthicons is shown in Fig. 45. Note that when a Wratten 6 filter is used with the tube, a spectral curve closely approximating that of the human eye is obtained.

From the standpoint of noise, the total television system can be represented as shown in Fig. 46, where I_s = signal current, I_n = total image-orthicon noise current, E_{nt} = thermal noise in R_1 , E_{ns} = shot noise in the input amplifier tube, R_1 = input load, C_1 = total input shunt capacitance, and R_i = shot-noise equivalent resistance of the input amplifier = $2.5/g_m$ for triode or cascode input = $[I_b/(I_b + I_c)][(2.5/g_m) + (20I_{c2}/g_m^2)]$ for pentode input, with g_m = transconductance of input tube or cascode combination, I_b = amplifier direct plate current, and I_c = amplifier direct screen-grid current.

The noise added per stage is

$$\Delta n = [\sigma / (\sigma - 1)]^{1/2}$$

where σ = stage gain in the multiplier. For a total multiplier noise figure to be directly usable, it must be referred to the first-dynode current; therefore, for five multiplier stages

$$\overline{\Delta N} = \Delta n^2 + \frac{\Delta n^2}{\sigma^2} + \frac{\Delta n^2}{\sigma^4} + \frac{\Delta n^2}{\sigma^6} + \frac{\Delta n^2}{\sigma^8}$$

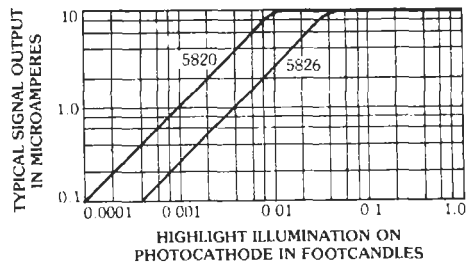


Fig. 44. Basic light-transfer characteristic for type 5820 and 5826 image orthicons. The curves are for small-area highlights illuminated by tungsten light, white fluorescent light, or daylight. (By permission of RCA, copyright proprietor.)

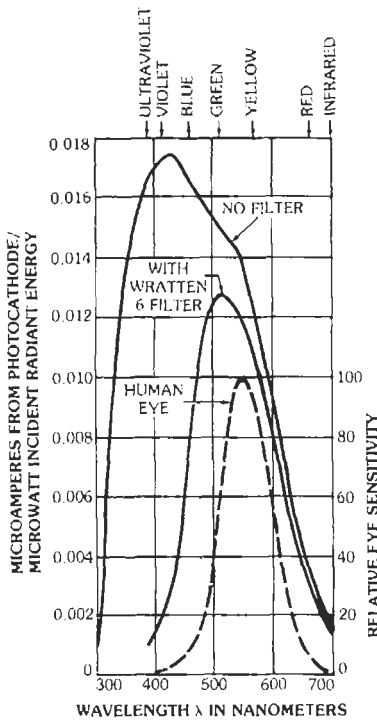


Fig. 45. Spectral sensitivity of image orthicon. (By permission of RCA, copyright proprietor.)

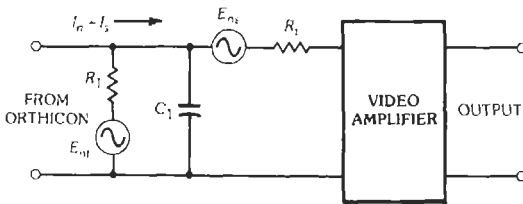


Fig. 46. Equivalent circuit for noise in orthicon and first amplifier stage.

where ΔN = electron-multiplier noise factor referred to multiplier input.

After combining all noise sources

$$\frac{S}{N} = \frac{I_s}{\left[F \left[2elk_m^2 + 4kT \left(\frac{1}{R_1} + \frac{R_1}{R_1^2} + \frac{\omega^2 C_1^2 R_1}{3} \right) \right] \right]^{1/2}}$$

where

- S/N = signal-to-noise ratio
- F = bandwidth in hertz
- e = electron charge = 1.6×10^{-19} coulomb
- I = image-orthicon beam current

- k_m = electron-multiplier noise factor, referred to multiplier output $m\Delta N$
- k = Boltzmann's constant = 1.38×10^{-23} joule/kelvin
- T = absolute temperature in kelvins
- $\omega = 2\pi f$ in hertz

The signal current is an alternating-current signal superimposed on a larger direct beam current. This can be thought of as a modulation of the beam current. Properly adjusted tubes obtain as much as 30% modulation.

$$I_s = mMI$$

where

- m = multiplier gain
- M = percentage modulation

If S/N is now rewritten,

$$\frac{S}{N} = \frac{I_s}{\left[4kTF \left(\frac{2el_s m \Delta N^2}{4kTM} + \frac{1}{R_1} + \frac{R_1}{R_1^2} + \frac{\omega^2 C_1^2 R_1}{3} \right) \right]^{1/2}}$$

In typical television operation, the thermal noise of the load resistor and the shot noise of the first amplifier can be neglected.

Focusing and Scanning Fields—The electron optics of the scanning section of the tube are quite complicated, and space does not permit the inclusion of the complete equations. A simple relationship between the strength of the magnetic focusing field and the magnetic deflection field is given below.

The image orthicon is usually operated with multiple-node focus in the scanning section. Working at a multiple-node focus not only demands more focus current but also more deflection current. Note the deflection path in Fig. 47. Let H = horizontal dimension of scanned area or target, L = effective length of horizontal deflection field, H_d = horizontal deflection field (peak-to-peak value), and H_f = focusing field. Then

$$H_d = H_f H/L$$

For the image orthicon, $H \approx 1.25$ inches, and $L \approx 4$ inches. Thus $H_f \approx 75$ gauss, and $H_d \approx 23$ gauss.

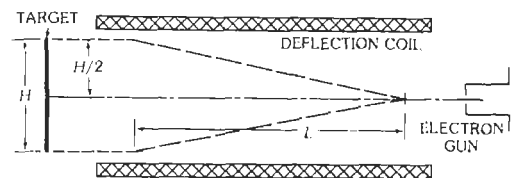


Fig. 47. Deflection in image orthicon.

Vidicons

The vidicon* is a small television camera tube that is used primarily for industrial television, space application, and studio film pickup because of its small size and simplicity.

As shown in Fig. 48, the tube consists of a signal electrode composed of a transparent conducting film on the inner surface of the faceplate, a thin layer (a few micrometers) of photoconductive material deposited on the signal electrode, a fine mesh screen (grid 4) located adjacent to the photoconductive layer, a focusing electrode (grid 3) connected to grid 4, and an electron gun.

Principle of Operation—Each elemental area of the photoconductor can be likened to a leaky capacitor with one plate electrically connected to the signal electrode that is at some positive voltage (usually about 20 volts) with respect to the thermionic cathode of the electron gun and the other plate floating except when commutated by the electron beam. Initially, the gun side of the photoconductive surface is charged to cathode potential by the electron gun, thus leaving a charge on each elemental capacitor. During the frame time, these capacitors discharge in accordance with the value of their leakage resistance, which is determined by the amount of light falling on each elemental area. Hence, there appears on the gun side of the photoconductive surface a positive-potential pattern corresponding to the pattern of light from the scene imaged on the opposite surface of the layer. Even those areas that are dark discharge slightly, since the dark resistivity of the material is not infinite.

The electron beam is focused at the surface of the photoconductive layer by the combined action of the uniform magnetic field and the electrostatic field of grid 3. Grid 4 serves to provide a uniform decelerating field between itself and the photoconductive layer such that the electron beam always approaches the surface normally and at a low velocity. When the beam scans

the surface, it deposits electrons where the potential of the elemental area is more positive than that of the electron-gun cathode. At this moment the electrical circuit is completed through the signal-electrode circuit to ground. The amount of signal current depends on the amount of discharge in the elemental capacitor, which in turn depends on the amount of light falling on this area.

Alignment of the beam is accomplished by a transverse magnetic field produced by external coils located at the base end of the focusing coil.

Deflection of the beam is accomplished by the transverse magnetic fields produced by external deflecting coils.

Signal and Noise—Because the vidicon acts as a constant-current generator as far as signal current is concerned, the value of the load resistor is determined by bandpass and noise considerations in the input circuit of the video amplifier. Unlike the image orthicon, the vidicon has the signal current removed at the target, and only that portion of the scan beam actually involved in the target discharge contributes shot noise. Moreover, electron-beam contributions to noise are minimal for low-light portions of the scene.

The primary noise associated with vidicon operation is seldom scan-beam shot noise. Where the signal current is less than 1 microampere and the bandpass is relatively wide, the principal noise in the system is contributed by the input circuit and the first stage of the video amplifier. To minimize the thermal noise of the load resistor, its resistance is made much higher than flat-bandpass considerations would indicate, as signal voltage increases directly and noise voltage increases as the square root. To correct for attenuation of the signal with increasing frequency, the amplitude response of the video amplifier frequently employs high-frequency boost of the following form, where C_1 and R_1 refer to Fig. 49:

$$G = G_0(1 + 4\pi^2 F^2 C_1^2 R_1^2)^{1/2} / R_1$$

A representative plot of amplitude response as a function of the number of television lines (per raster height) is shown in Fig. 50.

The vidicon has somewhat more lag or image persistence than the image orthicon. This is the result of two factors. To obtain high-sensitivity surfaces, the

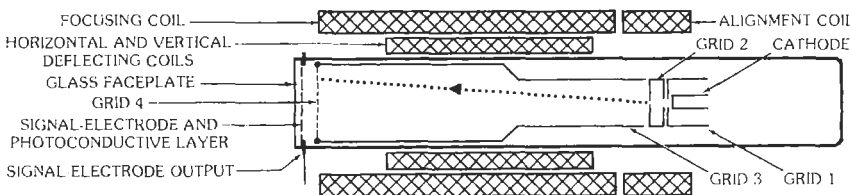


Fig. 48. Vidicon construction. (By permission of RCA, copyright proprietor.)

* B. H. Vinc, R. B. Janes, and F. S. Veith, "Performance of the vidicon—A small developmental camera tube," *RCA Review*, vol. 13, no. 1, pp. 3-10; March, 1952; P. Weimer, S. Forgue, and R. Goodrich, "The vidicon photoconductive camera tube," *Electronics*, vol. 23, no. 5, pp. 70-73, May, 1950.

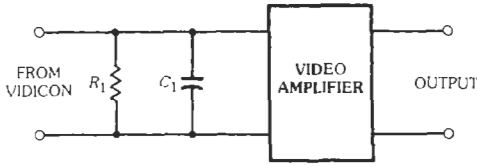


Fig. 49. Equivalent input circuit for first-stage amplifier in vidicon circuit.

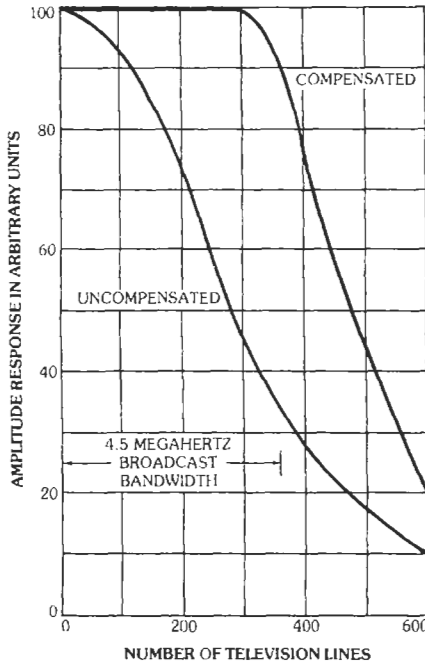


Fig. 50. Vidicon resolution, showing uncompensated and compensated horizontal responses. Highlight signal-electrode microamperes = 0.35; test pattern = transparent square-wave resolution wedge; 80 television lines = 1-megahertz bandwidth. (By permission of RCA, copyright proprietor.)

photoconductive decay time is made as long as tolerable, because quantum efficiency is limited by the ratio of effective carrier lifetime to carrier transit time across the photoconductor. A second source of lag is simply the RC time constant of the target recharging circuit, that is, the target capacitance and the beam impedance.

The spectral response of most commercial vidicons designated S-18 is more actinic than that of the human eye. Fig. 51 compares these responses with the spectrum of a 2854 K tungsten source.

Variations of the Vidicon

Interest in optical guidance and surveillance from aircraft and spacecraft has given rise to a wide variety

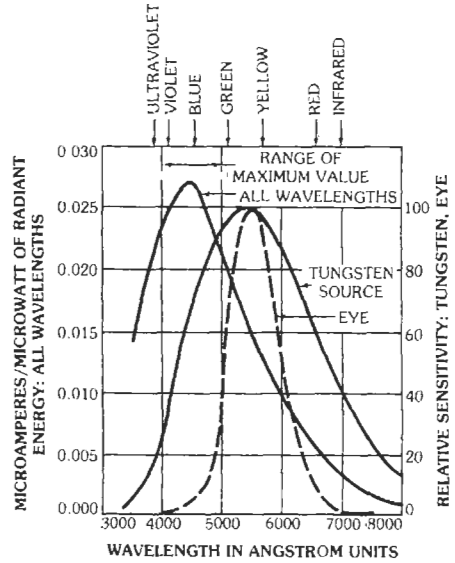


Fig. 51. Spectral response of vidicon. (By permission of RCA, copyright proprietor.)

of vidicon camera tubes. To treat these variations in detail becomes encyclopedic, but the following gives some indication of the choices now available to the user.

Effective Sensitivity—True photoconductive tubes now offer sensitivities of 150–200 nanoamperes for 1/2-footcandle illumination with 20 nanoamperes dark current. Improved methods of deposition of photoconductors have made possible higher-voltage operation without objectionable dark shading. Special devices using junction effects promise even better sensitivity.

Spectral Response—Available photoconductors, taken as a whole, provide sensitivity over the entire visible range with usual (7056) glass windows. Quartz-window tubes offer useful sensitivity to below 2000 angstroms. Numerous applications of direct excitation of photoconductors by X-radiation have been reported. High-velocity electron excitation (bombardment-induced conductivity) is also in use.

Size and Deflection—Vidicons are available in sizes ranging from 1/2 inch to 2 inches in diameter. Various combinations of deflection and focus are available.

Storage—A number of manufacturers have produced vidicons with long-storage characteristics. Many are merely long-lag tubes; however, a few rely on high-resistivity materials or on barrier layers to retain stored charge through minimal dark current. One such device, once exposed properly to a scene, regenerates the scene through readout over a period of the order of half an hour.

LIGHT-EMITTING TUBES

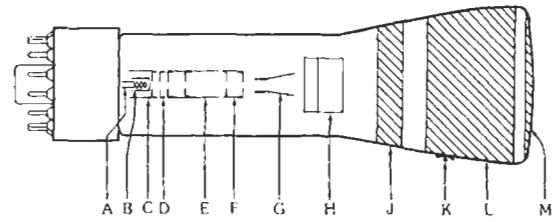
Cathode-Ray Tubes

A cathode-ray tube (CRT)* is a vacuum tube in which an electron beam, deflected by applied electric or magnetic fields, produces a trace on a fluorescent screen. Cathode-ray tubes have many uses; the most popular is the television (TV) display. Other uses include radar displays, oscillography, and visual displays for computer terminals.

Principle of Operation—The function of the cathode-ray tube is to convert an electrical signal into a visual display. The tube contains an electron-gun structure (to provide a narrow beam of electrons) and a phosphor screen. The electron beam is directed to the phosphor screen and strikes it, causing light to be emitted in a small area or spot in proportion to the intensity of the electron beam. The beam intensity varies as a function of the beam-control element in the electron gun. In a TV or computer-display CRT, the electrical signal that controls the beam intensity corresponds to the desired picture information and is referred to as the video signal. Although one spot or picture element is not enough to reproduce a picture, by moving the spot over the entire screen in a systematic manner, the complete picture can be reproduced.

For oscilloscope use, the deflection of the beam at the screen is proportional to the voltages applied to the deflection electrodes within the tube, and a visual picture of time-dependent waveforms can be produced on the screen in the same manner as a graph is drawn on paper.

Electrostatic Deflection—Electrostatic deflection is generally used in oscillography where high-frequency signals need to be analyzed. Electrostatic deflection requires little power, and deflectors can be designed to perform at extremely high frequencies. Electrodes placed within the tube (Fig. 52) form an electric field perpendicular to the electron path. Deflection is determined by the voltage across the electrodes, the accelerating voltage, and the distance to the screen. For high sensitivity, the accelerating voltage in the gun needs to be low, but for high brightness, the accelerating voltage at the screen needs to be high. Some form of post-deflection acceleration (PDA) is therefore required. The most common form of PDA uses a dome-shaped mesh at the end of the electron gun to form a scan-expansion lens. The mesh is maintained at gun voltage (2000 V), and the conductive coating in the glass envelope is maintained at a higher voltage (20 000 V). The resultant electric field between the mesh and this conducting surface causes an outward radial force on the electrons so that their angle from the gun is increased and scan magnification



A—Heater
B—Cathode
C—Control Electrode
D—Screen grid or preaccelerator
E—Focusing electrode
F—Accelerating electrode
G—Deflection-plate pair
H—Deflection-plate pair
J—Conductive coating connected to accelerating electrode
K—Intensifier-electrode terminal
L—Intensifier electrode (conductive coating on glass)
M—Fluorescent screen

Fig. 52. Electrode arrangement of typical electrostatic-deflection cathode-ray tube.

occurs. A typical mesh-type PDA CRT is shown in Fig. 53.

Magnetic Deflection—Magnetic deflection (Fig. 54) is generally used with TV, computer-display, and radar CRTs, where high resolution and brightness are required. Magnetic coils are placed in pairs on the outside of the CRT to provide horizontal and vertical magnetic fields perpendicular to the electron flow. Current in these coils causes deflection of the electrons perpendicular to the magnetic field and to the direction of the electrons.

Deflection bandwidth is limited by the high deflection power required of magnetic-deflection systems. Repetitive resonant circuits are normally used. Deflection is proportional to the flux or current in the coil and inversely proportional to the accelerating voltage.

Beam Focusing—Magnetic focusing can be achieved by placing an external magnetic coil in the form of a short solenoid on the outside of the CRT over the electron gun (Fig. 55). Because of the disadvantages of weight and the difficulty in aligning the coil to the beam, most modern CRTs use electrostatic focusing. Cylindrical electrodes at differing voltages form electron lenses that focus the beam to a fine spot at the phosphor screen.

Storage Cathode-Ray Tubes

The storage cathode-ray tube† produces a visual display of controllable duration. The tube has two electron guns, a phosphor viewing screen, and two or more fine-mesh screens. One of the electron guns is referred to as the writing gun and the other as the flood gun. The screen nearest the guns is the collector mesh. The other mesh is the storage mesh and is coated with

* K. R. Spangenberg, *Vacuum Tubes* (New York: McGraw-Hill Book Co., 1948).

† M. Knoll and B. Kazan, *Storage Tubes and Their Basic Principles* (New York: John Wiley & Sons, Inc., 1952).

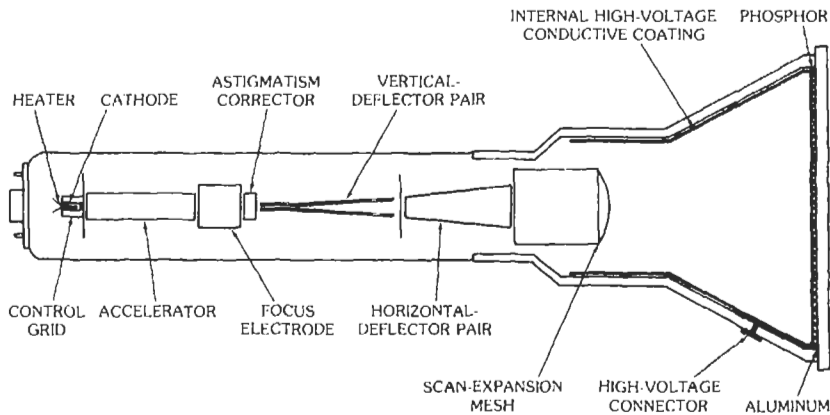


Fig. 53. CRT with mesh-type post-deflection acceleration.

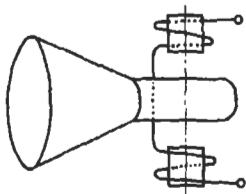


Fig. 54. Magnetic deflection.

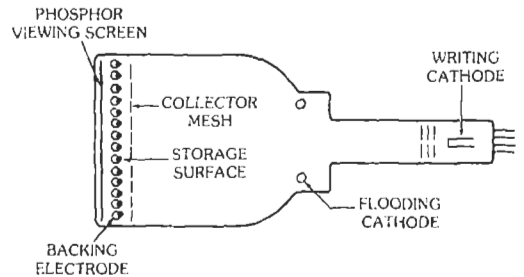


Fig. 56. Construction of storage cathode-ray tube.

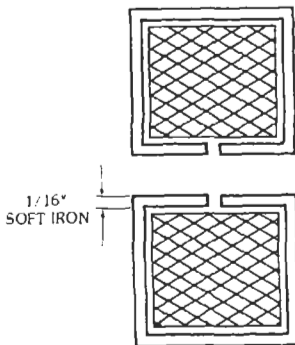


Fig. 55. Magnetic focusing.

a thin dielectric material to form a surface on which electrons store information. A typical storage tube is shown in Fig. 56.

The writing gun emits a pencil-like electron beam, which can be modulated by the information to be stored. This information is in the form of an electrical input signal that can be applied to the control grid for intensity modulation or to the deflection electrodes for spatial modulation. The storage surface is scanned by this high-resolution beam, which actually strikes the

surface. A positive-charge image corresponding to the input signal pattern is created on the storage surface by secondary-emission effects. The image remains on the storage surface until it decays from the neutralizing action of gas ions or is erased intentionally. The storage screen acts as an array of elemental electron guns, with each mesh hole acting as a control element of one of the guns. The desired information is stored on the storage mesh by the action of the writing gun, and the entire surface is flooded by electrons from the flood gun. The value of the positive charge deposited at each mesh aperture controls the amount of flood-gun current that can pass through the mesh aperture to the phosphor screen. The current that passes through the mesh strikes the phosphor, and light is emitted in proportion to the current density arriving and to the landing energy of the electrons. In other words, a gray scale is reproduced in the stored image. After the stored information has been observed or recorded, it is erased from the storage surface by fully writing the whole target first with high-energy flood electrons and then with low-energy electrons. The net negative charge so deposited on the storage target causes the target to be maintained at flood-gun cathode potential. The target is then prepared for storing a new image.

Increased storage sensitivity can be obtained by the addition of a third mesh coated with a dielectric material similar to that on the storage mesh and placed between the two meshes of the standard tube. This high-speed target is optimized to have very high sensitivity, but with consequently short retention (view time). The charge pattern written on this surface is quickly transferred to the front storage mesh, on which the coating has been optimized for long retention time and low sensitivity. Such CRTs are generally referred to as transfer storage tubes and are capable of storing signals with bandwidth in excess of 400 MHz.

Bistable Storage Tube

The bistable storage tube allows storage of charge in a CRT without the complication and expense of added meshes. Two guns are used as in the mesh storage tubes, but the phosphor screen is structured to have collector islands surrounded by a phosphor layer deposited on a transparent film (Fig. 57). The phosphor itself acts as both the dielectric storage surface and the light emitter. The potential of the phosphor can be maintained at two stable potentials by the action of secondary emission from the phosphor layer. To prepare the phosphor surface for storage, the collector electrode is pulsed with the waveform shown in Fig. 58. The flood electrons strike the whole target, and because the target potential is above first crossover, the target charges positively toward the collector potential. When the collector potential is suddenly dropped to 0 volts, the phosphor follows by capacitive coupling and is maintained at 0 volts while the collector potential slowly rises to its operational voltage. At 0 volts potential, few flood-gun electrons can land on the phosphor, and the light output is low. The target is now ready to store information. The writing-gun electrons have potential energy sufficiently high to charge the phosphor particles above first-crossover potential. Once the

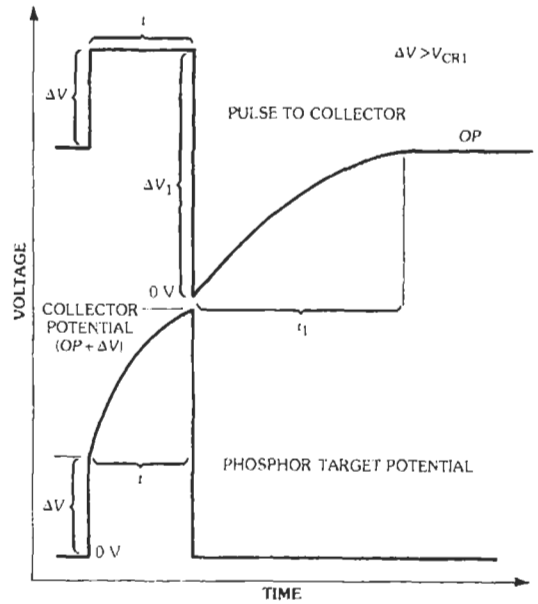


Fig. 58. Erasure waveforms.

particles are above this value, the flood-gun electrons can charge the written area to collector potential. The flood-gun electrons therefore can be made to maintain the phosphor surface at two stable potentials, flood-gun cathode and collector potential, and hence the term used is "bistable storage." At collector potential, the flood-gun electrons arrive with sufficient energy to excite the phosphor, and the stored image is clearly visible and remains as long as desired or until erased. These CRTs have enjoyed wide use in computer graphics terminals. Very high resolution can be obtained on large display tubes of 19 and 25 inches (diagonal). The 19-inch tube is capable of displaying

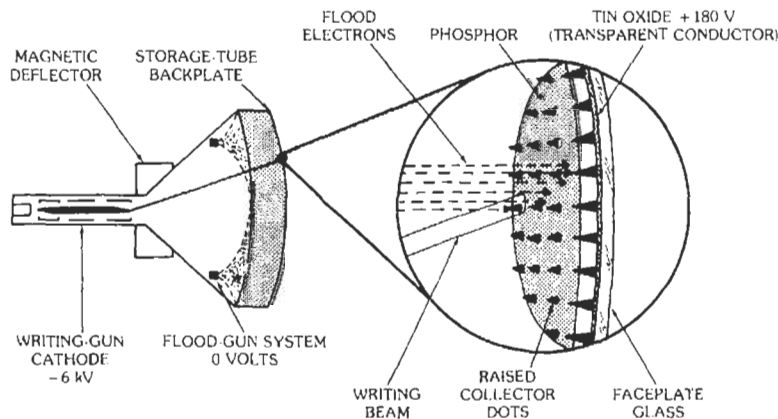


Fig 57. Typical bistable storage tube.

8500 characters, and the 25-inch CRT is capable of over 15 000 characters. Since these characters are actually stored on the face of the CRT until erased, no solid-state memory is required to refresh them. Another advantage of this type of CRT is that when a line is scanned with the writing gun at low beam current, the collector current will be a function of whether the area scanned is written or not. This signal current can be amplified and processed to produce an electronic image of the written areas. This image can be fed into a copy device for producing a hard copy of the display.

REFERENCES

1. A. S. Gilmour, Jr., *Microwave Tubes* (Artech House, Inc., Norwood, MA, 1986).
2. G. A. Haas, A. Shih, and R. E. Thomas, *Application of Surface Science* **2**, 293–321 (1979).
3. A. Shih and G. A. Haas, *Appl. Surf. Sci.* **2**, 164–172 (1979).
4. A. Shih and G. A. Haas, *Appl. Surf. Sci.* **2**, 275–283 (1979).
5. A. M. Shroff, "Review of dispenser cathodes," *Revue Technique Thomson-CSF* **23**, 947–1026 (1991).
6. C. R. K. Marrian, G. A. Haas, and A. Shih, *Appl. Surf. Sci.* **24**, 391 (1985).
7. M. C. Green, *Technical Report RADC-TR-81-211*, July 1981.
8. TriService/NASA Cathode Life Test Facility Annual Report, Naval Surface Warfare Center, Crane Division, Crane IN.
9. W. Mueller, *IEEE Transactions on Electron Devices* **36**, 180–187 (1989).
10. R. E. Thomas, J. W. Gibson, G. A. Haas, and R. H. Abrams, *IEEE Transactions on Electron Devices* **37**, 850–861 (1990).
11. J. Hasker, J. van Esdonk, and J. E. Crombeen, *Appl. Surf. Sci.* **26**, 173 (1986).
12. G. Gartner, P. Geittner, H. Lydtin, and A. Ritz, *Appl. Surf. Sci.* **111**, 11 (1997).
13. S. M. Sze, *Physics of Semiconductor Devices*, 2d ed. (Wiley, New York, 1981).
14. A. Modinos, *Field, Thermionic, and Secondary Electron Emission Spectroscopy* (Plenum, New York, 1984).
15. R. H. Fowler and L. W. Nordheim, *Proc. R. Soc. London Ser. A* **119**, 173 (1928).
16. E. L. Murphy and R. H. Good, *Phys. Rev.* **102**, 1464 (1956).
17. C. A. Spindt, I. Brodie, L. Humphrey, E. R. Westerber, *J. Appl. Phys.* **47**, 5248 (1976).
18. C. Bandis and B. B. Pate, *Appl. Phys. Lett* **69**, 366 (1996).
19. R. Stratton, *Phys. Rev.* **135**, A794 (1964).
20. C. A. Spindt, I. Brodie, L. Humphrey, E. R. Westerber, *J. Appl. Phys.* **47**, 5248 (1976).
21. S. T. Purcell, V. T. Binh, R. Baptist, "Nanoprotrusion model for field emission from integrated microtips," *J. Vac. Sci. Technol.* **B15**, 1666 (1997).
22. W. A. Mackie, T. Xie, and P. R. Davis, "Transition metal carbide field emitters for field-emitter array devices and high current applications," *J. Vac. Sci. Technol.* **B17**, 613 (1999).
23. W. Zhu, C. Bower, O. Zhou, G. Kochanski, and S. Jin, *Appl. Phys. Lett.* **75**, 875 (1999).
24. J. Robertson, "Field emission from carbon systems," *Mat. Res. Soc. Symp. Proc.* Vol. 621 (Materials Research Society, Warrendale, PA, 2000), R1.1.
25. D. R. Whaley, B. M. Gannon, C. R. Smith, C. M. Armstrong, and C. A. Spindt, "Application of field emitter arrays to microwave power amplifiers," *IEEE Trans. Plas. Sci.* **28**, 727 (2000).
26. K. L. Jensen, R. H. Abrams, and R. K. Parker, "Field emitter array development for high frequency applications," *J. Vac. Sci. Technol.* **B16**, 749 (1998).
27. C. M. Marrese, J. E. Polk, K. L. Jensen, A. D. Galimore, C. A. Spindt, R. L. Fink, and W. D. Palmer, "Performance of field emission cathodes in xenon electric propulsion system environments," chapter 11 in *Micropropulsion for Small Spacecraft* (Vol. 187 of Progress in Astronautics and Aeronautics), M. M. Micci, A. D. Ketsdever (eds.) (American Institute of Aeronautics and Astronautics, Reston, VA, 2000).
28. A. Ghis, R. Meyers, P. Rambaud, F. Levy, and T. Leroux, "Sealed vacuum devices—Fluorescent microtip displays," *IEEE Trans. Electron. Devices* **38**, 2320 (1991).
29. A. A. Talin, K. A. Dean, and J. E. Jaskie, "Field emission displays: A critical review," to appear in *Vacuum Microelectronics: A Special Issue of Solid State Electronics* (2001).
30. D. W. Jenkins, "Emission area of a field emitter array," *IEEE Trans. Electron. Devices* **40**, 666 (1993).
31. R. L. Hartman, W. A. Mackie, and P. R. Davis, "Use of boundary element methods in field emission computations," *J. Vac. Sci. Technol.* **B12**, 754 (1994).
32. K. L. Jensen, "Field emitter arrays for plasma and microwave source applications," *Phys. Plasmas* **6**, 2241 (1999).
33. C. Constancias and R. Baptist, "Emission observation of a microtip cathode array with an electrostatic-lens projector: Statistical approach," *J. Vac. Sci. Technol.* **B16**, 841 (1998).
34. R. Schlessler, M. T. McClure, B. L. McCarrson, and Z. Sitar, "Bias voltage dependent field-emission energy distribution analysis of wide band-gap field emitters," *J. Appl. Phys.* **82**, 5763 (1997).
35. C. A. Spindt, C. E. Holland, P. R. Schwoebel, and I. Brodie, "Field-emitter-array development for microwave applications," *J. Vac. Sci. Technol.* **B14**, 1986 (1996).

36. K. L. Jensen, R. H. Abrams, and R. K. Parker, "Field emitter array development for high frequency applications," *J. Vac. Sci. Technol.* **B16**, 749 (1998).
37. K. L. Jensen, "Field emitter arrays for plasma and microwave source applications," *Phys. Plasmas* **6**, 2241 (1999).
38. M. A. Kodis, K. L. Jensen, E. G. Zaidman, B. Goplen, and D. N. Smithe, "Optimization of field emission arrays for inductive output amplifiers," *J. Vac. Sci. Technol.* **B14**, 1990 (1996).
39. J. D. Levine, R. Meyer, R. Baptist, T. E. Fetler, and A. A. Talin, *J. Vac. Sci. Technol.* **B13**, 474 (1995).
40. J. A. Nation et al., *Proc. of the IEEE*, special issue "New Vista for Vacuum Electronics," vol. 87, no. 5, pp. 873-876, 1999.
41. M. Afif et al., *Applied Surface Science* **111**, 241 (1997).
42. W. C. Turner et al., "High-brightness high-current density cathode for induction LINAC FELs," *Lawrence Livermore Nat. Lab. Rep. UCRL-99042*, Sept. 1988.
43. Special issue of *Proc. of the IEEE*, vol. 87, 1999: R. M. Phillips and D. W. Sprehn, "High-Power Klystrons for the Next Linear Collider," *Proc. of the IEEE*, vol. 87, 1999; C. Smith, C. M. Armstrong and J. Duthic, "The Microwave Power Module—A Versatile RF Building Block for High Power Transmitters," *ibid.*; T. M. Antonsen, Jr., A. A. Mondelli, B. Levush, J. P. Verboncoeur, and C. K. Birdsall, "Advances in Modeling and Simulation of Vacuum Electric Devices," *ibid.*; J. A. Nation, L. Schachter, F. Mako, L. Len, W. Peter, C. M. Tang, and T. Srinivasan-Rao, "Advances in Non-Thermionic Cathode Physics and Technology," *ibid.*; and J. P. Calame and D. K. Abe, "Applications of Advanced Materials Technologies to Vacuum Electronics," *ibid.*
44. H. P. Freund and T. M. Antonsen, Jr., *Principles of Free-electron Lasers*, 2nd ed. (Chapman & Hall, London, 1996).
45. G. R. Neil et al., *Phys. Rev. Lett.* **84**, 662 (2000).
46. See NTIS Doc. No. DE98059292 (LCLS Design Group, "LCLS Design Report," April 1998). Copies may be ordered from the National Technical Information Service, Springfield, VA 22162.
47. B. Faatz et al., *Nucl. Instrum. Meth.* **A375**, 441 (1996).

17 Power Grid-Tube Circuits

Marvin Chodorow and Donald H. Preist

General Design 17-2

Graphic Design Methods 17-4

Class-C RF Amplifier or Oscillator

Class-B RF Amplifiers

Class-A and -AB AF Amplifiers

Class-AB and -B AF Amplifiers

Circuit Classification 17-8

RF Amplifier Circuits 17-8

Triodes

Tetrodes and Pentodes

UHF Operation

Klystrode™ Amplifiers

Circuits of Special Interest

Performance of Electron Power Tubes; Comparison With Solid-State Devices 17-13

Power Gain and Bandwidth of RF Amplifiers

Power Output Under Pulse Conditions

Life Expectancy of Tubes

Development Trends in Gridded Power Tubes

It is common practice to differentiate between types of vacuum-tube circuits, particularly amplifiers, on the basis of the operating regime of the tube.

Class-A: Grid bias and alternating grid voltages such that plate current flows continuously throughout electrical cycle ($\theta_p = 360^\circ$).

Class-AB: Grid bias and alternating grid voltages such that plate current flows appreciably more than half but less than entire electrical cycle ($360^\circ > \theta_p > 180^\circ$).

Class-B: Grid bias close to cutoff such that plate current flows only during approximately half of electrical cycle ($\theta_p \approx 180^\circ$).

Class-C: Grid bias appreciably greater than cutoff so that plate current flows for appreciably less than half of electrical cycle ($\theta_p < 180^\circ$).

A further classification between circuits in which positive grid current is conducted during some portion of the cycle and those in which it is not is denoted by subscripts 2 and 1, respectively. Thus a class-AB₂ amplifier operates with a positive swing of the alternating grid voltage such that positive electronic current is conducted and accordingly in-phase power is required to drive the tube.

GENERAL DESIGN

For quickly estimating the performance of a tube from catalog data, or for predicting the characteristics needed for a given application, the ratios given below may be used.

Table 1 gives correlating data for typical operation of tubes in the various amplifier classifications. If the maximum ratings of a tube are known, the maximum power output, currents, voltages, and corresponding load impedance may be estimated from the table. Take for example a type F-124-A water-cooled transmitting tube operated as a class-C radio-frequency power amplifier and oscillator (the constant-current characteristics are shown in Fig. 1). Published maximum ratings are as follows.

Dc plate voltage:

$$E_b = 20\,000 \text{ volts}$$

Dc grid voltage:

$$E_c = 3000 \text{ volts}$$

Dc plate current:

$$I_b = 7 \text{ amperes}$$

RF grid current:

$$I_g = 50 \text{ amperes}$$

Plate input:

$$P_i = 135\,000 \text{ watts}$$

Plate dissipation:

$$P_p = 40\,000 \text{ watts}$$

Maximum conditions may be estimated as follows. For $\eta = 75$ percent

$$P_i = 135\,000 \text{ watts}$$

$$E_b = 20\,000 \text{ volts}$$

Power output $P_o = \eta P_i = 100\,000$ watts.

Average dc plate current $I_b = P_i/E_b = 6.7$ amperes

From a tabulated typical ratio $^M I_b/I_b = 4$, the instantaneous peak plate current $^M I_b = 4I_b = 27$ amperes.*

The rms plate alternating-current component, taking the ratio $I_p/I_b = 1.2$, is

$$I_p = 1.2I_b = 8 \text{ amperes}$$

The rms value of the plate alternating-voltage component from the ratio $E_p/E_b = 0.6$ is $E_p = 0.6 E_b = 12\,000$ volts.

* In this discussion, the superscript M indicates the use of the maximum or peak value of the varying component; i.e., $^M I_b$ = maximum or peak value of the alternating component of the plate current.

TABLE 1. TYPICAL AMPLIFIER OPERATING DATA (MAXIMUM-SIGNAL CONDITIONS, PER TUBE)

Function	Class A	Class B AF (p-p)	Class B RF	Class C RF
Plate efficiency η (percent)	20-30	35-65	60-70	65-85
Peak instantaneous to dc plate-current ratio $^M I_b/I_b$	1.5-2	3.1	3.1	3.1-4.5
Rms alternating to dc plate-current ratio I_p/I_b	0.5-0.7	1.1	1.1	1.1-1.2
Rms alternating to dc plate-voltage ratio E_p/E_b	0.3-0.5	0.5-0.6	0.5-0.6	0.5-0.6
Dc to peak instantaneous grid-current $I_c/^M I_c$		0.1-0.25	0.1-0.25	0.1-0.25

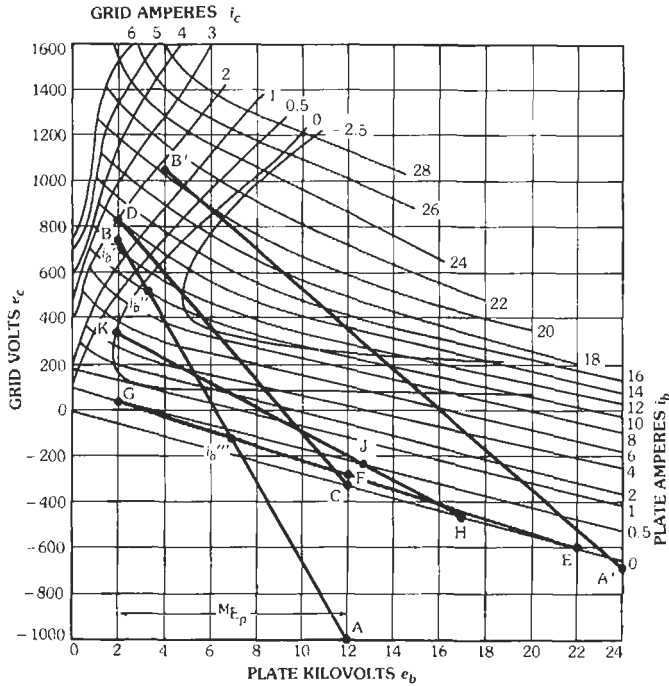


Fig. 1. Constant-current characteristics of type F-124-A tube with typical load lines: AB—class C, CD—class B, EFG—class A, HJK—class AB.

The approximate operating load resistance, R_l , is now found from

$$R_l = E_p / I_p = 1500 \text{ ohms}$$

An estimate of the grid drive power required may be obtained by reference to the constant-current characteristics of the tube and determination of the peak instantaneous positive grid current $M i_c$ and the corresponding instantaneous total grid voltage $M e_c$. If the value of grid bias for the given operating condition is E_c , the peak alternating grid drive voltage is

$$M E_g = (M e_c - E_c)$$

from which the peak instantaneous grid drive power can be determined:

$$M P_c = M E_g M i_c$$

An approximation to the average grid drive power, P_g , necessarily rough due to neglect of negative grid current, is obtained from the typical ratio of dc to peak value of grid current, $I_c / M i_c = 0.2$. The result is

$$P_g = I_c E_g = 0.2 M i_c E_g \text{ watt}$$

Plate dissipation P_p may be checked with published values since

$$P_p = P_i - P_o$$

It should be borne in mind that combinations of published maximum ratings as well as each individual maximum rating must be observed. Thus, for example in this case, the maximum dc plate operating voltage of 20 000 volts does not permit operation at the maximum dc plate current of 7 amperes since this exceeds the maximum plate input rating of 135 000 watts.

Plate load resistance R_l may be connected directly in the tube plate circuit as in the resistance-coupled amplifier, through impedance-matching elements as in audio-frequency transformer coupling, or effectively represented by a loaded parallel-resonant circuit as in most radio-frequency amplifiers. In any case, calculated values apply only to effectively resistive loads, such as are normally closely approximated in radio-frequency amplifiers. With appreciably reactive loads, operating currents and voltages will in general be quite different, and their precise calculation is quite difficult.

The physical load resistance present in any given setup may be measured by audio-frequency or radio-frequency bridge methods. In many cases, the proper value of R_l is ascertained experimentally as in radio-frequency amplifiers that are tuned to the proper minimum dc plate current. Conversely, if the circuit is to be matched to the tube, R_l is determined directly as in a resistance-coupled amplifier or as

$$R_l = N^2 R_o$$

in the case of a transformer-coupled stage, where N is the primary-to-secondary voltage transformation ratio.

In a parallel-resonant circuit in which the output resistance R_s is connected directly in one of the reactance legs

$$R_l = X^2/R_s = L/Cr_s = QX$$

where,

- X is the leg reactance at resonance (ohms),
- L and C are leg inductance in henrys and capacitance in farads, respectively,
- $Q = X/R_s$.

GRAPHIC DESIGN METHODS

When accurate operating data are required, more precise methods must be used. Because of the nonlinear nature of tube characteristics, graphic methods usually are most convenient and rapid. Examples of such methods are given below.

A comparison of the operating regimes of class A, AB, B, and C amplifiers is given in the constant-current-characteristics graph of Fig. 1. The lines corresponding to the different classes of operation are the loci of instantaneous grid voltage e_c and plate voltage e_b , corresponding to their respective load impedances.

For radio-frequency amplifiers and oscillators having tuned circuits that give an effectively resistive load, plate and grid tube and load alternating voltages are sinusoidal and in phase (disregarding transit time), and the loci become straight lines.

For amplifiers having nonresonant resistive loads, the loci are in general nonlinear except in the distortionless case of linear tube characteristics (constant r_p), for which they are again straight lines.

Thus, for determination of radio-frequency performance, the constant-current chart is convenient. For solution of audio-frequency problems, however, it is more convenient to use the i_b - e_c transfer characteristics of Fig. 2, on which a dynamic load line may be constructed.

Methods for calculation of the most important cases are given below.

Class-C RF Amplifier or Oscillator

Draw a straight line from A to B (Fig. 1) corresponding to the chosen dc operating plate and grid voltages, and to the desired peak alternating plate and grid voltage excursions. The projection of AB on the horizontal axis thus corresponds to ${}^M E_p$. Using Chaffee's 11-point method of harmonic analysis, lay out on AB points

$$e_p' = {}^M E_p$$

$$e_p'' = 0.866 {}^M E_p$$

$$e_p''' = 0.5 {}^M E_p$$

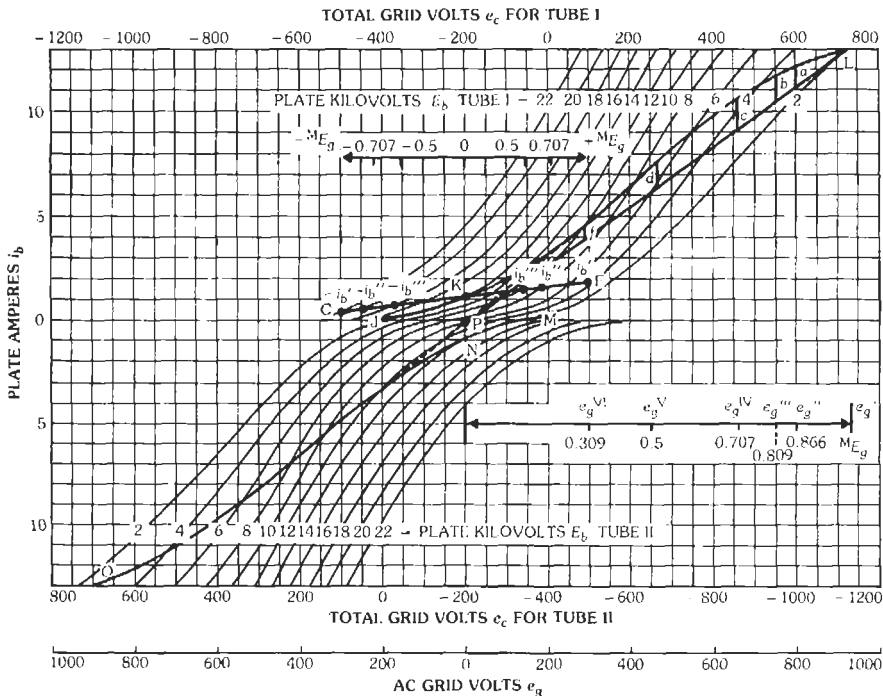


Fig. 2. Transfer characteristics i_b versus e_c with load lines: CKF—class A₂ and OPL—class B.

to each of which correspond instantaneous plate currents, i_b' , i_b'' , and i_b''' and instantaneous grid currents i_c' , i_c'' , and i_c''' . The operating currents are obtained from

$$I_b = [i_b' + 2i_b'' + 2i_b''']/12$$

$$I_c = [i_c' + 2i_c'' + 2i_c''']/12$$

$$^M I_p = [i_b' + 1.73i_b'' + i_b''']/6$$

$$^M I_g = [i_c' + 1.73i_c'' + i_c''']/6$$

Substitution of the above in the following gives the desired operating data.

$$\text{Power output } P_o = (^M E_p ^M I_p)/2$$

$$\text{Power input } P_i = E_b I_b$$

$$\text{Average grid excitation power } P_g = (^M E_g ^M I_g)/2$$

$$\text{Peak grid excitation power } ^M P_c = ^M E_g i_c'$$

$$\text{Plate load resistance } R_l = ^M E_p / ^M I_p$$

$$\text{Grid bias resistance } R_c = E_c / I_c$$

$$\text{Plate efficiency } \eta = P_o / P_i$$

$$\text{Plate dissipation } P_p = P_i - P_o$$

The above procedure may also be applied to plate-modulated class-C amplifiers. Take the above data as applying to carrier conditions, and repeat the analysis for $E_b^{\text{crest}} = 2E_b$ and $P_o^{\text{crest}} = 4P_o$, keeping R_l constant. After a cut-and-try method has given a peak solution, it will often be found that combination fixed and self grid biasing as well as grid modulation are indicated to obtain linear operation.

To illustrate the preceding exposition, a typical amplifier calculation is given below.

Operating requirements (carrier condition):

$$E_b = 12\,000 \text{ volts}$$

$$P_o = 25\,000 \text{ watts}$$

$$\eta = 75 \text{ percent}$$

Preliminary calculation (refer to Tables 1 and 2):

$$E_p/E_b = 0.6$$

$$E_p = 0.6 \times 12\,000 = 7200 \text{ volts}$$

$$^M E_p = 1.41 \times 7200 = 10\,000 \text{ volts}$$

$$I_p = P_o/E_p = 25\,000/7200 = 3.48 \text{ amperes}$$

$$^M I_p = 4.9 \text{ amperes}$$

$$I_p/I_b = 1.2$$

$$I_b = 3.48/1.2 = 2.9 \text{ amperes}$$

$$P_i = 12\,000 \times 2.9 = 35\,000 \text{ watts}$$

$$^M i_b/I_b = 4.5$$

$$^M i_b = 4.5 \times 2.9 = 13.0 \text{ amperes}$$

$$R_l = E_p/I_p = 7200/3.48 = 2060 \text{ ohms}$$

Complete Calculation: Lay out carrier operating line AB on the constant-current graph, Fig. 1, using values of E_b , $^M E_p$, and $^M i_b$ from the preliminary calculated data. Operating carrier bias voltage E_c is chosen somewhat greater than twice the cutoff value (1000 volts) to locate point A.

The following data are taken along AB.

$$i_b' = 13 \text{ amperes}$$

$$i_b'' = 10 \text{ amperes}$$

TABLE 2. CLASS-C RF AMPLIFIER DATA FOR 100-PERCENT PLATE MODULATION

Symbol	Preliminary Carrier	Detailed	
		Carrier	Crest
E_b (volts)	12 000	12 000	24 000
$^M E_p$ (volts)	10 000	10 000	20 000
E_c (volts)	—	-1 000	-700
$^M E_g$ (volts)	—	1 740	1 740
I_b (amperes)	2.9	2.8	6.4
$^M I_p$ (amperes)	4.9	5.1	10.2
I_c (amperes)	—	0.125	0.083
$^M I_g$ (amperes)	—	0.255	0.183
P_i (watts)	35 000	33 600	154 000
P_o (watts)	25 000	25 500	102 000
P_g (watts)	—	220	160
η (percent)	75	76	66
R_l (ohms)	2 060	1 960	1 960
R_c (ohms)	—	7 100	7 100
E_{cc} (volts)	—	-110	-110

$$i_b''' = 0.3 \text{ ampere}$$

$$i_c' = 1.7 \text{ amperes}$$

$$i_c'' = -0.1 \text{ ampere}$$

$$i_c''' = 0 \text{ ampere}$$

$$E_c = -1000 \text{ volts}$$

$$e_c' = 740 \text{ volts}$$

$${}^M E_p = 10\,000 \text{ volts}$$

From the equations, complete carrier data as follows are calculated.

$${}^M I_p = [13 + 1.73 \times 10 + 0.3]/6 = 5.1 \text{ amperes}$$

$$P_o = (10\,000 \times 5.1)/2 = 25\,500 \text{ watts}$$

$$I_b = [13 + 2 \times 10 + 2 \times 0.3]/12 = 2.8 \text{ amperes}$$

$$P_i = 12\,000 \times 2.8 = 33\,600 \text{ watts}$$

$$\eta = (25\,500/33\,600) \times 100 = 76 \text{ percent}$$

$$R_l = (10\,000/5.1) = 1960 \text{ ohms}$$

$$I_c = [1.7 + 2(-0.1)]/12 = 0.125 \text{ ampere}$$

$${}^M I_g = [1.7 + 1.7(-0.1)]/6 = 0.255 \text{ ampere}$$

$$P_g = (1740 \times 0.255)/2 = 220 \text{ watts}$$

Operating data at 100-percent positive modulation crests are now calculated based on the fact that here

$$E_b = 24\,000 \text{ volts}$$

$$R_l = 1960 \text{ ohms}$$

and for undistorted operation

$$P_o = 4 \times 25\,500 = 102\,000 \text{ watts}$$

$${}^M E_p = 20\,000 \text{ volts}$$

The crest operating line A'B' is now located by trial so as to satisfy the above conditions, by the use of the same equations and method as for the carrier condition.

It is seen that to obtain full-crest power output, in addition to doubling the alternating plate voltage, the peak plate current must be increased. This is accomplished by reducing the crest bias voltage with a resultant increase of current conduction period but lower plate efficiency.

The effect of grid secondary emission to lower the crest grid current is taken advantage of to obtain the reduced grid-resistance voltage drop required. By use of

combination fixed and grid-resistance bias, proper variation of the total bias is obtained. The value of grid resistance required is given by

$$R_c = -(E_c - {}^{\text{crest}} E_c)/(I_c - {}^{\text{crest}} I_c)$$

and the value of fixed bias by

$$E_{cc} = E_c - (I_c R_c)$$

Calculations at carrier and positive crest together with the condition of zero output at negative crest give sufficiently complete data for most purposes. If accurate calculation of audio-frequency harmonic distortion is necessary, the above method may be applied to the additional points required.

Class-B RF Amplifiers

A rapid approximate method is to determine by inspection from the tube i_b - e_b characteristics the instantaneous current, i_b' , and voltage, e_b' , corresponding to the peak alternating voltage swing from operating voltage E_b .

Ac plate current:

$${}^M I_p = i_b'/2$$

Dc plate current:

$$I_b = i_b'/\pi$$

Ac plate voltage:

$${}^M E_p = E_b - e_b'$$

Power output:

$$P_o = [(E_b - e_b')i_b']/4$$

Power input:

$$P_i = E_b i_b'/\pi$$

Plate efficiency:

$$\eta = (\pi/4)[1 - (e_b'/E_b)]$$

Thus $\eta \approx 0.6$ for the usual crest value of ${}^M E_p \approx 0.8 E_b$.

The same method of analysis used for the class-C amplifier may also be used in this case. The carrier and crest-condition calculations, however, are now made from the same E_b , the carrier condition corresponding to an alternating-voltage amplitude of ${}^M E_p/2$ such as to give the desired carrier power output.

For greater accuracy than the simple check of carrier and crest conditions, the radio-frequency plate currents ${}^M I_p'$, ${}^M I_p''$, ${}^M I_p'''$, ${}^M I_p^0$, $-{}^M I_p''$, $-{}^M I_p'$, and $-{}^M I_p'$ may be calculated for seven corresponding selected

points of the audio-frequency modulation envelope $+^M E_g$, $+0.707^M E_g$, $+0.5^M E_g$, 0 , $-0.5^M E_g$, $-0.707^M E_g$, and $-^M E_g$, where the negative signs denote values in the negative half of the modulation cycle. If the designations

$$S' = ^M I_p' - (-^M I_p')$$

$$D' = ^M I_p' + (-^M I_p') - 2^M I_p^0$$

are used, the fundamental and harmonic components of the output audio-frequency current are obtained as

$$^M I_{p1} = (S'/4) + [S''/2(2)^{1/2}] \text{ (fundamental)}$$

$$^M I_{p2} = (5D'/24) + (D''/4) - (D'''/3)$$

$$^M I_{p3} = (S'/6) - (S''/3)$$

$$^M I_{p4} = (D'/8) - (D''/4)$$

$$^M I_{p5} = (S'/12) - [S''/2(2)^{1/2}] + (S'''/3)$$

$$^M I_{p6} = (D'/24) - (D''/4) + (D'''/3)$$

This detailed method of calculation of audio-frequency harmonic distortion may, of course, also be applied to calculation of the class-C modulated amplifier, as well as to the class-A modulated amplifier.

Class-A and -AB AF Amplifiers

Approximate equations assuming linear tube characteristics:

Maximum undistorted power output

$$^M P_o = (^M E_p ^M I_p)/2$$

when plate load resistance

$$R_l = r_p \left[\frac{E_c}{(^M E_p/\mu) - E_c} - 1 \right]$$

and negative grid bias

$$E_c = (^M E_p/\mu)[(R_l + r_p)/(R_l + r_p)/(R_l + 2r_p)]$$

giving maximum plate efficiency

$$\eta = ^M E_p ^M I_p / 8 E_b I_b$$

Maximum maximum undistorted power output

$$^{MM} P_o = ^M E_p^2 / 16 r_p$$

when

$$R_l = 2r_p$$

$$E_c = \frac{3}{4} (^M E_p/\mu)$$

An exact analysis may be obtained by use of a dynamic load line laid out on the transfer characteristics of the tube. Such a line is CKF of Fig. 2, which is constructed about operating point K for a given load resistance r_l from

$$i_b^S = [(e_b^R - e_b^S)/R_l] + i_b^R$$

where R, S, etc., are successive conveniently spaced construction points.

Using the seven-point method of harmonic analysis, plot instantaneous plate currents $i_b^1, i_b^2, i_b^3, i_b^4, i_b^5, i_b^6, i_b^7$, and $-i_b^1, -i_b^2, -i_b^3, -i_b^4, -i_b^5, -i_b^6, -i_b^7$, corresponding to $+^M E_g, +0.707^M E_g, +0.5^M E_g, 0, -0.5^M E_g, -0.707^M E_g, \text{ and } -^M E_g$, where 0 corresponds to operating point K. In addition to the equations given under class-B radio-frequency amplifiers

$$I_b \text{ average} = I_b + (D'/8) + (D''/4)$$

from which complete data may be calculated.

Class-AB and -B AF Amplifiers

Approximate equations assuming linear tube characteristics give (referring to Fig. 1, line CD) for a class-B audio-frequency amplifier

$$^M I_p = i_b'$$

$$P_o = ^M E_p ^M I_p / 2$$

$$P_i = (2/\pi) E_b ^M I_p$$

$$\eta = (\pi/4) (^M E_p / E_b)$$

$$R_{pp} = 4(^M E_p / i_b') = 4R_l$$

An exact solution may be derived by use of dynamic load line JKL on the i_b-e_c characteristic of Fig. 2. This line is calculated about the operating point K for the given R_l (in the same way as for the class-A case). However, since two tubes operate in phase opposition in this case, an identical dynamic load line MNO represents the other half cycle, laid out about the operating bias abscissa point but in the opposite direction (see Fig. 2).

Algebraic addition of instantaneous current values of the two tubes at each value of e_c gives the composite dynamic characteristic OPL for the two tubes. Inasmuch as this curve is symmetrical about point P, it may be analyzed for harmonics along a single half-curve PL by use of the Mourontseff 5-point method. A straight line is drawn from P to L, and ordinate plate-current differences a, b, c, d, f between this line and the curve, corresponding to $e_g^2, e_g^3, e_g^{1V}, e_g^V, \text{ and } e_g^{VI}$, are measured. Ordinate distances measured upward from curve PL are taken positive.

Fundamental and harmonic current amplitudes and power are found from

$$M_{I_{p1}} = i_b' - M_{I_{p3}} + M_{I_{p5}} - M_{I_{p7}} + M_{I_{p9}} - M_{I_{p11}}$$

$$M_{I_{p3}} = 0.4475(b + f) + (d/3) - 0.578d - \frac{1}{2} M_{I_{p5}}$$

$$M_{I_{p5}} = 0.4(a - f)$$

$$M_{I_{p7}} = 0.4475(b + f) - M_{I_{p3}} + 0.5 M_{I_{p5}}$$

$$M_{I_{p9}} = M_{I_{p3}} - \frac{2}{3}d$$

$$M_{I_{p11}} = 0.707c - M_{I_{p3}} + M_{I_{p5}}$$

Even harmonics are not present due to the symmetry of the dynamic characteristic. The direct-current and power-input values are found by the 7-point analysis from curve PL and doubled for two tubes.

CIRCUIT CLASSIFICATION

The classification of amplifiers in classes A, B, and C is based on the operating conditions of the tube. Another classification can be used, based on the type of circuits associated with the tube.

A tube can be considered as a four-terminal network with two input terminals and two output terminals. One of the input terminals and one of the output terminals are usually common; this common junction or point is usually called "ground."

When the common point is connected to the filament or cathode of the tube, we can speak of a grounded-cathode circuit (the most conventional type of vacuum-tube circuit). When the common point is the grid, we can speak of a grounded-grid circuit; and when the common point is the plate or anode, we can speak of a grounded-anode circuit. This last type of circuit is most commonly known by the name "cathode-follower."

A fourth and most general class of circuit is obtained when the common point or ground is not directly connected to any of the three electrodes of the tube. This is the condition encountered at UHF where the series impedances of the internal tube leads make it impossible to ground any of them. It is also encountered in such special types of circuits as the phase-splitter, in which the impedance from plate to ground and the impedance from cathode to ground are made equal to obtain an output between plate and cathode balanced with respect to ground.

Design information for the first three classifications is given in Table 3, where

Z_2 = load impedance to which output terminals of amplifier are connected,

E_1 = phasor input voltage to amplifier,

E_2 = phasor output voltage across load impedance Z_2 ,

A = voltage gain of amplifier = E_2/E_1 ,

Y_1 = input admittance to input terminals of amplifier,

$$\omega = 2\pi \times (\text{frequency of excitation voltage } E_1),$$

$$j = (-1)^{1/2}.$$

RF AMPLIFIER CIRCUITS

The power grid tube requires external circuits. Examples are shown schematically below for each kind of tube.

Triodes

The triode has three electrodes: the thermionic cathode, which emits electrons; the control grid; and the anode, which collects most of the electrons. If the grid is "biased" to a sufficiently high negative potential (cutoff bias), no current flows. As the grid potential becomes less negative, more current flows to the anode. When the grid becomes positive with respect to the cathode, both grid and anode draw current. At some value of positive grid potential, the total space current starts to exceed the emitting capability of the cathode (cathode saturation) or the product of the grid current and grid-cathode voltage (grid dissipation) exceeds the limit above which the grid will emit electrons (primary emission). Excessive grid dissipation interferes with the desired operation of the tube or results in mechanical distortion due to excessive temperature. There is also a limit to the power dissipation of the anode, depending on the cooling method used.

In operation as an RF power amplifier, the triode must be either "neutralized" as in Fig. 3A or operated "grounded-grid" ("cathode-driven") as in Fig. 3B; otherwise, the internal capacitance between grid and anode produces positive feedback that may cause self-oscillation at a frequency close to the operating frequency. The triode may be operated as an efficient oscillator by optimizing the feedback through the addition of extra capacitance or by other means. Oscillators are used for RF heating of materials in industrial operations where precise control of frequency is not required.

Small planar triodes are used at UHF and microwave frequencies up to about 4 GHz, especially in pulse service where peak-to-average power ratios are 100–1000. The anode supply voltage is typically 1 kV or more.

Small cylindrical triodes are used mainly at VHF and UHF where CW power of a few hundred watts or pulse powers of tens of kilowatts are required. Modern triodes are designed with beam-forming cathode and control-grid geometry to allow the simplicity of design and circuit advantages of a triode with the gain of a tetrode.

Tetrodes and Pentodes

The tetrode and pentode have four or five electrodes, respectively. A tetrode has a cathode, a control grid, a screen grid, and an anode. The screen grid greatly reduces the capacitance between the anode and control

TABLE 3. DESIGN INFORMATION FOR THREE CLASSES OF AMPLIFIERS

Grounded-Cathode	Grounded-Grid	Grounded-Plate or Cathode-Follower
Circuit Schematic		
Equivalent Circuit, Alternating-Current Component, Class-A Operation		
Voltage Gain A for Output Load Impedance Z2; A = E2/E1		
$A = -\mu Z_2 / (r_p + Z_2)$ $= -g_m [r_p Z_2 / (r_p + Z_2)]$ Neglecting C_{gp} (Z_2 includes C_{pk})	$A = (1 + \mu) [Z_2 / (r_p + Z_2)]$ Neglecting C_{pk} (Z_2 includes C_{gp})	$A = \mu Z_2 / [r_p + (1 + \mu) Z_2]$ Neglecting C_{gk} (Z_2 includes C_{pk})
Input Admittance; $Y_1 = I_1/E_1$		
$Y_1 = j\omega [C_{gk} + (1 - A)C_{gp}]$	$Y_1 = j\omega [C_{gk} + (1 - A)C_{pk}]$ $+ [(1 + \mu) / (r_p + Z_2)]$	$Y_1 = j\omega [C_{gp} + (1 - A)C_{gk}]$
Equivalent Generator Seen by Load at Output Terminals		
Neglecting C_{gp}	Neglecting C_{pk}	Neglecting C_{gk}

grid and makes neutralization unnecessary or easy to accomplish. The pentode has an additional "suppressor grid" to control secondary electrons. Modern tetrodes accomplish this control in other ways and have almost completely displaced pentodes.

Two tetrode circuits for RF amplification are shown in Fig. 4. The arrangement in Fig. 4A is commonly used at frequencies below about 30 MHz for very large tubes and below about 400 MHz for small tubes. A power gain of 1000 (30 dB) is typical at lower frequencies.

At VHF and UHF, the cathode-lead inductance, L_k , of Fig. 4A produces excessive negative feedback, which

can be overcome by incorporating the inductance into the input resonant circuit as shown in Fig. 4B. In this arrangement, both the control grid and screen grid are maintained at RF ground potential by the bypassing capacitors (C_b). In this circuit, the power gain is reduced compared to that of Fig. 4A, but 10-15 dB is typical at frequencies below about 1000 MHz.

Tetrode tubes are available in various sizes for different power levels and frequencies. Applications are radio broadcasting (AM and FM), television (VHF and UHF), communications, radar, navigational aids, and high-energy physics, including particle accelerators and thermonuclear-fusion machines. Functions are RF power

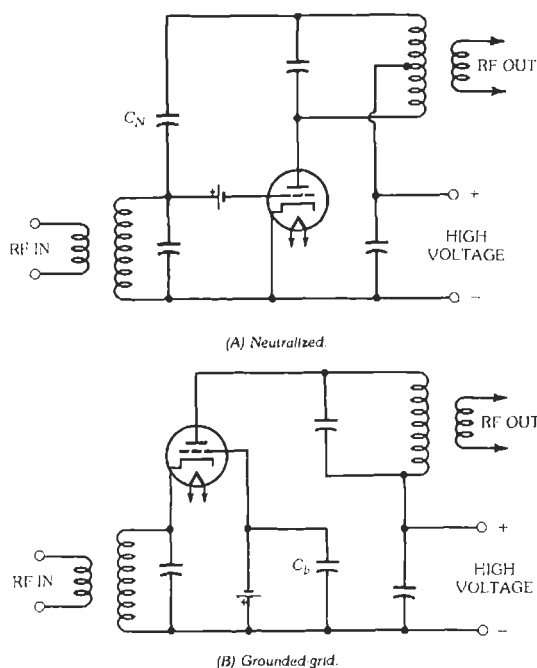


Fig. 3. Triode amplifier circuits.

generation, modulation (AM), and switching for pulse service.

UHF Operation

When the transit time of the electrons from cathode to anode is an appreciable fraction of one radio-frequency cycle:

(A) Input conductance due to reaction of electrons with the varying field from the grid becomes appreciable. This conductance, which increases as the square of the frequency, results in lowered gain, an increase in driving-power requirement, and loading of the input circuit.

(B) Grid-anode transit time introduces a phase lag between grid voltage and anode current. In oscillators, the problem of compensating for the phase lag by design and adjustment of a feedback circuit becomes difficult. Efficiency is reduced in both oscillators and amplifiers.

(C) Distortion of the current pulse in the grid-anode space increases the anode-current conduction angle and lowers the efficiency.

In amplifiers, the effect of cathode-lead inductance is to introduce a conductance component in the grid circuit. This effect is serious because the loading of the input circuit by the conductance current limits the gain of the stage. Cathode-grid and grid-anode capacitive reactances are of small magnitude at ultrahigh frequencies, and heavy currents flow as a result of these reactances. Tubes must be designed to carry these currents without serious loss. Coaxial cavities are used

to resonate with the tube reactances and to minimize resistive and radiation losses.

Klystrone™ Amplifiers

At frequencies between about 100 and 1000 MHz, a hybrid tube, part triode, part klystron, has promise. Called a Klystrone™, it was invented in 1939 and demonstrated at low power levels, but lay dormant for over 40 years. It now appears capable of high power, somewhat greater conversion efficiency, and considerably more power gain than the conventional tetrode. It is physically shorter than a klystron because of the absence of a multicavity buncher. The Klystrone requires a resonant output cavity; in this respect it resembles the klystron. At frequencies much above 1000 MHz, the power output falls off rapidly with increasing frequency, because the device is basically a "density-modulated" or grid-controlled tube limited by electron transit time and grid limitations. A diagram of a Klystrone is shown in Fig. 5, and expected power versus frequency is shown in Fig. 9.

Circuits of Special Interest

High-Efficiency Circuits for RF Amplification—

Increasing energy costs have led to attempts to improve the efficiency of conversion from DC to RF, especially in very high-power equipment. Most of these attempts are based on wave shaping. It can be shown that the power lost as heat at the anode of a tetrode, for example, can be reduced as the waveforms of voltage and current approach square waves rather than sine waves or portions of sine waves. One method due to Tyler* still in use in AM broadcast transmitters is to add harmonics in the correct phase and amplitude to the original sine functions. The harmonics are generated in tuned circuits that are resonant at the harmonic frequencies and placed in series with the circuits tuned to the fundamental. The current and voltage waveforms approach square waves.

Another approach not involving extra circuits is simply to restrict the time duration of the current (class-C operation). This, however, reduces the power output or increases the cathode current density.

For a power amplifier at 1.5 MHz, typical efficiencies are: Class B, 75%; Class C, 80%–85%; "Tyler," 93%.

Circuits That Provide High Linearity—In some applications, for example TV visual channels, it is necessary to achieve RF power amplification with very low distortion of the modulation content, measured either as total harmonic distortion or, more often, as intermodulation distortion. Present practice is to use triodes as class-A amplifiers or tetrodes as class-B or class-AB linear amplifiers. In class-B service, grid bias

* Reference 1.

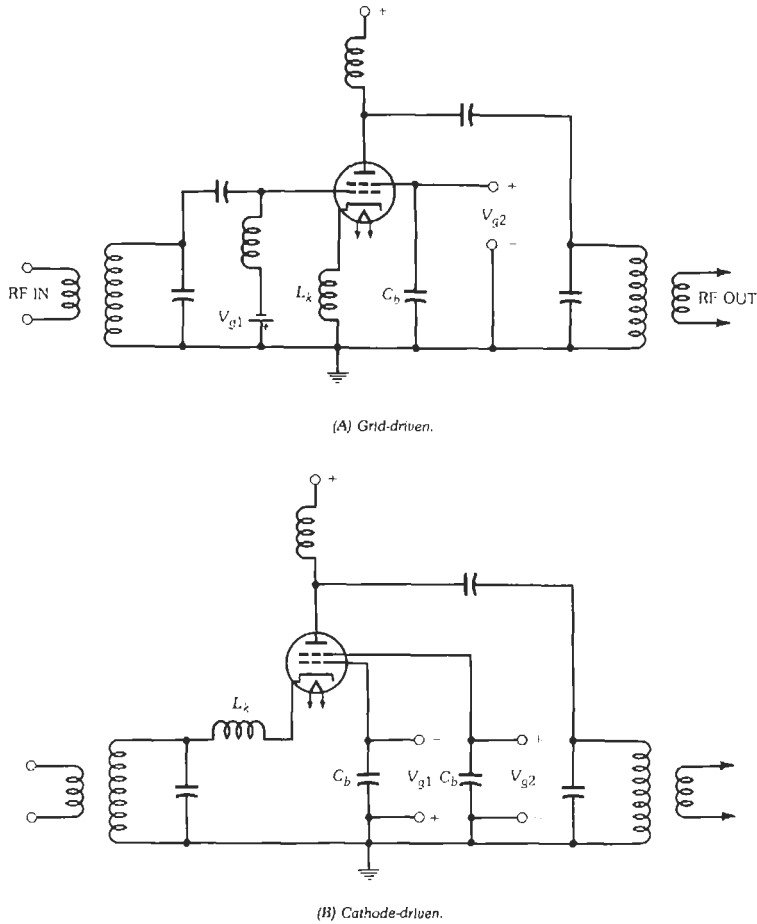


Fig. 4. Tetrode amplifier circuits.

is adjusted to the value that just reduces the anode current to zero. In class-AB, a small "idling" current is allowed to flow. By careful design of the tube and circuit combination, the modulated RF output power can be a faithful copy of the input over the bandwidth required for transmission of the information contained in the modulation. In this system, the average conversion efficiency over a period of time long compared to the modulation period or cycle can be reasonably high, because power input to the tube falls when the modulation level is low (the carrier is suppressed). In contrast, an amplifier operating in class A, such as a tetrode biased so that a sinusoidal anode current waveform is produced, would draw the same dc current whether the modulation were 100% or zero. The same is true of a klystron or TWT amplifier in which the beam current does not vary with the RF drive level.

In the linear amplifier with suppressed carrier, the efficiency varies as the square root of the modulation depth. For example, if with 100% modulation the efficiency is 50%, with 30% modulation the efficiency

would be $50\% \times \sqrt{0.3} = 27\%$. The corresponding klystron efficiency would be $50\% \times 0.3 = 15\%$.

Another tube-circuit combination of interest is the high-efficiency Doherty linear amplifier* shown schematically in Fig. 6. In this configuration, modulation is accomplished by causing two linear amplifiers, both driven by modulated RF grid-cathode voltages, to feed power into a common load impedance. The tubes are essentially in parallel. They are biased differently so that with no modulation the carrier power is supplied only by one (the "carrier" tube). The other ("peaking" tube) is cut off. At full modulation, both tubes deliver power. A quarter-wavelength transmission line is connected between the tubes. The impedance inverting properties of this line allow each tube to "see" a load impedance that varies with modulation depth. This arrangement is found in some AM broadcast transmitters in which the carrier and both sidebands are

* References 2 and 3.

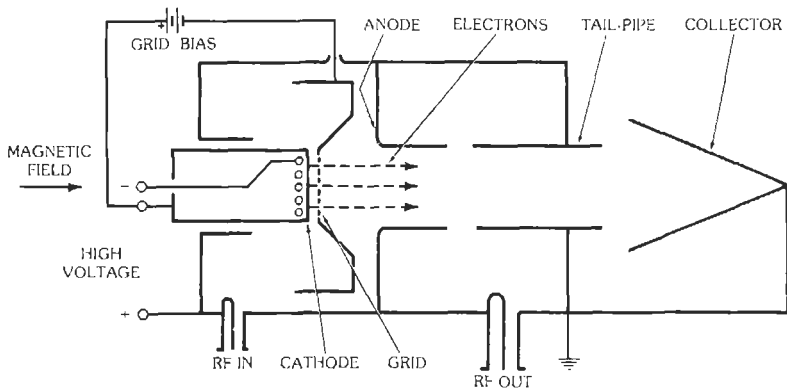


Fig. 5. KlystronTM amplifier schematic. (Courtesy Varian/Eimac)

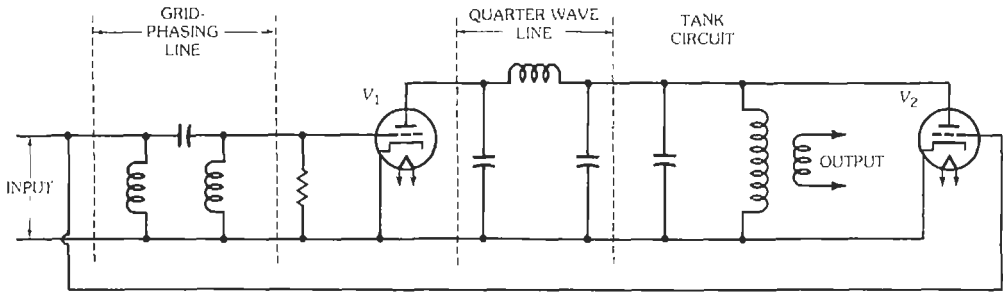
transmitted. This circuit has the advantage that efficiency tends to be high at all modulation levels, no high-level modulation transformer is required, and the peak voltage on the tubes is much lower than with conventional anode modulation.

Modulators—Modulation of the signal from a tetrode RF amplifier may be accomplished by:

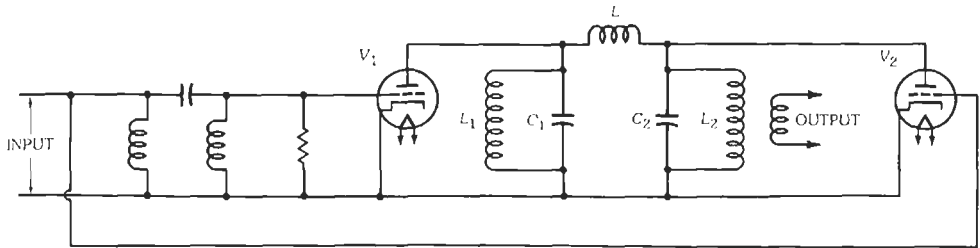
- A. Variation of RF drive, described above
- B. Variation of screen-grid voltage, keeping drive and anode voltage constant
- C. Variation of screen-grid and anode voltage together
- D. Variation of control-grid voltage only

Each of these systems has advantages and disadvantages. Most transmitters today use method A or C. Pulse transmission, where essentially square pulses are required, can be accomplished by method D, with variation of the grid-bias voltage between cutoff and the value required for efficient operation when current flows.

Power grid tubes are frequently used to pulse-modulate velocity-modulated microwave tubes such as klystrons, TWTs, magnetron oscillators, and cross-field amplifiers. Simple circuits can be used because a single gridded tube can perform the pulsing function. An interesting example at high power level is the "linear beam switch tube" series modulator. This tube is a



(A) Quarter-wave line shown separated.



(B) Quarter-wave line incorporated into tank circuits.

Fig. 6. Schematic diagram of high-efficiency linear amplifier. (After F. E. Terman, Radio Engineers Handbook, New York: McGraw-Hill Book Co., 1943; p. 456.)

form of tetrode in which the pseudo-screen grid is an apertured electrode and the final anode, or collector, is a reentrant chamber designed to minimize the production of secondary electrons, which would be accelerated back to the screen grid because of its higher potential. This arrangement is especially suitable for modulation of a high-power klystron that requires a high beam voltage negative in polarity with respect to ground. Fig. 7 shows this arrangement schematically. The switch tube itself uses a convergent electron beam from a concave spherical cathode and control grid in an electron gun that is similar to that of the klystron.

A relatively new modulation technique that is finding increasing use in AM broadcast transmitters is pulse duration modulation (PDM)*. An RF power amplifier is anode-modulated by a series modulator tube. The control grid of the modulator tube is pulsed by a train of rectangular pulses at a repetition frequency well above the highest modulation frequency required. The pulses are arranged to drive the modulator tube to its maximum anode current and to zero current alternately. The width of the pulses is varied at the modulation frequency. A filter tuned to the pulse repetition frequency is placed between the modulator tube and the RF amplifier tube so that only the required modulating signal is applied to the RF tube. The result is a relatively efficient system at all modulation levels because the modulator tube is either fully on or fully off, and because the efficiency of the RF tube is constant over a wide range of applied anode voltage.

Broad-Band Distributed Amplifiers—For some applications, an amplifier that covers several octaves of bandwidth is required. At microwave frequencies, a single TWT can be made to do this. At lower frequencies, from DC to several hundred megahertz, an arrangement of tetrode tubes and lumped-constant transmission lines known as a “distributed amplifier” can be used,

* Reference 4.

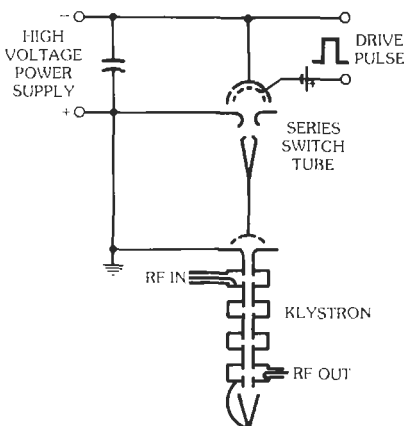


Fig. 7. Klystron modulated by linear beam switch tube. (Courtesy Varian)

as shown in Fig. 8. The transmission lines are terminated by load resistances of magnitude equal to their characteristic impedances. A growing wave of current is present on the output transmission line, each tube providing its contribution of current in the correct phase. Such amplifiers are quite inefficient. Typically eight to sixteen tubes are used. Tube requirements are high input impedance (there must be no grid current) and high anode dissipation capability, dictating the use of tetrodes with grounded cathodes. The upper cutoff frequency is limited mainly by cathode and grid lead inductances and grid-cathode capacitance of the tubes.*

PERFORMANCE OF ELECTRON POWER TUBES; COMPARISON WITH SOLID-STATE DEVICES

Power grid tubes are traditionally used for generation and amplification of power at radio frequencies; for modulation, switching, and rectification; and, formerly, for low-power-level applications including receivers and early computers. Solid-state devices and packages have eliminated tubes from all low-power-level applications where information processing is the objective, and they are becoming widely used where up to 2 kW of CW power is required for radio transmission below about 2 GHz. At much higher power levels, power tubes remain the economical choice and are likely to remain so for the foreseeable future.

Fig. 9 shows the RF power obtainable from various devices as a function frequency. Data are taken from manufacturers’ catalogs and other published information. It is clear that in terms of maximum CW power obtainable, a single power-tube device is many orders of magnitude more powerful than a single solid-state device over the whole frequency range. This situation is not likely to change. It exists because of the fundamental physical distinction between the properties of electrons moving in a vacuum and electrons moving in solid material, and the properties of the media themselves. This can be seen from the following discussion.

Radio-frequency generators are really converters of DC to RF power. The key elements are:

- A. A DC power supply
- B. An RF resonant circuit with Q greater than approximately 5
- C. A source of electrons
- D. A means of “bunching” electrons and accelerating the bunches
- E. A means for interaction between the bunches of electrons and the resonant circuit so that energy is extracted from the electrons and transferred via the circuit (which, like a flywheel, stores energy) to a useful load such as an antenna or a substance to be heated. (See Fig. 10.)

* Reference 5.

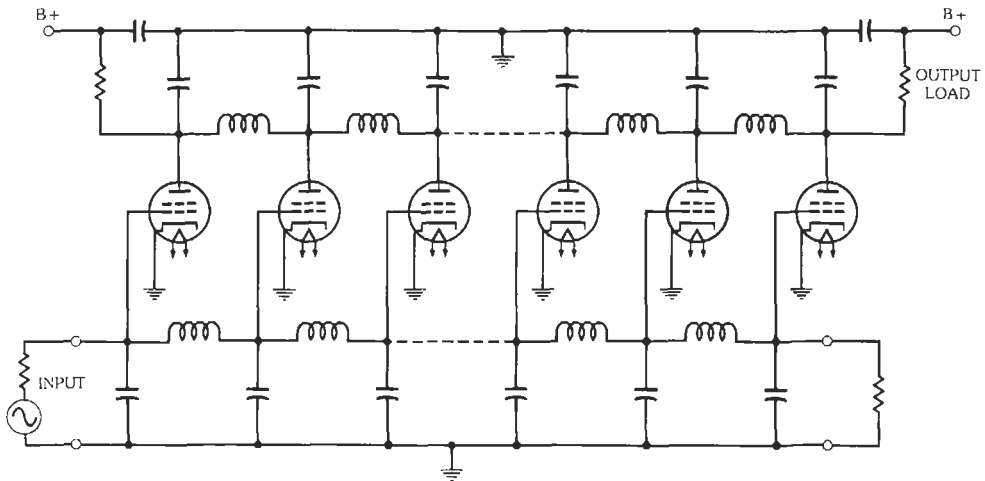


Fig. 8. Basic distributed amplifier circuit. (Courtesy IEEE)

More specifically, typical configurations are shown in Fig. 11. Simplified circuit diagrams are essentially the same for a triode tube and a transistor, both of which illustrate the principles involved.

Fig. 12, which illustrates current and voltage as functions of time, shows that the voltage appearing across the output region of the device is small when current is flowing, but peaks at a value roughly equal to twice the DC supply voltage half a period later. The device must withstand this voltage without internal arcing. The current waveform shown in Fig. 12 implies

that the time of transit of electrons between the source (cathode or emitter) and the final electrode (anode or collector) is very small compared to the period of the radio frequency. If the transit time is large, the current waveform will depart from that shown as the solid line in Fig. 12 and will tend to be shown by the dash line. The RF component of this current will decrease as the transit time increases, and the conversion efficiency of the device will fall, eventually to zero. In quantitative terms, it can be shown that, other factors being equal, the parameter v/fd (where v is electron velocity across

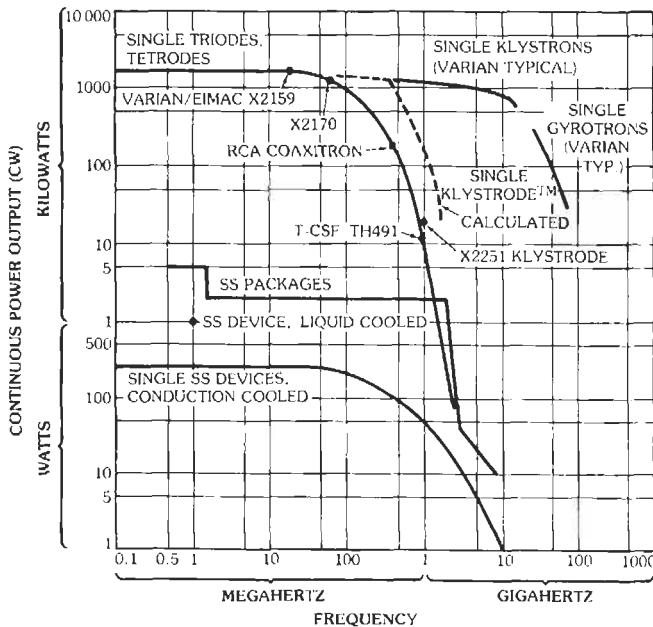


Fig. 9. Maximum CW power versus frequency from RF power sources, 1982. Available except where noted.

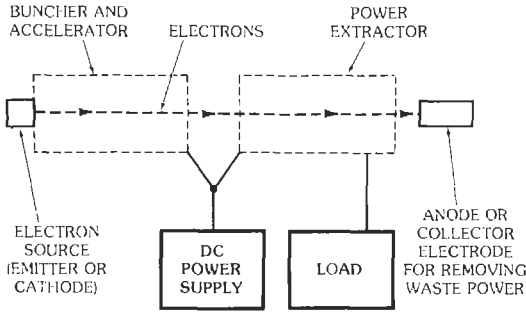


Fig. 10. Basic schematic diagram of RF amplifier using electron device.

the interelectrode gap, f is frequency, and d is gap length in the direction of electron flow) must be held constant for a given performance.

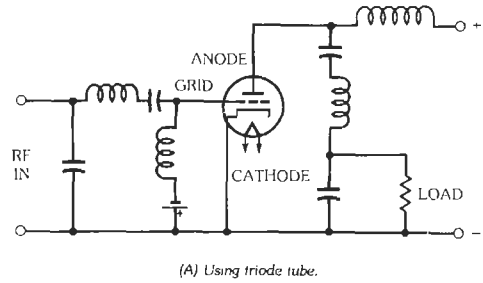
The above is true for electrons both in a vacuum and in the solid state, but if we now consider the actual values of v and d , we note a fundamental difference.

$$v \text{ in vacuum} \propto \phi^{1/2}$$

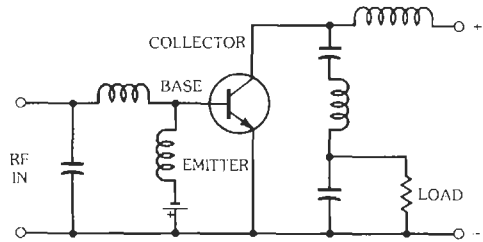
where ϕ is the potential difference through which the electron has passed. For example, if $\phi = 2000$ volts, as commonly encountered in a high-power tube, $v \approx 2 \times 10^9$ cm/s. On the other hand, in silicon, for example, v_{max} is a constant 6×10^6 cm/s.

Assume that for efficient energy extraction an electron must pass through the output gap in a time $t = (1/2\pi) \times (\text{period of radio frequency})$. (The output gap is grid to anode in a tube or base to collector in a transistor.) At a frequency of 100 MHz, for example, $t = (1/2\pi)10^{-8}$ s. Then the output gap length must be less than the product vt , or 10^{-2} cm for the transistor compared to 3 cm for the tube. At 1000 MHz, these dimensions would be 10^{-3} cm for the transistor and 0.3 cm for the tube.

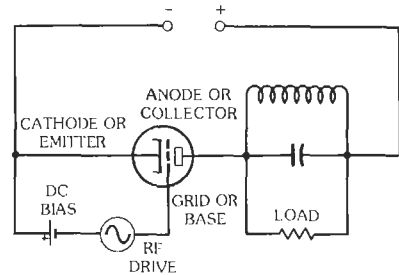
The maximum applied voltage is determined by the above dimensions and the dielectric strength of the medium. For silicon, the breakdown voltage is about 2×10^5 V/cm. In a vacuum, the situation is more complicated, but under typical conditions in a high-power tube the same number can be used. This leads to the conclusion that at 1000 MHz the maximum voltage across the tube could be $2 \times 10^5 \times 3 \times 10^{-1} = 60\,000$ V and the maximum voltage across the transistor could be 200 V. In both devices, when used in the circuits shown in Fig. 11, the peak applied voltage occurs during the half cycle when no current flows and will reach a value of about twice the DC power-supply voltage. Therefore, the DC voltage would ideally be limited, at 1000 MHz, to 100 volts for the transistor and 30 000 volts for the tube. In practice, the transistor DC voltage tends to be much lower, typically 28 volts, because the base-to-collector distance is in fact less than the distance limit determined by electron transit time (for reasons given below), and the maximum usable



(A) Using triode tube.



(B) Using transistor.



(C) Equivalent circuit for both types of device.

Fig. 11. Simplified RF amplifier circuit diagrams.

voltage remains constant as the frequency is reduced. In the tube, the voltage ideally could rise inversely with frequency, but it is limited in practice because of the cost of high-voltage power supplies and the problem of insulation outside the vacuum.

Summarizing, then, it appears that for reasons of breakdown voltage alone, the relative power levels of transistors and gridded power tubes would be in the ratio of 1 to 300, assuming the same current in both and no other limitations. In practice, the ratio is higher because the problem of heat removal in the transistor sets an upper limit to the voltage that tends to be constant at about 60 volts rather than 200 volts, so that power-supply voltages of 28 volts are typical. The heat-removal problem is severe because, unlike the case in a vacuum, heat is generated by the passage of current through the semiconductor due to collisions between electrons and atoms (the reason for the velocity limit referred to above) and also by the thermal vibration of

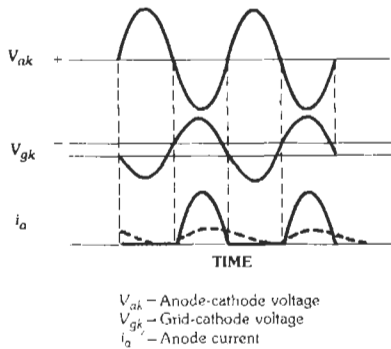


Fig. 12. Waveforms of voltage and current in class-B radio-frequency amplifier.

the atoms in the RF field (dielectric loss). In a vacuum, this is zero. The heat generated must be conducted through the semiconductor to the heat sink at the collector terminal, so the thermal conductivity of the material is critical. To add to the difficulties, semiconductor performance deteriorates as temperature increases. The heat-removal problem in fact sets an upper limit to power output well below the more basic limitation due to electron transit time.

The remaining factors that determine the power output of a device are its cross-sectional area and current density. Here again is a great distinction. The size of the solid-state device is limited by fabrication technology and cost to the order of 0.1 cm^2 . The gridded power tube is limited by the more fundamental fact that its efficiency falls off if one dimension perpendicular to electron flow exceeds approximately one-tenth the operating wavelength. The other dimension may be larger. A cylindrical tube at 1000 MHz (30 cm wavelength) may be $3 \times 10 = 30 \text{ cm}^2$ in area. At lower frequencies, it can be correspondingly larger, as will be seen later. The maximum current density on peaks of RF achievable in the transistor is about 200 A/cm^2 at UHF and 800 A/cm^2 at 100 MHz. In the tube, it is about 3 A/cm^2 . The current per device is the product of the current density and the area.

In summary, the power output per device is proportional to applied voltage times current density times cross-sectional area times conversion efficiency. If efficiency is neglected, Table 4 is a rough guide to relative performance. From the table, it is evident that large numbers of transistors must be combined to give

the power output obtainable from a single tube. Combination of transistors has been done to the extent indicated by Fig. 9, in which available power from combinations is shown.

It should be noted that the major advantage of the solid-state device, inherent long life under normal conditions, becomes less and less significant as the number of devices combined increases. This is to some extent overcome by techniques that allow complete failure of one or more devices to occur without failure of the system but with only a degradation in performance that can often be tolerated ("graceful degradation"). It is also not surprising that some performance criteria—bandwidth, for example—tend to fall off as the number of devices combined increases. This happens because the input and output impedances of the combination, which with simple paralleling are inversely proportional to the number of devices, become very low compared to the impedances of the circuits with which the combination is associated. Also, it should be noted that paralleling requires low-voltage regulated power supplies with high current capability. This is by no means easy to achieve and is very costly. One form of degradation, hardly graceful, is total failure of all devices due to a transient overcurrent or overvoltage phenomenon such as the electromagnetic pulse (emp) that can result from nearby lightning or a nuclear explosion. A less dramatic but equally fatal circumstance producing the same result is an unplanned increase in load impedance that causes the voltage to rise to the breakdown limit. It is interesting that the solid-state device and the tube device react differently to this event because of their different physical nature. In the solid-state device, the effect is a permanent breakdown, or short-circuit. In the vacuum device, the effect is a temporary breakdown followed by a resumption of full performance a short time after the transient has subsided.

In summary, then, it appears that both tubes and solid-state devices have fundamental limitations. Both have reached a stage of development where subsequent improvements will probably be small relative to those already made. Essentially, the transistor is a low-power device. The art of combining transistors may progress so that higher-power packages may become available, but at extra cost. In nonmilitary systems, it is ultimately the cost in terms of amortized capital cost plus operating and maintenance costs per unit time that governs the choice of device type.

TABLE 4. RELATIVE PERFORMANCE OF TUBE AND TRANSISTOR

Device	Volts	Area (cm^2)	Current Density (A/cm^2)	Product	Approx Ratio
Transistor at 1000 MHz	28	0.1	200	5.6×10^2	1
Transistor at 100 MHz	28	0.1	800	22.4×10^2	4
Tetrode at 1000 MHz	6000	30	3	5.4×10^5	1000
Tetrode at 100 MHz	18 000	1000	1	1.8×10^7	30 000

(It is of historical interest that before the advent of high-power velocity-modulated tubes, serious and successful attempts were made to combine triodes and tetrodes to increase RF power output at VHF and UHF, and many transmitters were built this way. As soon as tubes giving sufficient power from a single type for the application became available, these multitube equipments were superseded for economic reasons.)

Power Gain and Bandwidth of RF Amplifiers

Power gain and bandwidth of RF amplifiers depend on both device and circuit properties. As noted above, the power gain of a tetrode falls in the range between 30 dB at LF and MF and 10 dB at UHF. Power gain at UHF is limited by electron-transit-time effects and by the inevitable effect of cathode-lead inductance.

Corresponding bandwidths with conventional resonant circuits fall between about 1% (AM broadcast, UHF TV) and 10% (TV at 50 MHz) of the carrier frequency.

For a simplified ideal tube amplifier, the gain-bandwidth product tends to be constant and equal to G_m/C , where G_m is the mutual conductance (transconductance) and C is the output capacitance. This product is independent of frequency. In practice, the product is always less because of the presence of the external circuits that store energy and because the bandwidth of the input circuit may be smaller than that of the output circuit. Furthermore, the performance of the amplifier may be largely determined by the amount of feedback present. It is nevertheless true that the gain-bandwidth product will always increase with G_m/C .

This factor is related to interelectrode spacings and is

approximately inversely proportional to the grid-cathode spacing. The achievement of small spacing, which implies fine grid wires and structures of high dimensional stability with respect to change of temperature, has always been a major challenge in tube design.

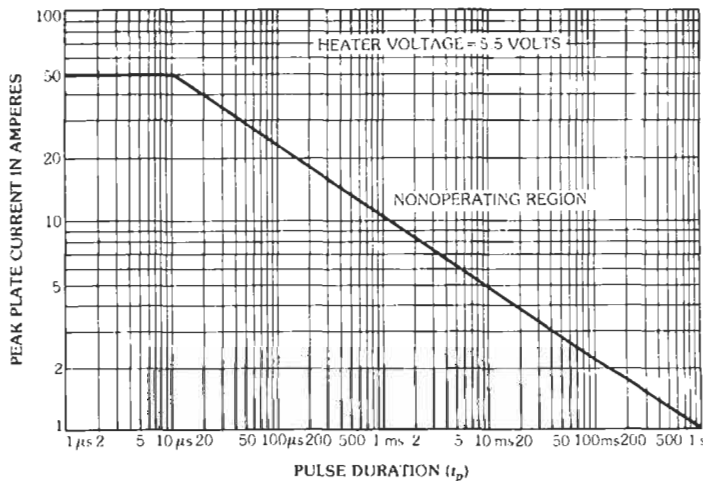
Power Output Under Pulse Conditions

Fig. 9 shows CW power available versus frequency. Under pulse conditions, gridded tubes can be made to give many times this power, depending on the pulse width and duty factor (ratio of pulse width to the interval between pulses). Tube manufacturers will supply data on the exact amount of power. A typical pulse derating chart is shown in Fig. 13. The power is limited by the thermal capacity of the grids, by the breakdown voltage between electrodes, and by the emission capability of the cathode.

Life Expectancy of Tubes

The principal factor governing tube life is the temperature of the thermionic cathode. This is related to emission and therefore to output power. In general, longer life results from lower cathode temperature. In a well designed tube, the available emission at the specified cathode heating power is many times that required for satisfactory operation. It is therefore often possible to reduce the heating power and increase the life of the tube.

Another factor affecting life is the relationship between maximum ratings and operating conditions. Exceeding ratings causes overheating, which tends to shorten life. Careful attention must be paid to manufacturers' ratings and recommendations.



Peak (pulse) plate current capability is dependent on pulse duration (t_p) and duty factor (D). Maximum peak plate current for a given value t_p is shown. Maximum D may then be derived from the relationship:

$$1.0 = i_p \times D$$

Fig. 13. Pulse-derating data for triode type 3CPX1500A7, pulse modulator or regulator service. (Courtesy Varian/Eimac)

Life expectancy for a highly stressed oxide-cathode tube at UHF may be 2500 hours. For a thoriated-tungsten-filament tube in broadcast service, 50 000 hours is not uncommon.

Development Trends in Gridded Power Tubes

Improved performance of power grid tubes, including longer life, can be anticipated from the incorporation of improved materials and fabrication techniques. For example, the use of pyrolytic graphite for grids is becoming standard practice. Impregnated tungsten dispenser cathodes commonly found in microwave velocity-modulated tubes may find their way into gridded power tubes.

In addition, improved structural forms such as modular construction may allow the manufacture of very-high-power tubes (capable of 10 MW of CW power output) having higher gain-bandwidth products because of smaller grid-cathode spacings.

At VHF and UHF, developments of the Klystrode™ amplifier may appear in applications where the combi-

nation of small size compared to a klystron; relatively high efficiency, especially with amplitude modulation; high power capability; and interesting power gain makes it more attractive than a klystron or a conventional tetrode.

REFERENCES

1. Tyler, V. J. "A New High-Efficiency Power Amplifier." *Marconi Review*, Vol. 21, No. 130 (1958), pp. 6-109.
2. Doherty, W. H. "A New High-Efficiency Power Amplifier for Modulated Waves." *Proc. IRE*, Vol. 24 (Sept. 1936), p. 1163.
3. Terman, F. E. *Radio Engineers Handbook*. New York: McGraw-Hill Book Co., 1943; pp. 455-458.
4. Swanson, H. "The Pulse Duration Modulator, a New Method of High Level Modulation in Broadcast Transmitters." *IEEE Trans. on Broadcasting*, Vol. BC17, No. 4 (Dec. 1971), pp. 89-92.
5. Ginzton, Hewlett, Jasberg, and Noe. "Distributed Amplifications." *Proc. IRE*, 1948, pp. 956-969.

18 Semiconductors and Transistors

Ben G. Streetman

Semiconductors	18-2
Semiconductor Materials	
Energy Bands and Charge Carriers	
Carrier Concentrations	
Drift of Carriers	
Excess Carriers; Diffusion	
pn Junctions	18-8
A Junction at Equilibrium	
Forward and Reverse Bias	
Reverse Breakdown	
Capacitance and Transient Behavior	
Diodes and Rectifiers	18-12
Transistors	18-13
Junction Field-Effect Transistors	
MOS Field-Effect Transistors	
Bipolar Junction Transistors	
Other Electronic Devices	18-21
Display Devices	
Microwave Devices	
Switching Devices	

The basic building blocks of any electronic circuit or system are electronic devices—the transistors, diodes, and other elements that collectively allow the system to perform its function. These devices may be individual transistors, diodes, and passive elements soldered in appropriate interconnection. Or, the system may be a collection of complex integrated circuits, each composed of thousands of active and passive elements on a silicon chip (see Chapter 20). The system may even involve information transfer by a modulated light beam and detector (see Chapter 21). In any case, a great variety of electronic elements is available for performing the desired function, and those who design or use electronic equipment have a wide range of devices from which to choose.

Most modern electronic devices are made with semiconductor materials. The unusual properties of these solids are responsible for a revolution in electronics since the invention of the transistor. This chapter* deals with semiconductors and their applications in pn junctions, transistors, and other devices. Circuit models of transistors are discussed in Chapter 19, and fabrication methods are discussed in Chapter 20.

SEMICONDUCTORS

Semiconductors are a group of materials that have electrical conductivities intermediate between those of metals and insulators. It is significant that the conductivity of these materials can be varied considerably by changes in temperature, optical excitation, and impurity content. This variability of electrical properties makes possible the wide range of modern electronic devices.

*Much of this discussion is from Ben G. Streetman, *Solid State Electronic Devices*, 3rd ed., © 1990. Reprinted by permission of Prentice-Hall, Inc., Englewood Cliffs, N.J.

Semiconductor Materials

Semiconductors are found in column IV and neighboring columns of the periodic table (Table 1). The column-IV semiconductors silicon (Si) and germanium (Ge) are called *elemental* semiconductors because they are composed of single species of atoms. In addition to the elemental materials, compounds of column-III and column-V atoms, as well as certain combinations from columns II and VI, make up the *intermetallic*, or *compound*, semiconductors.

As Table 1 indicates, there are numerous semiconductor materials. Among these, silicon is used for the majority of semiconductor devices; rectifiers, transistors, and integrated circuits are now made mostly in silicon, although GaAs is also used in high-speed integrated circuits. The compounds are most widely used in devices requiring the emission or absorption of light. For example, semiconductor light-emitting diodes commonly are made of such compounds as GaAs, GaP, and alloys such as GaAsP. An important microwave device, the Gunn diode, is usually made in GaAs. Thus, the wide range of semiconductor materials offers considerable variety in properties and provides experts in electronic circuits and systems with much flexibility in the design of electronic functions.

The electronic and optical properties of semiconductor materials are strongly affected by impurities, which may be added in precisely controlled amounts. Such impurities are used to vary the conductivities of semiconductors, and even to alter the nature of the conduction processes from conduction by negative charge carriers to conduction by positive charge carriers. For example, an impurity concentration of one part per million can change a sample of silicon from a poor conductor to a good conductor of electric current. This process of controlled addition of impurities, called *doping*, will be discussed below.

TABLE 1. COMMON SEMICONDUCTOR MATERIALS

The Portion of the Periodic Table Where Semiconductors Occur				
II	III	IV	V	VI
	B	C		
	Al	Si	P	S
Zn	Ga	Ge	As	Se
Cd	In		Sb	Te
Elemental and Compound Semiconductors				
Elemental	IV Compounds	III-V Compounds	II-VI Compounds	
Si	SiC	AlP	ZnS	
Ge	SiGe	AlAs	ZnSe	
		AlSb	ZnTe	
		GaP	CdS	
		GaAs	CdSe	
		GaSb	CdTe	
		InP		
		InAs		
		InSb		

Virtually all semiconductor devices require single crystals of extremely pure materials. Since small concentrations of impurities can radically alter their electrical properties, device-grade semiconductor crystals are grown with greater perfection than any other materials.

Energy Bands and Charge Carriers

The atomic arrangement in most semiconductor crystals is similar to the diamond lattice, in which each atom is surrounded by four nearest neighbors. The atomic bonding in the crystal is largely covalent (as in the H_2 molecule). That is, the sharing of electrons between adjacent atoms in covalent bonds holds the crystal together. In the column-IV semiconductor silicon, for example, each atom has four valence electrons shared with four nearest neighbors.

Electrons in isolated atoms are restricted to certain discrete energy levels, predictable generally from the Bohr model and more precisely from the results of quantum mechanics. In a similar fashion, electrons in solids are restricted to certain energies and are not allowed at other energies. The basic difference between the case of an electron in a solid and that of an electron in an isolated atom is that in the solid the electron has a *range*, or *band*, of available energies. General features of the energy bands for insulators, metals, and semiconductors are shown in Fig. 1. In an insulator or semiconductor, a lower (valence) band is filled with electrons at low temperatures, and an upper (conduction) band is empty of electrons. The separation between the valence and conduction bands is called the *band gap* E_g . Insulators have a wider band gap than do semiconductors, and metals have overlapping bands.

The variety of electrical conductivities of metals, insulators, and semiconductors arises primarily from their band structures (Fig. 1). If electrons are to experience acceleration in an electric field, they must be able to move into new energy states. In a completely filled band, there are no available empty states, and therefore no net charge transport (no current) can take place. The same is true for an empty band, which contains no electrons to move. As a result, a perfect insulator does not conduct current. Metals, on the other hand, have electrons mixed with unfilled allowable

energy states. Electrons are therefore available to participate in current flow, and there are plenty of energy states for the electrons to occupy in response to a field. This accounts for the high conductivity of metals.

Although semiconductors and insulators have similar band structures, there is an important difference: The band gap E_g of an insulator is several electronvolts wide,⁶ whereas E_g for a typical semiconductor is sufficiently small (about one eV) to allow excitations of electrons from the valence band to the conduction band. These excitations can result from the application of thermal or optical energy. For example, at room temperature in a perfect silicon crystal ($E_g = 1.1$ eV), the equilibrium concentration of thermally generated electrons in the conduction band is about 1.5×10^{10} per cubic centimeter (cm^{-3}). Of course, there is an equal number of empty energy states in the valence band. These unoccupied valence-band states are referred to as *holes*. Electron-hole pairs (EHP) can also be created by optical excitation; for example, a quantum of light (*photon*) with energy greater than the band gap E_g can be absorbed by the semiconductor crystal, exciting a valence-band electron to the conduction band. The result is a conduction-band electron and a valence-band hole. As discussed below, this optical generation of electron-hole pairs is basic to the operation of semiconductor photoconductors and photodiodes.

At thermal equilibrium, the excitation of electrons to the conduction band is balanced by electrons falling back to the valence band (recombination). This balance of electron-hole generation and recombination results in a predictable concentration of EHP at a given temperature. Similarly, a steady optical excitation causes a new steady-state balance between generation and recombination.

It is clear that electrons which have been excited to the conduction band of a semiconductor can participate in current conduction, since ample unoccupied energy states are available in that band. A less obvious but important feature of semiconductors is charge transport involving holes in the valence band. The presence of an empty state (hole) in the otherwise filled valence band allows electrons in the band to move, with a resulting net motion of charge, or current. The approach for calculating this charge transport is to sum the contributions of all electrons in a filled valence band (which results in zero net current) and then subtract the contribution of the missing electron. The result is equivalent to considering the hole as a positive charge carrier (i.e., a particle with charge $+q$, where $-q$ is the charge on an electron).⁷ For example, a semiconductor

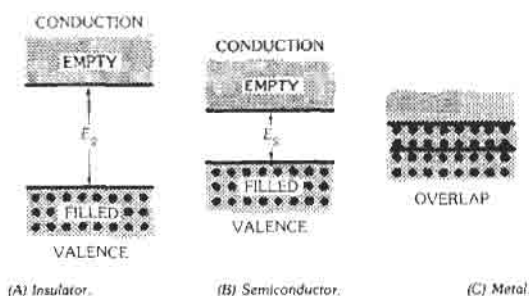


Fig. 1. Energy bands for typical materials.

⁶The electronvolt (eV) is a particularly convenient unit of measure for the energy of an electron. One eV is defined as the energy acquired by an electron moving through a potential of one volt. Thus, $1 \text{ eV} = 1.6 \times 10^{-19}$ joule.

⁷Table 2 lists the magnitude of the electronic charge and other information.

TABLE 2.* PHYSICAL CONSTANTS AND CONVERSION FACTORS

Avogadro's number	$N_A = 6.02 \times 10^{23}$ molecules/mole
Boltzmann's constant	$k = 1.38 \times 10^{-23}$ J/K $= 8.62 \times 10^{-5}$ eV/K
Electronic charge (magnitude)	$q = 1.60 \times 10^{-19}$ C
Electronic rest mass	$m_0 = 9.11 \times 10^{-31}$ kg
Permittivity of free space	$\epsilon_0 = 8.85 \times 10^{-14}$ F/cm $= 8.85 \times 10^{-12}$ F/m
Planck's constant	$h = 6.63 \times 10^{-34}$ J s $= 4.14 \times 10^{-15}$ eV s
Room temperature value of kT	$kT = 0.0259$ eV
Speed of light	$c = 2.998 \times 10^{10}$ cm/s
1 Å (angstrom) = 10^{-8} cm	Prefixes:
1 μm (micron) = 10^{-4} cm	milli-, m- = 10^{-3}
1 mil = 10^{-3} in	micro-, μ- = 10^{-6}
2.54 cm = 1 in	nano-, n- = 10^{-9}
1 eV = 1.6×10^{-19} J	pico-, p- = 10^{-12}
	kilo-, k- = 10^3
	mega-, M- = 10^6
	giga-, G- = 10^9
A wavelength λ of 1 μm corresponds to a photon energy of 1.24 eV.	

* From Ben G. Streetman, *Solid State Electronic Devices*, 3rd ed., © 1990. Reprinted by permission of Prentice-Hall, Inc., Englewood Cliffs, N.J.

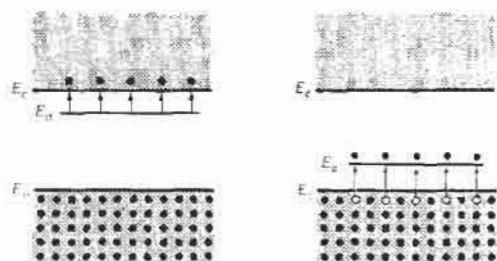
with 10^{10} electron-hole pairs per cubic centimeter contains a concentration of 10^{10} cm $^{-3}$ negatively charged electrons in the conduction band, which drift opposite to an applied electric field, and 10^{10} cm $^{-3}$ positively charged holes in the valence band, which drift in the direction of the field. The two current components add in the field direction to give the total current. Both electrons and holes are important charge carriers in semiconductors.

In addition to the direct excitation of electrons and holes in pairs, each type of charge carrier can be introduced by appropriate doping of the semiconductor with impurities. Impurities can be added to the crystal in the growth process or introduced later during device fabrication. By doping, a crystal can be altered so that it has a predominance of either electrons or holes. Thus there are two types of doped semiconductors, n-type (mostly electrons) and p-type (mostly holes). The designation n and p is chosen to reflect the dominance of negative or positive charge carriers.

When impurities are introduced into an otherwise perfect crystal, additional levels are created in the energy-band structure, usually within the band gap. For example, an impurity from column V of the periodic table (phosphorus, P; arsenic, As; antimony, Sb) introduces an energy level very near the conduction band in germanium or silicon. This level is filled with electrons at 0 K, and very little thermal energy is required to excite these electrons to the conduction band (Fig. 2A). Thus at room temperature virtually all of the electrons in the impurity level are "donated" to the conduction band. Such an impurity level is called a *donor* level, and the column-V impurities in germanium or silicon are called donor impurities. Material doped with donor

impurities can have a considerable number of electrons in the conduction band, even when the temperature is too low for the thermally generated EHPs to be appreciable. Thus for semiconductors doped with a significant number of donor atoms, the electron concentration n will be much greater than the hole concentration p . This is n-type material.

Atoms from column III (boron, B; aluminum, Al; gallium, Ga; indium, In) introduce impurity levels in germanium or silicon near the valence band (Fig. 2B). These levels are empty of electrons at 0 K. At higher temperatures, enough thermal energy is available to excite electrons from the valence band into the impurity level, leaving behind holes in the valence band. Since this type of impurity level "accepts" electrons from the valence band, it is called an *acceptor* level, and the column-III impurities are acceptor impurities in germa-



(A) Electrons are thermally excited from a donor level to the conduction band.

(B) Valence-band electrons are excited to E_a , leaving holes in the valence band.

Fig. 2. Donor and acceptor levels (E_d and E_a) in a semiconductor.

nium and silicon. Doping with acceptor impurities can create a semiconductor with a hole density much greater than the conduction-band electron density (p-type material).

A semiconductor without doping impurities is called *intrinsic* material, and the *intrinsic carrier concentration* n_i is the concentration of thermally generated electron-hole pairs. In silicon, for example, the intrinsic concentration n_i is about 10^{10} cm^{-3} at room temperature. If silicon is doped with 10^{15} antimony atoms/cm³, the conduction electron concentration n increases from 10^{10} to 10^{15} cm^{-3} , a change of five orders of magnitude. As a result, the resistivity of silicon changes from about $2 \times 10^5 \text{ ohm-cm}$ to 5 ohm-cm with this doping.

When a semiconductor is doped n-type or p-type, one type of carrier dominates. In the example given above, the conduction-band electrons outnumber the holes in the valence band by many orders of magnitude. We refer to the small number of holes in n-type material as *minority carriers* and the relatively large number of conduction electrons as *majority carriers*. Similarly, electrons are the minority carriers in p-type material and holes are the majority carriers.

Carrier Concentrations

In calculating semiconductor electrical properties and analyzing device behavior, it is often necessary to know the concentration of charge carriers in the material. The majority-carrier concentration is usually obvious in heavily doped material, since one majority carrier is obtained for each impurity atom. The concentration of minority carriers is not obvious, however, unless we consider details of the electron distribution in the solid. An important result of semiconductor statistics is that the *electron-hole product* np at equilibrium is a constant for a given material at a given temperature, whether the material is doped or not:

$$np = n_i^2 \tag{Eq. 1}$$

where n and p are the equilibrium electron and hole concentrations, respectively. Thus, if we know the majority-carrier concentration (e.g., from the doping density) we can find the minority-carrier concentration from Eq. 1. For example, if a silicon sample ($n_i = 1.5 \times 10^{10} \text{ cm}^{-3}$) is doped with 10^{15} donors/cm³, the electron concentration is essentially 10^{15} and the hole concentration is only $2.25 \times 10^5 \text{ cm}^{-3}$ at room temperature. It is interesting to note from Eq. 1 that as the majority-carrier concentration increases, the minority-carrier concentration must decrease.

If a semiconductor contains both donors and acceptors, the exact relationship among the electron, hole, donor, and acceptor concentrations can be obtained by considering the requirements for *space-charge neutrality*. If the material is to remain electrostatically neutral, the sum of the positive charges (holes and ionized donor atoms N_d^+) must balance the sum of the negative charges (electrons and ionized acceptor atoms N_a^-):

$$p + N_d^+ = n + N_a^- \tag{Eq. 2}$$

If the material is strongly n-type ($n \gg p$) and all of the impurities are ionized, we can approximate Eq. 2 by $n \cong N_d - N_a$.

Drift of Carriers

The electrons and holes in a semiconductor are in constant motion due to their thermal energy. At thermal equilibrium, the movement is random, and carriers are scattered from lattice atoms, impurities, and defects. Since the scattering is completely random, there is no net motion of the group of carriers over a period of time. On the other hand, if an electric field is applied to the sample, a net drift of carriers is superimposed on the random thermal motion. Thus, a field \mathcal{E}_x applied in the direction results in a net force $q\mathcal{E}_x$ on each hole and $-q\mathcal{E}_x$ on each electron. The current resulting from this net drift is just the number of carriers crossing a given area A per unit time (number of carriers times velocity) multiplied by the charge of the carrier:

$$I = qA p v_p - qA n (-v_n) \tag{Eq. 3}$$

In this equation, v_p is the average drift velocity of holes in the x direction, and $-v_n$ is the average electron velocity, which is in the opposite direction. The resulting electrical conductivity is usually written in terms of the *mobility* of each charge carrier. This quantity measures the ease with which carriers drift in a given material. The hole mobility μ_p is the average drift velocity per unit applied field v_p/\mathcal{E} , and similarly for electron mobility μ_n . Thus the current is

$$I = qA (\mu_p p + \mu_n n) \mathcal{E} \equiv \sigma \mathcal{E} A \tag{Eq. 4}$$

where the conductivity of the sample σ has the units (ohm-cm)⁻¹. Conductivity is determined by electron and hole concentration (n, p), and therefore by the doping of the sample, and also by the carrier mobilities.

If a semiconductor bar contains both types of carriers, holes drift as a group in the direction of the electric field, and electrons drift in the opposite direction. Both components of current are in the direction of the \mathcal{E} field, however, since conventional current is positive in the direction of hole flow and opposite to the direction of electron flow.

If a magnetic field is applied perpendicular to the direction in which holes drift in a p-type bar, the path of the holes tends to be deflected (Fig. 3). In vector notation, the total force on a single hole due to the combined electric and magnetic fields is

$$\mathbf{F} = q(\mathbf{E} + \mathbf{v} \times \mathbf{B}) \tag{Eq. 5}$$

In the y direction the force is

$$F_y = q(\mathcal{E}_y - v_x \mathcal{B}_z) \tag{Eq. 6}$$

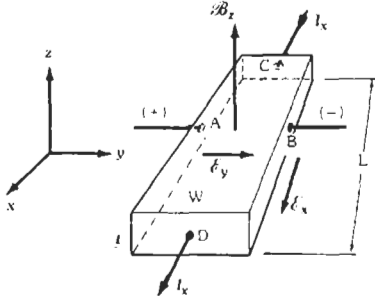


Fig. 3. A p-type semiconductor bar in a magnetic field, resulting in the production of a Hall-effect voltage.

An electric field \mathcal{E}_y builds up in the y direction as holes are displaced laterally, and in steady state this field is sufficient to just balance the force due to the magnetic field ($\mathcal{E}_y = v_x B_z$). The resulting voltage between A and B can be measured with a high-impedance voltmeter. This is called the Hall effect, and V_{AB} is called the Hall voltage. This effect can be used to measure magnetic fields and in a variety of other applications. An important use for the Hall effect in semiconductor research is that the Hall voltage can be easily related to the concentration of majority carriers in the sample. Although the discussion here relates to p-type material, similar results obtain for n-type samples, in which the majority carriers are electrons. By combining Hall measurements with conductivity measurements, the majority-carrier concentration and mobility can be obtained for a semiconductor sample.

Excess Carriers; Diffusion

To this point, we have discussed electron and hole carrier concentrations at thermal equilibrium or with steady fields applied. In each case, the carrier concentrations have been the equilibrium values (let us call them n_o and p_o). However, most semiconductor devices operate by the creation of carrier concentrations greater than the equilibrium values. In such cases, we can refer to the *excess* electron concentration as δn , such that the total electron concentration n is

$$n = n_o + \delta n \quad (\text{Eq. 7A})$$

and similarly

$$p = p_o + \delta p \quad (\text{Eq. 7B})$$

for holes.

The excess carrier concentrations denoted by δn and δp can arise from several causes, most commonly from *optical excitation* or from *carrier injection*, which is characteristic of pn junctions. Optical excitation is the basis of semiconductor light detectors.

When excess carriers are created at a rate g_{op} by a steady light, a steady-state balance is established between generation and recombination. That is, in steady

state the EHPs recombine as fast as they are generated; otherwise, there would be a buildup of EHPs with time. Thus, the recombination rate equals the generation rate in steady state. However, the electron and hole concentrations in the illuminated sample are greater than the equilibrium values. In particular,

$$\delta n = g_{op} \tau_n, \quad \delta p = g_{op} \tau_p \quad (\text{Eq. 8})$$

where τ_n and τ_p are called the electron and hole *lifetimes*, respectively. The lifetime of a carrier is the average time it spends in its respective band before recombination. Thus $\delta n/\tau_n$ is the electron recombination rate, and $\delta p/\tau_p$ is the rate at which holes are generated and recombine in steady state.

The complementary mechanism to optical absorption is *radiative recombination*. In some semiconductors, the energy lost by an electron in recombining with a hole is given off by a photon of light. If the recombination takes place directly, without involving an intermediate state in which a carrier is temporarily captured, the energy of the emitted photon is equal to the band gap (see Table 3). For example, recombination of excess carriers in GaAs results in light emission with photon energies equal to the band-gap energy of about 1.4 eV. Such emission is called *luminescence*. Recombination in many semiconductors (such as silicon and germanium) takes place indirectly; for example, an electron is trapped at an impurity or lattice defect, and then a hole is captured to complete the recombination process. In such cases, the energy of the electron is given up as heat to the lattice, and no light is emitted.

If a current is passed through an illuminated semiconductor sample, we find that the conductivity has increased due to the presence of the excess carriers. From Eqs. 4 and 7, the conductivity is

$$\begin{aligned} \sigma &= q(\mu_p p + \mu_n n) \\ &= q[\mu_p(p_o + \delta p) + \mu_n(n_o + \delta n)] \quad (\text{Eq. 9}) \end{aligned}$$

The change in conductivity between its dark value and its value with the generation rate g_{op} is

$$\begin{aligned} \Delta\sigma &= q(\mu_p \delta p + \mu_n \delta n) \\ &= q g_{op}(\tau_p \mu_p + \tau_n \mu_n) \quad (\text{Eq. 10}) \end{aligned}$$

This increase in conductivity for the illuminated sample is called *photoconductivity*. We notice from Eq. 10 that a sensitive photoconductive detector should have reasonably long carrier lifetimes and high carrier mobilities. Carrier lifetime is determined for most semiconductors by crystal quality, doping, and other material properties.

If the excess carrier concentrations δn and δp vary with position in the sample, *diffusion* occurs. Diffusion is the familiar process by which particles migrate from regions where their concentration is high to regions of lower concentration. A simple example of this process is the opening of a perfume bottle in a closed room with

TABLE 3.* PROPERTIES OF SEMICONDUCTOR MATERIALS (300 K)

	E_g (eV)	μ_n (cm ² /Vs)	μ_p (cm ² /Vs)	Lattice	a (Å)	ϵ_r	Density (g/cm ³)	Melting point (°C)
Si	1.11	1350	480	D	5.43	11.8	2.33	1415
Ge	0.67	3900	1900	D	5.65	16	5.32	936
SiC(α)	2.86	500		W	3.08	10.2	3.21	2830
AlP	2.45	80		Z	5.46	9.8	2.40	2000
AlAs	2.16	180		Z	5.66	10.9	3.60	1740
AlSb	1.6	200	300	Z	6.14	11	4.26	1080
GaP	2.26	300	150	Z	5.45	11.1	4.13	1467
GaAs	1.43	8500	400	Z	5.65	13.2	5.31	1238
GaSb	0.7	5000	1000	Z	6.09	15.7	5.61	712
InP	1.35	4000	100	Z	5.87	12.4	4.79	1070
InAs	0.36	22600	200	Z	6.06	14.6	5.67	943
InSb	0.18	10 ⁵	1700	Z	6.48	17.7	5.78	525
ZnS	3.6	110		Z,W	5.409	8.9	4.09	
ZnSe	2.7	600		Z	5.671	9.2	5.65	
ZnTe	2.25		100	Z	6.101	10.4	5.51	
CdS	2.42	250	15	W,Z	4.137	8.9	4.82	1475
CdSe	1.73	650		W	4.30	10.2	5.81	1258
CdTe	1.58	1050	100	Z	6.482	10.2	6.20	1098
PbS	0.37	575	200	H	5.936	161	7.6	1119
PbSe	0.27	1000	1000	H	6.147	280	8.73	1081
PbTe	0.29	1600	700	H	6.452	360	8.16	925

Definitions of symbols: a is lattice constant; D is diamond; Z is zinc blende; W is wurtzite; H is halite (NaCl). Values of mobility are for material of available purity. These values are considered approximate.

* From Ben G. Streetman, *Solid State Electronic Devices*, 3rd ed., © 1990. Reprinted by permission of Prentice-Hall, Inc., Englewood Cliffs, N.J.

very still air. The scented air molecules soon disperse throughout the room, even if no air currents exist. The diffusion is due simply to the random thermal motion of the air molecules. The diffusion continues until the particles are evenly distributed in space. This process is described mathematically by *Fick's first law of diffusion*

$$\phi(x) = -D[dN(x)/dx] \quad (\text{Eq. 11})$$

for a one-dimensional problem. This relation states that the rate of particle flow in the x direction [the particle flux density $\phi(x)$] is proportional to the negative gradient of the particle concentration $N(x)$ at each point x . The proportionality constant D is called the *diffusion coefficient* (cm²/s). Particles diffuse in the direction of decreasing particle concentration; the rate of flow depends on how steeply the concentration profile is graded; and no external driving force besides random thermal motion is involved.

The same diffusion mechanism applies to charge carriers if their concentration varies spatially, and the result is a diffusion current

$$\begin{aligned} J_n(\text{diff.}) &= -(-q) D_n [dn(x)/dx] \\ &= +q D_n [dn(x)/dx] \end{aligned} \quad (\text{Eq. 12A})$$

$$\begin{aligned} J_p(\text{diff.}) &= -(+q) D_p [dp(x)/dx] \\ &= -q D_p [dp(x)/dx] \end{aligned} \quad (\text{Eq. 12B})$$

In these equations $J(\text{diff.})$ is the current density (A/cm²) due to diffusion, and the n and p subscripts refer to electrons and holes, respectively. The current density J is related to the particle flux density ϕ by the charge on the carrier. As a result of the opposite charge on the two types of carriers, Eqs. 12 indicate that electrons and holes diffusing together in a particle gradient give rise to currents in opposite directions.

In most cases, electric fields are present along with gradients in the carrier concentrations. We can include drift in the field along with diffusion:

$$J_n(x) = q\mu_n n(x)\mathcal{E}(x) + qD_n [dn(x)/dx] \quad (\text{Eq. 13A})$$

drift diffusion

$$J_p(x) = q\mu_p p(x)\mathcal{E}(x) - qD_p [dp(x)/dx] \quad (\text{Eq. 13B})$$

The total current density is the sum of the electron and hole components:

$$J(x) = J_n(x) + J_p(x) \quad (\text{Eq. 14})$$

At equilibrium, there is no net current in a semiconductor. Thus, any fluctuation that would begin a diffusion current sets up an electric field that redistributes carriers by drift. An examination of the requirements for equilibrium indicates that the diffusion coefficient and mobility must be related by

$$D/\mu = kT/q \quad (\text{Eq. 15})$$

for either carrier type. This important equation is called the *Einstein relation*. It allows us to measure either D or μ and calculate the other.

In many carrier-diffusion problems, a steady-state distribution is maintained by a constant generation of excess carriers at some point in the semiconductor. For example, in Fig. 4, a steady excess hole concentration Δp is maintained by the injection of holes at $x = 0$ in a long semiconductor bar. Clearly, these excess holes will diffuse into the semiconductor; but they will also recombine with a characteristic lifetime τ_p . Thus the excess hole concentration $\delta p(x)$ decreases with distance into the semiconductor. This process is governed by the *steady-state diffusion equations*:

$$d^2 \delta n / dx^2 = \delta n / D_n \tau_n \equiv \delta n / L_n^2 \quad (\text{Eq. 16A})$$

$$d^2 \delta p / dx^2 = \delta p / D_p \tau_p \equiv \delta p / L_p^2 \quad (\text{Eq. 16B})$$

where,

$L_n \equiv \sqrt{D_n \tau_n}$ is called the *electron diffusion length*, L_p is the diffusion length for holes.

For the example of Fig. 4, the solution to Eq. 16B is

$$\delta p(x) = C_1 e^{x/L_p} + C_2 e^{-x/L_p} \quad (\text{Eq. 17})$$

We can evaluate C_1 and C_2 from the boundary conditions. Since recombination must reduce $\delta p(x)$ to zero for large values of x , $\delta p = 0$ at $x = \infty$, and therefore $C_1 = 0$. Similarly, the condition $\delta p = \Delta p$ at $x = 0$ gives $C_2 = \Delta p$, and the solution is

$$\delta p(x) = \Delta p e^{-x/L_p} \quad (\text{Eq. 18})$$

for the steady-state hole distribution. The injected excess hole concentration dies out exponentially in x due to recombination, and the diffusion length L_p represents the average distance a hole diffuses before recombining.

The steady-state distribution of excess holes causes diffusion, and therefore a hole current, in the direction of decreasing concentration. From Eqs. 12B and 18 we have

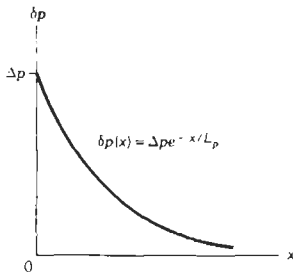


Fig. 4. Excess holes injected into a long n-type semiconductor diffuse and recombine, giving an exponential distribution.

$$\begin{aligned} J_p(x) &= -qD_p(dp/dx) = -qD_p(d\delta p/dx) \\ &= q(D_p/L_p)\Delta p e^{-x/L_p} = q(D_p/L_p)\delta p(x) \end{aligned} \quad (\text{Eq. 19})$$

Since $p(x) = p_o + \delta p(x)$, the space derivative involves only the excess holes. Notice that since $\delta p(x)$ is proportional to its derivative for an exponential distribution, the diffusion current at any x is just proportional to the excess density δp at that position.

Although this example seems rather restricted, it is typical of carrier injection in pn junctions. The injection of minority carriers across a junction often leads to exponential distributions as in Eq. 18, with the resulting diffusion current of Eq. 19.

pn JUNCTIONS

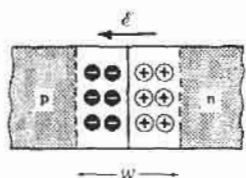
Most semiconductor devices contain at least one junction between p-type and n-type material. Such junctions are responsible for the injection and collection of charge carriers necessary for the operation of diodes, transistors, and other devices. Junctions are typically formed by crystal growth processes, by alloying, or by ion implantation or diffusion of doping impurities into crystals of the opposite conductivity type. We will consider here an abrupt junction between a uniformly doped p region and a uniformly doped n region of a single crystal semiconductor.

A Junction at Equilibrium

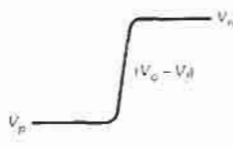
Let us consider separate regions of p- and n-type semiconductor material, brought together to form a junction (Fig. 5). Since the n material has a large concentration of electrons and few holes, and conversely for the p material, we expect considerable diffusion of carriers across the junction. Thus holes diffuse from the p side into the n side, and electrons diffuse from n to p.

Uncharged particles would diffuse throughout the material, resulting in a homogeneous distribution. This cannot occur in the case of the charged particles in a pn junction because of the development of space charge and the electric field \mathcal{E} in the neighborhood of the junction. Electrons diffusing from n to p leave behind uncompensated* donor ions (N_d^+) in the n material, and holes leaving the p region leave behind uncompensated acceptors (N_a^-); thus a region of positive space charge develops near the n side of the junction, and negative charge develops near the p side. The resulting

*Neutrality is maintained in the bulk materials of Fig. 5 by the presence of one electron for each ionized donor ($n = N_d^+$) in the n material, and one hole for each ionized acceptor ($p = N_a^-$) in the p material (neglecting minority carriers). Thus, if electrons leave n, some of the positive donor ions near the junction are left uncompensated. The donors and acceptors are fixed in the lattice, in contrast to the mobile electrons and holes.



(A) At equilibrium.



(B) With bias.

Fig. 5. A pn junction.

electric field is directed from the positive charge toward the negative charge. The influence of \mathcal{E} for each type of carrier is in the direction opposite to that of diffusion current. Therefore the field creates a drift component of current from n to p, opposing the diffusion current.

Since no net flow of electrons or holes across the junction can take place at equilibrium, the current due to the drift of each type of carrier in the \mathcal{E} field must exactly cancel the diffusion current in Eq. 13. Therefore, the electric field \mathcal{E} builds up to the point that the net current is zero at equilibrium. The electric field appears in some region W about the junction, and there is an equilibrium potential difference $V_0 = V_n - V_p$ between the two. The region W is called the *transition region* (also called the *depletion region*), and the potential V_0 is called the *contact potential*.

If we consider the junction to be made up of material with N_a acceptors per cubic centimeter on the p side

and N_d donors on the n side, we can show that the contact potential is

$$V_0 = (kT/q) \ln N_a N_d / n_i^2 \quad (\text{Eq. 20})$$

Forward and Reverse Bias

One of the useful features of a pn junction is the fact that current flows quite freely in the p-to-n direction when the p region has a positive external voltage bias relative to n (forward bias and forward current), whereas virtually no current flows when p is made negative relative to n (reverse bias and reverse current). This current asymmetry makes the pn junction diode very useful as a rectifier of ac signals, and forms the basis for many other applications of junctions.

Assume an applied voltage bias V appears across the transition region of the junction rather than in the neutral n and p regions. In typical devices, the neutral regions have low resistance, and the space charge regions about the junction take up most of the applied voltage. Since an applied voltage changes the electrostatic potential barrier and thus the electric field within the transition region, we would expect changes in the various components of current at the junction.

The *electrostatic potential barrier* at the junction is lowered by a forward bias, V_f , from the equilibrium contact potential, V_0 , to the smaller value $V_0 - V_f$. This lowering of the potential barrier occurs because a forward bias (p positive with respect to n) raises the electrostatic potential on the p side relative to the n side (Fig. 5B). For a reverse bias ($V = -V_r$) the opposite occurs; the electrostatic potential of the p side is depressed relative to the n side, and the potential barrier at the junction becomes larger ($V_0 + V_r$).

The *electric field* within the transition region decreases with forward bias, since the applied electric field opposes the built-in field. With reverse bias, the field at the junction is increased by the applied field.

The change in electric field at the junction calls for a change in the *transition region width*, W , since it is necessary that a proper number of positive and negative charges (in the form of uncompensated donor and acceptor ions) be exposed for a given value of the \mathcal{E} field. Thus, we expect the width, W , to decrease under forward bias (smaller \mathcal{E} , fewer uncompensated charges) and to increase under reverse bias.

The *diffusion current* is composed of majority-carrier electrons on the n side surmounting the potential-energy barrier to diffuse the p side, and holes surmounting their barrier from p to n. With forward bias, the barrier is lowered from V_0 to $V_0 - V_f$, and many more electrons in the n-side conduction band have sufficient energy to diffuse from n to p over the smaller barrier. Therefore, the electron diffusion current can be quite large with forward bias. Similarly, more holes can diffuse from p to n under forward bias because of the lowered barrier. For reverse bias, the barrier becomes so large ($V_0 + V_r$) that virtually no electrons in the n-side conduction band or holes in the p-side valence

band have enough energy to surmount it. Therefore, the diffusion current is usually negligible for reverse bias.

The *drift current* is relatively insensitive to the height of the potential barrier. The reason for this is that the drift current is limited *not* by *how fast* carriers are swept down the barrier, *but* rather *how often*. For example, minority-carrier electrons on the p side which wander into the transition region will be swept down the barrier by the \mathcal{E} field, giving rise to the electron component of drift current. However, this current is small not because of the size of the barrier, but because there are very few minority electrons in the p side to participate. Every p-side electron that diffuses to the transition region will be swept down the potential energy hill, whether the hill is large or small. Similar comments apply regarding the drift of minority holes from the n side to the p side of the junction. To a good approximation, therefore, the electron and hole drift currents at the junction are independent of the applied voltage.

The supply of minority carriers on each side of the junction required to participate in the drift component of current is generated by thermal excitation of electron-hole-pairs (EHPs). For example, an EHP created near the junction on the p side provides a minority electron in the p material. If the EHP is generated within a diffusion length L_n of the transition region, this electron can diffuse to the junction and be swept down the barrier to the n side. The resulting current is commonly called the *generation current* since its magnitude depends entirely upon the rate of generation of EHP. This generation current can be greatly increased by optical excitation of EHPs near the junction (the pn junction *photodiode*).

The *total current* crossing the junction is composed of the sum of the diffusion and drift components. The electron and hole diffusion currents are both directed from p to n (although the particle flow directions are opposite to each other), and the drift currents are from n to p. The *net* current crossing the junction is zero at equilibrium, since the drift and diffusion components cancel for each type of carrier. Under reverse bias, both diffusion components are negligible because of the large barrier at the junction, and the only current is the relatively small (and essentially voltage-independent) generation current from n to p.

We expect from Fig. 4 that injection of excess holes from p to n will produce a *distribution* of excess holes in the n material. As the holes diffuse deeper into the n region, they recombine with electrons, and the resulting excess hole distribution is obtained as a solution of the diffusion equation, Eq. 16B. If the n region is long compared with the hole diffusion length, L_p , the solution is exponential, as in Fig. 4. Similarly, the injected electrons in the p material diffuse and recombine, giving an exponential distribution of excess electrons. These excess carrier distributions must be accounted for in switching a pn junction from forward bias to reverse bias.

By evaluating the diffusion currents of electrons and holes across the junction, it can be shown that the current-voltage characteristic of the diode is given by

$$I = qA[(D_p/L_p)p_n + (D_n/L_n)n_p](e^{qV_r/kT} - 1) \quad (\text{Eq. 21})$$

Equation 21 is called the *diode equation*, which describes the total current through the diode for either forward or reverse bias. We can calculate the current with reverse bias (commonly called the reverse saturation current) by letting $V = -V_r$,

$$I = qA[(D_p/L_p)p_n + (D_n/L_n)n_p](e^{-qV_r/kT} - 1) \quad (\text{Eq. 22})$$

If V_r is larger than a few kT/q , the total current is just the reverse saturation current

$$I = -qA[(D_p/L_p)p_n + (D_n/L_n)n_p] \quad (\text{Eq. 23})$$

The *I-V* characteristic for the pn junction is shown in Fig. 6. It exhibits conduction for forward bias (V positive), negligible current for moderate reverse bias, and avalanche breakdown at a large reverse bias (V_{br}).

Reverse Breakdown

The small reverse saturation current typical of reverse bias in a junction is valid until a critical reverse bias is reached, for which *reverse breakdown* occurs (Fig. 6). At this critical voltage (V_{br}), the reverse current through the diode increases sharply, and relatively large currents can flow with little further increase in voltage. The

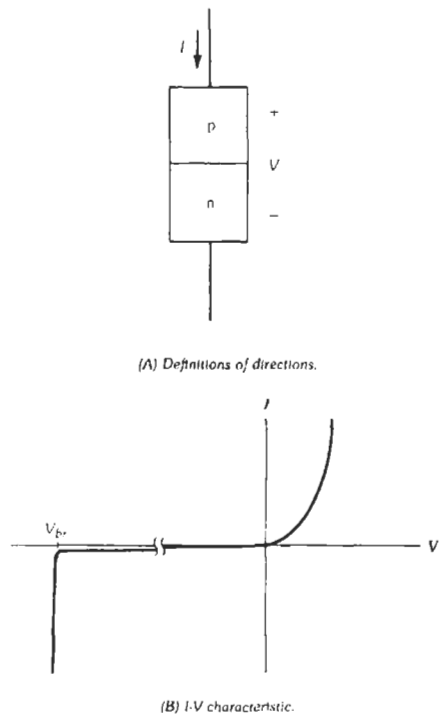


Fig. 6. Current and voltage in a pn junction.

existence of a critical breakdown voltage introduces almost a right-angle appearance to the reverse characteristics.

There is nothing inherently destructive about reverse breakdown. If the current is limited to a reasonable value by the external circuit, the pn junction can be operated safely in the reverse-breakdown condition. Useful devices called *breakdown diodes* are designed to operate in the reverse-breakdown regime of their characteristics.

Reverse breakdown can occur by two mechanisms, each of which requires a critical electric field in the junction transition region. The first mechanism, called the *zener effect*, is operative at low voltages (up to a few volts reverse bias). If the breakdown occurs at higher voltages (from a few volts to hundreds of volts), the mechanism is *avalanche breakdown*.

Zener breakdown occurs in heavily doped junctions in which the transition from the p side to the n side is very abrupt. Under these conditions, the potential barrier separating carriers on the two sides of the junction is very narrow, and a quantum-mechanical process called *tunneling* can take place. In effect, electrons from the p-side valence band are able to tunnel across the barrier into the empty states in the n-side conduction band when a small reverse bias is applied. The result is a strong current from n to p in the diode, causing zener breakdown. We shall discuss tunneling again in relation to tunnel diodes in a later section.

For diodes with lighter doping or with a graded junction, the avalanche mechanism is the operative breakdown process. In this effect, the electric field in the transition region becomes strong enough under reverse bias to accelerate carriers to quite high kinetic energy. For example, if a minority hole is generated thermally on the n side and wanders into the transition region, it is swept by the junction field to the p side. If the field is high enough, the hole can collide with the lattice with sufficient energy to create an electron-hole pair. As a result, the newly created hole is swept to the p side and the electron is swept to the n side. With one carrier entering the transition region, we in effect have three carriers leaving. This *carrier multiplication* adds to the reverse current. Furthermore, if on the average each carrier (including carriers created by ionizing collisions) can create one EHP during its transit across the transition region, a runaway (avalanche) process develops. Avalanche breakdown usually occurs at a sharply defined breakdown voltage which corresponds to the establishment of a critical field in the transition region. This critical field is that necessary to cause a series of ionizing collisions by carriers in transit across the transition region.

Capacitance and Transient Behavior

Since many pn-junction devices are used in ac circuits or for switching applications, it is important to consider their time-dependent properties. The transient

behavior of a junction is influenced strongly by the presence of stored charge. For example, the charge (Q) due to excess minority carriers in a forward-biased junction must be changed as the current varies. Since this charge cannot be altered instantaneously, current and voltage may vary in a complex fashion for the transient case.

The problem of stored charge is particularly important in applications that call for switching a diode from forward conduction to the reverse-biased state. In this case, the distribution of excess minority carriers is swept across the junction, giving rise to a brief pulse of reverse current after the applied voltage has been switched negative (Fig. 7). The time this reverse current ($-I_r$) lasts is called the *storage delay time*, t_{sd} . This time is proportional to the minority-carrier lifetime, τ , in the region where the charge is stored. After the charge distribution has properly rearranged to that expected for a reverse-biased junction, the small reverse saturation current is reached. However, the relatively large reverse current during switching must be accounted for in circuit designs. Special switching devices can be constructed which store very little charge in forward bias, and therefore are able to switch from the forward to the reverse state with a minimum of time delay or reverse current.

When ac signals are applied to a junction, the necessity for buildup and depletion of stored charge causes the voltage to lag behind the current. The resulting capacitance is an important junction property. For a reverse-biased junction, the capacitance is due to the changing charge in the depletion region (Fig. 5). As the reverse bias increases, W must increase to accommodate the increased charge on each side of the junction. We calculate the junction capacitance, C_j , from the rate of change of charge within W as the voltage is varied. The result is equivalent to the usual parallel-plate capacitor formula

$$C_j = \epsilon A/W \tag{Eq. 24}$$

where,

- W is the depletion width,
- ϵ is the permittivity of the semiconductor.

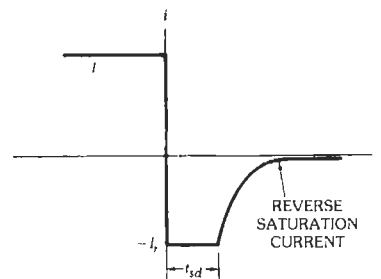


Fig. 7. Reverse recovery transient for a diode switched from forward conduction (I) to reverse bias.

Since W depends on the reverse bias, C_j is a *voltage-variable capacitance*. This property can be used in tuning circuits, as discussed below.

For forward bias, ac variations impressed on a steady forward current I (dc) call for variations in stored charge of the minority carrier distribution (Fig. 4). A forward-biased pn diode in which hole injection dominates responds to ac signals with an equivalent RC time constant τ_p required for redistribution of the charge due to holes stored in the n region. The current in a forward-biased junction ($V \gg kT/q$) is

$$I = I_s e^{qV/kT} \quad (\text{Eq. 25})$$

where I_s is the magnitude of the saturation current, described in Eq. 23.

If we take the derivative with respect to V , we find the ac conductance to be

$$G_{ac} = (q/kT) I_s e^{qV/kT} = (q/kT) I \quad (\text{Eq. 26})$$

Thus charge-storage capacitance $C_s = \tau_p G$ is

$$C_s = (q \tau_p / kT) I \quad (\text{Eq. 27})$$

In summary, the reverse-bias junction capacitance, C_j , varies with the reverse bias voltage, V_r , and the forward bias capacitance, C_s , varies with the forward bias current, I . The storage delay time is an important limitation in switching circuits, and the forward-bias capacitance is important in high-frequency ac circuits. In either case, the effects of charge storage can be minimized by reducing the carrier lifetime in the minority-carrier storage region. One approach used in silicon devices involves doping the junction with gold, which introduces a high density of recombination centers, thereby reducing carrier lifetime.

DIODES AND RECTIFIERS

The most obvious property of a pn junction is its *unilateral* nature; that is, to a good approximation it conducts current in only one direction. We can think of an *ideal diode* as a short circuit when forward biased and an open circuit when reverse biased (Fig. 8A). The pn junction diode does not quite fit this description, but the I - V characteristics of many junctions can be approximated by the ideal diode in series with other circuit elements to form an equivalent circuit. For example, most forward-biased diodes exhibit an *offset voltage*, E_o , which can be approximated in a circuit model by a battery in series with the ideal diode (Fig. 8B). The series battery in the model keeps the ideal diode turned off for applied voltages less than E_o . In some cases, the approximation to the actual diode characteristic is improved by adding a series resistor, R . The values of E_o and R depend on the particular diode, and can be obtained from the actual I - V characteristic of the device. The circuit approximations illustrated in Fig. 8B are called *piecewise-linear equivalents*, since the

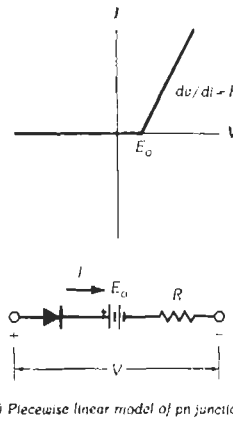
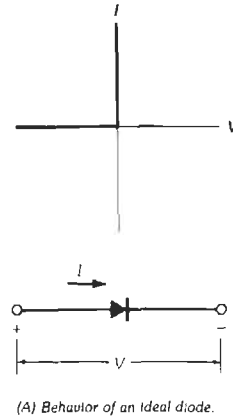


Fig. 8. Semiconductor diode.

approximate characteristics are linear over specific ranges of voltage and current.

An ideal diode can be placed in series with an ac voltage source to provide *rectification* of the signal. Since current can flow only in the forward direction through the diode, only the positive half-cycles of an input sine wave are passed to the load. During the negative half-cycles, the diode is nonconducting, and the resulting voltage across a load resistor is zero. The output voltage is a *half-rectified sine wave*. The rectified signal has an average value, and therefore contains a dc component. By appropriate filtering, this dc level can be extracted from the rectified signal.

The unilateral nature of diodes is useful for many other circuit applications that require *waveshaping*. This involves alteration of ac signals by passing only certain portions of the signal while blocking other portions.

Other properties of pn junctions can be exploited for electronic-device applications. For example, the *breakdown diode* (sometimes called *zener diode*) makes use of the abrupt reverse breakdown shown in Fig. 6. The breakdown voltage is strongly dependent on doping, and junctions can be made with accurately chosen values of V_{br} . Breakdown diodes can be used as *voltage*

regulators in circuits with varying inputs. Such a device can also be used as a *reference diode*; since the breakdown voltage of a particular diode is known, the voltage across it during breakdown can be used as a reference in circuits that require a known value of voltage.

The junction capacitance described by Eq. 24 can be put to use in a device called a *varactor*. Since C_j is inversely proportional to the depletion width, W , and W varies with the reverse bias, V_r , the junction is a *voltage-variable capacitor*. For an abrupt junction, C_j is proportional to $V_r^{-1/2}$. Thus it is possible to vary C_j by choosing values of the applied bias. A set of varactors can be used, for example, to replace variable-plate capacitors in the tuning section of a radio receiver. By proper adjustment of doping profiles in the junction, the voltage dependence of C_j can be tailored for the specific application.

In abrupt pn junctions that are very heavily doped on both sides, current can be carried by the quantum-mechanical process called tunneling. In reverse bias, a tunnel diode passes current at very small voltage, and the same occurs for a fraction of a volt of forward bias (Fig. 9). However, at a critical forward voltage (V_p), the tunneling decreases, and eventually the usual forward diode current dominates. As Fig. 9 indicates, a range of the I - V characteristic exists in which the tunneling current decreases with increasing voltage, giving rise to a *negative resistance* region. Such a negative resistance can be used in various switching, oscillation, amplification, and other circuit functions.

When light shines on a pn junction, electron-hole pairs are generated by optical excitation. The collection of these optically generated charge carriers across the junction results in a net excess current from n to p (i.e., I - V curve A in Fig. 10 is moved down to curve B). The shift of the I - V curve is proportional to the optical generation rate of excess electrons and holes, and therefore to the intensity of the light. In the third quadrant (III), the junction can be used as a *photodiode*, in which the reverse current at a fixed reverse bias voltage is proportional to the intensity of the light. When the junction operates in the fourth quadrant of the I - V curve (IV), power is delivered from the junction to

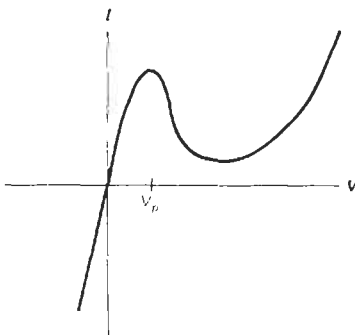


Fig. 9. I - V characteristic for a tunnel diode.

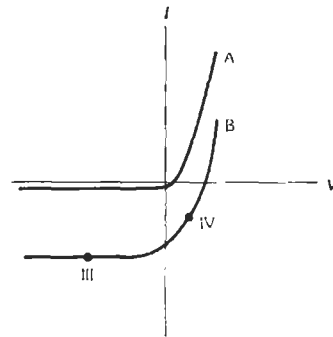


Fig. 10. I - V characteristics for a photodiode in the dark (A) and with light applied (B).

the external load. (Notice that at IV, the diode current is negative while the junction voltage is positive. Thus, the current in this case flows from minus to plus relative to the voltage, as in a battery.) In this case, the pn junction can be used as a *solar cell* (or *photovoltaic junction*). Such solar cells are widely used as power sources for satellites and in certain terrestrial applications.

In many compound semiconductors (e.g., GaAs, InP, and alloys such as InGaAsP), electrons and holes recombine with each other at a forward-biased pn junction, giving off light. The wavelength of the resulting light depends on the material used, generally ranging from the infrared (e.g., GaAs) to the green (GaP). By proper choice of materials, *light-emitting diodes* (LEDs) can be made that cover a wide range of wavelengths. Such LEDs are useful as visible indicator lamps and in alphanumeric displays. In addition, infrared LEDs can be used in optical communication systems in which light signals are sent over optical fibers and are detected at the other end with a photodiode. Such optoelectronic systems can be made very efficient if the compound semiconductor diode is made in the form of a laser and the photodetector operates in the reverse-bias avalanche condition. For further discussion, see Chapters 21 and 22.

TRANSISTORS

The basic building block of modern electronics is the transistor. This device, which has replaced the vacuum tube in all but a few special applications, is used in most circuits requiring amplification or switching. Complex integrated circuits are also based on the transistor. This section deals with the operation of three transistor types—the junction field-effect transistor, the MOS field-effect transistor, and the bipolar junction transistor. We concentrate here on the internal operation of these devices, reserving treatments of their circuit operation for later chapters. Transistor fabrication methods are discussed in Chapter 20.

As an amplifying device, the transistor converts weak time-varying signals into strong signals. As a switching element, it can be changed from a conducting state to a

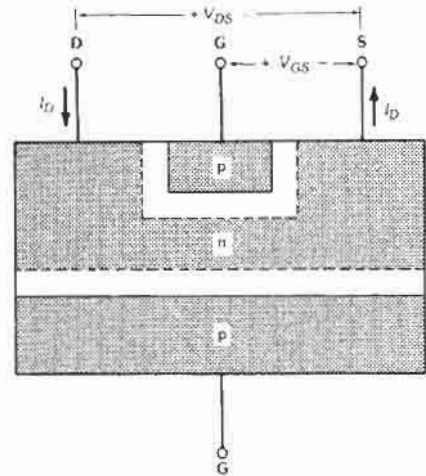
nonconducting state quickly with the application of very little control power. In the field-effect transistor, current through two terminals is varied by voltage applied to a third terminal. In the bipolar transistor, the current through two terminals is controlled by a small current applied to a third terminal.

Junction Field-Effect Transistors

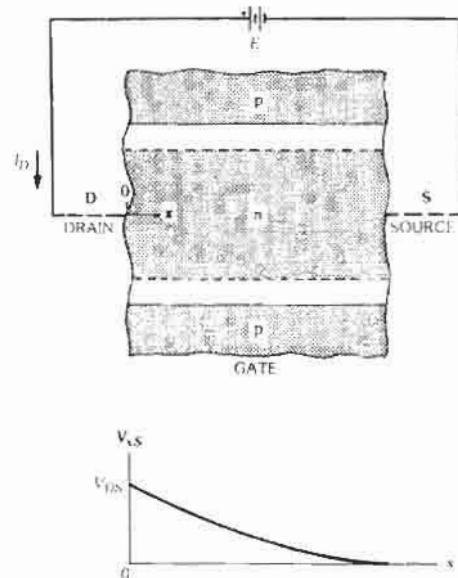
In Fig. 5, the width of the depletion region (W) under reverse bias is controlled by the voltage applied to the junction. This property is the basis for the operation of the *junction field-effect transistor (JFET)*. In a JFET, the voltage-variable depletion-region width of a junction is used to control the effective cross-sectional area of a conducting *channel*. In the device of Fig. 11, current I_D flows through an n-type channel between two p regions. A reverse bias between these p regions and the channel causes the depletion regions to intrude into the n material, and therefore the effective width of the channel can be restricted. Since the resistivity of the channel region is fixed by its doping, the channel resistance varies with changes in the effective cross-sectional area.

In Fig. 11, electrons in the n-type channel drift from right to left, opposite to the current. The end of the channel from which electrons flow is called the *source*, and the end toward which they flow is called the *drain*. The p regions are called *gates*. If the channel were p-type, holes would flow from source to drain, in the same direction as the current, and the gate regions would be n. Voltage V_{GS} refers to the potential from each gate region, G, to the source, S. Since the conductivity of the heavily doped gate regions is high, we can assume that the potential is uniform throughout each gate. In the lightly doped channel material, however, the potential varies with position (Fig. 11B). For low values of current, we can consider the channel to be a distributed resistor, and assume that voltage V_{xS} varies linearly from V_{DS} at the drain end to zero at the source end of the channel.

In Fig. 12, assume the gates are short-circuited to the source ($V_{GS} = 0$), such that the potential at $x = L$ is the same as the potential everywhere in the gate regions. For very small currents, the widths of the depletion regions are close to the equilibrium values (Fig. 11A). As current I_D is increased, however, it becomes important that V_{xS} is large near the drain end and small near the source end of the channel. Since the reverse bias across each point in the gate-to-channel junction ($-V_{Gx}$) is simply V_{xS} when V_{GS} is zero, we can estimate the shape of the depletion regions as in Fig. 12A. The reverse bias is relatively large near the drain ($-V_{GD} = V_{DS}$) and decreases toward zero near the source. As a result, the depletion region intrudes into the channel near the drain, and the effective channel area is constricted. Since the resistance of the constricted channel is higher, the I - V plot for the channel begins to depart from the straight line that was valid at low current



(A) Transistor geometry

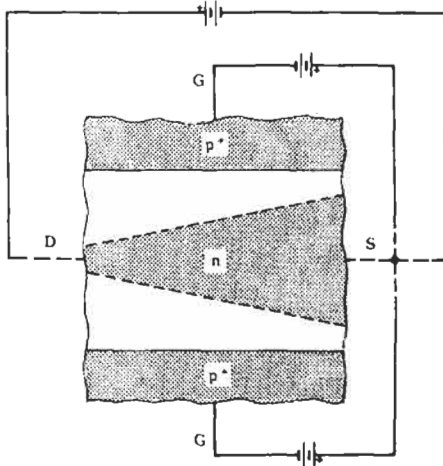


(B) Detail of channel and voltage variation along channel with $V_{DS} = 0$ and small I_D .

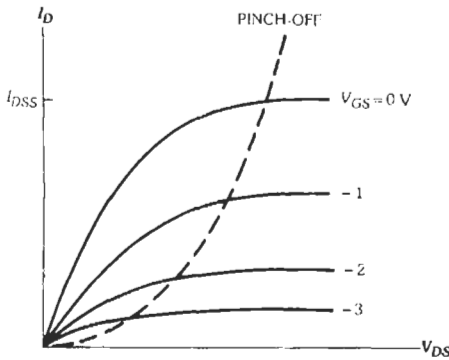
Fig. 11. Simplified cross-sectional view of a junction FET.

levels. As V_{DS} is increased, there must be some bias voltage at which the depletion regions meet near the drain and essentially *pinch off* the channel. When this happens, current I_D cannot increase significantly with further increase in V_{DS} . For higher voltages, the current is *saturated* approximately at its value at pinch-off.

The effect of a negative gate bias, $-V_{GS}$, is to increase the resistance of the channel and induce pinch-off at a lower value of current (Fig. 12B). Since the depletion regions are larger with V_{GS} negative, the slopes of the I_D -versus- V_{DS} curves below pinch-off become smaller. As current I_D increases, the pinch-off condition is reached at a lower drain-to-source voltage,



(A) Variation of depletion-region widths.



(B) Family of current-voltage curves for the channel as V_{GS} is varied.

Fig. 12. Effects of negative gate bias in the JFET.

and the saturation current is lower than for the case of zero gate bias. As V_{GS} is varied, a family of curves is obtained for the I - V characteristic of the channel, as in Fig. 12B. By varying the gate bias, we can obtain amplification of an ac signal or switch the device from its off ($I_D = 0$) to its on condition. Since the input control voltage, V_{GS} , appears across the reverse-biased gate junctions, the input impedance of the device is relatively high.

The JFET is generally used in the saturation region of its I - V characteristic (i.e., beyond pinch-off). The value of I_D at saturation (called I_{DSS}) depends on V_{GS} , as Fig. 12B indicates. To a good approximation, I_{DSS} beyond pinch-off can be written as follows:

$$I_{DSS} \approx I_{DSS}(1 + V_{GS}/V_p)^2 \quad (V_{GS} \text{ negative}) \quad (\text{Eq. 28})$$

In this expression, I_{DSS} is the saturation value with the gate shorted to the source (i.e., with $V_{GS} = 0$). The

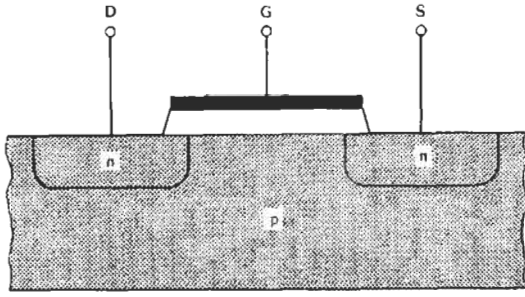
term V_p is the pinch-off voltage defined as the gate bias required to pinch off the channel with $I_D = 0$. The quantity V_p is dictated by the device geometry and doping. For an n-channel device, V_{GS} is a negative number. Thus, Eq. 28 indicates that the saturation drain current, I_{DSS} , is greatest (I_{DSS}) when V_{GS} is zero, and decreases to zero when V_{GS} becomes as negative as the pinch-off voltage V_p is positive.

MOS Field-Effect Transistors

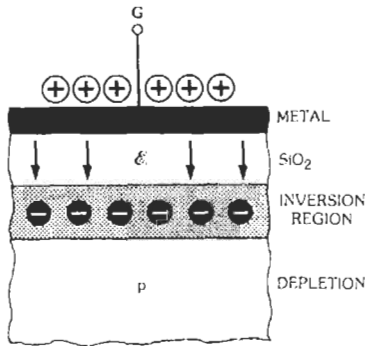
In the second class of FETs, the channel current is controlled by a voltage at a gate electrode that is isolated from the channel by an insulator. The resulting device is called an *insulated-gate field-effect transistor (IGFET)* (Fig. 13). Since the channel exists at the semiconductor surface, such a device is also called a *surface FET*. In the most common configuration, an oxide layer is grown or deposited on the semiconductor surface, and the metal gate electrode is deposited onto this oxide layer. This structure is commonly called a *metal-oxide-semiconductor transistor (MOST)*.

The n source and drain regions of Fig. 13 are diffused into a high-resistivity p substrate. The channel region may be a thin diffused n layer, or more commonly an *induced inversion region*. If an n-type diffused channel is included between source and drain, the effect of the field is to raise or lower the conductance of the channel by either depleting or enhancing the electron density in the channel. If the gate voltage is positive, the conductivity of the channel is enhanced, whereas a negative gate voltage tends to deplete the channel of electrons. Thus a diffused-channel MOST can be operated in either the *depletion* or *enhancement* mode.

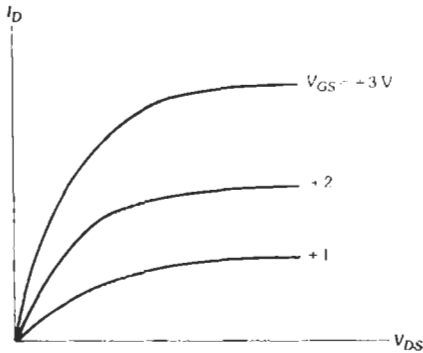
We shall concentrate on the *induced-channel MOST*, in which no diffused n-type region exists between source and drain at equilibrium. As with the JFET of the previous section, we shall consider an n-channel device; the corresponding p-channel case can be deduced from this discussion. When a positive gate voltage is applied to this structure, a depletion region is formed in the p material, and a thin layer of mobile electrons is drawn from the source and drain into the channel (Fig. 14A). Where the mobile electrons dominate, the material is effectively n-type. This is called an *inversion layer*, since that material was originally p-type. Once an inversion layer is formed near the semiconductor surface, a conducting channel exists from source to drain. The operation of the device is then quite similar to that of the JFET. The channel conductance is controlled by field \mathcal{E}_i in the insulator, but the magnitude of this field varies along the channel as V_{Gx} varies from $V_{GS} - V_{DS}$ at the drain to V_{GS} at the source (Fig. 14C); this corresponds qualitatively to the variation in junction reverse bias along the channel for the JFET. Since a positive voltage is required between the gate and each point x in the channel to maintain inversion, a large enough value of V_{DS} can cause \mathcal{E}_i to go to zero at the drain. This corresponds to the pinch-off condition. Once pinch-off is reached, the



(A) Cross-sectional view.



(B) Expanded view of the channel with positive polarity of gate bias.



(C) I-V characteristics as gate bias is varied.

Fig. 13. The induced-channel MOST.

saturation current remains essentially constant, as in the JFET.

The saturation drain current for the MOST is described by an equation of the same form as Eq. 28 for the JFET. For the MOS transistor with an induced channel,

$$I_{DS} = (\mu_n C_o / 2L^2) (V_{GS} - V_T)^2 \quad (\text{Eq. 29})$$

where,

μ_n is the electron mobility in the channel,
 C_o is the capacitance between gate and channel,

L is the length of the channel from source to drain, V_T is the *threshold voltage*, the gate bias that must be applied to induce a channel in the underlying semiconductor.

The value of the threshold voltage is dependent on several properties of the MOS structure. For example, differences in work function between the metal gate and the semiconductor must be balanced by the applied gate voltage. Furthermore, trapped charges at the semiconductor-oxide interface and in the oxide layer may induce image charges in the channel region that must also be overcome by the applied gate voltage. In many cases, several volts must be applied to the gate to create an inversion region (and therefore a conducting channel) in the semiconductor. The existence of such a threshold voltage can be a serious limitation in circuits that are designed to operate with low-voltage batteries. Therefore, much effort has been devoted to designing MOS structures with low values of V_T . Techniques used for this purpose include matching the work function of the gate and the underlying silicon by employing heavily doped polycrystalline silicon as a gate electrode, rather than the more common aluminum metal gate. Silicon crystals oriented along certain crystallographic directions that minimize surface-charge effects are also used to reduce V_T . One of the most successful methods for controlling the threshold voltage is the use of ion implantation in the channel region. For example, in a p-channel MOST, the source and drain regions are diffused with boron impurities, a thin gate oxide layer is grown, and then boron ions are implanted through the oxide region into the channel. Since the concentration of implanted impurities can be precisely controlled, a sufficient number of boron atoms can be implanted into the channel such that the gate bias required to create an inversion region is greatly reduced. Such low-threshold transistors can be used in circuits that are driven by quite low-voltage battery sources.

The MOS transistor has the advantage of extremely high input impedance between the gate and source electrodes, since these terminals are separated by an oxide layer. The input impedance of a MOST can be on the order of 10^{14} ohms, making this device very useful in amplifying signals delivered by circuits that are sensitive to loading.

The MOST is useful in integrated circuits utilizing silicon planar technology. For many years, it was difficult to control the surface charge and other effects in MOS structures, but these devices can now be manufactured in large numbers with good reproducibility, as discussed in Chapter 20.

Bipolar Junction Transistors

The BJT is a current-controlled device, in which the current through two terminals is controlled by a relatively small current in a third terminal. We shall begin the discussion by considering the reverse-biased pn junction. As discussed earlier, the reverse saturation current through a junction depends on the rate at which

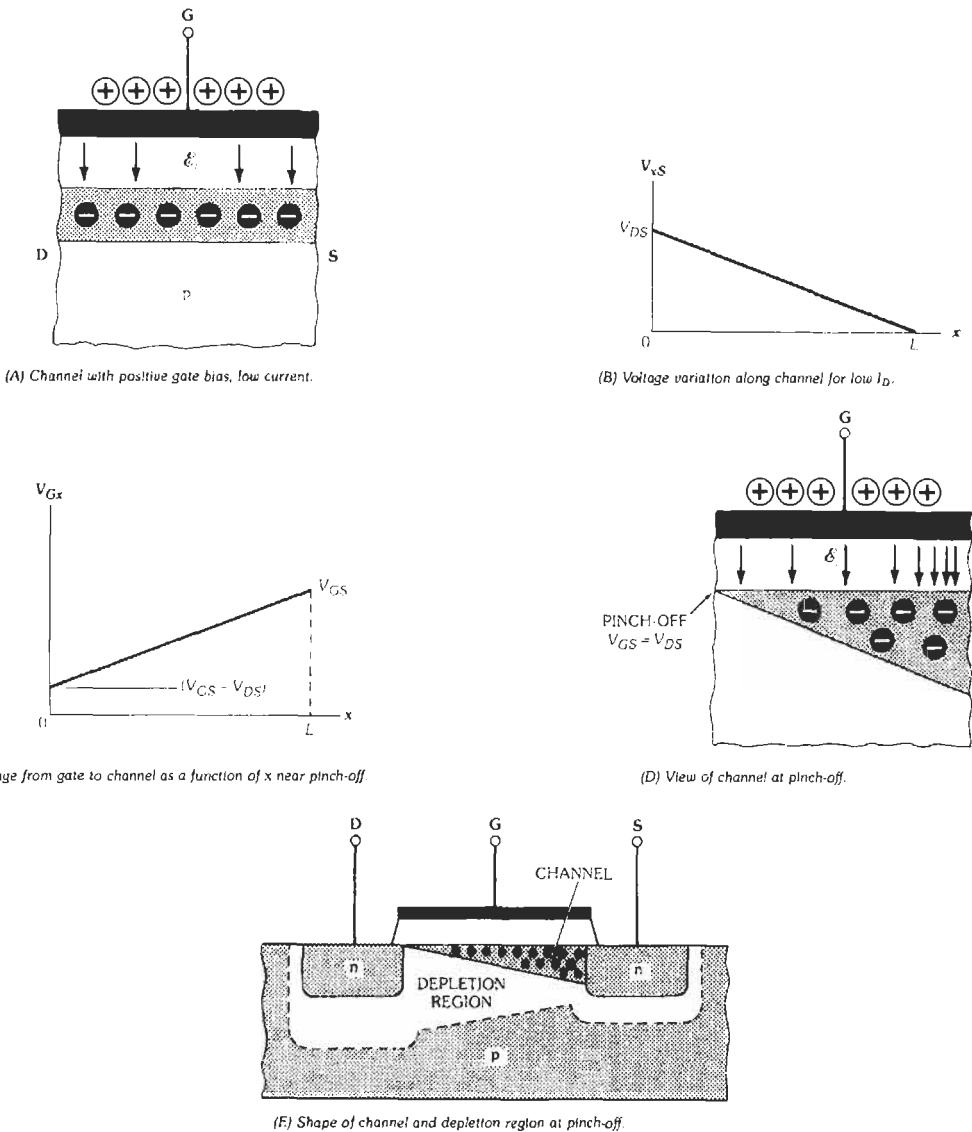


Fig. 14. Properties of an induced-channel MOST.

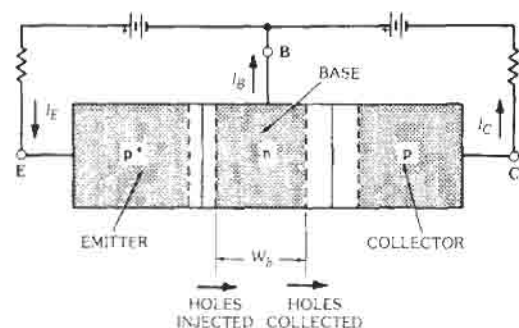
minority carriers are generated in the neighborhood of the junction. For example, the reverse current due to holes being swept from n to p is essentially independent of the size of the junction \mathcal{E} field and hence independent of the reverse bias. The reason for this is that the hole current depends on how often minority holes are generated by EHP creation within a diffusion length of the junction, and not upon how fast a particular hole is swept across the depletion layer by the field. As a result, it is possible to increase the reverse current through the diode by increasing the rate of EHP generation. One convenient method for accomplishing this is optical

excitation of EHPs with light, as in the photodiode. With steady photoexcitation, the reverse current will still be essentially independent of bias voltage, and if the dark saturation current is negligible, the reverse current is directly proportional to the optical generation rate, g_{op} (Fig. 10).

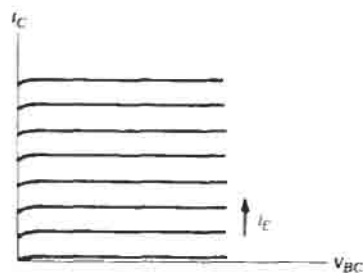
The example of external control of current through a junction by optical generation raises an interesting question: Is it possible to inject minority carriers into the neighborhood of the junction *electrically* instead of optically? If so, we could control the junction reverse current by simply varying the rate of minority-carrier

injection. For example, let us consider a hypothetical hole-injection device that can inject holes at a predetermined rate into the n side of the junction. The effect on the junction current will resemble the effects of optical generation. The current from n to p will depend on the hole injection rate (similar to the reverse current in Fig. 10) and will be essentially independent of the bias voltage. There are several obvious advantages to such external control of a current; for example, the current through the reverse-biased junction would change very little if the load resistance were altered, since the magnitude of the junction voltage is relatively unimportant. Therefore, such an arrangement should be a good approximation to a controllable constant-current source.

A convenient hole-injection device is a forward-biased p^+n junction (p^+ refers to a very heavily doped p region). The current in such a junction is due primarily to holes injected from the p^+ region into the n material. If we make the n side of the forward-biased junction the same as the n side of the reverse-biased junction, the p^+np structure of Fig. 15 results. With this configuration, injection of holes from the p^+n junction into the center n region supplies the minority-carrier holes to participate in the current through the reverse-biased np junction. Of course, it is important that the injected holes do not recombine in the n region before they can diffuse to the depletion layer of the reverse-biased junction. Thus we must make the n region narrow compared with a hole diffusion length.



(A) Schematic representation of PNP device with forward biased emitter junction and reverse biased collector junction.



(B) I_C characteristics of the reverse-biased np junction as a function of emitter current.

Fig. 15. A PNP transistor.

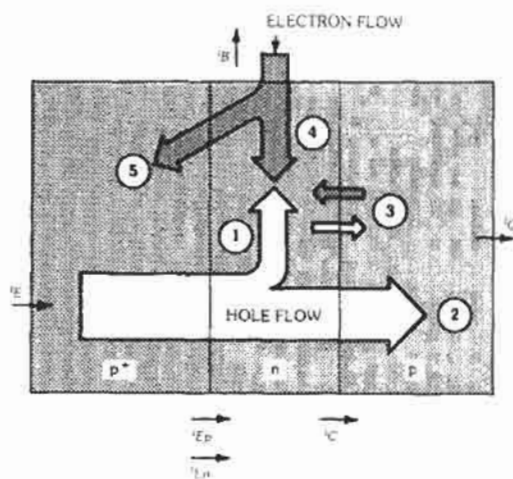
The structure we have described is a PNP bipolar junction transistor. The forward-biased junction that injects holes into the center n region is called the *emitter junction*, and the reverse-biased junction that collects the injected holes is called the *collector junction*. The p^+ region that serves as the source of injected holes is called the *emitter*, and the p region into which the holes are swept by the reverse-biased junction is called the *collector*. The center n region is called the *base*. The biasing arrangement of Fig. 15 is called the *common-base* configuration, since the base electrode (B) is common to the emitter and collector circuits.

To have a good PNP transistor, almost all of the holes injected by the emitter into the base should be collected. Thus the n-type base region should be narrow, and the hole lifetime, τ_p , should be long. This requirement is summed up by specifying $W_b \ll L_p$, where W_b is the length of the *neutral* n material of the base (measured between the depletion regions of the emitter and collector junctions) and L_p is the diffusion length for holes in the base, $(D_p \tau_p)^{1/2}$. With this requirement satisfied, an average hole injected at the emitter junction will diffuse to the depletion region of the collector junction without recombination in the base. A second requirement is that current I_E crossing the emitter junction should be composed almost entirely of holes injected into the base, rather than electrons crossing from base to emitter. This requirement is satisfied by doping the base region lightly compared with the emitter, so that the p^+n emitter junction of Fig. 15 results.

It is clear that current I_E flows into the emitter of a properly biased PNP transistor and that I_C flows out at the collector, since the direction of hole flow is from emitter to collector. However, base current I_B requires a bit more thought; I_B flows out since this current supplies electrons to the base region. In a good transistor, the base current will be very small since I_E is essentially hole current and the collected hole current, I_C , is almost equal to I_E . There must be some base current, however, due to requirements of electron flow into the n-type base region (Fig. 16). We can account for I_B physically by three dominant mechanisms:

- There must be some recombination of injected holes with electrons in the base, even with $W_b \ll L_p$. The electrons lost to recombination must be resupplied through the base contact.
- Some electrons will be injected from n to p in the forward-biased emitter junction, even if the emitter is heavily doped compared to the base. These electrons must also be supplied by I_B .
- Some electrons are swept into the base at the reverse-biased collector junction due to thermal generation in the collector. This small current reduces I_B by supplying electrons to the base.

The dominant mechanism in the base current is usually recombination, and we can often approximate the base current by calculating the recombination rate in



- 1 Injected holes lost to recombination in the base
- 2 Holes reaching the reverse biased collector junction
- 3 Thermally generated electrons and holes making up the reverse saturation current of the collector junction
- 4 Electrons supplied by the base contact for recombination with holes
- 5 Electrons injected across the forward biased emitter junction

Fig. 16. Hole and electron flow in a PNP transistor with normal biasing.

the base. In a well-designed transistor, i_B will be a very small fraction (perhaps one-hundredth) of i_E .

In an NPN transistor, the three current directions are reversed, since electrons flow from emitter to collector and holes must be supplied to the base.* The physical mechanism for operation of the NPN device can be understood simply by reversing the roles of electrons and holes in the PNP discussion.

The BJT is useful in amplifiers because the currents at the emitter and collector are controllable by the relatively small base current. The essential mechanisms are easy to understand if various secondary effects are neglected. We shall use total current (dc plus ac) in this discussion, with the understanding that the simple analysis applies only to dc and to small-signal ac at low frequencies. We can relate the terminal currents of the transistor (i_E , i_B , and i_C) by several important factors. In this introduction, we shall neglect the saturation current at the collector (3 in Fig. 16) and such effects as recombination in the transition regions. Under these assumptions, the collector current is made up entirely of those holes injected at the emitter that are not lost to recombination in the base. Thus i_C is proportional to the hole component of the emitter current i_{Ep} .

$$i_C = Bi_{Ep} \quad (\text{Eq. 30})$$

*In the present discussion of device operation, we show the directions of the currents as they actually flow in normal operation. In circuit analysis, however, it is common to define all currents as flowing into the transistor, and incorporate minus signs where they are needed.

Proportionality factor B is simply the fraction of injected holes that make it across the base to the collector; B is called the *base transport factor*. The total emitter current, i_E , is made up of the hole component, i_{Ep} , and an electron component, i_{En} , due to electrons injected from base to emitter (5 in Fig. 16). The *emitter injection efficiency*, γ , is

$$\gamma = i_{Ep}/(i_{En} + i_{Ep}) \quad (\text{Eq. 31})$$

For an efficient transistor, B and γ should each be very near unity; that is, the emitter current should be due mostly to holes ($\gamma = 1$), and most of the injected holes should eventually participate in the collector current ($B = 1$). The relation between the collector and emitter currents is

$$i_C/i_E = Bi_{Ep}/(i_{En} + i_{Ep}) = B\gamma \quad (\text{Eq. 32})$$

The product $B\gamma$ is defined as the factor α , which represents the emitter-to-collector current amplification. There is no real amplification between these currents, since α is smaller than unity. On the other hand, the relation between i_C and i_B is more promising for amplification.

In accounting for the base current, we must include the rates at which electrons are lost from the base by injection across the emitter junction (i_{En}) and the rate of electron recombination with holes in the base. In each case, the lost electrons must be resupplied through the base current, i_B . If the fraction of injected holes making it across the base *without* recombination is B , then it follows that $(1 - B)$ is the fraction *recombining* in the base. Thus the base current is

$$i_B = i_{En} + (1 - B)i_{Ep} \quad (\text{Eq. 33})$$

neglecting the collector saturation current. The relation between the collector and base currents is found from Eq. 30 and Eq. 33:

$$\begin{aligned} i_C/i_B &= Bi_{Ep}/[i_{En} + (1 - B)i_{Ep}] \\ &= B[i_{Ep}/(i_{En} + i_{Ep})]/\{1 - B[i_{Ep}/(i_{En} + i_{Ep})]\} \\ &= B\gamma/(1 - B\gamma) = \alpha/(1 - \alpha) \equiv \beta \quad (\text{Eq. 34}) \end{aligned}$$

The factor β relating the collector current to the base current is the *base-to-collector current-amplification factor*. Since α is near unity, it is clear that β can be large for a good transistor, and the collector current is large compared with the base current.

It remains to be shown that collector current i_C can be controlled by variations in the small current i_B . In the discussion up to this point, we have indicated the control of i_C by emitter current i_E , with the base current characterized as a small side effect. In fact, we can show from space-charge-neutrality arguments that i_B can indeed be used to determine the magnitude of i_C . Let us consider the transistor of Fig. 17A, in which i_B is determined by a biasing circuit. For simplicity, we shall

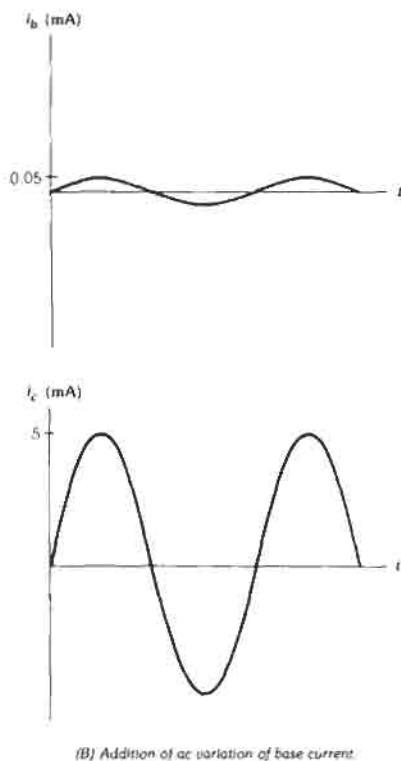
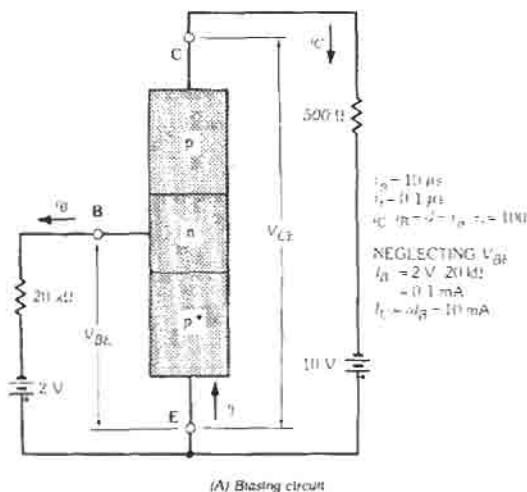


Fig. 17. Example of amplification in a common-emitter transistor circuit.

assume unity emitter injection efficiency and negligible collector saturation current. Since the n-type base region is electrostatically neutral between the two transition regions, the presence of excess holes in transit from emitter to collector calls for compensating excess electrons from the base contact. However, there is an important difference in the times which electrons and

holes spend in the base. The average hole spends a time τ_i , defined as the *transit time* from emitter to collector. Since the base width W_b is made small compared with L_p , this transit time is much less than the average hole lifetime, τ_p , in the base. On the other hand, an average excess electron supplied from the base contact spends τ_p seconds in the base (for simple recombination and equal excess carrier densities, τ_n and τ_p are equal). While the average electron waits τ_p seconds for recombination, many holes can enter and leave the base region, each with an average transit time τ_i . In particular, for each electron entering from the base contact, τ_p/τ_i holes can pass from emitter to collector while maintaining space-charge neutrality. Thus the ratio of collector current to base current is simply

$$i_C/i_B = \beta = \tau_p/\tau_i \quad (\text{Eq. 35})$$

for $\gamma = 1$ and negligible collector saturation current.

If the electron supply to the base (i_B) is restricted, the traffic of holes from emitter to base is correspondingly reduced. This can be argued in a simple way by supposing that the hole injection does continue despite a reduction of electrons from the base contact. The result would be a net buildup of positive charge in the base and a loss of forward bias (and therefore a loss of hole injection) at the emitter junction. Clearly, the supply of electrons through i_B can be used to raise or lower the hole flow from emitter to collector.

The base current is controlled independently in Fig. 17A. This is called a *common-emitter* circuit, since the emitter electrode is common to the base and collector circuits. The emitter junction is clearly forward biased by the battery in the base circuit. The voltage drop in the forward-biased emitter junction is small, however, so that almost all of the voltage from collector to emitter appears across the reverse-biased collector junction. Since v_{BE} is small for the forward-biased junction, we can neglect it and approximate the base current as $2\text{ V}/20\text{ k}\Omega = 0.1\text{ mA}$. If we assume $\tau_i = 0.1\ \mu\text{s}$ and $\tau_p = 10\ \mu\text{s}$, β for the transistor is 100, and the collector current i_C is 10 mA. It is important to notice that i_C is determined by β and the base current, rather than by the battery and resistor in the collector circuit (as long as these are of reasonable values to maintain a reverse-biased collector junction). In this example 5 V of the collector-circuit battery voltage appears across the 500 Ω resistor, and 5 V serves to reverse-bias the collector junction.

If a small ac current, i_b , is superimposed on the steady-state base current of Fig. 17, a corresponding ac current, i_c , appears in the collector circuit (Fig. 17B). The time-varying portion of the collector current will be larger than i_b by the factor β , and current gain results.

In a switching operation, a transistor is usually controlled in two conduction states that can be loosely referred to as the *on* state and the *off* state. Ideally, a switch should appear as a short circuit when turned on and an open circuit when turned off. Furthermore, it is desirable to switch the device from one state to the other with no lost time in between. Transistors do not fit this

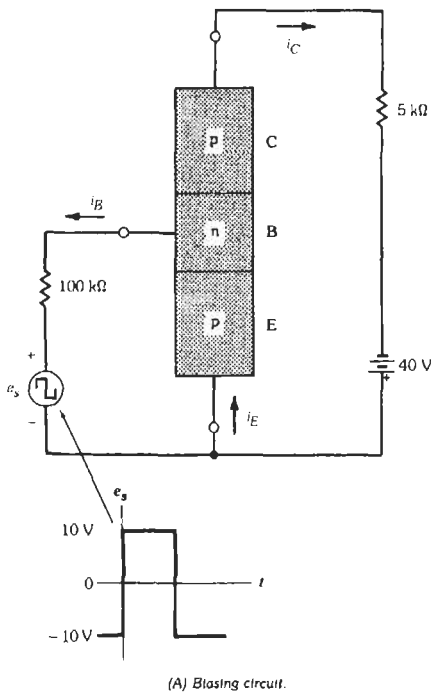
ideal description of a switch, but they can serve as a useful approximation in practical electronic circuits. The two states of a transistor in switching can be seen in the simple common-emitter example of Fig. 18. In this figure, collector current i_C is controlled by base current i_B over most of the family of characteristic curves. The load line specifies the locus of allowable ($i_C, -v_{CE}$) points for the circuit. The load line is obtained by writing a loop equation around the collector circuit, which results in a linear equation with an intercept on the $-V_{CE}$ axis ($I_C = 0$) and an intercept on the i_C axis ($V_{CE} = 0$).

If I_B is such that the operating point lies somewhere between the two end points of the load line (Fig. 18B), the transistor operates in the *normal active mode*. That

is, the emitter junction is forward biased and the collector is reverse biased, with a reasonable value of I_B flowing out of the base. On the other hand, if the base current is zero or negative, the cutoff point (C in Fig. 18B) is reached at the bottom end of the load line, and the collector current is negligible. This is the *off* state of the transistor, and the device is said to be operating in the *cutoff* regime. If the base current is positive and sufficiently large, the device is driven to the *saturation* regime, marked S. This is the *on* state of the transistor, in which a large value of I_C flows with only a very small voltage drop V_{CE} . The saturation regime corresponds to the loss of reverse bias across the collector junction. In a typical switching operation, the base current swings from positive to negative, thereby driving the device from saturation to cutoff and vice versa.

The response of a BJT to switching signals is not instantaneous. Junction capacitance and charge-storage effects, discussed earlier, cause delays in switching response that must be taken into account in circuit design. A typical time response of collector current to the input square wave of Fig. 18 is shown in Fig. 19. An initial delay, t_d , is required for charging of the junction space-charge regions. Then a rise time, t_r , is required before the minority-carrier distribution is established in the base region. Similarly, a fall time, t_f , characterizes the decay of the minority-carrier distribution from the conducting state to the cutoff condition (in which the base is essentially depleted of minority carriers). Some switching transistors utilize gold doping to reduce carrier recombination times, thereby minimizing these switching delay times.

Circuit models of transistors and applications in electronic systems are discussed in Chapter 19.



(A) Blasing circuit.
(B) Collector characteristics and load line for circuit.

Fig. 18. Simple common-emitter switching circuit for a transistor.

OTHER ELECTRONIC DEVICES

The basic devices discussed thus far—diodes, transistors, and their applications in integrated circuits—form the foundation of modern electronics. However, there are a host of other devices that play important roles in electronic circuit and system design. Symbols for many solid-state devices are shown in Table 4. In

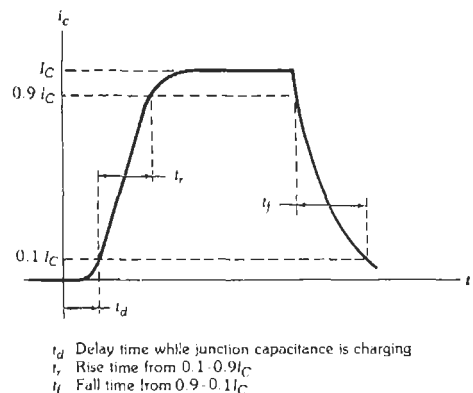











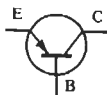

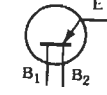

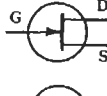
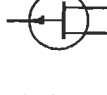













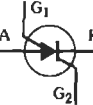
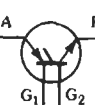
Fig. 19. Collector current during switching transients.

TABLE 4.* COMMON CIRCUIT SYMBOLS FOR SOLID STATE DEVICES†

<i>Diodes</i>			
pn diode, PIN, IMPATT, Schottky barrier diode			
Breakdown ("zener") diode			
Bidirectional zener diode			
Varactor			
Tunnel diode			
Photodiode, solar cell			
Light-emitting diode (LED)			
<i>Transistors</i>			
Bipolar PNP			
Bipolar NPN			
Unijunction (n-base)			
Unijunction (p-base)			
JFET (n-channel)			
JFET (p-channel)			
IGFET (n-channel)	Depletion		
	Enhancement		

Continued on next page.

TABLE 4* (CONT). COMMON CIRCUIT SYMBOLS FOR SOLID STATE DEVICES†

<i>Transistors (Continued)</i>			
IGFET (p-channel)			
<i>Four-Layer and Related Devices</i>			
PNPN (Shockley) diode			
SCR			
Bilateral switch			
SCS			

* From Ben G. Streetman, *Solid State Electronic Devices*, 3rd ed., © 1990. Reprinted by permission of Prentice-Hall, Inc., Englewood Cliffs, N.J.

† This table gives many of the commonly used symbols for semiconductor devices. Standardization is incomplete, however, and other symbols are often used in the electronics literature.

this section, we will briefly consider a few devices that are particularly important in many electronic systems.

Display Devices

Many electronic systems require some means of displaying information. Such display may take the form of a printout, a cathode-ray-tube image, or various types of digital (*numeric*) or letter and number (*alpha-numeric*) displays. Even analog displays, such as the indicator needle on a meter scale, are designed to present information from a circuit or system to the viewer. Most of these display methods are familiar and need no further elaboration. We will discuss here only a few of the newer displays that are currently being used.

Liquid crystal (LC) materials are made up of molecules that have considerable freedom of movement and can be aligned in various configurations by application of an electric field. Depending on the molecular alignment, light can be passed through the LC or can be scattered within it. A typical display cell is made up of the liquid crystal held between two glass plates with conductive coatings. The top plate is transparent, and the bottom plate is either reflecting or absorbing. By the application of appropriate voltages across the plates, incident light may be either scattered within the LC or

transmitted through the cell and absorbed or reflected from the back plate. As a result, an array of LC cells can be made to reflect light in certain segments but not others. Typically, the array is in the form of a series of bars that can be selectively addressed to display numbers or letters.

Liquid-crystal displays require little power for their operation and can be made in rather large arrays. Certain LC systems can be made to reflect colors selectively, so that colored displays are possible. Disadvantages of the LC devices include the requirement of incident light, sensitivity to temperature variations, and relatively slow response time (milliseconds) to the control signal.

One of the most versatile display systems is the *light-emitting diode* (LED). As discussed earlier, carriers injected across a forward-biased pn junction recombine within approximately one diffusion length of the junction. If the semiconductor is a luminescent material, the energy given up when an electron-hole pair recombines results in a photon of light. For example, when an electron in the conduction band of GaAs recombines with a hole in the valence band, a photon is emitted with energy equal to the GaAs band gap (1.4 eV). A photon with this energy is in the infrared portion of the spectrum. Thus in a properly constructed GaAs

junction, forward bias results in considerable carrier injection, recombination, and the resulting emission of infrared light. The intensity of the light can be varied (with nanosecond response time) by changing the diode current; thus electrical signals can be converted into optical signals by the LED. Used in conjunction with a photodiode, the LED can send information from one point to another optically. Another application of the LED-photodiode is in an *isolator*. Since the signal transmission is optical, complete electrical isolation can be achieved. Such an isolator pair can be mounted on an insulating substrate and packaged to form a device that allows signal transfer between input and output while maintaining electrical isolation.

If the information receiver is the human eye, the display device must emit light in the visible part of the spectrum (approximately 4000 to 7000 Å in wavelength, or 3.1 to 1.8 eV in photon energy). Visible LEDs are built that cover the range from red to green. For example, the compound GaAsP has a band gap that can be varied from that of GaAs (1.4 eV) to that of GaP (2.3 eV) by choosing the appropriate mixture. When the phosphorus content is about 40% of the column V constituent, the band gap is 2 eV, in the red. This is the most common material for LED display fabrication. By using segments made up of GaAsP strips, a numeric display can be made that emits light in the appropriate segments to display numbers or letters. Such an LED display operates at very low voltage, has a long operating lifetime, and is highly reliable.

Microwave Devices

The use of transistors at high frequencies is generally limited to the range of a few gigahertz by capacitance and transit-time effects. Therefore, generation and amplification of microwave signals usually depend on special devices that can deliver high-frequency ac power to a resonant cavity or waveguide. Power is delivered to an ac signal if there is an increase in the motion of charge through a region where the field is changing such as to retard such motion. For charges in a solid or gas, this is called *negative differential conductivity*. As an example, suppose holes are drifting down a typical semiconductor bar in the direction of an electric field. If the field varies with time, the holes speed up or slow down as the field changes, and power is extracted from the field. However, if the arrival of holes at a certain point were to increase during a time interval in which the field at that point decreases, power would be delivered to the field, and the apparent conductivity $dI/d\mathcal{E}$ would be negative. The object of microwave oscillators is to modulate the rate of arrival of electrons or holes at a point where the electric field is varying, such that the arrival coincides with the retarding half-cycle of the ac field.

A number of solid-state microwave devices have been introduced in recent years. Most of these devices are variations of the *IMPATT diode* and the *gunn diode*. In each of these devices, a negative conductance allows

current oscillations that can be used to generate microwave signals.

The negative conductance in an IMPATT diode results from avalanche multiplication and transit-time effects. A simple version of the IMPATT structure is the N^+PIP^+ device shown in Fig. 20A. In this diode, the n^+ and p^+ regions are heavily doped, and the i region is essentially intrinsic. Assume the diode is biased such that avalanche multiplication at the n^+p junction begins at $t = 0$ in the voltage cycle (Fig. 20B). Then a positive ac voltage superimposed on the dc bias causes more multiplication. The avalanche-generated holes drift to the right and into the i region in the diagram of Fig. 20C. Since the hole pulse grows as long as the critical field, \mathcal{E}_c , is exceeded, the pulse reaches its maximum at $\omega t = \pi$ (i.e., the hole pulse grows during the entire

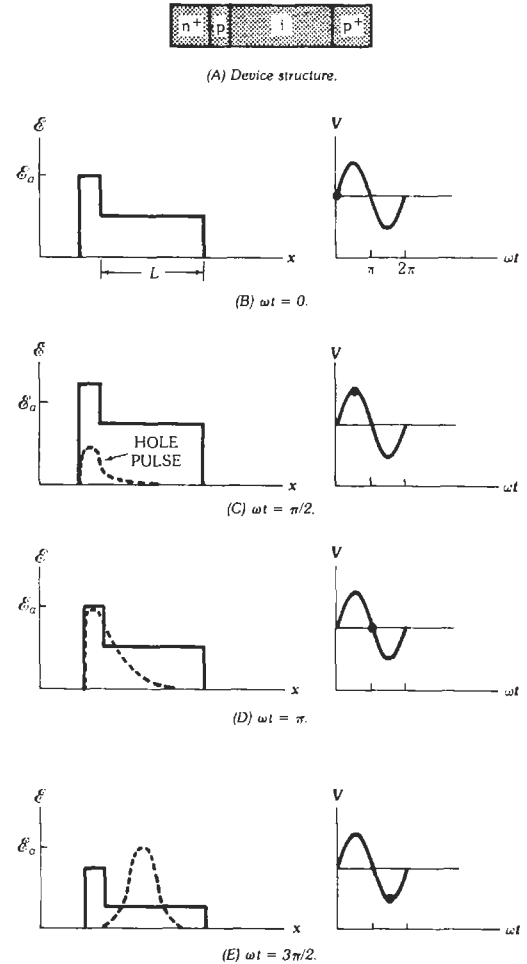


Fig. 20. Time dependence of the growth and drift of holes during a cycle of applied voltage for the N^+PIP^+ IMPATT diode.

positive half-cycle of the voltage). Then as v enters its negative half-cycle, avalanche ceases at the n^+p junction (Fig. 20D). However, the hole pulse drifts through the i region from left to right while the ac terminal voltage is negative. Since the holes are collected during the negative half-cycle of voltage, the ac conductance is negative during this period. If length L is chosen such that the pulse drifts through the i region for the full negative half-cycle of voltage, negative conductance will be obtained for the entire half-cycle. The IMPATT device can be placed in a resonant cavity tuned to the appropriate frequency, and microwave generation results.

Current pulses are created in a GaAs gunn diode by a mechanism that transfers electrons from one region of

the conduction band to another. The conduction electrons in GaAs normally reside in a band 1.43 eV above the valence band. In addition, a subsidiary conduction band lies about 0.3 eV above the first. Of course, this higher-lying band is generally of little interest, since conduction-band electrons are usually found only in the lower-energy band. It is possible, however, to excite electrons into the upper conduction band by applying a sufficiently large electric field. The reason this transfer is of interest is that the mobility of electrons in the upper band is much smaller than the usual mobility of electrons in the lower conduction band. Thus, when electrons are transferred to the low-mobility subsidiary band, they actually *slow down*. At a critical threshold field, \mathcal{E}_{th} , electrons begin to transfer to the upper band,

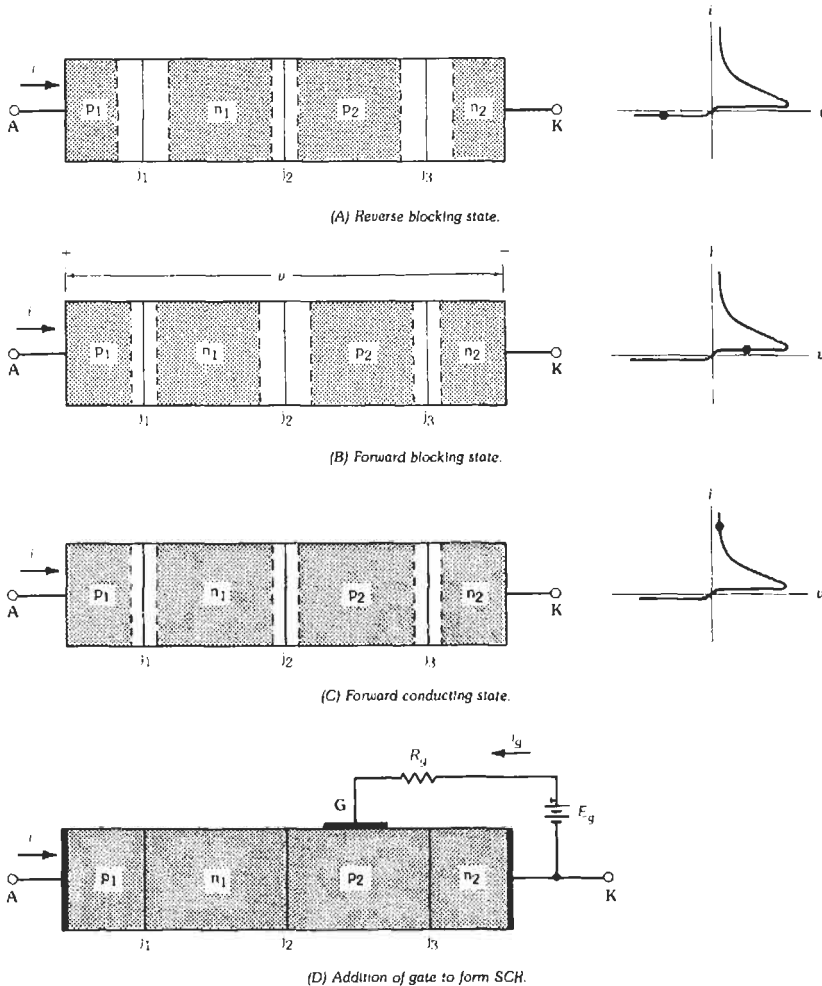


Fig. 21. Biasing of a PNP device.

where their velocity is smaller. This transfer results in a negative differential conductance that can be used in microwave generation.

Switching Devices

Many electronic applications call for a device that can be switched from a nonconducting "off" state to a conducting "on" state. Several devices can be used in switching applications, and selection of the appropriate device depends on requirements of power level, switching time, and other factors. A common electronic switch is the bipolar transistor, which can be driven from cutoff to saturation by controlling the small base current, or the FET, which can be switched by controlling the gate voltage.

The most widely used controllable switch for large currents is the *semiconductor controlled rectifier (SCR)*. This is a four-layer PNP device (Fig. 21) with terminals attached to the anode (A) and the cathode (K) and a third terminal attached to one of the central regions. This third terminal is called the gate (G). When the gate is left open, the SCR has the characteristic of Figs. 21A-C. With a negative applied voltage, junctions j_1 and j_3 are reverse biased (Fig. 21A), and current through the device is effectively blocked, even at high reverse voltage. With a positive voltage (A positive with respect to K), junctions j_1 and j_3 are forward biased, while j_2 is reverse biased (Fig. 21B). Initially, the device current is restricted to the small saturation current of j_2 . This is called the *forward blocking state* and corresponds to the "off" condition of the switch. The SCR can be thought of as two coupled transistors ($p_1n_1p_2$ and $n_1p_2n_2$) with a common collector junction, j_2 . If transistor action is initiated, holes injected into n_1 from the forward-biased emitter junction, j_1 , can be transported across the base of the PNP structure into p_2 . Such transistor action in effect feeds holes into the base of the NPN structure, thereby increasing electron injection from n_2 into p_2 . Such injected electrons can then be collected across j_2 into n_1 by transistor action, and the process continues. The result is a combination of two saturated transistors, typified by high current at low voltage (the *forward conducting state*), as shown in Fig. 21C.

The initiation of transistor action (and therefore switching) can occur as a result of raising the bias to a critical value (V_p). At this voltage, avalanche multiplication at j_2 and base-width narrowing* in n_1 and p_2 combine to initiate transport of minority carriers across the two base regions. Alternatively, a small gate current i_g can supply sufficient base current to the NPN structure to initiate transistor action (Fig. 21D). The latter switching method is the most common type for an SCR. In this mode of operation, the device is nonconducting at

forward voltages until a small pulse of current is applied to the gate. Such a pulse initiates transistor action in the device and switches it into the conducting state. The SCR remains in the conducting state until current I is dropped below a value called the *holding current* required to maintain transistor action. In addition, some SCR devices can be turned off by applying a negative current to the gate, thereby extracting carriers and terminating the transistor action within the device.

REFERENCES

Semiconductor Physics

- Blakemore, J. S. *Semiconductor Statistics*. New York: Dover Publications, 1987.
- Ferry, D. K. *Semiconductors*. New York: Macmillan, 1991.
- Hess, K. *Advanced Theory of Semiconductor Devices*. Englewood Cliffs, N.J.: Prentice Hall, 1988.
- Shur, M. *GaAs Devices and Circuits*. New York: Plenum Press, 1987.
- Shur, M. *Physics of Semiconductor Devices*. Englewood Cliffs, N.J.: Prentice Hall, 1990.
- Wang, S. *Fundamentals of Semiconductor Theory and Device Physics*. Englewood Cliffs, N.J.: Prentice Hall, 1989.
- Wolfe, C. M., Holonyak, N., and Stillman, G. E. *Physical Properties of Semiconductors*. Englewood Cliffs, N.J.: Prentice Hall, 1989.

Junctions and Transistors

- Capasso, F., and Margaritondo, G. eds. *Heterojunction Band Discontinuities: Device Physics and Applications*. New York: North Holland, 1987.
- Jaeger, R. C. *The Modular Series on Solid State Devices*, Vol. V: *Introduction to Microelectronic Fabrication*. Reading, Mass.: Addison-Wesley, 1988.
- Muller, R. S., and Kamins, T. I. *Device Electronics for Integrated Circuits*. 2d ed. New York: John Wiley, 1986.
- Roulston, D. J. *Bipolar Semiconductor Devices*. New York: McGraw-Hill, 1990.
- Sze, S. M. *High-Speed Semiconductor Devices*. New York: John Wiley, 1990.
- Sze, S. M. *Semiconductor Devices: Physics and Technology*. New York: John Wiley, 1985.
- Wang, C. T. *Semiconductor Technology: GaAs and Related Compounds*. New York: John Wiley, 1990.
- Warner, R. M., and Grung, B. L. *Semiconductor-Device Electronics*. Philadelphia: Holt, Rinehart and Winston, 1991.
- Wolf, S., and Tauber, R. N. *Silicon Processing for the VLSI Era*. Sunset Beach, Calif.: Lattice Press, 1986.

*Base-width narrowing occurs as the depletion region about j_2 grows due to the increased reverse bias. As the effective widths of n_1 and p_2 become smaller, base transport factor B increases.

Yang, E. S. *Microelectronic Devices*. New York: McGraw-Hill, 1988.

Optoelectronic Devices

- Agrawal, G. P., and Dutta, N. K. *Long-Wavelength Semiconductor Lasers*. New York: Van Nostrand Reinhold, 1989.
- Cheo, P. K. *Fiber Optics and Optoelectronics*. Englewood Cliffs, N.J.: Prentice Hall, 1990.
- Coutts, T. J., and Meakin, J. D. *Current Topics in Photovoltaics*. New York: Academic, 1985.
- Gowar, J. *Optical Communication Systems*. Englewood Cliffs, N.J.: Prentice Hall, 1984.
- Miller, S. E., and Kaminow, I. P. eds. *Optical Fiber Telecommunications II*. San Diego: Academic Press, 1988.
- Nishihara, H., Haruna, M., and Suhara, T. *Optical Integrated Circuits*. New York: McGraw-Hill, 1989.
- Senior, J. *Optical Fiber Communications*. Englewood Cliffs, N.J.: Prentice Hall International, 1985.
- Wilson, J., and Hawkes, J. F. B. *Optoelectronics*. 2nd ed. Cambridge: Prentice Hall International, 1989.

Microwave Devices

- Blakey, P. A., East, J. R., and Haddad, G. I. "Impact of Submicron Technology on Microwave and Millimeter-wave Devices." *VLSI Electronics: Microstructure Science*. Vol. 2. New York: Academic Press, 1981. Ch. 4, pp. 105-147.
- Bosch, B. G., and Engelmann, R. W. N. *Gunn Effect Electronics*. New York: John Wiley-Halstead Press, 1975.
- Bulman, P. J., Hobson, G. S., and Taylor, B. C. *Transferred Electron Devices*. New York: Academic Press, 1972.
- Capasso, F., and Margaritondo, G., eds. *Heterojunction Band Discontinuities: Device Physics and Applications*. New York: North Holland, 1987.
- Ferry, D. K., ed. *Gallium Arsenide Technology*. Indianapolis: Howard W. Sams & Co., 1985.
- Gibbons, G. *Avalanche-Diode Microwave Oscillators*. Oxford: Clarendon Press, 1973.
- Hobson, G. S. *The Gunn Effect*. Oxford: Clarendon Press, 1974.

Liao, S. Y. *Microwave Devices and Circuits*. Englewood Cliffs, N.J.: Prentice Hall, 1980.

Sze, S. M. *High-Speed Semiconductor Devices*. New York: John Wiley, 1990.

Semiconductor Controlled Rectifiers

- Blicher, A. *Thyristor Physics*. New York: Springer-Verlag, 1976.
- Gentry, F. E., Gutzwiller, F. W., Holonyak, N., Jr., and Von Zastrow, E. E. *Semiconductor Controlled Rectifiers: Principles and Application of p-n-p-n Devices*. Englewood Cliffs, N.J.: Prentice Hall, 1964.
- Ghandi, S. K. *Semiconductor Power Devices*, New York: John Wiley, 1977.
- Special Issue on High-Power Semiconductor Devices, *IEEE Transactions on Electron Devices* ED-23, no. 8 (August 1976).
- Laster, C. *Thyristor Theory and Application*. Blue Ridge Summit, Pa.: Tab Books, 1986.
- Sze, S. M. *Physics of Semiconductor Devices*. New York: John Wiley, 1981.
- Taylor, P. D. *Thyristor Design and Realization*. Chichester: John Wiley, 1987.

General

- Allison, J. *Electronic Engineering Semiconductors and Devices*. London: McGraw-Hill, 1990.
- Beadle, W. E., Tsai, J. C. C., and Plummer, R. D., eds. *Quick Reference Manual for Silicon Integrated Circuit Technology*. New York: John Wiley-Interscience, 1986.
- Dalven, R. *Introduction to Applied Solid State Physics*. 2nd ed. New York: Plenum, 1990.
- Kittel, C. *Introduction to Solid State Physics*. 6th ed. New York: John Wiley, 1986.
- Pierret, R. F., and Neudeck, G. W. *Modular Series on Solid State Devices*. Reading, Mass.: Addison-Wesley Publishing Co., Inc., 1983.
- Streetman, B. G. *Solid State Electronic Devices*. 3rd ed. Englewood Cliffs, N.J.: Prentice Hall, Inc., 1990.
- Sze, S. M. *Physics of Semiconductor Devices*. 2nd ed. New York: John Wiley & Sons, Inc., 1981.

19 Transistor Circuits

Horace G. Jackson

Device Models and Equations	19-2
Bias Techniques	19-2
Bipolar Transistor	
Field-Effect Transistor	
Small-Signal Models and Equations	19-5
Single-Stage Amplifiers	19-8
Small-Signal Characteristics	
Frequency Response	
Large-Signal Characteristics	
Output Stages	
Harmonic Distortion	
Differential Amplifiers	19-17
Current Sources	19-17
Current Mirrors	
Active Loads	
Feedback Amplifiers	19-18
Basic Properties	
Basic Feedback Circuit Topologies	
Bandpass Amplifiers	19-22
Single-Tuned Interstage	
Double-Tuned Interstage	
Sinusoidal Oscillators	19-24
Pulse Circuits	19-26
Pulse Shaping	
Multivibrators	

This chapter gives condensed descriptions of many types of circuits in which transistors are used. Also presented is design information that makes possible the determination of the various circuit parameters. In accordance with the accepted practice, upper-case variables with upper-case subscripts (V_{CE}) are used to indicate the static, or large-signal, quantities, and lower-case variables with lower-case subscripts (v_{ce}) are used to indicate the dynamic, or small-signal, values.

The overwhelming majority of transistor circuits are made with silicon devices; therefore, this is the assumed technology in this chapter. However, germanium and gallium-arsenide devices are also available to a limited degree.

DEVICE MODELS AND EQUATIONS*

This section presents the basic large-signal circuit models and equations for three common semiconductor active devices, namely:

- (1) Bipolar junction transistor (BJT, or simply transistor)
- (2) Junction field-effect transistor (JFET or FET)
- (3) Metal-oxide-semiconductor field-effect transistor (MOSFET or MOS)

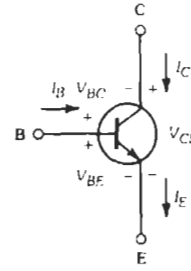
The first two are readily found in both discrete and integrated transistor circuits. Except for some high-power applications, the third is restricted to integrated circuits only.

The symbol and large-signal model for an npn transistor are illustrated in Fig. 1. For a pnp transistor, the polarities of the terminal voltages, V_{BC} , V_{BE} , and V_{CE} , must be reversed; the direction of the junction diodes must be reversed; and the direction of all the currents must also be reversed. The model equations are listed in Chart 1. Typical values for the device parameters are given in Chart 2.

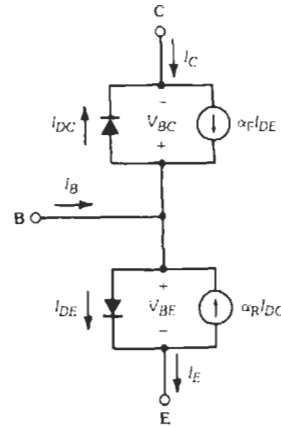
The symbol and large-signal model for an n-channel JFET are given in Fig. 2. For a p-channel device, the polarities of the terminal voltages, V_{GD} , V_{GS} , and V_{DS} ; the direction of the junction diodes; and the direction of all the currents must be reversed. The model equations, given in Chart 3, assume that the JFET approximates a square-law device. Typical parameter values are listed in Chart 4.

The symbol and large-signal model for an n-channel MOSFET are given in Fig. 3. For a p-channel device, the polarities of the five terminal voltages, V_{GS} , V_{GD} , V_{DS} , V_{BS} , and V_{BD} ; the two junction diodes; and the currents must all be reversed. The model equations and parameter values are given in Chart 5 and Chart 6, respectively.

- Note that a MOSFET may be one of two types:
- A. Enhancement type: with $V_{GS} = 0$ V there is no conducting channel, and the drain current is zero.



(A) Symbol and nomenclature.



(B) Large-signal model.

Fig. 1. Symbol, nomenclature, and large-signal model for npn transistor (BJT).

- B. Depletion type: with $V_{GS} = 0$ V there is a conducting channel, and the drain current is finite.

The FET terms I_{DSS} and V_T are graphically defined in the transfer characteristic of Fig. 4.

A summary of the operating modes for JFETs and MOSFETs is presented in Table I.

BIAS TECHNIQUES†

Bipolar Transistor

As an amplifier, the BJT is normally operated in the forward-active region. That is, the base-emitter junction is forward-biased, and the base-collector junction is reverse-biased. Thus, with $V_{BE} \gg 4V_T$ and $V_{BC} \ll -4V_T$ the equations in Chart 1 reduce to

$$I_E = (\beta_F + 1)(I_B + I_{CO})$$

$$I_C = \beta_F I_B + (\beta_F + 1)I_{CO}$$

where

$$I_{CO} = (1 - \alpha_F \alpha_R)I_{CS}$$

* References 1, 13, and 17.

† References 6 and 14.

CHART 1. MODEL EQUATIONS FOR AN NPN TRANSISTOR

$$I_E = I_{ES}(e^{V_{BE}/V_T} - 1)(1 + V_{BC}/V_A) - \alpha_R I_{CS}(e^{V_{BC}/V_T} - 1)(1 + V_{BE}/V_B)$$

$$I_C = \alpha_F I_{ES}(e^{V_{BE}/V_T} - 1)(1 + V_{BC}/V_A) - I_{CS}(e^{V_{BC}/V_T} - 1)(1 + V_{BE}/V_B)$$

$$I_B = I_E - I_C$$

Also,

$$\alpha_F I_{ES} = \alpha_R I_{CS}$$

where,

α_F = Forward current gain,

α_R = Reverse current gain,

I_{ES} = Emitter junction saturation current,

I_{CS} = Collector junction saturation current,

V_A = Forward-mode basewidth modulation factor,

V_B = Reverse-mode basewidth modulation factor,

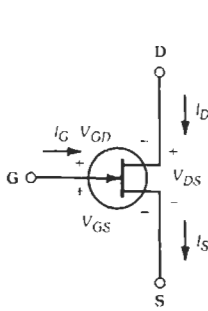
V_T = Thermal voltage* = kT/q ,

k = Boltzmann's constant = 1.38×10^{-23} J/K,

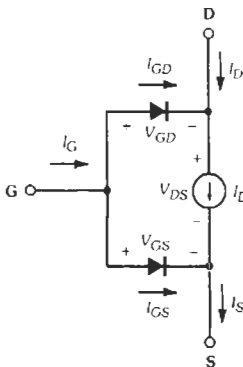
T = Absolute temperature = $[273 + \text{temp } (^{\circ}\text{C})]$ K,

q = Electronic charge = 1.60×10^{-19} C.

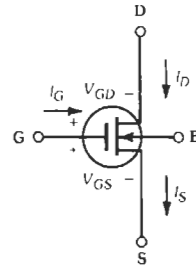
* The thermal voltage, used with bipolar devices, is not to be confused with the threshold voltage (V_T) used with field-effect devices.



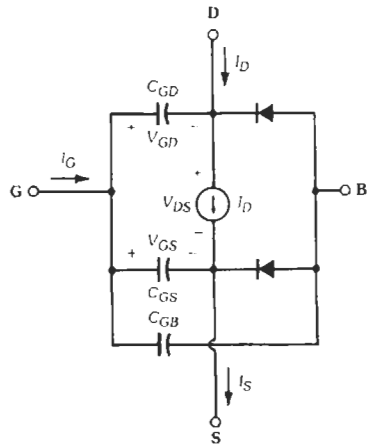
(A) Symbol and nomenclature.



(B) Large-signal model.



(A) Symbol and nomenclature.



(B) Large-signal model.

Fig. 2. Symbol and large-signal model for n-channel junction field-effect transistor (JFET).

Fig. 3. Symbol and large-signal model for n-channel MOSFET.

CHART 2. TYPICAL VALUES OF THE DEVICE PARAMETERS FOR A LOW-POWER SILICON BJT

$\alpha_F = 0.98$	$\beta_F = \alpha_F / (1 - \alpha_F) = 50$
$\alpha_R = 0.49$	$\beta_R = \alpha_R / (1 - \alpha_R) = 1$
$I_{ES} = 1 \times 10^{-14}$ A	$V_A = 20$ to 200 V
$I_{CS} = 2 \times 10^{-14}$ A	$V_B = 10$ to 100 V
	$V_T = 26$ mV at 27 °C

Good design of the bias circuit for a BJT requires that two deficiencies of the transistor be overcome. They are:

1. The transistor is a temperature sensitive device, in particular with respect to V_{BE} and β_F .

$$\Delta V_{BE}(T) \approx -2.5 \text{ mV}/^\circ\text{C}$$

$$\Delta \beta_F(T) \approx +0.7\%/^\circ\text{C}$$

Also note that I_{CO} doubles for every 8 °C rise in temperature. But with silicon devices I_{CO} is usually so small that it may be neglected.

2. The parameters of a transistor are subject to process variation, in particular with respect to β_F . Typically $\beta_F = 30$ to 300 (or greater).

Generally, a constant collector current I_C and constant collector-emitter voltage V_{CE} are required. That is, the operating, or quiescent, point—the Q-point—of the transistor is designed to be insensitive to temperature and process variations. A self-bias, or emitter bias, circuit is illustrated in Fig. 5A. The circuit in Fig. 5B is the equivalent.

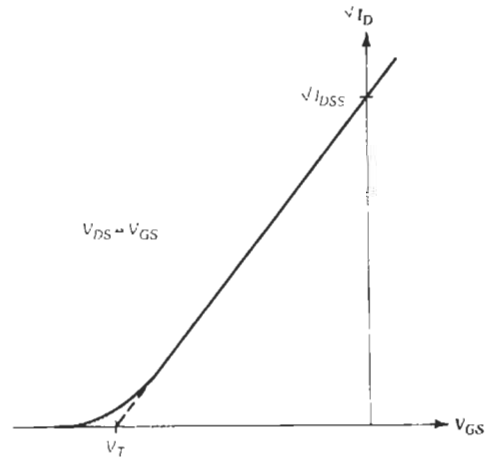


Fig. 4. Typical JFET transfer characteristic.

Given I_C , V_{CE} , and V_{BE} , the design procedure is as follows:

- A. Choose R_E , where $R_E \geq 5\Delta V_{BE}/I_E$. Here, ΔV_{BE} is the change in V_{BE} over the temperature range of interest. This tends to desensitize the operating point to temperature variations of V_{BE} .
- B. Determine R_B , where $R_B \leq \beta_F R_E / 5$. Here, β_F is the nominal value. This tends to make the operating point independent of variations in the value of β_F .
- C. Find R_C , from $R_C = [(V_{CC} - V_{CE})/I_C] - R_E$.
- D. Using nominal values for V_{BE} and β_F , solve for

$$V_{BB} \approx V_{BE} + I_C [(R_B + \beta_F R_E) / \beta_F]$$

In an amplifier design, R_E is usually paralleled by a capacitor or by a series resistor-capacitor combination.

CHART 3. MODEL EQUATIONS FOR AN N-CHANNEL JFET

Nonsaturation region,* $V_{DS} < (V_{GS} - V_T)$:

$$I_D = \frac{I_{DSS}}{V_T^2} [2(V_{GS} - V_T) V_{DS} - V_{DS}^2] (1 + \lambda V_{DS})$$

Saturation region,** $V_{DS} \geq (V_{GS} - V_T)$:

$$I_D = \frac{I_{DSS}}{V_T^2} [V_{GS} - V_T]^2 (1 + \lambda V_{DS})$$

where,

I_{DSS} = Zero-bias saturation drain current, i.e., $V_{GS} = 0$ V, $V_{DS} > (V_{GS} - V_T)$.

V_T = Threshold voltage, †

λ = Channel-length modulation factor.

* Also termed the linear, ohmic, or triode region.

** Also the constant current, pinchoff, or pentode region.

† Also the pinchoff voltage (V_p) or $V_{GS(off)}$. Not to be confused with the thermal voltage (V_T) used with bipolar devices.

CHART 4. TYPICAL DEVICE PARAMETERS FOR A LOW-POWER SILICON JFET

	N-Channel	P-Channel
$ I_{DSS} $	1 to 10 mA	1 to 10 mA
V_T	-2 to -6 V	+2 to +6 V
λ	0.1 to 0.01 V ⁻¹	0.1 to 0.01 V ⁻¹

Field-Effect Transistor

The JFET is also a temperature sensitive device.

$$\Delta I_{DSS}(T) \approx -0.5\%/^{\circ}C$$

$$\Delta V_T(T) \approx -2 \text{ mV}/^{\circ}C$$

Due to process variations, I_{DSS} may range 3:1 and V_T may vary 2:1.

For a MOSFET, the parameters k and V_T have temperature and process dependence similar to those of I_{DSS} and V_T for a JFET.

In designing for a constant I_D and constant V_{DS} in the saturation region, the self-bias circuit shown in Fig. 6A can be used with JFETs or depletion-type MOSFETs. Since the gate current is negligible

$$R_S = -V_{GS}/I_D$$

$$R_D = [(V_{DD} - V_{DS})/I_D] - R_S$$

The circuit of Fig. 6A is unsuitable for biasing an enhancement-type MOSFET, since for an n-channel device V_T is positive and it is required that $V_{GS} > V_T$. Two alternative circuits for biasing an enhancement-type MOSFET are included in Fig. 6. In Fig. 6B, with negligible gate current $V_{GS} = V_{DS}$ so that the transistor is operating in the saturation region.

$$R_D = (V_{DD} - V_{DS})/I_D$$

$$R_G \geq 10R_D$$

When it is required that $V_{GS} \neq V_{DS}$, the circuit in Fig. 6C may be used.

As in the BJT case, in most amplifier applications R_S is paralleled by a capacitor.

SMALL-SIGNAL MODELS AND EQUATIONS*

For a BJT biased in the forward-active region, the basic small-signal model for the device is illustrated in Fig. 7. Included in the model are the bulk ohmic effects of the neutral base, collector, and emitter regions, designated r_b , r_c , and r_e . Typically, these are about 100 ohms, 50 ohms, and 1 ohm, respectively. The forward-biased base-emitter junction is modeled by r_{π} and C_{π} .

* References 2, 4, and 5.

CHART 5. MODEL EQUATIONS FOR AN N-CHANNEL MOSFET

Nonsaturation region, $V_{DS} < (V_{GS} - V_T)$:

$$I_D = \frac{k}{2} [2(V_{GS} - V_T) V_{DS} - V_{DS}^2] (1 + \lambda V_{DS})$$

Saturation region, $V_{DS} \geq (V_{GS} - V_T)$:

$$I_D = \frac{k}{2} [V_{GS} - V_T]^2 (1 + \lambda V_{DS})$$

where

k = Conduction factor = $(\mu \epsilon_{ox} / t_{ox}) (W/L)$,

μ = Mobility of the conducting carriers,

ϵ_{ox} = Permittivity of the gate oxide,

t_{ox} = Thickness of the gate oxide,

W = Width of the conducting channel,

L = Length of the conducting channel,

V_T = Threshold voltage* = $V_{T0} + \gamma[\sqrt{|2\phi_F| + |V_{DS}|} - \sqrt{|2\phi_F|}]$,

V_{T0} = Zero-bias threshold voltage, i.e., $V_{BS} = 0 \text{ V}$,

$2\phi_F$ = Surface potential $\approx 0.6 \text{ V}$,

γ = Body-effect factor,

λ = Channel-length modulation factor.

* Not to be confused with the thermal voltage (V_T) used with bipolar devices.

TABLE 1. SUMMARY OF OPERATING MODES FOR FET DEVICES

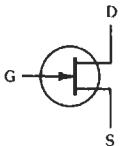
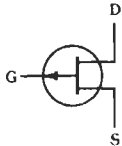
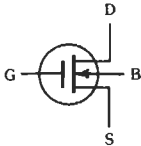
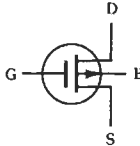
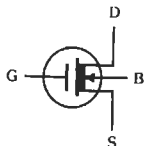
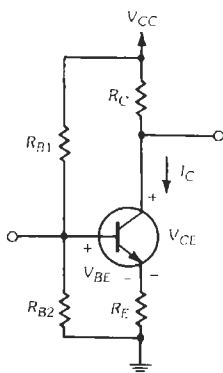
	JFET		MOSFET		
	N-Channel	P-Channel	N-Channel	P-Channel	N-Channel Depletion
Symbol					
V_T Polarity	Negative	Positive	Positive	Negative	Negative
V_{GS} Polarity	Negative	Positive	Positive	Negative	Neg or Pos
V_{DS} Polarity	Positive	Negative	Positive	Negative	Positive
Cutoff Region	$V_{GS} < V_T$	$V_{GS} > V_T$	$V_{GS} < V_T$	$V_{GS} > V_T$	$V_{GS} < V_T$
Nonsaturation Region	$V_{DS} < (V_{GS} - V_T)$	$V_{DS} > (V_{GS} - V_T)$	$V_{DS} < (V_{GS} - V_T)$	$V_{DS} > (V_{GS} - V_T)$	$V_{DS} < (V_{GS} - V_T)$
Saturation Region	$V_{DS} \leq (V_{GS} - V_T)$	$V_{DS} \leq (V_{GS} - V_T)$	$V_{DS} \geq (V_{GS} - V_T)$	$V_{DS} \leq (V_{GS} - V_T)$	$V_{DS} \geq (V_{GS} - V_T)$

CHART 6. TYPICAL DEVICE PARAMETERS FOR A LOW-POWER SILICON MOSFET

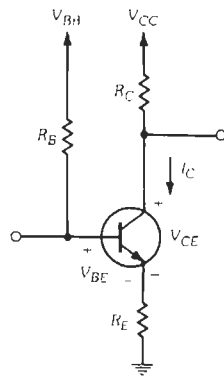
	N-Channel	P-Channel
k	$20(W/L)\mu A/V^2$	$7(W/L)\mu A/V^2$
V_{TEO}	+1 to +3 V	-1 to -3 V
V_{TDO}	-1 to -3 V	+1 to +3 V
$ 2\phi_F $	-0.6 to -0.8 V	+0.6 to +0.8 V
γ	+0.3 to +0.6 V	-0.3 to -0.6 V
λ	0.1 to 0.01 V^{-1}	0.1 to 0.01 V^{-1}

The reverse-biased base-collector junction is modeled by C_{μ} , and r_{μ} is due to the basewidth modulation effect. That the collector current is a function of the base-emitter voltage is modeled by current source $g_m v_{\pi}$ with output resistance r_{ce} . How these small-signal parameters relate to the operating point of the transistor is indicated in Chart 7.

The small-signal model for a JFET biased in the saturation region is given in Fig. 8. The ohmic resistances of the neutral gate, drain, and source regions are modeled by the linear resistors r_g , r_d , and r_s . Typically,



(A) Single voltage supply.



(B) Two-voltage-supply equivalent.

$$V_{B3} = \frac{R_{B2}}{R_{B1} + R_{B2}} V_{CC}$$

$$R_3 = R_{B1} \parallel R_{B2} = \frac{R_{B1} R_{B2}}{R_{B1} + R_{B2}}$$

Fig. 5. Bias circuits for bipolar transistors.

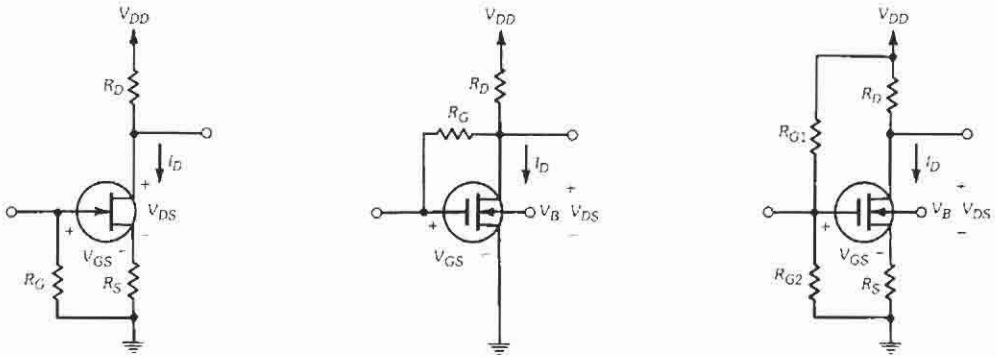


Fig. 6. Bias circuits for FETs.

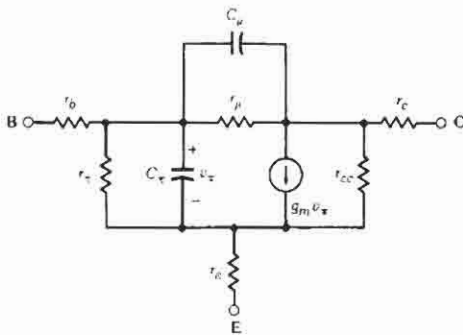


Fig. 7. Small-signal model for BJT.

CHART 7. SMALL-SIGNAL PARAMETERS FOR A BJT

$$g_m = i_c/v_{\pi} = \left. \frac{dI_C}{dV_{BE}} \right|_{OP} = I_C/V_T$$

$$r_{\pi} = v_{\pi}/i_b = \left. \frac{dV_{BE}}{dI_B} \right|_{OP} = \beta_F V_T/I_C = \beta_F/g_m$$

$$r_{ce} = v_{ce}/i_c = \left. \frac{dV_{CE}}{dI_C} \right|_{OP} = (V_A + V_{CE})/I_C$$

$$r_{\mu} = v_{cb}/i_b = \left. \frac{dV_{CB}}{dI_B} \right|_{OP} = (V_A + V_{CB})/I_B$$

$$C_{\pi} = (g_m/2\pi f_T) - C_{\mu}$$

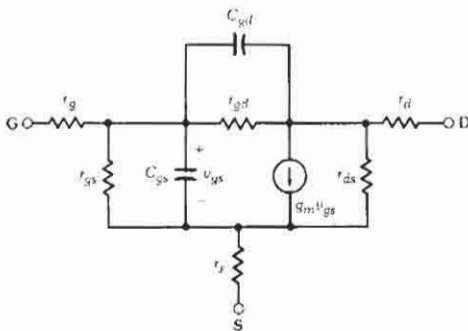


Fig. 8. Small-signal model for JFET.

these are 1 ohm, 10 ohms, and 10 ohms, respectively. In most applications they may be neglected.

The reverse-biased gate-source and gate-drain junctions of the JFET are modeled by $r_{gs}-C_{gs}$ and $r_{gd}-C_{gd}$, respectively. The voltage-dependent drain-current generator is modeled by $g_m v_{gs}$ with an output resistance of

r_{ds} . The equations for the latter two parameters are found in Chart 8.

For low-frequency applications, the charge-storage elements, C_{gs} and C_{gd} , may be neglected. Because r_{gs} and r_{gd} are generally very large, they also may be neglected in many applications.

The small-signal model for a MOSFET biased in the saturation region is given in Fig. 9. The ohmic effects of the gate, drain, source, and body regions are modeled by the linear resistors r_g , r_d , r_s , and r_b . Typical values are 1 ohm, 1 ohm, 1 ohm, and 10 ohms, respectively. In most applications they may be neglected.

CHART 8. SMALL-SIGNAL PARAMETERS FOR JFET

$$g_m = i_d/v_{gs} = \left. \frac{dI_D}{dV_{GS}} \right|_{OP} = (2I_{DSS}/V_T)(1 - V_{GS}/V_T)$$

$$= (2/V_T)(I_D I_{DSS})^{1/2}$$

$$r_{ds} = v_{ds}/i_d = \left. \frac{dV_{DS}}{dI_D} \right|_{OP} = (V_A + V_{DS})/I_D$$

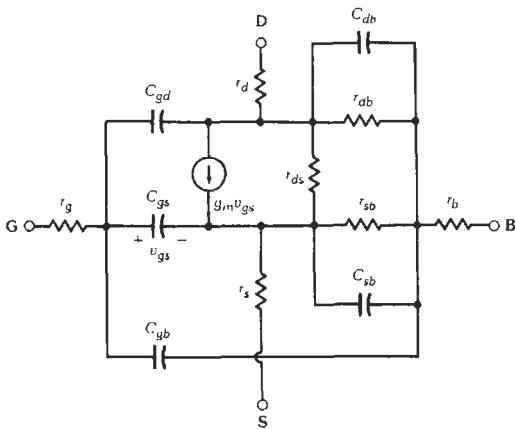


Fig. 9. Small-signal model for MOSFET.

Charge stored between the insulated gate and the source, drain, and body regions is modeled by the capacitors C_{gs} , C_{gd} , and C_{gb} , respectively. The reverse-biased source-body and drain-body junctions are included in the model by $r_{sb}-C_{sb}$ and $r_{db}-C_{db}$. However, generally the resistance values are so large that their effect is usually neglected.

As with the JFET, the voltage-dependent current source is modeled by $g_m v_{gs}$ with an output resistance of r_{ds} . The equations for these parameters are found in Chart 9.

The general ranges of the small-signal parameters for each of the devices considered in this section are compared in Table 2.

SINGLE-STAGE AMPLIFIERS*

Small-Signal Characteristics

The three possible circuit configurations for a single-transistor amplifier are illustrated in Fig. 10. For clarity, the bias networks have been omitted. It is assumed that the bipolar transistor is operating in the

CHART 9. SMALL-SIGNAL PARAMETERS FOR MOSFET

$$g_m = i_c/v_{gs} = \frac{dI_C}{dV_{GS}} \Big|_{OP} = k(V_{GS} - V_T) = (2kI_D)^{1/2}$$

$$r_{ds} = v_{ds}/i_d = \frac{dV_{DS}}{dI_D} \Big|_{OP} = (V_A + V_{DS})/I_D$$

forward-active region and the field-effect transistor is operating in the saturation region. The small-signal design equations listed in Table 3 are for the input resistance (r_i), output resistance (r_o), voltage gain ($a_v = v_o/v_i$), and current gain ($a_i = i_o/i_i$). The small-signal transducer power gain is obtained from $a_p = a_v a_i$. The second-order terms due to the bulk ohmic effects have been omitted. The equations for the FET apply for both the JFET and the MOSFET. Numerical values for each of the amplifier configurations are compared in Table 4.

Frequency Response

In Figs. 11 through 13 are shown simple circuit schematics for each of the single-stage amplifier configurations. Included in each figure are the small-signal equivalent circuits appropriate for the three ranges of frequencies, low, middle, and high. The corresponding small-signal transfer functions derived from the equivalent circuits are listed in Tables 5 through 7.

For the common-emitter and common-source amplifiers, the low cutoff frequency (-3 dB) is due to the greater of ω_{t1} or ω_{t2} , provided the ratio of the two is greater than 10. In the high-frequency region, normally $\omega_{zh} \gg \omega_{h2} \gg \omega_{h1}$, so that ω_{h1} is the high cutoff frequency (-3 dB).

For the emitter-follower and the source-follower, there is only one low break frequency. In the high-frequency region, typically ω_{zh} is only slightly greater than ω_{h1} , though ω_{h2} is much greater than either. Hence, the high cutoff frequency is difficult to determine accurately. However, it is greater than ω_{h1} , which is approximately equal to ω_T of the transistor.

For the common-base amplifier, the determination

* References 2, 8, and 14.

TABLE 2. COMPARISON OF SMALL-SIGNAL PARAMETERS

BJT	JFET	MOSFET
$g_m = 1$ to 100 mA/V	$g_m = 0.1$ to 10 mA/V	$g_m = 0.1$ to 10 mA/V
$r_{\pi} = 1$ to 100 k Ω	$r_{gs} \geq 10^9 \Omega$	$r_{gs} \geq 10^{10} \Omega$
$r_{ce} = 0.01$ to 1 M Ω	$r_{ds} = 0.01$ to 1 M Ω	$r_{ds} = 0.01$ to 1 M Ω
$C_{\pi} = 10$ to 100 pF	$C_{gs} = 1$ to 10 pF	$C_{gs} = 1$ to 10 pF
$C_{\mu} = 1$ to 10 pF	$C_{gd} = 1$ to 10 pF	$C_{gd} = 1$ to 10 pF

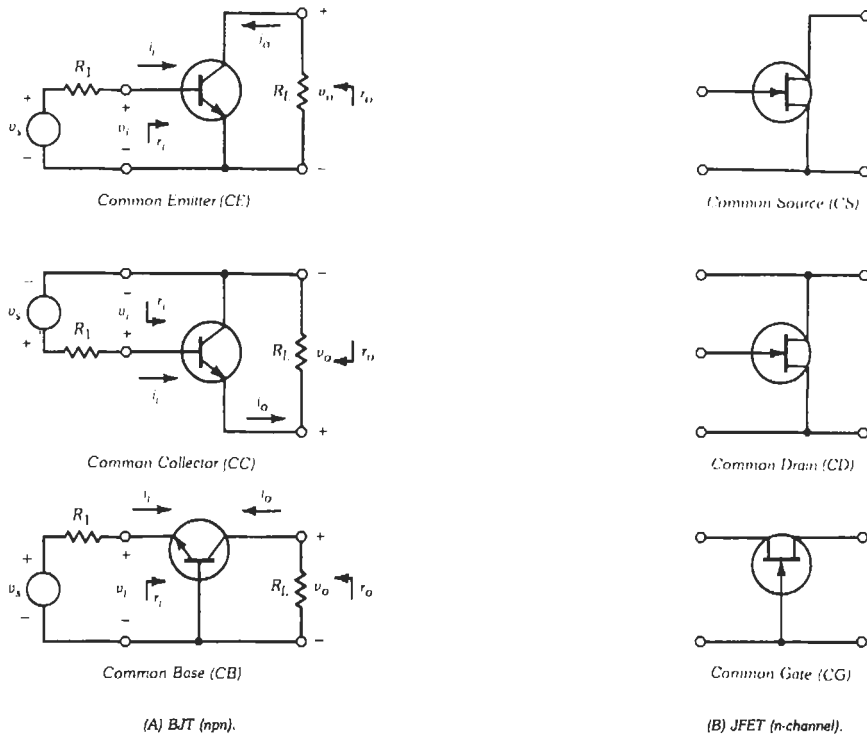
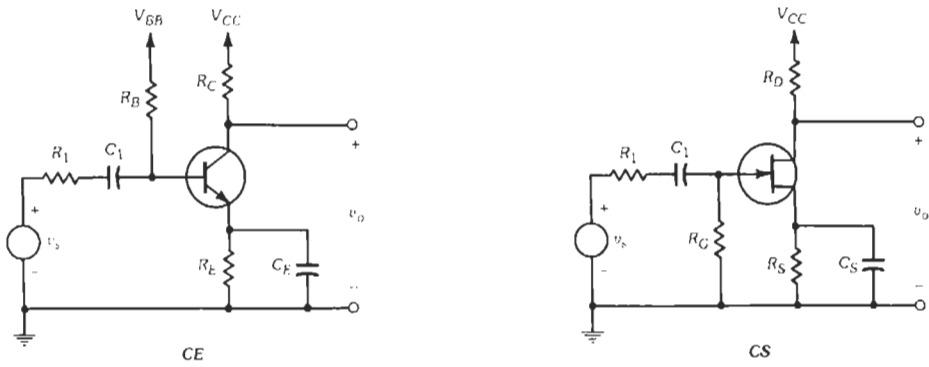


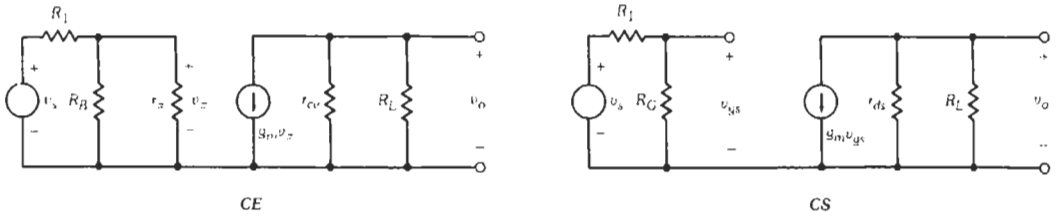
Fig. 10. Single-transistor amplifier configurations.

TABLE 3. SMALL-SIGNAL DESIGN EQUATIONS

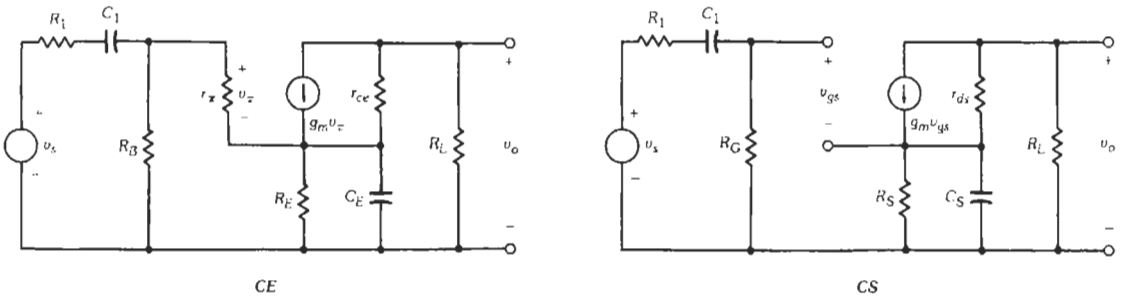
Bipolar Transistor	Field-Effect Transistor
Common-Emitter	Common-Source
$r_i = r_\pi$	$r_i = r_{gs}$
$r_o = R_L \parallel r_{ce}$	$r_o = R_L \parallel r_{ds}$
$a_v = -g_m(R_L \parallel r_{ce})$	$a_v = -g_m(R_L \parallel r_{ds})$
$a_i = \beta_0$	$a_i \rightarrow \infty$
Common-Collector (Emitter-Follower)	Common-Drain (Source-Follower)
$r_i = r_\pi + (\beta_0 + 1)R_L$	$r_i = r_{gs}$
$r_o = R_L \parallel [(R_1 + r_\pi)/(\beta_0 + 1)]$	$r_o = R_L \parallel 1/g_m$
$a_v = \frac{(\beta_0 + 1)R_L}{r_\pi + (\beta_0 + 1)R_L}$	$a_v = R_L/(1/g_m + R_L)$
$a_i = \beta_0 + 1$	$a_i \rightarrow \infty$
Common-Base	Common-Gate
$r_i = r_\pi/(\beta_0 + 1) \approx 1/g_m$	$r_i = 1/g_m$
$r_o = R_L \parallel r_{ce}(1 + g_m R_1)$	$r_o = R_L \parallel r_{ds}(1 + g_m R_1)$
$a_v = a_0 \frac{[R_L \parallel r_{ce}(1 + g_m R_1)]}{r_\pi/(\beta_0 + 1)}$	$a_v = g_m[R_L \parallel r_{ds}(1 + g_m R_1)]$
$a_i = a_0$	$a_i \rightarrow \infty$



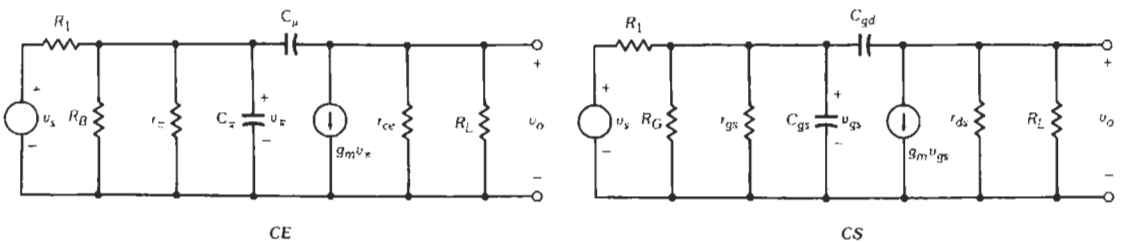
(A) Circuit schematic.



(B) Midband equivalent circuit.

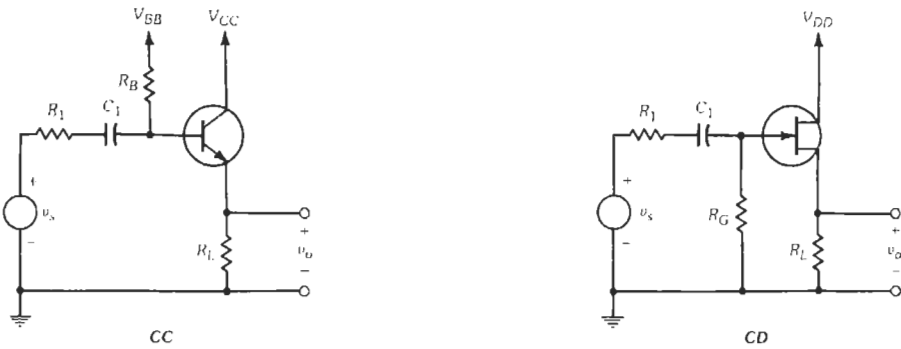


(C) Low-frequency equivalent circuit.

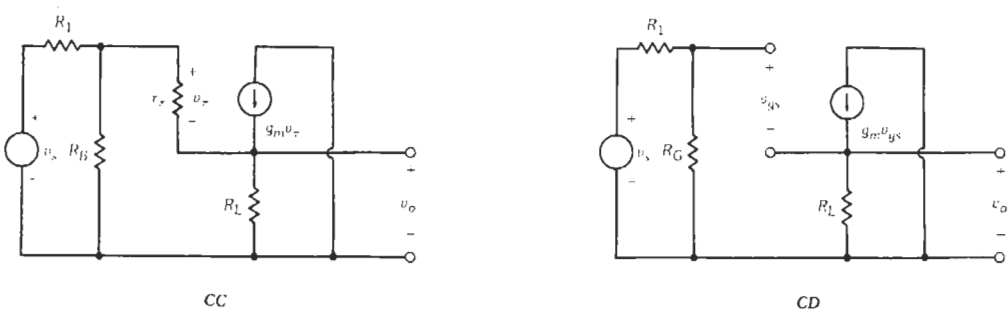


(D) High-frequency equivalent circuit.

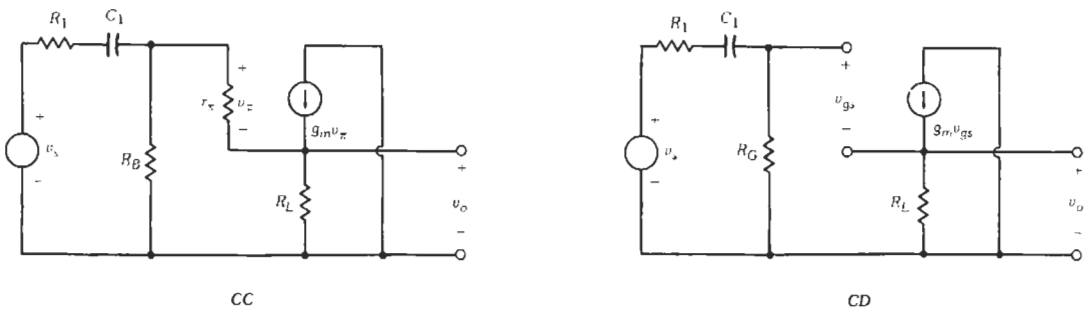
Fig. 11. Circuit schematic and small-signal equivalent circuits for common-emitter and common-source amplifier.



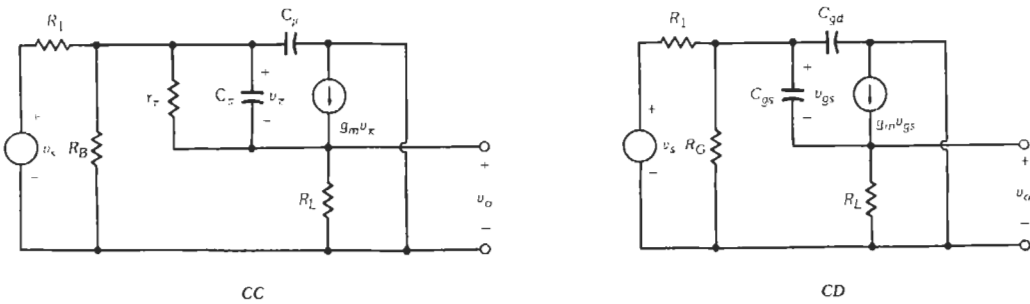
(A) Circuit schematic.



(B) Midfrequency equivalent circuit.



(C) Low-frequency equivalent circuit.



(D) High-frequency equivalent circuit.

Fig. 12. Circuit schematic and small-signal equivalent circuits for common-collector (emitter-follower) and common-drain (source-follower) amplifier.

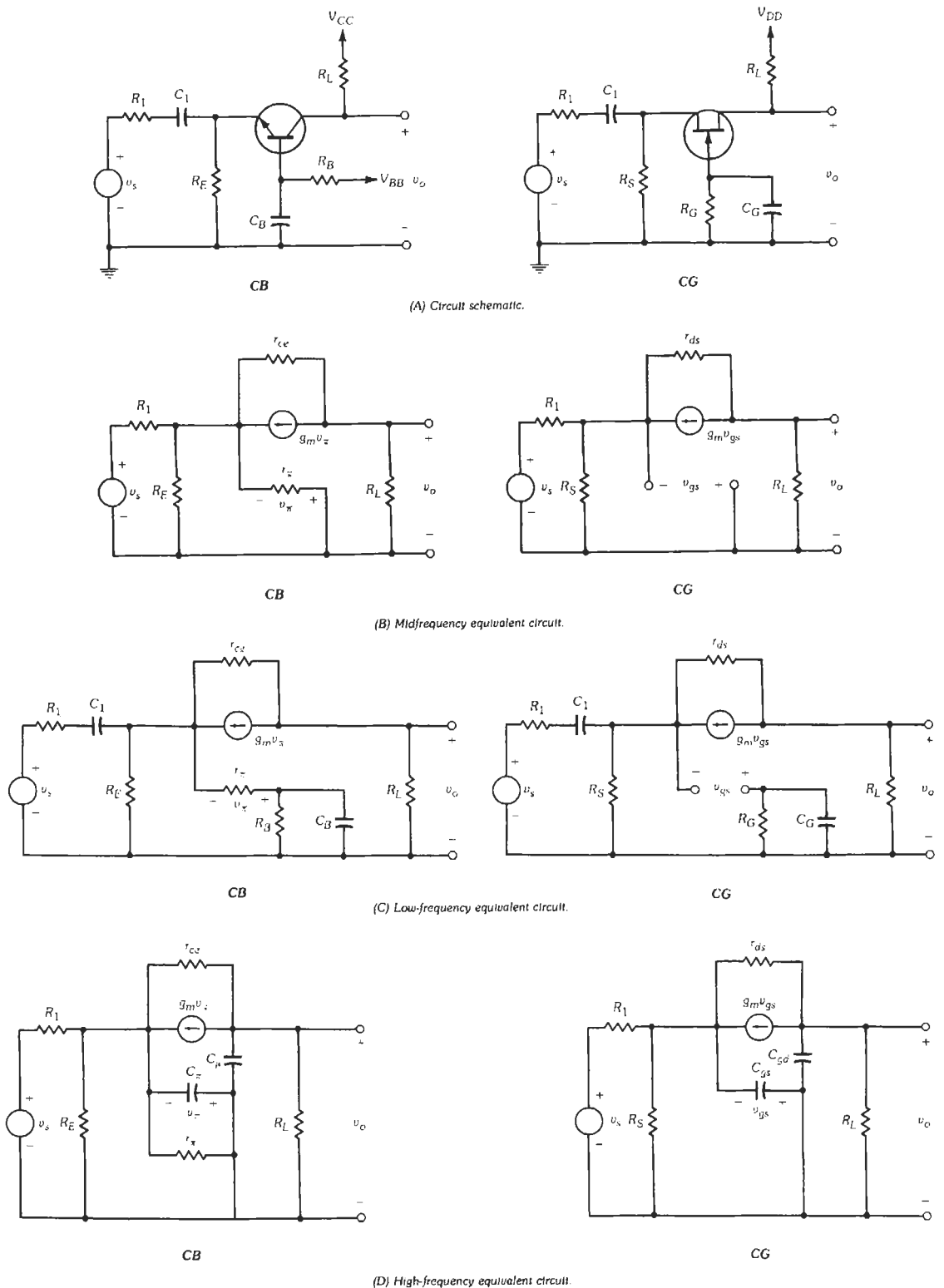


Fig. 13. Circuit schematic and small-signal equivalent circuits for common-base and common-gate amplifier.

TABLE 4. COMPARISON OF SMALL-SIGNAL AMPLIFIER PARAMETERS

	BJT			FET		
	CE	CC	CB	CS	CD	CG
r_i	1.3 kΩ	256 kΩ	26 Ω	10 ⁹ Ω	10 ⁹ Ω	224 Ω
r_o	5 kΩ	120 Ω	5 kΩ	5 kΩ	214 Ω	5 kΩ
a_v	-192	0.99	188	-22	0.96	22
a_i	50	51	0.98	—	—	1.0
Circuit Data	$R_1 = 5 \text{ k}\Omega, R_C = 5 \text{ k}\Omega$ $I_C = 1 \text{ mA}$ $\beta_0 = 50, V_A = 50 \text{ V}$			$R_1 = 5 \text{ k}\Omega, R_D = 5 \text{ k}\Omega$ $I_D = 1 \text{ mA}$ $I_{DSS} = 5 \text{ mA}, V_T = -1 \text{ V}, V_A = 50 \text{ V}$		

TABLE 5. SMALL-SIGNAL TRANSFER FUNCTIONS: COMMON-EMITTER AND COMMON-SOURCE

Bipolar Transistor	Field-Effect Transistor
<p>Midfrequency</p> $a_0 = v_o/v_s = -g_m(R_L \parallel r_{ce})[r_\pi/(R_1 + r_\pi)]$ <p>Low Frequency</p> <p>1. Due to C_1 (assume that C_E is a short-circuit)</p> $a_v(j\omega) = a_0/(1 - j\omega/\omega_{l1})$ <p>where,</p> $\omega_{l1} \approx 1/(R_1 + r_\pi)C_1$ <p>2. Due to C_E (assume that C_1 is a short-circuit)</p> $a_v(j\omega) = a_0(1 + j\omega/\omega_{i1})/(1 + j\omega/\omega_{l2})$ <p>where,</p> $\omega_{i1} \approx 1/R_E C_E$ $\omega_{l2} \approx 1/[R_E \parallel ((R_1 + r_\pi)/(\beta_0 + 1))]C_E$ <p>High Frequency</p> $a_v(j\omega) = a_0 \frac{(1 + j\omega/\omega_{zh})}{(1 + j\omega/\omega_{h1})(1 + j\omega/\omega_{h2})}$ <p>where,</p> $\omega_{zh} \approx 1/(C_\mu/g_m)$ $\omega_{h1} \approx 1/C_T(R_1 \parallel r_\pi)$ $C_T = C_\pi + C_\mu[1 + g_m(R_L \parallel r_{ce})]$ $\omega_{h2} \approx g_m/C_\pi + 1/C_\pi(R_1 \parallel R_1) + 1/C_\mu R_L$	<p>Midfrequency</p> $a_0 = v_o/v_s = -g_m(R_L \parallel r_{ds})[R_G/(R_1 + R_G)]$ <p>Low Frequency</p> <p>1. Due to C_1 (assume that C_S is a short-circuit)</p> $a_v(j\omega) = a_0/(1 - j\omega/\omega_{l1})$ <p>where,</p> $\omega_{l1} \approx 1/(R_1 + R_G)C_1$ <p>2. Due to C_S (assume that C_1 is a short-circuit)</p> $a_v(j\omega) = a_0(1 + j\omega/\omega_{i1})/(1 + j\omega/\omega_{l2})$ <p>where,</p> $\omega_{i1} \approx 1/R_S C_S$ $\omega_{l2} \approx 1/[R_S \parallel (1/g_m)]C_S$ <p>High Frequency</p> $a_v(j\omega) = a_0 \frac{(1 + j\omega/\omega_{zh})}{(1 + j\omega/\omega_{h1})(1 + j\omega/\omega_{h2})}$ <p>where,</p> $\omega_{zh} \approx 1/(C_{gd}/g_m)$ $\omega_{h1} \approx 1/C_T(R_1 \parallel R_G)$ $C_T = C_{gs} + C_{gd}[1 + g_m(R_L \parallel r_{ds})]$ $\omega_{h2} \approx g_m/C_{gs} + 1/C_{gd}(R_1 \parallel R_L) + 1/C_{gd}R_L$

Note: In this analysis it is assumed
 1. For the BJT circuit that $R_1 \ll R_B \gg r_\pi$.
 2. For the FET circuit that $R_1 \ll R_G \ll r_{gs}$.

TABLE 6. SMALL-SIGNAL TRANSFER FUNCTIONS: COMMON-COLLECTOR AND COMMON-DRAIN

Bipolar Transistor	Field-Effect Transistor
<p>Midfrequency</p> $a_0 = v_o/v_s = [(\beta_0 + 1)R_L]/[R_1 + r_\pi + (\beta_0 + 1)R_L]$ <p>Low Frequency</p> $a_v(j\omega) = a_0/(1 - j\omega/\omega_{h1})$ <p>where,</p> $\omega_{h1} \approx 1/[R_1 + r_\pi + (\beta_0 + 1)R_L]C_1$ <p>High Frequency</p> $a_v(j\omega) = a_0 \frac{(1 + j\omega/\omega_{zh})}{(1 + j\omega/\omega_{h1})(1 + j\omega/\omega_{h2})}$ <p>where,</p> $\omega_{zh} \approx 1/(C_\pi/g_m)$ $\omega_{h1} \approx 1/[(R_1 + R_L)/(1 + g_m R_L)]C_\pi + R_1 C_\mu\}$ $\omega_{h2} \approx g_m/C_\pi + 1/(R_1 \parallel R_L)C_\mu$	<p>Midfrequency</p> $a_0 = v_o/v_s = R_L/(1/g_m + R_L)$ <p>Low Frequency</p> $a_v(j\omega) = a_0/(1 - j\omega/\omega_{h1})$ <p>where,</p> $\omega_{h1} \approx 1/(R_1 + R_G)C_1$ <p>High Frequency</p> $a_v(j\omega) = a_0 \frac{(1 + j\omega/\omega_{zh})}{(1 + j\omega/\omega_{h1})(1 + j\omega/\omega_{h2})}$ <p>where,</p> $\omega_{zh} \approx 1/(C_{gs}/g_m)$ $\omega_{h1} \approx 1/[(R_1 + R_L)/(1 + g_m R_L)]C_{gs} + R_1 C_{gd}\}$ $\omega_{h2} \approx g_m/C_{gs} + 1/(R_1 \parallel R_L)C_{gd}$

Note: In this analysis it is assumed
 1. For the BJT circuit that $R_1 \ll R_B \gg r_\pi$.
 2. For the FET circuit that $R_1 \ll R_G \ll r_{gs}$.

TABLE 7. SMALL-SIGNAL TRANSFER FUNCTIONS: COMMON-BASE AND COMMON-GATE

Bipolar Transistor	Field-Effect Transistor
<p>Midfrequency</p> $a_0 = v_o/v_s = \alpha_0 \frac{R_L \parallel r_{ce} [1 + gm(R_1)]}{R_1 + r_\pi / (\beta_0 + 1)}$ <p>Low Frequency</p> <p>1. Due to C_1 (assume that C_B is a short-circuit)</p> $a_v(j\omega) = a_0/(1 - j\omega/\omega_{h1})$ <p>where,</p> $\omega_{h1} \approx 1/(R_1 + 1/g_m)C_1$ <p>2. Due to C_B (assume that C_1 is a short-circuit)</p> $a_v(j\omega) = a_0(1 + j\omega/\omega_{i1})/(1 + j\omega/\omega_{i2})$ <p>where,</p> $\omega_{i1} \approx 1/R_B C_B$ $\omega_{i2} \approx 1/[R_B \parallel [r_\pi + (\beta_0 + 1)R_1]]C_B$ <p>High Frequency</p> $a_v(j\omega) = a_0/(1 + j\omega/\omega_{h1})(1 + j\omega/\omega_{h2})$ <p>where,</p> $\omega_{h1} \approx 1/(C_\mu/g_m + R_L C_\mu)$ $\omega_{h2} \approx g_m/C_\pi + 1/R_L C_\mu$	<p>Midfrequency</p> $a_0 = v_o/v_s = \frac{R_L \parallel r_{ds} [1 + gm(R_1)]}{R_1 + 1/g_m}$ <p>Low Frequency</p> <p>1. Due to C_1 (assume that C_G is a short-circuit)</p> $a_v(j\omega) = a_0/(1 - j\omega/\omega_{h1})$ <p>where,</p> $\omega_{h1} \approx 1/(R_1 + 1/g_m)C_1$ <p>2. Due to C_G (assume that C_1 is a short-circuit)</p> $a_v(j\omega) = a_0$ <p>High Frequency</p> $a_v(j\omega) = a_0/(1 + j\omega/\omega_{h1})(1 + j\omega/\omega_{h2})$ <p>where,</p> $\omega_{h1} \approx 1/(C_{gs}/g_m + R_L C_{gd})$ $\omega_{h2} \approx g_m/C_{gs} + 1/R_L C_{gd}$

Note: In this analysis it is assumed
 1. For the BJT circuit that $R_1 \ll R_E \gg r_\pi$.
 2. For the FET circuit that $R_1 \ll R_S \ll r_{gs}$.

of the low cutoff frequency is similar to the case of the common-emitter amplifier. However, note that for the common-gate amplifier there is only one low break frequency. In the high-frequency region, typically $\omega_{h2} \gg \omega_{h1}$ so that ω_{h1} is the high cutoff frequency.

Large-Signal Characteristics

For large-signal operation, both bipolar and field-effect transistor amplifiers are classified as one of the following.

Class A—The transistor is conducting at all times. Essentially this describes a transistor operating in the linear region of the characteristic curves.

Class B—The transistor is conducting only one-half of the operating cycle, *i.e.*, for 180° of a sine-wave input. This operation is used in output stages to obtain a higher power-conversion efficiency than class A provides.

Class AB—Operation is similar to class B, but in the absence of an input signal the transistor is conducting a small quiescent current.

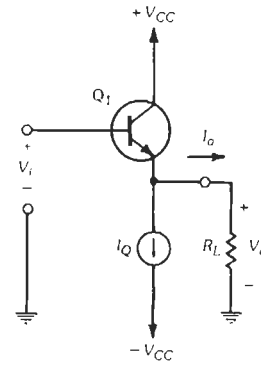
Class C—The transistor is conducting less than 180° of a sine-wave input. Class-C operation is used mostly in high-power amplifiers and oscillators.

Output Stages

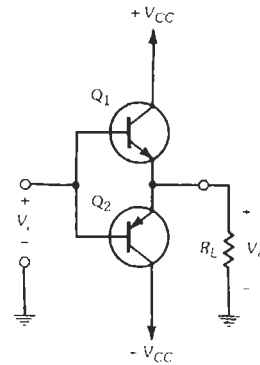
The simple emitter-follower of Fig. 14A illustrates a class-A output stage. From the maximum undistorted output voltage swing, where $|V_{max}| = |V_{min}|$, and the maximum output current swing, the maximum output power is $V_{rms}I_{rms} = V_{max}I_{max}/2$. For the ideal case ($V_{CE(sat)} = 0$ V for the output transistor and the current source is effective with 0 V across it) $V_{max} = V_{CC+}$ and $V_{min} = V_{CC-}$. With $|V_{CC+}| = |V_{CC-}| = V_{CC}$, the maximum output power is as listed in Table 8. Also listed in Table 8 is the average supply power for the output stage. The ratio of the output power to the supply power gives the power-conversion efficiency of the circuit. For a class-A amplifier this is 25% maximum, but a typical figure is about 15%.

A class-B output stage is shown in Fig. 14B. For this circuit, in the quiescent state neither transistor is conducting. On the positive half-cycle of an input sine wave, the npn transistor conducts and acts as a source of current to the load. On the negative half-cycle, the pnp transistor conducts and sinks current from the load. Hence, this configuration is known as a push-pull output stage. Since each transistor conducts only during one-half cycle of the input sine wave, the average supply power is as given in Table 8. For the class-B amplifier, the maximum power-conversion efficiency is 79%, but a typical figure is about 65%.

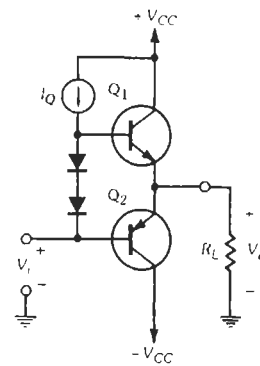
A problem with the class-B push-pull output stage is that each base-emitter junction must be forward biased before current can flow to the load. This gives rise to



(A) Emitter-follower (class A).



(B) Push-pull (class B).



(C) Push-pull (class AB).

Fig. 14. Output stages.

crossover distortion—a deadband of about $2V_{BE}$ around the zero axis of the output waveform. This problem can be avoided with the class-AB stage illustrated in Fig. 14C. The addition of a current source and two diodes permits a controlled quiescent current to flow in the two output transistors. The crossover distortion is improved at the cost of a small decrease in power-conversion efficiency.

TABLE 8. OUTPUT CHARACTERISTICS FOR LARGE-SIGNAL AMPLIFIERS

	Class A	Class B
Output voltage swing	$V_{\max} - V_{\min}$	$V_{\max} - V_{\min}$
Output current swing	$I_{\max} - I_{\min}$	$I_{\max} - I_{\min}$
Maximum output power	$V_{\max} I_{\max} / 2$ $= V_{CC} I_Q / 2$	$V_{\max} I_{\max} / 2$ $= V_{CC} I_Q / 2$
Average supply power	$2V_{CC} I_Q$	$2V_{CC} I_Q / \pi$
Maximum conversion efficiency	25%	79%

Note: This analysis assumes that $|V_{CC+}| = |V_{CC-}| = V_{CC}$

Harmonic Distortion

For a common-emitter amplifier, the small-signal output voltage is given as

$$v_o = -g_m R_L v_i$$

But for large-signal class-A operation (Fig. 15),

$$V_o = -R_L (I_S e^{V_i/V_T} - I_Q)$$

where,

$$I_Q = I_S e^{V_{BE}/V_T}$$

and

$$V_i = V_s + V_{BE}$$

Therefore

$$V_o = -R_L I_Q (e^{V_s/V_T} - 1)$$

Expanding this in a power series

$$V_o = -R_L I_Q [(V_s/V_T) + (V_s/V_T)^2/2 + (V_s/V_T)^3/6 + \dots]$$

$$= a_1 V_s + a_2 V_s^2 + a_3 V_s^3 \dots$$

where,

$$a_1 = -R_L I_Q / V_T$$

$$a_2 = -R_L I_Q / 2V_T^2$$

$$a_3 = -R_L I_Q / 6V_T^3$$

Note that in this equation, with $V_s/V_T \ll 1$ the first term dominates and the circuit is essentially linear.

For a sine-wave input

$$V_s = \hat{V}_s \sin \omega t$$

Then

$$V_o = a_1 \hat{V}_s \sin \omega t + a_2 \hat{V}_s^2 \sin^2 \omega t + a_3 \hat{V}_s^3 \sin^3 \omega t + \dots$$

$$= a_1 \hat{V}_s \sin \omega t + (a_2 \hat{V}_s^2 / 2)(1 - \cos 2\omega t) + (a_3 \hat{V}_s^3 / 4)(3 \sin \omega t - \sin 3\omega t) + \dots$$

Hence, with large-signal operation, harmonics are introduced into the output voltage waveform due to the exponential relationship of the transfer characteristic for the bipolar transistor.

Equations for the second (HD_2) and third (HD_3) harmonic distortion for each of the bipolar transistor configurations are given in Table 9. Note, for example, that for the CE stage with $\hat{V}_s = 10$ mV, $HD_2 = 10\%$ and $HD_3 = 0.62\%$.

In an ideal push-pull amplifier stage, the positive half-cycle of the output waveform exactly matches the negative half-cycle. The result is that all even harmonic terms are balanced out, leaving the third harmonic as the prime source of harmonic distortion.

The transfer characteristic of an FET has a square-law relationship. Hence, ideally, only second-harmonic distortion appears at the output of a field-effect transistor amplifier. In practice, all FETs deviate somewhat from the ideal transfer characteristic and exhibit a third-order term. However, this term is usually very small. For a common-source amplifier

$$HD_2 = (V_4)(\hat{V}_s/V_T)(I_{DSS}/I_Q)^{1/2}$$

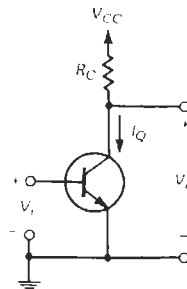


Fig. 15. Large-signal operation of common-emitter amplifier circuit.

TABLE 9. HARMONIC DISTORTION IN BIPOLAR TRANSISTOR AMPLIFIERS

Parameter	CE	CC	CB
HD_2	$\frac{1}{4}(\hat{V}_s/V_T)$	$\frac{1}{4}(\hat{V}_s/V_T)/(1 + g_m R_L)$	$\frac{1}{4}(\hat{V}_s/V_T)/(1 + g_m R_E)$
HD_3	$\frac{1}{24}(\hat{V}_s/V_T)^2$	$\frac{1}{24}(\hat{V}_s/V_T)^2 \frac{1 - 3g_m R_L/(1 + g_m R_L)}{1 + g_m R_L}$	$\frac{1}{24}(\hat{V}_s/V_T)^2 \frac{1 - 3g_m R_E/(1 + g_m R_E)}{1 + g_m R_E}$

The second harmonic in other FET single amplifier configurations follows from Table 9.

DIFFERENTIAL AMPLIFIERS*

A simple differential amplifier is exemplified by the BJT emitter-coupled pair illustrated in Fig. 16. A source-coupled FET pair could readily replace the bipolar transistors. Fig. 17 shows the small-signal difference and common-mode input voltages of the circuit. The relationship between the difference and common-mode operation is given in Table 10. The design equations are listed in Table 11.

CURRENT SOURCES†

A useful component in many circuits, but especially with differential amplifiers, is a transistor current source.

Current Mirrors

The concept of current mirrors is widely used in linear integrated circuits, where the matching and tracking of transistor characteristics are well controlled. The idea may also be used with discrete dual transistors that have the same matching and tracking characteristics as in ICs.

Four simple, but effective, current sources are shown in Fig. 18. Their characteristics are listed in Table 12.

In Fig. 18A, with $V_{BE1} = V_{BE2}$ and the characteristics of Q_1 closely matching those of Q_2 , to a first order $I_{C1} = I_{C2}$. The effect of finite base currents is included in the equations given in Table 12, where, provided $\beta_F \gg 2$, $I_{C2} = I_R$. The effects of base current are even further reduced in the improved circuit of Fig. 18B.

The output resistance of a transistor current source can be improved, as in the Widlar current source of Fig. 18C. For this circuit, note that typically $r_{\pi 2} > (1/g_{m1} + R_2)$; then the output resistance is simply given as $r_o = r_{ce} (1 + g_{m2}R_2)$. The solution of the transcendental equation for I_{C2} can be done readily by trial and error with one or two iterations, or with a programmable calculator. An even higher output resistance can generally be obtained with the Wilson current source of Fig. 18D, where, due to negative feedback, $r_o \approx \beta_0 r_{ce2}/2$.

The Widlar circuit is useful for obtaining small output currents, and the Wilson circuit is useful for obtaining high output resistance and low sensitivity to transistor base currents.

Each of these circuits is also viable with matching dual JFETs or MOSFETs.

Active Loads

Another useful application for current sources is as active loads in amplifier circuits. The output resistance of the current source provides a high collector resistance

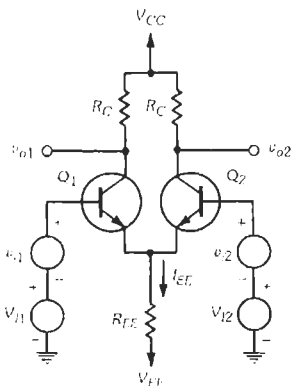


Fig. 16. Differential amplifier pair.

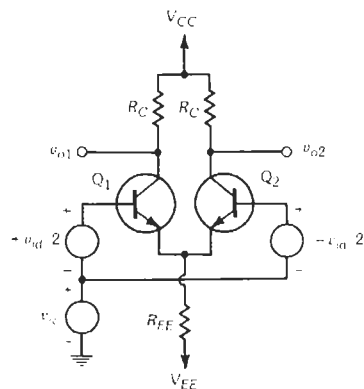


Fig. 17. Differential amplifier pair showing small-signal difference and common-mode input voltages.

* References 8 and 10.

† References 9 and 11.

TABLE 10. RELATIONSHIP FOR DIFFERENCE AND COMMON-MODE OPERATION

	Difference Mode	Common Mode
Input	$v_{id} = v_{i1} - v_{i2}$	$v_{ic} = (v_{i1} + v_{i2})/2$
Output	$v_{od} = v_{o1} - v_{o2}$	$v_{oc} = (v_{o1} + v_{o2})/2$
Gain	$a_{dm} = v_{od}/v_{id}$	$a_{cm} = v_{oc}/v_{ic}$
Input	$v_{i1} = v_{ic} + v_{id}/2$ $v_{i2} = v_{ic} - v_{id}/2$	
Output	$v_{o1} = a_{cm}v_{ic} + a_{dm}v_{id}/2$ $v_{o2} = a_{cm}v_{ic} - a_{dm}v_{id}/2$	

TABLE 11. DIFFERENTIAL AMPLIFIERS: SMALL-SIGNAL DESIGN EQUATIONS

Difference Mode	Common Mode
$a_{dm} = v_{od}/v_{id} = -g_m R_C$	$a_{cm} = v_{oc}/v_{ic} = -g_m R_C r_{\pi} / [r_{\pi} + (\beta_0 + 1)2R_{EE}]$
$r_{id} = v_{id}/i_{id} = 2r_{\pi}$	$r_{ic} = v_{ic}/i_{ic} = r_{\pi} + (\beta_0 + 1)2R_{EE}$
$r_{od} = 2R_C$	$r_{oc} = 2R_C$
CMRR (dB) = $20 \log a_{dm}/a_{cm} $	

TABLE 12. CURRENT SOURCES

Type	Reference Current (I_R)	Output Current (I_{C2})	Output Resistance (r_o)
Simple	$(V_{CC} - V_{BE1})/R$	$I_R/(1 + 2/\beta_F)$	r_{ce2}
Improved	$(V_{CC} - V_{BE3} - V_{BE1})/R$	$I_R/[1 + 2/(\beta_F^2 + \beta_F)]$	r_{ce2}
Widlar	$(V_{CC} - V_{BE1})/R$	$(V_T/R_2) \ln (I_{C1}/I_{C2})$	$r_{ce2}[1 + \beta_0 R_2/(1/g_{m1} + r_{\pi 2} + R_2)]$
Wilson	$(V_{CC} - V_{BE2} - V_{BE3})/R$	$I_R[1 - 2/(\beta_F^2 + 2\beta_F + 2)]$	$\beta_0 r_{ce2}/2$

with a relatively small voltage drop. A basic circuit is illustrated in Fig. 19A.

Neglecting base currents, in this circuit

$$I_{C1} = I_{C2} = I_{C3} = I_{C4} = I_{C5} = I_{EE}/2$$

The effective load resistance for differential input signals is

$$R_{L1} = R_{L2} = r_{ce1} \parallel r_{ce3}$$

The output resistance of the current source may be further increased by adding resistances, $R_{E3} = R_{E4}$, in the emitter leads of Q_3 and Q_4 , respectively. The effective load resistance then becomes

$$R_{L1} = R_{L2} = r_{ce1} \parallel [r_{ce3} (1 + g_{m3} R_{E3})]$$

Fig. 19B shows an amplifier with a differential input signal but a single-ended output signal. In this circuit, not only does Q_4 effectively provide a large collector resistance for Q_2 , but also since the collector current of

Q_4 mirrors that of Q_3 , the small-signal change in the collector current of Q_1 is added to that of Q_2 to provide the full differential voltage gain of the stage at the single-ended output, v_o .

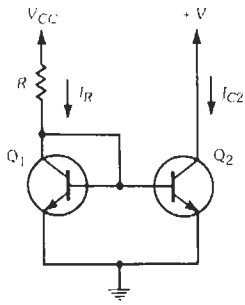
Each of these circuits may be designed with JFETs or MOSFETs. Also, many times it is advantageous to use JFETs as the amplifier devices and BJTs for active loads.

FEEDBACK AMPLIFIERS*

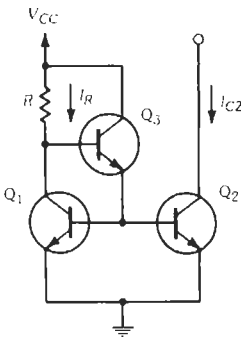
Negative feedback is widely used in amplifier design because it produces several benefits. Among these are:

- Desensitivity against parameter changes
- Improved input/output characteristics
- Reduction in harmonic distortion
- Increased frequency response

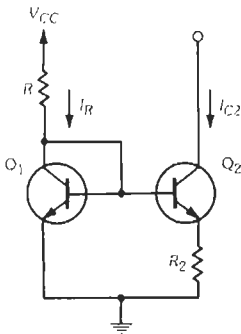
* References 7 and 15.



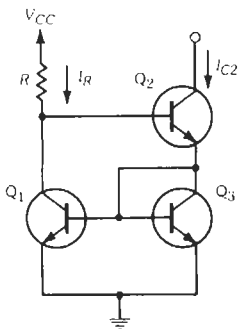
(A) Simple.



(B) Improved.



(C) Widlar.



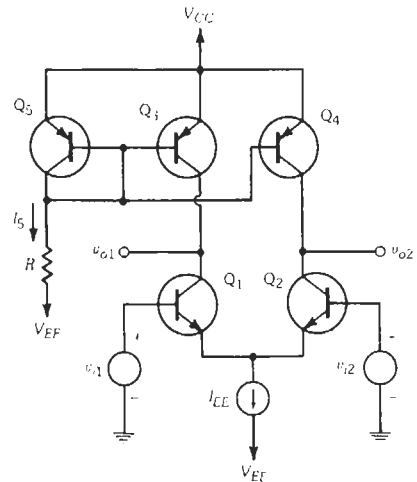
(D) Wilson.

Fig. 18. Current-mirror circuits.

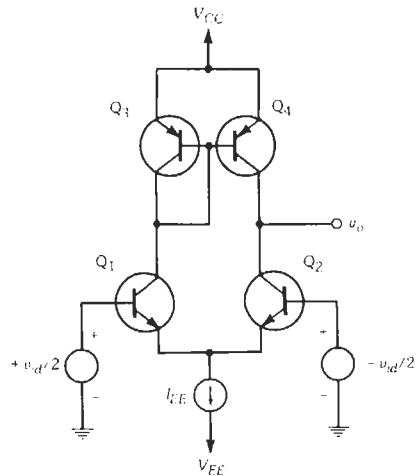
Unfortunately there are some disadvantages:
 Reduced gain
 Instability problems

Basic Properties

The basic properties of negative feedback are presented with the aid of the block diagram in Fig. 20. In this diagram, the gain of the basic amplifier, a , is controlled by undependable transistor parameters that are subject to considerable variation. It is assumed that gain a is much larger than needed. It is also assumed in the block diagram that there is no loading of one block by another block and the signal path is only in the direction of the arrows; *i.e.*, the gain of each block is unilateral.



(A) Differential amplifier pair with active load.



(B) Modified active load for differential amplifier pair.

Fig. 19. Active loads.

With positive feedback:

$$s_i = s_s + s_f$$

That is, the signals add and $s_i > s_s$.

With negative feedback:

$$s_i = s_s - s_f$$

That is, the signals subtract and $s_i < s_s$.

Also from the block diagram:

$$s_o = a s_i$$

$$s_f = f s_o$$

Chart 10 presents the basic equation for negative feedback, along with other important definitions. Note from the basic equation that with negative feedback the open-loop gain is reduced by the factor $1 + T$. Further, if $T \gg 1$ the closed-loop gain becomes $1/f$. It is common for the feedback network to be made up of stable, high-precision, passive components; the value of f is then well defined, and so is the overall amplifier gain.

Basic Feedback Circuit Topologies

Preliminary to the presentation of the feedback topologies, consider the circuit diagrams of four basic classes of amplifier (Fig. 21). The classification is based on the magnitude of the input and output impedance relative to the source and load impedance. The ideal characteristics for each class of amplifier are given in Table 13.

A feedback amplifier is described in terms of the way in which the feedback network is connected to the basic amplifier. There are four basic feedback circuit topologies. These are illustrated in the block diagrams of Fig. 22. Clearly, in Fig. 22A the feedback network is connected in series with the input terminals of the basic amplifier and in shunt with the output terminals. Notice:

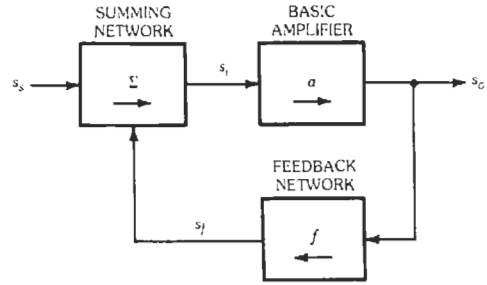
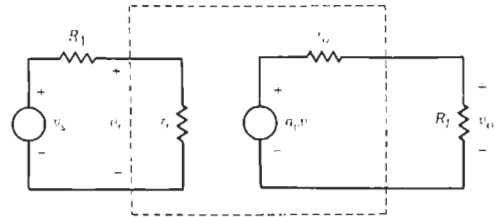
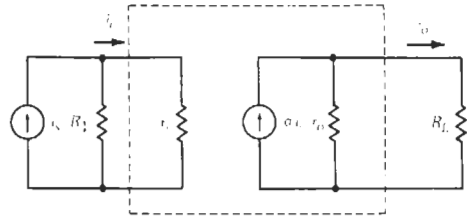


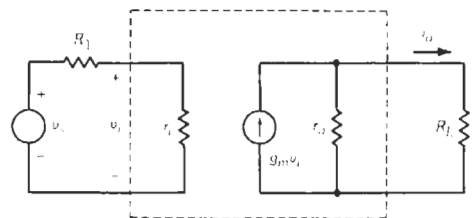
Fig. 20. Ideal feedback block diagram.



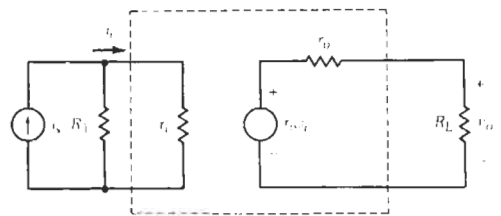
(A) Voltage amplifier ($v_{in} \rightarrow v_{out}$).



(B) Current amplifier ($i_{in} \rightarrow i_{out}$).



(C) Transconductance amplifier ($v_{in} \rightarrow i_{out}$).



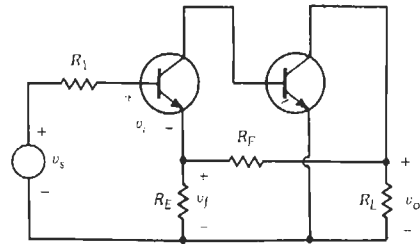
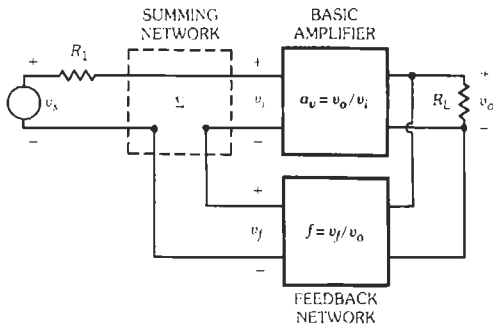
(D) Transresistance amplifier ($i_{in} \rightarrow v_{out}$).

CHART 10. BASIC EQUATIONS FOR NEGATIVE FEEDBACK

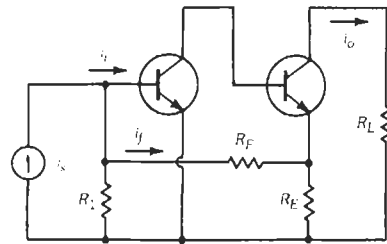
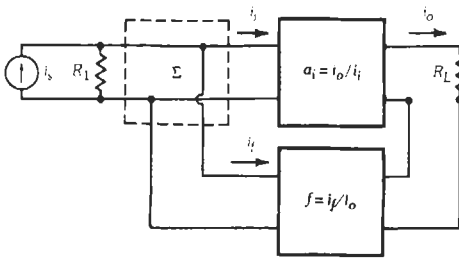
A	$= a/(1 + af) = a/(1 + T)$	$= a/D_s$
A	$=$ closed-loop gain	$= s_o/s_s$
a	$=$ open-loop gain	$= s_o/s_i$
f	$=$ feedback factor	$= s_f/s_o$
T	$=$ loop gain $= af$	$= s_f/s_i$
D_s	$=$ desensitvity factor	$= 1 + T$

Note: The basic equation for negative feedback is often written as $A_{CL} = A_{OL}/(1 + A_{OL}\beta)$

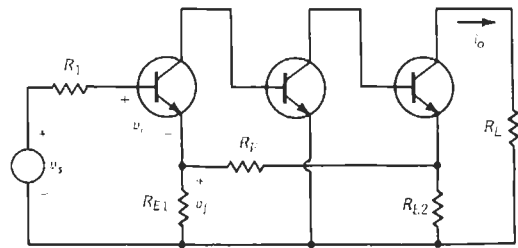
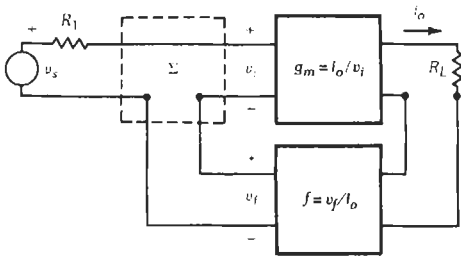
Fig. 21. Amplifier classifications.



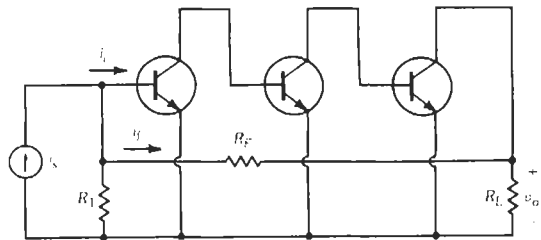
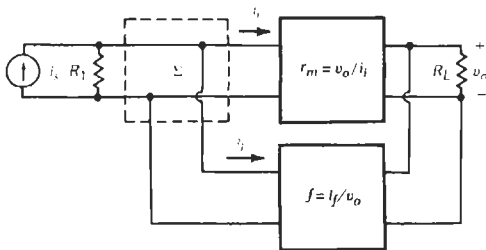
(A) Series-shunt feedback configuration (a voltage amplifier).



(B) Shunt-series feedback configuration (a current amplifier).



(C) Series-series feedback configuration (a transconductance amplifier).



(D) Shunt-shunt feedback configuration (a transresistance amplifier).

Fig. 22. Feedback-amplifier topologies.

TABLE 13. IDEAL AMPLIFIER CHARACTERISTICS

Parameter	Voltage Amplifier	Current Amplifier	Trans-conductance Amplifier	Trans-resistance Amplifier
Input resistance (r_i)	$\rightarrow \infty (>> R_i)$	$\rightarrow 0 (<< R_i)$	$\rightarrow \infty (>> R_i)$	$\rightarrow 0 (<< R_i)$
Output resistance (r_o)	$\rightarrow 0 (<< R_L)$	$\rightarrow \infty (>> R_L)$	$\rightarrow \infty (>> R_L)$	$\rightarrow 0 (<< R_L)$
Transfer characteristic (a)	$a_v = v_o/v_i$	$a_i = i_o/i_s$	$g_m = i_o/v_s$	$r_m = v_o/i_s$

Note: Transconductance (g_m) used here refers to an amplifier parameter, not necessarily just a device parameter.

With series feedback at the input, voltages v_s and v_f are algebraically summed.

With shunt feedback at the input, currents i_s and i_f are algebraically summed.

With series feedback at the output, a current i_o is sampled.

With shunt feedback at the output, a voltage v_o is sampled.

Included in Fig. 22 are simple examples of each feedback connection, implemented with bipolar transistors. Especially note the correspondence between each circuit schematic and the related block diagram. To avoid complexity, all biasing resistors have been omitted from the circuit diagrams, but it is assumed that all transistors are biased in the forward-active region to yield a high-gain amplifier.

Presented below is a method of analysis of feedback amplifiers. The design of a feedback amplifier would follow a similar procedure.

- Identify the feedback topology:
 - Is feedback signal s_f applied in series (v_f) or in shunt (i_f) with the signal source s_s ?
 - Is sampled signal s_o obtained at the output node (v_o) or from the output loop (i_o)?
- Draw the basic amplifier circuit with the feedback set to zero; that is
 - For the correct input circuit:
 - With shunt sampling, short-circuit the output nodes to set $v_o = 0$.
 - With series sampling, open-circuit the output loop to set $i_o = 0$.
 - For the correct output circuit:
 - With series summing, open-circuit the input loop to set $v_f = 0$.
 - With shunt summing, short-circuit the input nodes to set $i_f = 0$.
- Indicate s_f and s_o , and solve for the feedback factor ($f = s_f/s_o$).
- Evaluate the open-loop gain function (a).
- From a and f , find T , D_s , A , R_i , and R_o .

Information to aid in the analysis and design of feedback amplifiers is summarized in Table 14. Notice that the effect of negative feedback is to modify the open-loop parameters of an amplifier so that the closed-loop performance approaches the ideal characteristics as listed in Table 13.

BANDPASS AMPLIFIERS*

A bandpass amplifier selectively amplifies a narrow band of frequencies around a center frequency. The selectivity is indicated by

$$Q = \omega_0/\Delta\omega = \omega_0/(\omega_h - \omega_l)$$

where,

ω_0 is the center frequency,

$\Delta\omega$ is the bandwidth,

ω_h and ω_l are the high and low cutoff frequencies (-3 dB).

For a bandpass amplifier, typically $Q > 10$.

Single-Tuned Interstage

A single-tuned interstage is modeled in Fig. 23. Resistor R is the total shunt resistance at the output of an amplifying stage including the input resistance of the following stage. Similarly, capacitor C is the total shunt capacitance.

For a single-tuned interstage the general transfer function is

$$a_v(j\omega) = v_o/v_i = -g_m R/[1 + jQ(\omega/\omega_0 - \omega_0/\omega)]$$

and

$$\Delta\omega = \omega_h - \omega_l = \omega_0/Q = 1/RC$$

At resonance:

$$\omega_0 = 1/(LC)^{1/2}$$

* References 18 and 20.

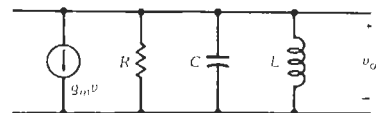


Fig. 23. Single-tuned interstage.

TABLE 14. FEEDBACK-AMPLIFIER ANALYSIS

	Series-Shunt	Shunt-Series	Series-Series	Shunt-Shunt
Input signal (s_i)	v_i	i_i	v_i	i_i
Feedback signal (s_f)	v_f	i_f	v_f	i_f
Output Signal (s_o)	v_o	i_o	i_o	v_o
To calculate loading of feedback network				
At input	Short output node	Open output loop	Open output loop	Short output node
At output	Open input loop	Open input loop	Short input node	Short input node
To calculate feedback factor	Drive feedback network with a voltage and calculate open-circuit voltage v_f	Drive feedback network with a current and calculate short-circuit current i_f	Drive feedback network with a current and calculate open-circuit voltage v_f	Drive feedback network with a voltage and calculate short-circuit current i_f
Feedback factor (f)	v_f/v_o	i_f/i_o	v_f/i_o	i_f/v_o
Open-loop gain (a)	$a_v = v_o/v_i$	$a_i = i_o/i_i$	$g_m = i_o/v_i$	$r_m = v_o/i_i$
Loop gain (T)	$a_v f$	$a_i f$	$g_m f$	$r_m f$
Closed-loop				
Gain (A)	$A_v = a_v/(1 + T)$	$A_i = a_i/(1 + T)$	$G_m = g_m/(1 + T)$	$R_m = r_m/(1 + T)$
Input resistance (R_i)	$r_i(1 + T)$	$r_i/(1 + T)$	$r_i(1 + T)$	$r_i/(1 + T)$
Output resistance (R_o)	$r_o/(1 + T)$	$r_o(1 + T)$	$r_o(1 + T)$	$r_o/(1 + T)$

and the center-frequency voltage gain is

$$a_v(j\omega_o) = -g_m R$$

With $Q \gg 1$, the high and low cutoff frequencies may be determined from

$$\omega_h = \omega_o (1 + 1/2Q)$$

$$\omega_l = \omega_o (1 - 1/2Q)$$

For a cascade of n synchronous single-tuned stages that are unilateral and noninteracting, the magnitude of the center-frequency voltage gain is

$$|A_v| = a_v^n = (g_m R)^n$$

and the overall bandwidth is

$$\Delta \omega_n = \Delta \omega (2^{1/n} - 1)^{1/2}$$

Double-Tuned Interstage

A double-tuned interstage is modeled in Fig. 24. Here, R_1 and C_1 are the total shunt resistance and capacitance at the output of an amplifying stage, and R_2 and C_2 are the total shunt resistance and capacitance at the input of the following stage. The transformer is loosely coupled so that the coefficient of coupling $k \ll 1.0$, typically < 0.1 .

For a double-tuned interstage, the general transfer function is

$$a_v(j\omega) = v_o/v_i =$$

$$-g_m(R_1 R_2)^{1/2} \frac{s}{[(1 + s^2)^2 - 2(s^2 - b/2)s^2 + a^4]^{1/2}}$$

where,

$$s = k(Q_1 Q_2)^{1/2}$$

$$p = Q_1/Q_2$$

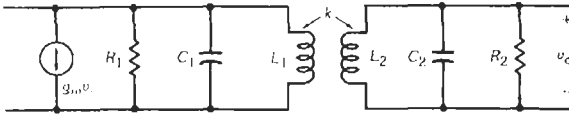


Fig. 24. Double-tuned interstage.

$$b = (Q_1/Q_2 + Q_2/Q_1)$$

$$a = (Q_1 Q_2)^{1/2} (\omega/\omega_0 - \omega_0/\omega)$$

At resonance:

$$\omega_0 = 1/(L_1 C_1)^{1/2} = 1/(L_2 C_2)^{1/2}$$

The gain of the interstage is at a maximum with $s = 1$; the transformer is then critically coupled, and

$$a_v(j\omega_0) = -g_m(R_1 R_2)^{1/2}/2$$

with

$$k_c = 1/(Q_1 Q_2)^{1/2}$$

In practice, it is customary to use not the value of k which gives the maximum gain, but that value which gives the flattest selectivity curve. The circuit is then said to be transitionally coupled.

$$k_t = [(1/Q_1^2 + 1/Q_2^2)/2]^{1/2}$$

For a transitionally coupled circuit

$$a_v(j\omega_0) = -g_m(R_1 R_2)^{1/2} [2p(1 + p^2)]^{1/2} / (1 + p)^2$$

and the bandwidth is

$$\Delta\omega = [(1 + p)/\sqrt{2} p] \omega_0 / Q_2$$

With equal Q 's, $Q_1 = Q_2 = Q_m$, it follows that $k_t = k_c = k_m$ and

$$a_v(j\omega_0) = -g_m(R_1 R_2)^{1/2} / 2$$

$$\Delta\omega = \sqrt{2} \omega_0 / Q_m$$

and

$$k_m = \Delta\omega / \sqrt{2} \omega_0 = 1/Q_m$$

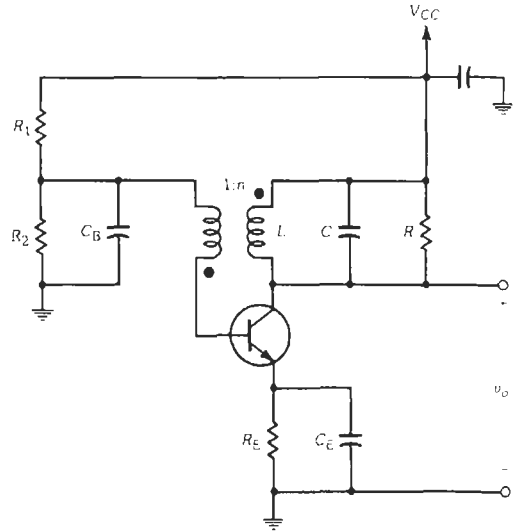
Comparing a double-tuned interstage (with $n = 1$) and a single-tuned interstage yields the same center-frequency gain, but the double-tuned circuit provides a squarer response over the passband.

For a cascade of n synchronous double-tuned stages that are unilateral and noninteracting, the magnitude of the center-frequency voltage gain is

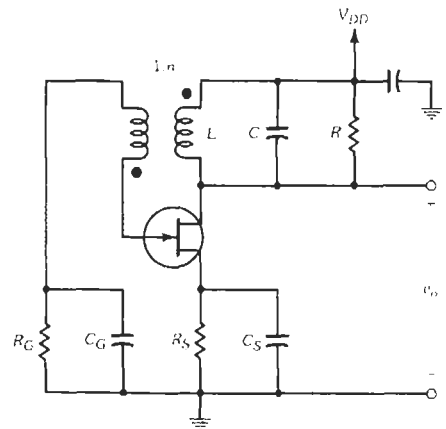
$$|A_v| = a_v^n = [g_m(R_1 R_2)^{1/2} / 2]^n$$

and the overall bandwidth is

$$\Delta\omega_n = \Delta\omega(2^{1/n} - 1)^{1/4}$$



(A) Bipolar junction transistor.



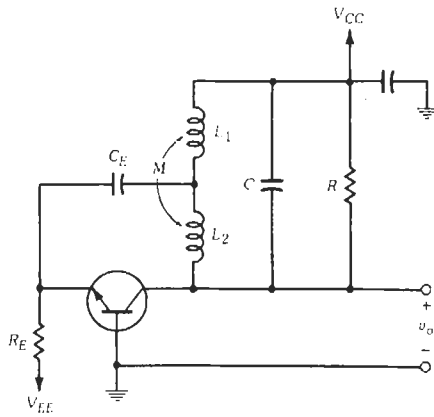
(B) Junction field-effect transistor.

Fig. 25. Transformer-coupled feedback oscillators.

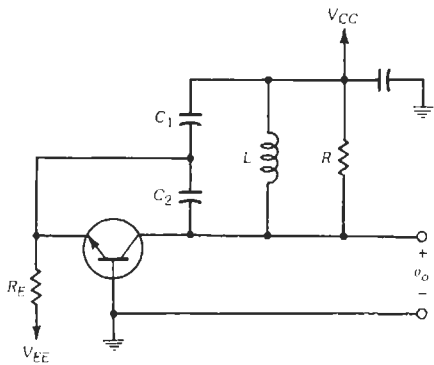
SINUSOIDAL OSCILLATORS*

In the section on feedback amplifiers, *negative* feedback was used to advantage to produce near-ideal amplifiers. In this section, *positive* feedback is used to produce near-ideal sinusoidal oscillators. Fig. 25 illustrates two transformer-coupled feedback oscillator circuits. In Fig. 26, the popular Colpitts and Hartley circuits are shown. These are LC tuned-circuit feedback oscillators. Two RC oscillator circuits are presented in Fig. 27. Equations for the oscillation frequency and conditions for oscillation are listed in Table 15.

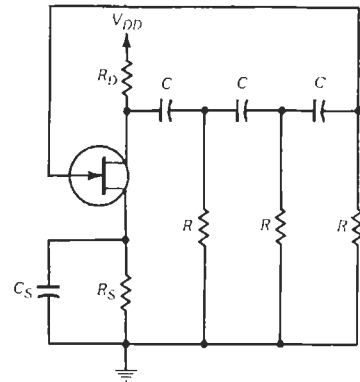
* References 3 and 19.



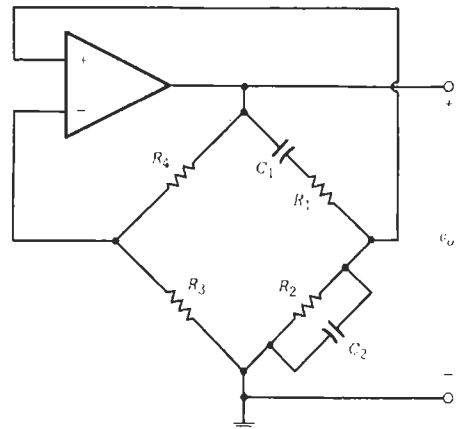
(A) Hartley.



(B) Colpitts.



(A) Phase-shift.



(B) Wien bridge.

Fig. 26. Tuned-circuit feedback oscillators.

Tight coupling is used in the transformer-coupled oscillators so that coefficient of coupling k is close to 1.0. To sustain oscillations, but to prevent gross distortion in the output waveform, the ratio $g_m R/n$ is typically chosen to be about 3. The turns ratio for the transformer is not critical; a value of about 10 is appropriate. For the Hartley and Colpitts oscillators, "turns ratio" n is typically about 3.

Fig. 27. RC oscillators.

For the RC oscillators, loop gain T is equal to 1 at ω_0 , provided that the voltage gain of the basic amplifier is as listed in Table 15. The RC oscillators are used at low frequencies, where the size of the inductor in an LC oscillator would be too large.

TABLE 15. OSCILLATOR EQUATIONS

Circuit	Oscillation Frequency	Condition for Oscillation
Transformer Coupled	$\omega_o = 1/(LC)^{1/2}$	$g_m R/n \geq 1$
Hartley	$\omega_o = 1/(LC)^{1/2}$ $L = L_1 + L_2$	$n \geq 1/\alpha_0$ $n = (L_1 + L_2)/L_1$
Colpitts	$\omega_o = 1/(LC)^{1/2}$ $C = C_1 C_2 / (C_1 + C_2)$	$n \geq 1/\alpha_0$ $n = (C_1 + C_2)/C_2$
Phase-Shift	$\omega_o = 1/RC$	$a_v(j\omega_o) \geq 29$
Wien Bridge	$\omega_o = 1/R_1 C_1 = 1/R_2 C_2$	$a_v(j\omega_o) = 3$ i.e., $R_4 = 2R_3$

Prime requirements of sinusoidal oscillators are (1) good amplitude stability, (2) good frequency stability, and (3) a low harmonic content in the output waveform.

A stable voltage supply and stable passive components are basic to good amplitude stability. Stable passive components are required also for a stable output frequency. In addition, the oscillation frequency should be much less than the effective f_T of the active device; the tank-circuit capacitance should be large compared with C_{μ} or C_{gd} ; and, where the feedback is to the emitter or source, a small resistance ($R \approx 10/g_m$) is useful in the feedback path. The harmonic content at the output is minimized with a low L/C ratio. Also it should be noted that compared to a common-emitter or common-source connection a common-base or common-gate circuit gives better frequency stability and is generally less critical in adjustment.

For the best sinusoidal output, the collector current and load resistor R should be chosen so that the transistor just cuts off but does not saturate on the peaks

of the waveform. However, this class-A operation does lead to a low power-conversion efficiency. An oscillator operated near class A has good frequency stability and a low harmonic content, but for a high power-conversion efficiency class-C operation is preferred.

The RC oscillators are normally used at frequencies below 1 MHz. The Hartley is a good general-purpose oscillator to about 10 MHz. The Colpitts oscillator has good stability and a low harmonic output. It is used almost exclusively at frequencies above 10 MHz.

PULSE CIRCUITS†

Pulse Shaping

The shaping of pulses by means of passive RC and RL circuits is common. Fig. 28 shows the response of

† References 12 and 16.

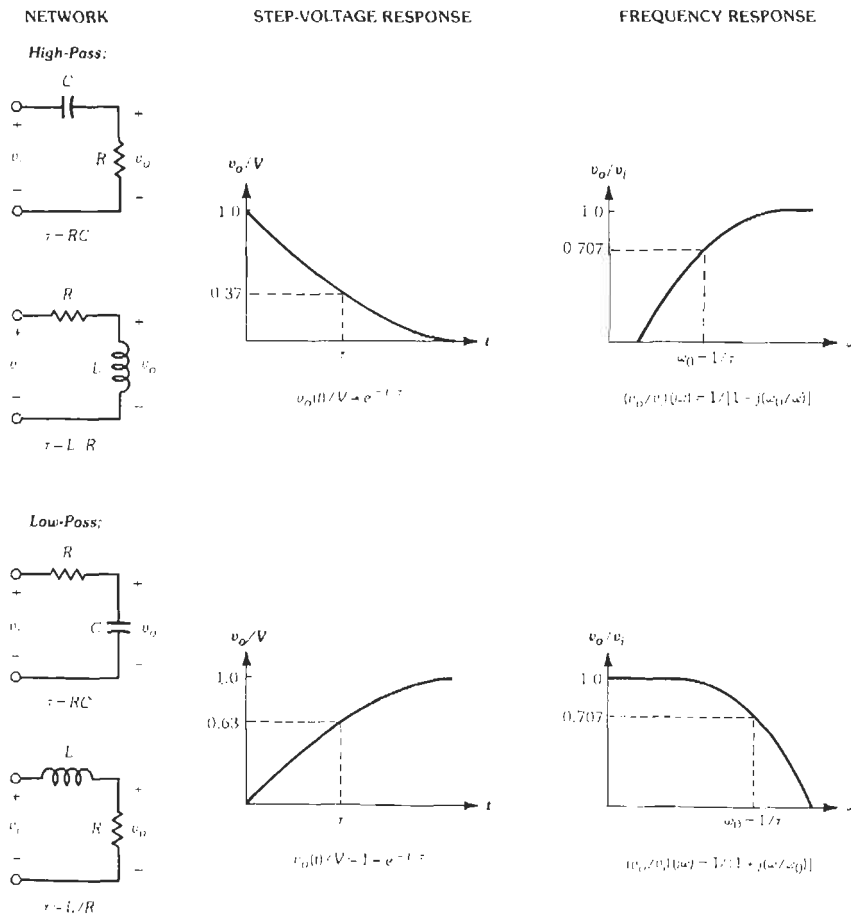
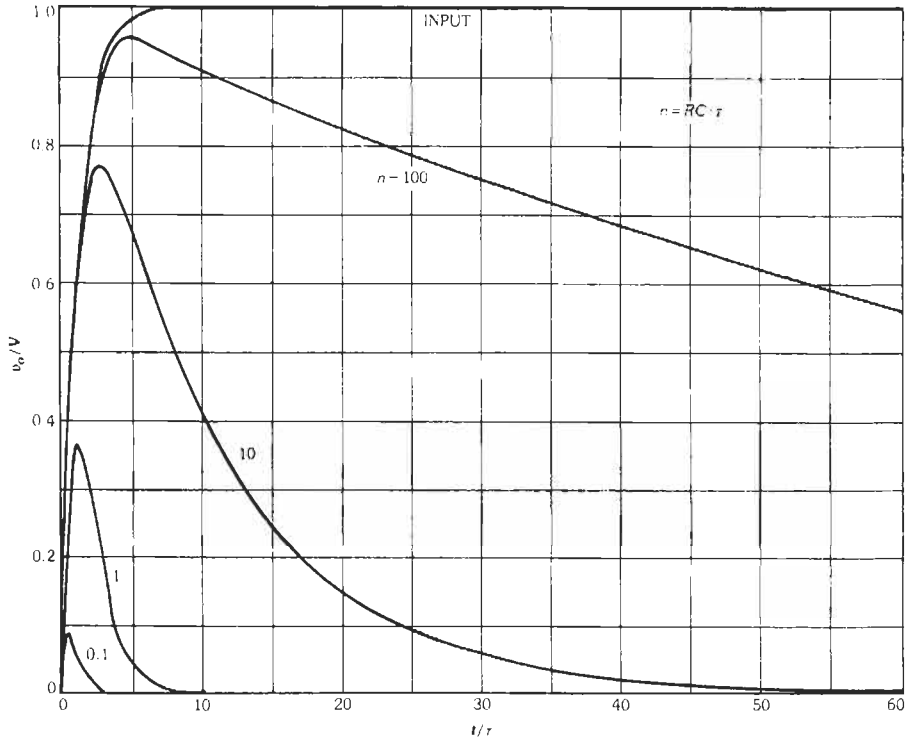
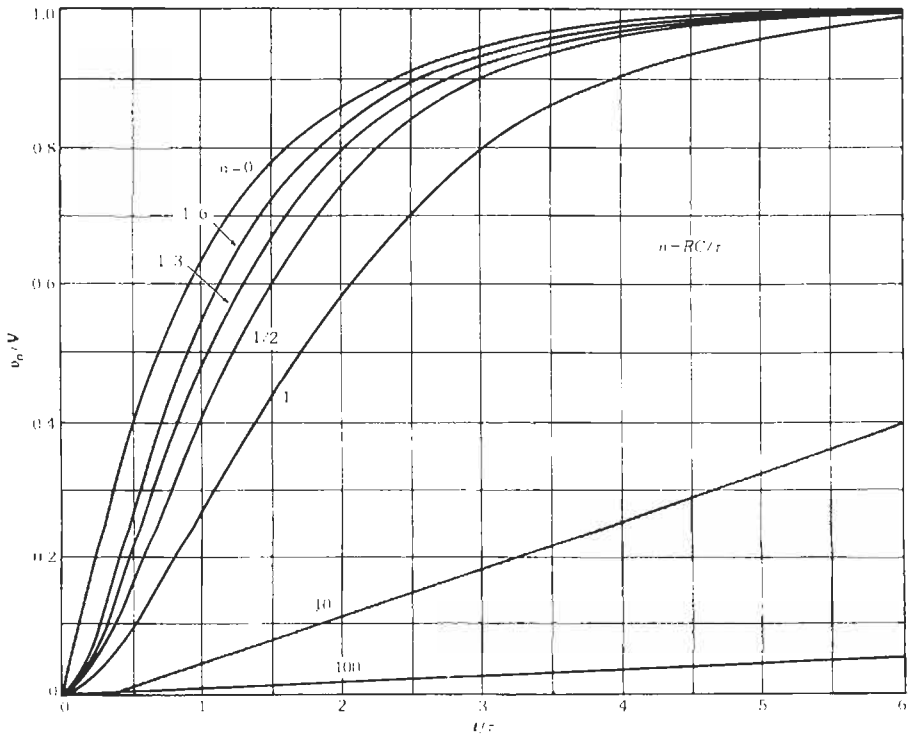


Fig. 28. Response of single RC and RL networks.

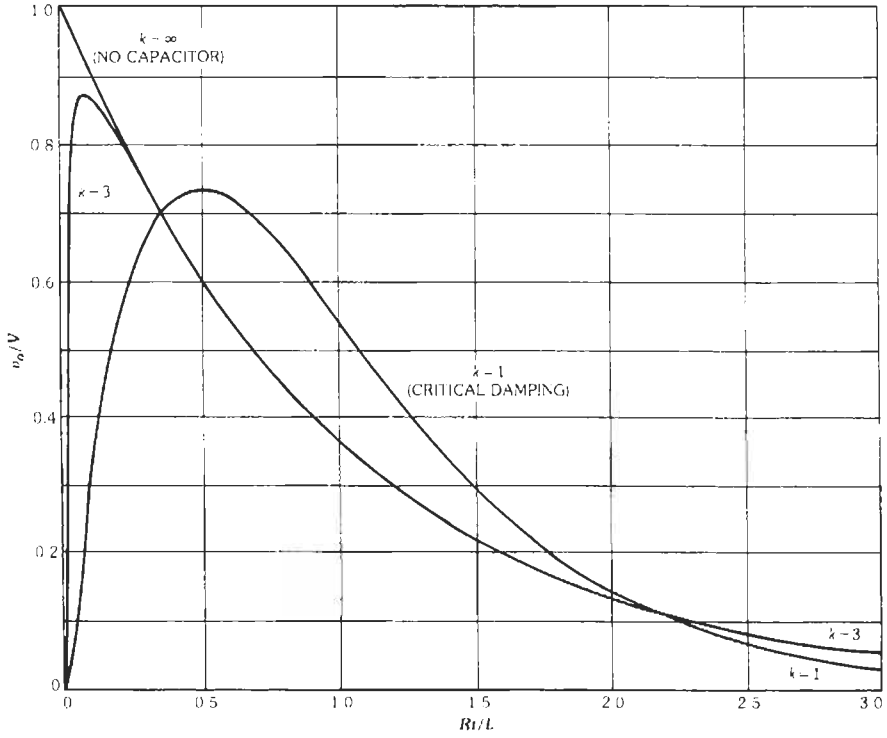


(A) Response of a high-pass RC circuit to an exponential input.

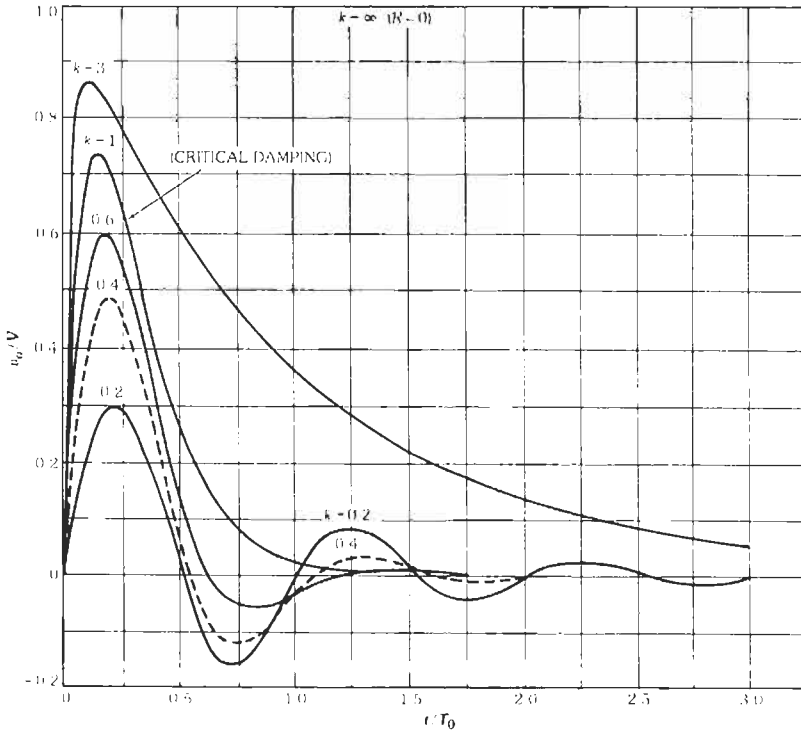


(B) Response of a low-pass RC circuit to an exponential input.

Fig. 29. Response curves for RC networks.



(A) Critically damped and overdamped cases for fixed values of R and L



(B) Fixed values of L and C

Fig. 30. Response curves for the RLC circuit of Fig. 31.

each of these simple networks to a step-voltage input and a sine-wave input. For a step input, $v_i = V$. With a step-voltage at the input of a simple high-pass network, the fall time between the 0.9 and 0.1 points on the output waveform is $t_f = 2.2\tau$. Similarly, for a simple low-pass network the rise time between the 0.1 and 0.9 points of the output waveform is $t_r = 2.2\tau$. The high-pass network also serves as a voltage differentiator, since $v_o(t) = \tau(dv_i/dt)$. The low-pass network serves as a voltage integrator, since $v_o(t) = (1/\tau) \int v_i dt$. The response of these shaping networks to an exponential voltage input, $v_i = V(1 - e^{-t/\tau})$, is illustrated in Fig. 29.

Fig. 30 shows the response of an RLC circuit (Fig. 31) to a step-voltage input. Parameter k in these curves is related to capacitance C or resistance R by the following equation:

$$k = (1/2R) \sqrt{L/C}$$

For the critically damped case, $k = 1$ and

$$v_o(t)/V = (4Rt/L)e^{-2Rt/L}$$

In the overdamped case $k > 1$, and with $T_0 \equiv 2\pi \sqrt{LC}$

$$v_o(t)/V = e^{-\pi k t/T_0} - e^{-4\pi k t/T_0} \approx e^{-Rt/L}$$

In the underdamped case, $k < 1$ and

$$v_o(t)/V = [2k/(1 - k^2)^{1/2}] e^{-2\pi t/T_0} \sin 2\pi(1 - k^2)^{1/2} t/T_0$$

Multivibrators

The three forms of a multivibrator circuit are bi-stable, monostable, and astable. Each of these operations is best performed with digital integrated circuits. However, discrete transistor versions of these circuits may be used, and some simple examples using saturating bipolar transistors are shown in Figs. 32-34.

For the bistable circuit of Fig. 32, typically $R_B \approx 10R_C$. Included in the figure is the truth table for this

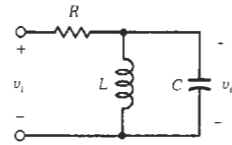


Fig. 31. RLC network.

bistable latch circuit. Note that since \bar{Q} is defined as the inverse of Q , the $R = S = 1, Q = \bar{Q} = 0$ state is not allowed.

In the monostable circuit of Fig. 33, a rectangular pulse is generated at v_{out} due to the trigger pulse at v_{in} . The output pulse width is given by

$$PW \approx 0.7R_x C_x$$

The width of the input pulse should be much less than the output pulse width.

The feedback biasing scheme in the astable circuit of Fig. 34 ensures that the astable will oscillate when power is initially applied to the circuit. For this circuit

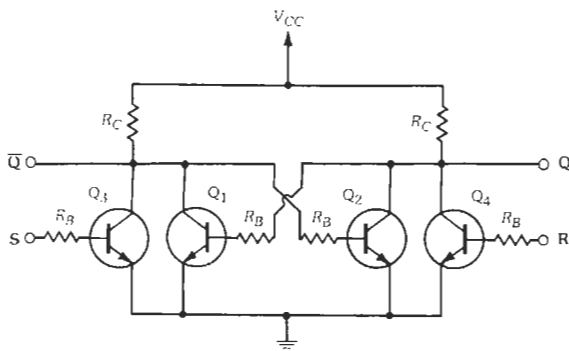
$$PW_1 \approx 0.7C_1(R_1 + R_3)$$

$$PW_2 \approx 0.7C_2(R_2 + R_4)$$

With $PW_1 = PW_2$ the output at each collector of the astable is a square wave with a frequency of $1/2PW_1$.

The timing of these circuits may be improved if the transistors are prevented from saturating by connecting a Schottky-barrier diode, or other fast-recovery low forward-voltage diode, from the base to the collector. The anode of the diode is connected to the base terminal.

Finally, each of these multivibrator circuits may be implemented with field-effect transistors, especially complementary MOS(CMOS).



(A) Circuit.

S	R	Q	\bar{Q}
0	0	No change	
0	1	0	1
1	0	1	0
1	1	0*	0*

*Not allowed

(B) Truth table.

Fig. 32. Bistable multivibrator.

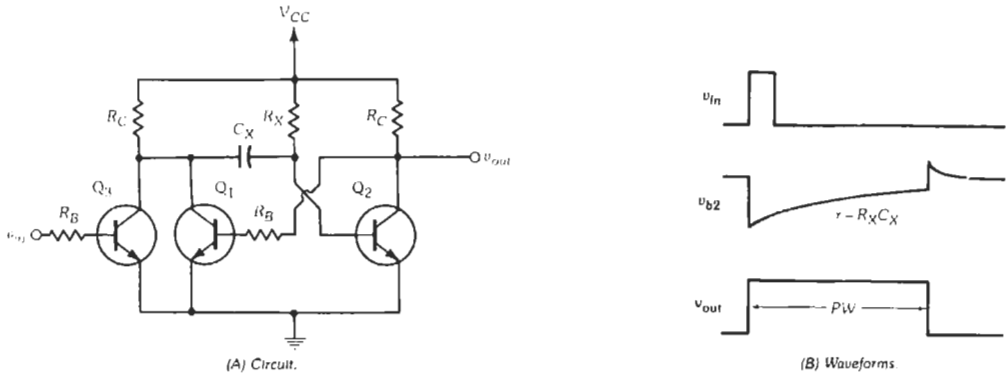


Fig. 33. Monostable multivibrator.

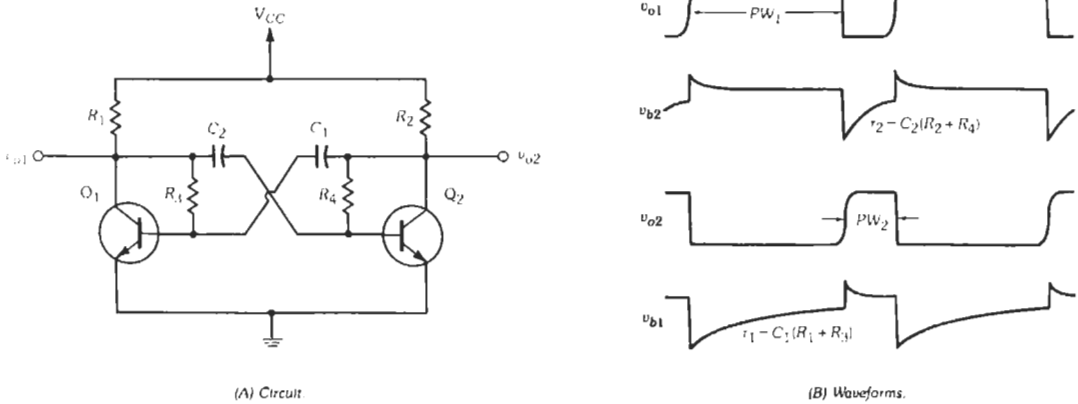


Fig. 34. Astable multivibrator.

REFERENCES

1. Carr, W. N., and Mize, J. P. *MOS/LSI Design and Application*, Chapter 1. New York: McGraw-Hill Book Co., 1972.
2. *Ibid.*, Chapter 9.
3. Chirlian, P. M. *Analysis and Design of Integrated Circuits*, Chapter 18. New York: Harper and Row Publishers, Inc., 1981.
4. Getreu, I. *Modeling the Bipolar Transistor*. Beaverton, Ore.: Tektronix Inc., 1976.
5. Gray, P. E., and Searle, C. L. *Electronic Principles—Physics, Models, and Circuits*, Chapter 11. New York: John Wiley and Sons, Inc., 1969.
6. *Ibid.*, Chapter 13.
7. *Ibid.*, Chapter 18.
8. Gray, P. R., and Meyer, R. G. *Analysis and Design of Analog Integrated Circuits*, Chapter 3. New York: John Wiley and Sons, Inc., 1977.
9. *Ibid.*, Chapter 4.
10. Grinich, V. H., and Jackson, H. G. *Introduction to Integrated Circuits*, Chapter 8. New York: McGraw-Hill Book Co., 1975.
11. Hamilton, D. J., and Howard, W. G. *Basic Integrated Circuit Engineering*, Chapter 9. New York: McGraw-Hill Book Co., 1975.
12. Hodges, D. A., and Jackson, H. G. *Analysis and Design of Digital Integrated Circuits*, Chapter 8. New York: McGraw-Hill Book Co., 1983.
13. Holt, C. A. *Electronic Circuits—Digital and Analog*, Chapters 2 and 3. New York: John Wiley and Sons, Inc., 1978.
14. Millman, J. *Microelectronics*, Chapter 11. New York: McGraw-Hill Book Co., 1979.
15. *Ibid.*, Chapter 12.
16. Millman, J., and Taub, H. *Pulse Digital and Switching Waveforms*, Chapter 2. New York: McGraw-Hill Book Co., 1965.
17. Muller, R. S., and Kamins, T. I. *Device Electronics for Integrated Circuits*, Chapters 3 and 5. New York: John Wiley and Sons, Inc., 1977.
18. Pederson, D. O. *Electronic Circuits*, Chapter 12. New York: McGraw-Hill Book Co., 1965.
19. *Ibid.*, Chapter 14.
20. Valley, G. E., and Wallman, H. *Vacuum Tube Amplifiers*, Chapters 4 and 5. New York: McGraw-Hill Book Co., 1948.

20

Integrated Circuits

*Pallab K. Chatterjee, Satwinder D. S. Malhi, Michiel deWit,
Harold Hosack, Mark G. Harward, and Mehrdad M. Moslehi
Revised by Harold Hosack*

Introduction	20-3
Definitions and Terminology	20-3
IC Fabrication	20-11
IC Processing	20-11
Substrate-Cleaning Procedures	
Thin Films	
Lithography	
Etching	
Substrate Doping Procedures	
An IC Process Example—NMOS Process	
Yield Statistics	
IC Design	20-26
Modeling and Simulation	
Topological Realization of Circuits	
Packaging Integrated Circuits	20-28
Plastic Packages	
Ceramic Packages	
Ceramic Flat Package	
Leadless Ceramic Chip Carrier	
Flip-Chip Package	
Thermal Resistance	
Digital Integrated Circuits	20-30
Bipolar Logic Families	
MOS Logic Families	
Microprocessors and Minicomputers	

- Memory Integrated Circuits 20-39
 - Read-Only Memories (ROM)
 - Nonvolatile Read-Only Memories
 - Random-Access Memories (RAM)
 - New Developments

- Linear Integrated Circuits 20-42
 - Differential Amplifier
 - Operational Voltage Amplifier
 - Other Amplifiers
 - Voltage Reference
 - Voltage Regulators
 - Comparators
 - Special-Purpose Linear Integrated Circuits
 - Miscellaneous Linear ICs

- Trends in Integrated Circuits 20-45
 - Scaling and Miniaturization
 - Image-Sensing ICs
 - Speech-Synthesizer IC
 - System in a Package (SIP) and System on a Chip (SOC)
 - Micro-Electro-Mechanical Systems (MEMS)

- Application Specific Integrated Circuits 20-49
 - Gate Arrays
 - Future Trends in ASICs

- Single-Wafer Integrated Device Processing 20-56
 - Single-Wafer Semiconductor Manufacturing
 - Integrated Device Manufacturing
 - RTP for Integrated Device Processing
 - Single-Wafer Surface Cleaning
 - Single-Wafer Epitaxy
 - Single-Wafer Plasma Processing
 - SWP-Based Process and Device Integration

- Silicon-on-Insulator Technology and Devices 20-102
 - SOI Material Options
 - Device Structures Using SOI Materials
 - Devices on SOI
 - Nonfully Depleted Devices
 - Fully Depleted Devices
 - SOI Circuits
 - Three-Dimensional ICs
 - Conclusion

INTRODUCTION

Progress in the complexity and usefulness of integrated circuits (ICs) has been extremely rapid over the past four decades. In fact, the pace of introduction of new ICs and new IC applications has continued to increase over this time. Currently ICs are a pervasive element in essentially all parts of commercial, scientific, and military systems. They have allowed major advances in such areas as computers, entertainment electronics, video cameras, global positioning systems, and innumerable other applications. The impact of ICs on society in the past forty years has been comparable to the impact of the industrial revolution. The rapid rate of change experienced during the recent past is expected to continue well into the next millennium. This chapter describes some of the basic concepts of integrated circuits, integrated circuit fabrication, designs, and a few applications.

DEFINITIONS AND TERMINOLOGY

access time: 1. Address-to-read output access time—The time delay in random-access memories from the application of appropriate address signal levels to the presence of valid data signal levels at the output. The signal to enable the memory chip enable is present. 2. Chip-enable-to-read output access time—With an appropriate address signal level present, the time delay in random-access memories from the application of appropriate signal levels to the enable inputs to the presence of valid data signal levels at the output.

adder: Switching circuit that combines binary bits to generate the sum and carry of these bits. Takes the bits from two binary numbers to be added (addend and augend), plus the carry from the preceding less-significant bit, and generates the sum and carry.

APCVD: Atmospheric pressure chemical vapor deposition.

ARC: Antireflection coating used in resist processing.

ash: Isotropic dry etching of resist with O_2 .

avalanche breakdown: Reverse-voltage breakdown due to electrons gaining sufficient speed to dislodge valence electrons and thus create more current carriers.

background: Process of removal of material from back side of wafers by grinding.

ball bond: Connecting wire from package to IC chip.

barrier layer: A thin layer deposited on contacts before metallization that prevents metal spikes

from penetrating through thin junctions and causing unwanted shorts.

beta ratio: For all standard p-channel and n-channel MOS structures, the coefficient $(\mu\epsilon_{ox})/(2t_{ox})$ for the driver device is equal to that of the load device. In this equation, μ = majority carrier mobility, ϵ_{ox} = permittivity of the oxide, and t_{ox} = thickness of the oxide over the channel. Hence, the beta ratio simply becomes the ratio of $(W/L)_{driver}$ to $(W/L)_{load}$, where W = width of the channel and L = length of the channel in the direction of the current. Some engineers use the symbol K_R to denote this same ratio.

bias, forward: An external voltage applied in the conducting direction of a pn junction. This is accomplished by connecting the positive terminal of the source to the p-type region and the negative terminal to the n-type region.

bias, reverse: An external voltage applied in the nonconducting direction of a pn junction. The connections are opposite to those for forward bias.

binary coded decimal (BCD): A binary numbering system for coding decimal numbers in groups of four bits. The binary value of these four-bit groups ranges from 0000 to 1001 and codes the decimal digits 0 through 9. To count to 9 takes four bits; to count to 99 takes two groups of four bits; to count to 999 takes three groups of four bits; etc.

binary logic: Digital logic elements that operate with two distinct states. The two states are variously called true and false, high and low, on and off, or 1 and 0. In computers they are represented by two different voltage levels. The level that is more positive (or less negative) than the other is called the high level; the other is called the low level. If the true (1) level is the most positive voltage, such logic is referred to as “positive true” or “positive logic.”

bipolar technology: Semiconductor fabrication involving an active semiconductor device with two or more pn junctions. Conduction is by a flow of electrons (in n-type material) or positive holes (in p-type material) across the junction of the two materials.

bird beak: Edge feature produced during local oxidation.

bistable element: Another name for flip-flop. A circuit in which the output has two stable states (output levels 0 and 1) and can be caused to go to either of these states by input signals, but remains in that state permanently after the input signals are removed. This differentiates the bistable element from a gate, which also has two output states but

which requires the retention of the input signals to stay in a given state. The characteristic of two stable states also differentiates the bistable element from a monostable element, which keeps returning to a specific state, and an astable element, which keeps changing from one state to the other.

bit: Binary digit.

boat: Wafer holder generally made of quartz or polysilicon used during furnace processes.

BTS: Bias temperature stress.

buffer memory: A memory system, usually of small capacity compared to a main-frame memory, that provides a buffer function between two digital activities.

bus: Long signal-carrying line.

byte: Eight bits.

CAD/CAE/CAM: Computer-aided design/computer-aided engineering/computer-aided manufacture.

CDI: Collector diffusion isolation. A bipolar semiconductor fabrication process that uses an epitaxial layer for the base structure (eliminating a base-diffusion step) and combines the deep collector diffusion contact and isolation steps.

Cell: A collection of devices that is repeated in a circuit. Mostly used to describe unit memory structures.

character generation: A design technique for integrated circuits that utilizes a fixed program in a storage element as the means for generating a dot matrix to represent alphanumeric characters on a CRT display.

charge-coupled device (CCD): A type of semiconductor device in which the presence or absence of charge represents the information. Charge-coupled devices store minority-carrier charges in potential wells created at the silicon surface and transport these charge packets along the surface by moving the potential wells.

chip enable: The control signal that activates a complete storage element. No reading or writing of data occurs unless this signal activates the storage element.

class N: Measure of particulate contamination in clean rooms. N particles of size greater than $0.3 \mu\text{m}$ per cubic foot.

clock: A timing control signal required by storage elements and memory systems.

CMOS: See complementary MOS.

C_o : Gate capacitance corresponding to the oxide layer over the channel area; $C_o = A(\epsilon_{ox}/t_{ox})$. $A =$

channel area, ϵ_{ox} = permittivity of oxide, t_{ox} = thickness of oxide over channel.

column decoder: Circuits within the storage elements designed to route data from a column of storage cells.

complementary MOS: An MOS fabrication process that combines both p-channel and n-channel transistors on the same substrate. The n-channel transistor is usually the driver device, and the p-channel transistor is the load. Only one transistor of the pair is normally on, except during the switching mode.

complexity: The number of equivalent gates in a circuit.

Cu-low k: Description of advanced IC metallization that is composed of thin films of copper separated by a very low dielectric constant medium.

current, forward: The net current that flows across the forward-biased pn junction.

CVD: Chemical vapor deposition of thin films.

cycle time: 1. Read cycle time—The total time required between the application of address information for reading data number 1 and the application of address information for reading data number 2 with the memory in the read mode. 2. Write cycle time—Same as read cycle time, but for writing of data and memory in the write mode. 3. System cycle time—The time between memory cycle initiations. It usually is determined by the longest of the read or write access or read or write cycle times.

Czochralski: Crystal-pulling process from a melt using a single crystal seed.

Damascene processing: Metal application to ICs using a process where metal is deposited on a surface having grooves where metal is desired, and then the metal is polished from the wafer, leaving metal in the grooves.

decoder: A conversion circuit that accepts digital input information (in the memory case, binary address information) that appears as a small number of lines and selects and activates one of a large number of output lines.

deglaze: Removal of thin layer of SiO_2 by the use of weak HF.

depletion region: The region in a semiconductor in which the mobile-carrier charge density is insufficient to neutralize the net charge density of the fixed donor and acceptor ions.

descum: Anisotropic plasma removal of resist residues after development.

die-by-die alignment: Alignment mechanism in which each die is aligned separately.

diffusion current: The current produced when charges move by diffusion.

diffusion length: The average distance excess minority carriers diffuse between injection and recombination.

diode isolation: A method of producing a high electrical resistance between an integrated-circuit element and the substrate by surrounding the element with a reverse-biased pn junction. The method is also called junction isolation (JI).

diode, semiconductor: A two-electrode semiconductor device that conducts current more easily in one direction than in the other.

DIP: Dual in-line package.

direct mapping: A design technique, used in virtual memory design, that determines the way data are organized and transferred from primary memory to the buffer memory. Data are mapped directly as blocks.

dope: Introduction of different atoms (dopant) in a semiconductor structure, Q .

dose: Number of implanted or diffused ions per unit area.

dot matrix: A matrix of dots that is used to identify alphanumeric characters.

dRAM/sRAM: Dynamic (refreshed storage cells), static (V_{cc} only) RAM.

drift current: A current that is produced when the carriers move under the influence of an applied voltage, i.e., owing to a voltage gradient.

drive-in: Deeper thermal diffusion of dopants introduced during predeposition.

DSW: Direct step on wafer for step and repeat projection optical lithography.

dynamic storage elements: Storage elements that contain storage cells that must be refreshed at appropriate time intervals to prevent the loss of information content.

EAROM: Electrically alterable ROM. A fixed-program semiconductor storage element whose program can be altered by the application of external electrical or optical means.

E-beam: Electron beam (exposes resist to pattern masks or wafers).

E_C : Conduction-band energy level.

ECL circuits: Bipolar emitter-coupled logic; also called current-mode logic circuits.

ECL storage cell: A type of bipolar storage cell that is like the multiple-emitter cell but is used for ECL interface.

E_F : Fermi energy level.

E_g : Band-gap energy = $E_C - E_V$.

E_i : Intrinsic energy level; assumed to be at the center of the band gap.

electromigration: Metal line failure due to mass transport of metal under large dc electric field.

electron-beam lithography: Use of focused electron beam to delineate patterns in resist. Commonly used for mask making.

epitaxial growth: The deposition of a monocrystalline layer of material onto a substrate material such that the layer thus formed has the same crystal orientation as the substrate.

EPROM, EEPROM: Electrically programmable, electrically erasable and programmable ROM.

etch filament: Thin residue left along steps during anisotropic dry etching.

etch selectivity: Ability of an etchant to remove layers of certain material with minimal attack on others.

E_v : Valence-band energy level.

excess minority carriers: The number of minority carriers that exceeds the normal equilibrium number in a semiconductor.

factory-programmed ROM: A fixed-program semiconductor storage element that has been programmed at the factory with a unique bit pattern.

failure rate: The number of components that fail in a given amount of time, usually expressed as a percentage per 1000 hours.

fall time: A measure of the time required for the output voltage of a circuit to change from a high-voltage level to a low-voltage level once a level change has started.

fan-out: The number of loads connected to the output of a logic stage. (A load normally consists of the input impedance of a logic circuit.)

fiducial: Markers on reticles for aid in alignment in projection printers.

filming: Lifting of resist caused by lack of adhesion.

fixed-program storage (read-only memory or read-only store): A special application of random-access storage in which storage is fixed after programming.

flip-flop: A storage element—a circuit having two stable states and the capability of changing from

- one state to another with the application of a control signal and remaining in that state after removal of signals. (See bistable element.)
- flip-flop, D*: The letter D stands for delay. A flip-flop whose output is a function of the input that appeared one pulse earlier; for example, if a 1 appeared at the input, the output after the next clock pulse would be a 1.
- flip-flop, JK*: A flip-flop having two inputs, designated J and K. At the application of a clock pulse, a 1 on the J input and a 0 on the K input will set the flip-flop to the 1 state; a 1 on the K input and a 0 on the J input will reset it to the 0 state; and 1 simultaneously on both inputs will cause it to change state regardless of the previous state. J = 0 and K = 0 will prevent change.
- flip-flop, RS*: A flip-flop consisting of two cross-coupled NAND gates and having two inputs designated R and S. A 1 on the S input and 0 on the R input will reset (clear) the flip-flop to the 0 state; 1 on the R input and 0 on the S input will set it to the 1 state. It is assumed that 0 will never appear simultaneously at both inputs. If both inputs have 1, the flip-flop will stay as it was. A 1 is considered nonactivating. A similar circuit can be formed with NOR gates.
- flip-flop, RST*: A flip-flop having three inputs designated R, S, and T. This unit works like the RS flip-flop, except that the T input is used to cause the flip-flop to change states.
- flip-flop, T*: A flip-flop having only one input. A pulse appearing on the input will cause the flip-flop to change states. Used in ripple counters.
- float zone*: Crystal growth process by repetitive melting and slow solidification.
- fusible metallization*: Utilization of a fusible metal, such as Nichrome, in the fabrication of semiconductor memory circuits. The metal interconnection is severed by a high current to program the storage element.
- gate, AND*: A logic circuit in which all inputs must have 1-level signals at the input to produce a 1-level output (assuming positive logic).
- gate, NAND*: A logic circuit in which all inputs must have 1-level signals at the input to produce a 0-level output (assuming positive logic).
- gate, NOR*: A logic circuit in which any one input or more than one input having a 1-level signal will produce a 0-level output (assuming positive logic).
- gate, OR*: A logic circuit in which any one input or more than one input having a 1-level signal will produce a 1-level output (assuming positive logic).
- gate, XOR*: A logic circuit in which any odd number of inputs having a 1-level signal will produce a 1-level output (assuming positive logic).
- glass spray*: Quartz particles generated by friction of boat with furnace tube during loading and unloading.
- global alignment*: Alignment mechanism in which only two geometries per wafer are aligned automatically, assuring requisite alignment of others.
- g_m : Transconductance = $\delta I_D / \delta V_{GS} / V_D$.
- GOI*: Gate oxide integrity against applied voltage.
- hillocks*: Protrusions in deposited thin films.
- Hilton haze*: Hazy look on thin films caused by departure from stoichiometry, with appearance akin to Hilton Hotel symbol.
- hold time*: The time that address information and data information must be maintained after write-enable to guarantee successful writing of data in the memory.
- hydrophilic*: Surface readily wetted by water.
- hydrophobic*: Surface not wetted by water.
- IGFET*: Insulated-gate field-effect transistor.
- input loading factor (ILF)*: The load that an input line presents to a driver. It may be expressed as a current, voltage, or impedance.
- insulator*: A material in which the outer (valence) electrons are tightly bound to the atom and are not free to move. No current can flow when a voltage less than breakdown is applied across the material.
- integrated circuit*: "The physical realization of a number of electrical elements inseparably associated on or within a continuous body of semiconductor material to perform the functions of a circuit." (EIA definition)
- intrinsic concentration*: The number of free electrons (or holes) per cubic centimeter in an undoped semiconductor at thermal equilibrium.
- ion-beam lithography*: Use of focused ion beams for pattern delineation.
- ion implantation*: An MOS semiconductor fabrication process often used to adjust threshold voltage values by implantation of dopant ions in the gate region, after source and drain formation. The implanted doping level is controlled by the ion accelerator beam current and implant time.
- ion milling*: Micromachining of thin films using mechanical motion of energetic ions.

isoplanar: A bipolar semiconductor fabrication process in which the p-diffused isolation regions are replaced by selectively grown oxide isolation.

junction avalanche: Utilization of an avalanche junction in the fabrication of fixed-program (read-only) semiconductor circuits. The avalanche junction forms a connection by high voltage and current to program the storage element.

junction barrier: The opposition to the diffusion of majority carriers across a pn junction caused by the charge of fixed donor and acceptor ions.

junction capacitor: A capacitor utilizing the capacitance of a reverse-biased pn junction.

junction transistor: An active semiconductor device with a base electrode and two or more junction electrodes.

KOOI effect: Inhibition of subsequent oxidation of silicon underlying Si_3N_4 caused by NH_3 formed during wet local oxidation.

L: Effective channel length (in direction of current) in a MOSFET.

large-scale integration (LSI): The simultaneous realization of large-area chips and optimum component packing density, resulting in cost reduction by maximizing the number of system connections done at the chip level. Circuit complexity above 100 gates.

laser anneal: Use of high-energy laser beam for local melting and recrystallization of semiconductors.

lifetime: The average time interval between the introduction and recombination of minority carriers.

loading factors: Specifically used here for memory systems. A numerical measure of the load that must be supplied to drive lines in a memory system.

load line: A line drawn on the family of collector characteristic curves of a transistor showing how the transistor collector voltage changes as the current through the transistor and load resistance changes.

LOCOS: Local oxidation of silicon. This refers to the use of silicon nitride to protect against oxidation.

logic swing: The voltage difference between the two logic levels, 1 and 0.

LPCVD: Low pressure chemical vapor deposition.

mainframe memory: The main memory of the digital system.

mechanically programmable semiconductor ROM: A fixed-program (read-only) semiconductor storage element that can be programmed by breaking interconnection by mechanical means (wiping metal away).

medium-scale integration (MSI): The realization of circuit complexities between 12 and 100 equivalent gates.

MEMS (micro-electro-mechanical systems): Integrated mechanical elements, sensors, actuators, and electronics that are fabricated using micro-fabrication technology similar to that used for integrated circuits.

MESFET: Metal gate Schottky field-effect transistor.

metallization: A thin-film pattern of conductive material (usually aluminum) deposited on a substrate to interconnect electronic components or to provide conductive contacts to which interconnecting wires may be bonded.

MINIMOS: MOSFET simulation program (available from University of Vienna).

mobile charge: Alkali ions that move through oxide under the influence of applied electric field.

mobility: The average velocity attained by a charge carrier under the influence of a unit electric field.

monolithic integrated circuit: An electronic circuit that has been fabricated as an inseparable assembly of circuit elements in a single structure that cannot be divided without permanently destroying its intended electronic function.

MOS capacitor: A capacitor formed by depositing a silicon-oxide dielectric layer and then a metal top electrode on the surface of a semiconductor region that forms the bottom electrode. The use of silicon as one of the capacitor plates makes the capacitance a function of applied voltage.

MOSFET (metal-oxide-semiconductor field-effect transistor): An active semiconductor device in which a conducting channel is induced in the region between two electrodes by a voltage applied to an insulated electrode on the surface of the region.

MOS, MIS: Metal-oxide-semiconductor, metal-insulator-semiconductor.

MTBF: Mean time between failures of an electronic system.

multilevel oxide: Insulating layer between two conductive layers.

multiple-emitter cell: A type of bipolar storage cell that uses a multiple-emitter transistor control for coupling to the bit lines.

N_A , N_D , N : Doping levels. (A = acceptor, D = donor.)

negative logic: Logic in which the more negative voltage represents the 1 state and the less negative voltage represents the 0 state. (See binary logic.)

nibble: Four binary digits.

NMOS, *PMOS*: Designations for n-channel, p-channel MOSFET.

noise immunity: A measure of the insensitivity of a logic circuit to triggering or reacting to spurious or undesirable electrical signals or noise, largely determined by the signal swing of the logic. Noise can occur in either of two directions, positive or negative.

nondestructive readout: Semiconductor memory designed so that readout does not affect the content stored. It is not necessary to perform a write after every read operation.

NOT: A Boolean logic operation indicating negation. Actually an inverter. If input is 1, output is *NOT* 1 = 0; if input is 0, output is *NOT* 0 = 1. Graphically represented by a bar over a Boolean symbol: \bar{A} . \bar{A} means "when A is not 1."

n-type semiconductor: A semiconductor in which electric conduction is due to the presence of more free electrons than holes.

NVRAM: Nonvolatile RAM.

OC curve: Operating characteristic curve. For a particular sampling plan, the graph of the probability of acceptance for all values of percent defective. Completely describes the risks involved in using the plan.

ODE: Orientation dependent etch.

ohmic contact: A resistive contact area that permits aluminum to be used as interconnecting metal from one high-resistivity n-type region to another.

OR: Oxide removal (for contact openings).

oxidation: A process that converts the surface of a silicon wafer to silicon dioxide. This is accomplished by subjecting the wafer to an oxygen or steam atmosphere at very high temperatures.

P: Heavy p-type diffusion.

parallel: The technique for handling a binary data word that has more than one bit. All bits are acted upon simultaneously. It is like the line of a football team; upon a signal, all linemen act.

parallel operation: The organization of data manipulation within computer circuitry wherein all the digits of a word are transmitted simultaneously

on separate lines in order to speed up operation, as opposed to serial operation.

parasitics: Stray components associated with the desired components diffused into an integrated circuit. Such parasitics may consist of capacitances, resistances, diodes, or transistors effectively in series or in shunt with the diffused components. They tend to limit the performance of the desired components in a circuit unless compensated for in device and circuit design.

passivation: Protection against penetration by impurity atoms. A silicon surface is passivated by covering it with a thin layer of variously doped films. This layer cannot easily be penetrated by impurities at normal processing and operating temperatures and, therefore, provides the necessary protection.

passive elements: Electronic components, such as resistors and capacitors, that simply introduce resistance or reactance into an electrical circuit but cannot change the waveform of an applied sine wave.

photomasking: A semiconductor-fabrication process in which a photographic negative is used to delineate selective chemical change to portions of the semiconductor surface.

pinholes: Weak spots in thin films due to defects or contamination.

Pirahna: A cleaning agent consisting of H_2O_2 and H_2SO_4 .

PLA: Programmable logic array. An integrated circuit that employs ROM matrices to combine sum and product terms of logic networks.

planar transistor: A diffused-junction transistor in which the emitter, base, and collector regions all come to the same plane surface, with the junctions between the regions protected at the surface by a layer of material such as silicon oxide.

plasma: Ionized gas in which concentrations of positive and negative charge carriers are almost equal.

plasma etch: Use of a plasma of reactive ions to etch thin layers.

plug bar: Assortment of test structures added on IC mask for model extraction and diagnostic purposes.

pn junction: The region of transition between p-type and n-type semiconductor materials.

poly: Polycrystalline silicon.

potential barrier: The difference in potential across a pn junction.

precharge time (reset time): The timing pulse width within a memory cycle that is used for charging node capacitances to particular starting-point voltage levels.

PREDEP: Predeposition of dopants.

process flow: The detailed step-by-step sequence of a fabrication schedule.

propagation delay: The time required for a change in logic level to be transmitted through an element or a chain of elements.

PSG: Phosphosilicate glass.

p-type semiconductor: A semiconductor in which electric conduction is due to the presence of more holes than free electrons.

purge: Extended exposure to only a desired gas or a mixture of gases.

ϕ_F : Fermi potential; the amount the Fermi level is displaced from the intrinsic level or the center of the gap (as measured in the bulk). Units are volts.

ϕ_s : Surface potential; the amount the intrinsic Fermi level, at the surface, has been shifted with respect to the bulk Fermi level.

q : Electronic charge; 1.6×10^{-19} coulomb.

random-access memory (RAM): A memory from which information can be obtained at the output with approximately the same time delay by choosing an address randomly and without first searching through a vast amount of irrelevant data.

range: Depth of peak of implanted ion distribution.

rapid thermal processing (RTP): Process steps implemented by very sharp rises in temperature, short process times, and very steep cool down cycles. The complete process sequence is completed in a few tens of seconds.

ratio inverter: An inverter whose logic swing is determined by the beta ratio (β_R) of the load and driver devices.

ratioless inverter: An inverter whose logic swing closely approximates the power-supply voltage. ("Ratioless" because the load and driver device do not conduct simultaneously.)

ratioless-type shift register: Current does not flow through the inverter when the clock and data inputs are simultaneously at the logic 1 level.

ratio-type shift register: Current flows through the inverter when the clock and data inputs are simultaneously at logic 1.

reactive ion etching: Etching of films using plasma and mechanical motion of reactive ions.

reflow: High-temperature treatment of phosphosilicate glass that causes its mass flow.

refresh: Method that restores capacitance charge that deteriorates because of leakage.

register: Temporary storage for digital data.

reox: Oxide removal from heavily implanted region and its thermal oxidation for sake of passivation.

repair frequency: The rate at which an electronic system must be repaired, i.e., once/day, once/week, once/four weeks, etc.

reset: Also called "clear." Similar to set except it is the input through which the Q output can be made to go to 0.

resist: Photosensitive organic or inorganic resins.

resistivity: The (volume) resistivity ρ is the electric field E required to produce a unit current density J . That is, $\rho = E/J$.

rinse cascade: A series of containers with running deionized water used for wafer rinsing after chemical treatment.

ROM: Read-only memory.

row decoder: Circuits within the storage elements designed to route data from a row of storage cells.

SAG (self-aligned gate): An MOS fabrication process using a self-aligning gate formed from deposited silicon.

SAMPLE: Lithography and etch simulation program available from the University of Berkeley.

SATO (self-aligned thick oxide) process: An MOS fabrication process using nitride as the self-aligning gate material. Thus, the gate metallization does not have to withstand diffusion temperature.

Schmitt trigger: An input circuit with hysteresis. There is a higher threshold for positive-going inputs than for negative-going inputs.

Schottky barrier diode: Metal-semiconductor barrier diode.

segregation: Preferential accumulation of dopants in either silicon or oxide during thermal oxidation.

self-refresh: A circuit-design technique that incorporates the refresh method in the storage-element circuitry so that external refresh circuitry is not required.

SEM: Scanning electron microscope.

semiconductor: A material with conductivity roughly midway between that of conductors and insulators, and in which the conductivity

- increases with temperature over a certain temperature range.
- sense amplifier*: A sensitive amplifier accepting linear voltage or current signals and producing logic-level outputs.
- sequentially accessed memory*: A memory from which information is received at the output in varying time delays from a reference point depending on the position of the data in a time sequence.
- serial accumulator*: A register that receives data bits in sequence and temporarily holds the data for future use.
- serial operation*: The organization of data manipulation within computer circuitry wherein the digits of a word are transmitted one at a time along a single line. The serial mode of operation is slower than parallel operation but utilizes less complex circuitry.
- sheet resistance*: The resistance per square of a sheet of material.
- Silicon On Insulator (SOI)*: A wafer substrate composed of a thin layer of silicon (typically one micron or less) over a layer of insulating material (typically silicon dioxide, one half micron or less) over a silicon substrate. ICs are fabricated in the thin top layer of silicon.
- SIM*: Secondary ion microscope.
- single crystal*: A piece of material in which all the basic groups of atoms have the same crystallographic orientation.
- sinter*: Annealing treatment after metal etch to promote ohmic contacts.
- slice*: A single wafer cut from a silicon ingot, forming a thin substrate on which all active and passive elements for multiple integrated circuits have been fabricated by semiconductor epitaxial growth, diffusion, passivation, masking, photoresist, and metallization technologies. A completed slice generally contains hundreds of individual circuits, called chips or bars.
- SOI*: Silicon on insulator.
- solid-state diffusion*: The introduction of atoms of an impurity element into the surface regions of a solid semiconductor wafer.
- SOS*: Silicon on sapphire.
- SPICE 2*: Version of circuit simulation computer program (University of Berkeley).
- state*: The condition of an input or output of a circuit as to whether it is logic 1 or logic 0. The state of a circuit (gate or flip-flop) refers to its output. A flip-flop is said to be in the 1 state when its Q output is 1. A gate is in the 1 state when its output is 1.
- static storage elements*: Storage elements that contain storage cells that retain their information as long as power is applied unless the information is altered by external excitation.
- step coverage*: Ability of thin films to maintain thickness when going over feature steps.
- stored-charge programmable semiconductor ROM*: Utilizing a charge stored on a floating gate or a dielectric as a means of programming an MOS fixed-program (read-only) semiconductor storage element.
- straggle*: Standard deviation around the range of implanted ions.
- substrate*: The physical material on which an integrated circuit is fabricated. Its primary function is mechanical support, but it may serve some electrical function also.
- SUPREM*: Process simulation program available from Stanford University.
- System in a Package (SIP)*: A single package that hosts several different computation modules (i.e., microprocessor, controllers, DSPs, I/O unit, RF modules ...) together with memory to perform a complete system function.
- System on a Chip (SOC)*: A single chip that hosts several different computation modules (i.e., microprocessor, controllers, DSPs, I/O unit, RF modules ...) together with memory to perform a complete system function.
- TEM*: Transmission electron microscope.
- thermal compression bond*: A commonly used method for attaching a very fine wire to a point (usually a bonding pad) on an integrated-circuit chip.
- thermal generation*: The creation of a hole and a free electron by freeing a bound electron through the addition of heat energy.
- three-state output*: An output condition that has a low on impedance for driving to the high state, a low on impedance for driving to the low state, and an intermediate high-impedance off state.
- throughput*: Rate of material processed per hour.
- TIRAM*: Taper isolated random-access memory, an advanced dynamic RAM cell.
- toggle*: To switch between two states, as in a flip-flop.
- t_{ox} : Oxide thickness.

trigger: A timing pulse used to initiate the transmission of logic signals through the appropriate circuit signal paths.

truth table: A chart that tabulates and summarizes all the combinations of possible states of the inputs and outputs of a circuit. It tabulates what will happen at the output for a given input combination.

TTL: Transistor-transistor logic multiple emitter bipolar semiconductor circuit. Bipolar semiconductor transistor-transistor coupled logic circuits.

two-level main memory: A memory system featuring two separate memories. One memory is a buffer store, or cache memory. The other memory is the primary storage. Coupled together, these form a virtual memory with the capacity of the primary memory and the speed of the buffer memory.

two-level metallization: A semiconductor fabrication process in which there are two levels of interconnecting metal on the surface of the integrated circuit.

ULSI: Ultralarge-scale integration; circuit complexity above 1 million transistors.

vacancy: Unoccupied position in bond structure receptive to substitutional occupancy.

V_{BG} : Back-gate bias.

V_{CC}, V_{DD}, V_{SS} : High, intermediate, low supply-voltage levels.

vector generation: A design technique for generating the coordinates for positioning the beam in a CRT display that uses only the changes in coordinates of the beam rather than the absolute coordinates.

VFB: Flatband voltage.

V_{GG} : Gate supply voltage.

V_{GS} : Gate-to-source voltage.

VHPIIC: Very-high-performance integrated circuits (British program).

VHSIC: Very-high-speed integrated circuits.

VLSI: Very-large-scale integration; circuit complexity above 100 000 transistors.

V_{TD} : Threshold voltage of the driver device.

V_{TL} : Threshold voltage of the load device.

$V_T (V_{BG})$: Threshold voltage as a function of back-gate bias.

W: Effective channel width (perpendicular to current) in a MOSFET.

wafer flat: A flat portion of an otherwise circular wafer for orientation identification.

wafer stepper: Resist exposure system in which each die is focused and exposed separately.

word: Sixteen bits.

write enable: Also called read/write or R/W. The control signal to a storage element or a memory that activates the write mode or operation. When the device is not in the write mode, the read mode is active.

write time: 1. Address-to-write time—The time delay in random-access memories from the application of appropriate address signal levels until the write-mode control signal of an appropriate level is applied. 2. Chip-enable-to-write time—The time delay, with appropriate address signal levels present, from the application of appropriate signal levels to the enable inputs until the write-mode control signal of an appropriate level is applied. 3. The time that the appropriate level must be maintained on the write-enable line and that data must be present to guarantee successful writing of data in the memory.

yield: The percentage of acceptable circuits (chips) produced by a particular process, process step, inspection, or test.

zener: Breakdown in diodes due to tunneling.

IC FABRICATION

The process of producing ICs has three main components—design, processing, and assembly and test. Design is the step necessary to translate the desired electrical functions of the IC to the description of the physical structure required to produce those functions. Processing is the step needed to transform the description provided by design into a physical structure that is capable of providing the desired electrical functionality. Assembly is the step necessary to place the physical structure in a package that will protect it from environmental stresses and provide wires to connect it to the outside world, and test is the step that assures that the finished IC meets all of the desired functions. The subjects of IC design, processing, and assembly and test are extensive. In this chapter we will only describe the fundamentals of each of these areas. The reader is directed to the extensive literature on each of these areas for further details.*

IC PROCESSING

IC Processing is the step needed to transform the description provided by design into a physical structure that is capable of providing the desired electrical

* For example, References 1, 2, and 3.

functions. The discipline of IC processing is one that has been continuously and rapidly evolving, especially in the past decade. This section discusses processes that are currently, or are about to be, in widespread production use. New processes continue to be invented and developed, and to keep abreast of current understanding the reader is well advised to refer to pertinent technical journals. Since silicon-based circuits are the mainstream IC technology, the following discussion is set in the context of silicon IC processing. However, the concepts presented apply equally well to processing of other materials, such as gallium arsenide.

The silicon wafers used in IC processing are thin slices of a large silicon crystal. The wafers are round like a plate and have a small "flat" ground on them to allow indexing in IC fabrication equipment. These wafers are polished to a very fine mirrorlike finish to allow very small geometries to be printed on them by lithographic reproduction. Current production wafers are typically 200 mm in diameter. The mainstay of production is expected to be 300-mm-diameter wafers by the year 2005. The silicon wafers are cut to a typical thickness of a few hundred microns. The ICs usually occupy only the top few microns of the silicon wafer. The remaining thickness is necessary for mechanical stability during the wafer-processing steps. After the wafer has completed processing it is usually ground to a thickness of only a few tens of microns before the individual ICs are separated and assembled. The wafers are interchangeably called the "substrate" for IC fabrication.

Typical ICs are rectangular and are one centimeter or less on a side. They are produced by a process that is typically several hundred individual steps of repeated application of specific patterns and intervening process steps to produce structures in the silicon wafer in the form of the patterns. Since the individual IC is less than one square centimeter in area, and the wafers are several hundred square centimeters in area, many ICs are produced on a single wafer. This "mass production" is the key ingredient to being able to continuously decrease the cost of functionality available

from ICs. After the wafer is completed with processing, it appears as a disc with a clearly repetitive pattern of small rectangles, each of which is an individual IC.

The following sections describe some of the key processes used in IC processing and provide a simple example.

Substrate-Cleaning Procedures

The performance characteristics of semiconductor devices are sensitive to cleanliness in processing. The removal of unwanted impurities from the wafer surface is important because such impurities may diffuse into the semiconductor during high-temperature processing, altering its bulk and surface properties. Many of the spurious impurities may be either donor or acceptor dopants, directly affecting the device characteristics. Other impurities may cause surface or bulk defects such as traps, stacking faults, or dislocations. Surface contaminants such as oil, grease, or other organic matter may lead to poor film adhesion. This requires a careful chemical cleaning of wafers at the initiation of the process and appropriate cleaning at various steps during processing.

The initial cleanup generally starts with wafer scrubbing to remove loose particulate contaminants. This is followed by treatment with organic chemicals to get rid of any possible organic impurities such as hydrocarbons and greases that may be remnants from the wafer-grinding process. Organic solvents such as methanol and ethanol are suitable for this purpose. The final cleanup consists of a variety of inorganic chemicals to remove heavy metals, etc. Most of these chemical mixtures are strong oxidants,* forming a thin oxide at the wafer surface. This oxide is then stripped so that impurities absorbed therein are removed. Table 1 lists some cleaning reagents commonly used in IC processing.

* Reference 4.

TABLE 1. CLEANING REAGENTS USED IN IC PROCESSING

Cleaning Agent	Boiling Point (°C)	Purpose
Ethanol	78.3	
Trichloroethylene	87.2	Remove organic contaminants
Acetone	56.2	
p-Xylene	138.4	
5 H ₂ O: 1 H ₂ O ₂ : 1 NH ₄ OH	—	Remove inorganic contaminants (heavy metals)
6 H ₂ O: 1 H ₂ O ₂ : 1 HCl	—	
1 H ₂ SO ₄ : 1 H ₂ O ₂	—	
1 HCl: 1 HNO ₃	—	

Thin Films

All IC processes involve the use of a number of thin films—insulators, semiconductors, and conductors. These films may form an essential part of a particular structure or just be of use in the implementation of a process step. The following means of achieving thin films are of widespread use.

Thermal Growth—This technique involves heating the substrate in a furnace at precisely controlled temperature and gas ambient. Generally, a high temperature, in the range from 800 °C to 1200 °C for silicon processing, is used to promote chemical reaction between the ambient gases and the substrate. The prominent example is the growth of SiO₂ on silicon in O₂ ambient. It is also possible to grow Si₃N₄ in an N₂ or, preferably, NH₃ atmosphere. It is possible to exercise extremely good control on absolute thickness as well as thickness uniformity across the wafer. Figs. 1 and 2 show the relationship between oxidation time and resultant oxide thickness in dry O₂ and steam, respectively.*

Chemical Vapor Deposition—The gas-phase reduction of highly reactive chemicals under low pressure forms a convenient way of obtaining very uniform thin films. This method is widely used for deposition of oxides, nitrides, and polycrystalline semiconductors. A conformal deposition around sharp edges is an important attribute of this technique. Table 2 lists constituent gases and suitable temperatures for the deposition of some common films.‡

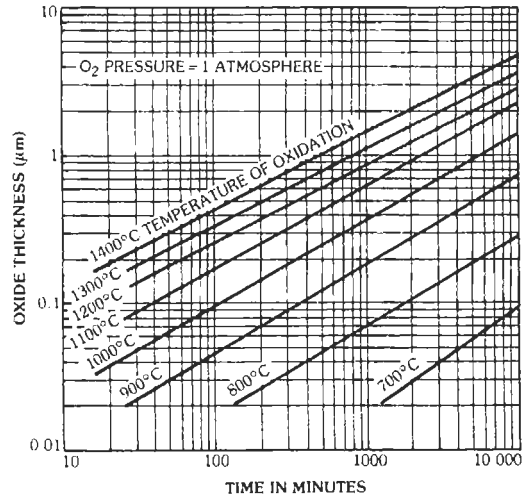


Fig. 1. Silicon dioxide growth rate in dry oxygen. (From A. M. Smith, "Experimental Measurements,"† in Burger and Donovan, Eds., *Fundamentals of Silicon Integrated Device Technology*, Vol. 1. Englewood Cliffs, NJ: Prentice-Hall, Inc.; Fig. 6-42.)

Plasma Deposition—The production of thin films by electric discharge or plasma depends on the capability of a nonequilibrium but sustained plasma to generate chemically reactive species at low temperature. The plasma is typically sustained at 0.1 to several torr and exhibits free electron temperatures of tens of thousands of kelvins, while the temperature of the translational or rotational modes of atoms, radicals, or molecules is only hundreds of kelvins. Thus, deposition can be made at very low temperatures. The deposition scheme is in use for oxides, nitrides, carbides, and amorphous semiconductors. Table 3 lists gases used.§

Evaporation—Evaporation involves film deposition by vaporizing the material on heating it past its melting point under vacuum to produce enough vapor

* Reference 5.

† References: B. I. Boltaks and H. Shih-yin, "Diffusion, Solubility and the Effect of Silver Impurities on Electrical Properties of Silicon," *Soviet Phys. Solid State* 2, May 1961, p. 2303. F. A. Trumbore, "Solid Solubilities of Impurity Elements in Germanium and Silicon," *Bell System Tech. J.* 39, January 1960, pp. 205-233.

‡ Reference 6.

§ Reference 7.

TABLE 2. GASES AND TEMPERATURES FOR DEPOSITION OF FILMS

Film	Gases	Temperature (°C)
Polysilicon (Si)	SiH ₄ , N ₂	650-700
	SiH ₄ , H ₂	850-950
Silicon nitride (Si ₃ N ₄)	SiH ₄ , NH ₃ , N ₂	750-800
	SiH ₂ Cl ₂ , NH ₃ , N ₂	750-900
Silicon dioxide (SiO ₂)	SiH ₄ , N ₂ O, N ₂	750-850
	SiH ₄ , CO ₂ , H ₂	950-1000

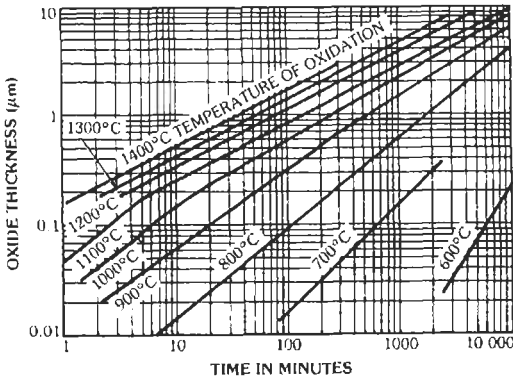


Fig. 2. Silicon dioxide growth rate in steam. (From A. M. Smith, "Experimental Measurements,"* in Burger and Donovan, Eds., *Fundamentals of Silicon Integrated Device Technology*, Vol. I. Englewood Cliffs, NJ: Prentice-Hall, Inc.; Fig. 6-42.)

TABLE 3. GASES USED FOR PLASMA DEPOSITIONS

Film	Gases
Amorphous silicon (Si)	SiH ₄ , Ar
Silicon dioxide (SiO ₂)	SiH ₄ , N ₂ O
Silicon nitride (Si ₃ N ₄)	SiH ₄ , NH ₃

pressure. Either resistive heating or E-beam heating is used to bring about melting. This technique is used mainly to deposit metals such as aluminum. A list of temperatures and support materials used to evaporate various elements is in Table 4.†

Sputter Deposition—In this scheme, the material to be deposited is bombarded with positive inert ions with kinetic energy far exceeding the heat of sublimation of the target material. This results in dislodging of target atoms and their ejection into the gas phase succeeded by deposition of the substrate, which may or may not be biased negatively. In widespread use are sputter deposition of metals and metal silicides.

Spin-On Deposition—The material to be deposited is mixed with a suitable solvent and spun coated on the substrate. The resulting thickness is a function of spin speed and viscosity of the solution used. Subsequently, an oven-bake drives out the solvent, leaving behind a stable layer. This procedure is used for deposition of layers of resists and doped or undoped oxides.

Rapid Thermal Processing (RTP)—The drive to smaller and smaller geometries for the individual transistors in ICs has produced a need for very thin depos-

ited layers and very careful control of the location of impurities in the silicon wafers. This need has recently led to a new processing technique called Rapid Thermal Processing (RTP). Instead of the usual diffusion furnaces and other large equipment, where many wafers are processed at one time, this technique uses a bank of very high intensity lamps to rapidly heat a single wafer to temperatures of 500 °C to over 1000 °C in just a few tens of seconds. The wafers are held at the high temperature for only a short time and then cooled. This rapid heat treatment allows very thin layers to be grown with extreme precision and effective activation of implanted impurities with minimal diffusion.

Lithography

Lithography is the aspect of the IC fabrication process that deals with transferring onto a substrate the detailed features associated with individual components that collectively comprise an integrated circuit. The following lithographic steps must be implemented for each level of an IC fabrication schedule.

Mask Generation—A mask, or reticle, is a flat plate or a membrane that features a geometrical pattern with areas that are selectively transparent or opaque to a wavelength or a band of wavelengths used in a particular lithographic system. The pattern dimensions may be the same as the final size required in the circuit, or they may be larger by a factor *N* in an *N*:1 exposure system, where generally *N* = 1, 5, or 10. If *N* > 1, then the resolution requirements during the process of mask generation are relaxed, and less exacting equipment and techniques may be utilized to reduce costs.

The first step during mask making is the circuit layout. The goals of a good layout are to transform all designed components associated with a circuit into a geometrical layout that achieves the required packing density while keeping the parasitics small. The parasitics are those electrical circuit components that are not designed in but that inevitably originate owing to some features of a particular layout. Also, the performance of certain designs critically depends on well-matched components; therefore, such circuit elements must be laid out identically and close together.

First, the layout rules are established with the capabilities of a particular technology kept in mind. The layout may be accomplished by means of a variety of computer-aided design methodologies. All schemes utilize advanced pattern-generating capabilities. Various levels are denoted by differing colors. The geometries may be laid out explicitly or implicitly with a symbolic representation that is later converted into required geometries.‡ After the layout is completed, this information is fed into the mask-generation system. The mask generation in itself must utilize various lithographic steps, which are discussed in detail below.

* Boltaks and Shih-yin, *loc. cit.*

† Reference 8.

‡ Reference 9.

TABLE 4. TEMPERATURES AND SUPPORT MATERIALS USED TO EVAPORATE VARIOUS ELEMENTS*

Element and Predominant Vapor Species	Temp (°C)		Support Materials		Remarks
	mp	$p^* = 10^{-2}$ Torr	Wire, Foil	Crucible	
Aluminum (Al)	659	1220	W	C, BN TiB ₂ , BN	Wets all materials readily and tends to creep out of containers. Alloys with W and reacts with carbon. Nitride crucibles preferred.
Antimony (Sb ₄ , Sb ₂)	630	530	Mo, Ta, Ni	Oxides, BN, metals, C	Polyatomic vapor, $\alpha v = 0.2$. Requires temperatures above mp. Toxic.
Arsenic (As ₄ , As ₂)	820	300	—	Oxides, C	Polyatomic vapor, $\alpha v = 5 \times 10^{-5} - 5 \times 10^{-2}$. Sublimates but requires temperatures above 300 °C. Toxic.
Barium (Ba)	710	610	W, Mo, Ta, Ni, Fe	Metals	Wets refractory metals without alloying. Reacts with most oxides at elevated temperatures.
Beryllium (Be)	1283	1230	W, Mo, Ta	C, refractory oxides	Wets refractory metals. Toxic, particularly BeO dust.
Bismuth (Bi, Bi ₂)	271	670	W, Mo, Ta, Ni	Oxides, C, metals	Vapors are toxic.
Boron (B)	2100 ± 100	2000	—	C	Deposits from carbon supports are probably not pure boron.
Cadmium (Cd)	321	265	W, Mo, Ta, Fe, Ni	Oxides, metals	Film condensation requires high supersaturation. Sublimates. Wall deposits of Cd spoil vacuum system.
Calcium (Ca)	850	600	W	Al ₂ O ₃	
Carbon (C ₃ , C ₁ , C ₂)	3700	2600	—	—	Carbon-arc or electron-bombardment evaporation. $\alpha v < 1$
Chromium (Cr)	1900	1400	W, Ta	—	High evaporation rates without melting. Sublimation from radiation-heated Cr rods preferred. Cr electrodeposits are likely to release hydrogen.
Cobalt (Co)	1495	1520	W	Al ₂ O ₃ , BeO	Alloys with W; charge should not weigh more than 30% of filament to limit destruction. Small sublimation rates possible.
Copper (Cu)	1084	1260	W, Mo, Ta	Mo, C, Al ₂ O ₃	Practically no interaction with refractory materials. Mo preferred for crucibles because it can be machined and conducts heat well.
Gallium (Ga)	30	1130	—	BeO, Al ₂ O ₃	Alloys with refractory metals. The oxides are attacked above 1000 °C.
Germanium (Ge)	940	1400	W, Mo, Ta	W, C, Al ₂ O ₃	Wets refractory metals but low solubility in W. Purest films by electron-gun evaporation.
Gold (Au)	1063	1400	W, Mo	Mo, C	Reacts with Ta, wets W and Mo. Mo crucibles last for several evaporations.

(Continued on next page)

TABLE 4 (CONT). TEMPERATURES AND SUPPORT MATERIALS USED TO EVAPORATE VARIOUS ELEMENTS*

Element and Predominant Vapor Species	Temp (°C)		Support Materials		Remarks
	mp	$p^* = 10^{-2}$ Torr	Wire, Foil	Crucible	
Indium (In)	156	950	W, Mo	Mo, C	Mo boats preferred.
Iron (Fe)	1536	1480	W	BeO, Al ₂ O ₃ , ZrO ₂	Alloys with all refractory metals. Charges should not weigh more than 30% of W filament to limit destruction. Small sublimation rates possible.
Lead (Pb)	328	715	W, Mo, Ni, Fe	Metals	Does not wet refractory metals. Toxic.
Magnesium (Mg)	650	440	W, Mo, Ta, Ni	Fe, C	Sublimes.
Manganese (Mn)	1244	940	W, Mo, Ta	Al ₂ O ₃	Wets refractory metals.
Molybdenum (Mo)	2620	2530	—	—	Small rates by sublimation from Mo foils. Electron-gun evaporation preferred.
Nickel (Ni)	1450	1530	W, W foil lined with Al ₂ O ₃	Refractory oxides	Alloys with refractory metals; hence charge must be limited. Small rates by sublimation from Ni foil or wire. Electron-gun evaporation preferred.
Palladium (Pd)	1550	1460	W, W foil lined with Al ₂ O ₃	Al ₂ O ₃	Alloys with refractory metals. Small sublimation rates possible.
Platinum (Pt)	1770	2100	W	ThO ₂ , ZrO ₂	Alloys with refractory metals. Multistrand W wire offers short evaporation times. Electron-gun evaporation preferred.
Rhodium (Rh)	1966	2040	W	ThO ₂ , ZrO ₂	Small rates by sublimation from Rh foils. Electron-gun evaporation preferred.
Selenium (Se ₂ , Se _n ; $n = 1-8$)	217	240	Mo, Ta, stainless steel 304	Mo, Ta, C, Al ₂ O ₃	Wets all support materials. Wall deposits spoil vacuum system. Toxic. $\alpha v = 1$
Silicon (Si)	1410	1350	—	BeO, ZrO ₂ , ThO ₂ , C	Refractory oxide crucibles are attacked by molten Si, and films are contaminated by SiO. Small rates by sublimation from Si filaments. Electron-gun evaporation gives purest films.
Silver (Ag)	961	1030	Mo, Ta	Mo, C	Does not wet W. Mo crucibles are very durable sources.
Strontium (Sr)	770	540	W, Mo, Ta	Mo, Ta, C	Wets all refractory metals without alloying.
Tantalum (Ta)	3000	3060	—	—	Evaporation by resistance heating of touching Ta wires, or by drawing an arc between Ta rods. Electron-gun evaporation preferred.
Tellurium (Te ₂)	450	375	W, Mo, Ta	Mo, Ta, C, Al ₂ O ₃	Wets all refractory metals without alloying. Contaminates vacuum system. Toxic. $\alpha v = 0.4$

(Continued on next page)

TABLE 4 (CONT). TEMPERATURES AND SUPPORT MATERIALS USED TO EVAPORATE VARIOUS ELEMENTS*

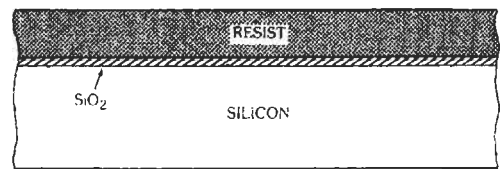
Element and Predominant Vapor Species	Temp (°C)		Support Materials		Remarks
	mp	$p^* = 10^{-2}$ Torr	Wire, Foil	Crucible	
Tin (Sn)	232	1250	W, Ta	C, Al ₂ O ₃	Wets and attacks Mo.
Titanium (Ti)	1700	1750	W, Ta	C, ThO ₂	Reacts with refractory metals. Small sublimation rates from resistance-heated rods or wires. Electron-gun evaporation preferred.
Tungsten (W)	3380	3230	—	—	Evaporation by resistance heating of touching W wires, or by drawing an arc between W rods. Electron-gun evaporation preferred.
Vanadium (V)	1920	1850	Mo, W	Mo	Wets Mo without alloying. Alloys slightly with W. Small sublimation rates possible.
Zinc (Zn)	420	345	W, Ta, Ni	Fe, Al ₂ O ₃ , C, Mo	High sublimation rates. Wets refractory metals without alloying. Wall deposits spoil vacuum system.
Zirconium (Zr)	1850	2400	W	—	Wets and slightly alloys with W. Electron-gun evaporation preferred.

* From Maissel and Glang, *Handbook of Thin Film Technology*. New York: McGraw-Hill Book Co., 1970; pp. 1-37 and 1-38.

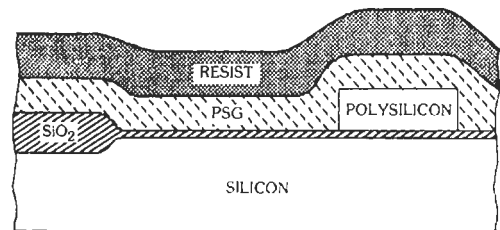
Resist Casting—The pattern on a mask is replicated on a substrate through the use of resists. Resists are organic or inorganic resins that are sensitive to a wavelength or a band of wavelengths used in a lithographic exposure system and, when exposed, undergo a chemical transformation so that selective removal can be accomplished. A resist is referred to as a positive resist if when exposed it is removed during development; it is referred to as a negative resist if when not exposed it is removable on development.

The resist application involves a thorough cleaning of the substrate to rid it of undesirable contaminants. After careful drying, the substrate is covered by a layer of resist, typically less than 3 μm thick, by spin-coating, spraying, or immersion.

The resist thickness uniformity is critical to obtaining high resolution in a lithographic system. With the advent of very-large-scale integration (VLSI), a tight line-width control has become of supreme importance. Often, the requirements of obtaining good line-width control, high resolution, and good step coverage are difficult to meet simultaneously. Good step coverage demands a thicker resist, whereas a thinner resist is necessary for good resolution. This is true for both positive and negative, as well as organic and inorganic, resists. The problem is shown schematically in Fig. 3. The resist thickness diminishes when the resist goes over steps. The effect becomes more pronounced in



(A) The resist is uniform when the wafer does not have any topology.



(B) Resist thins as it goes over steps in the presence of surface features.

Fig. 3. Schematic cross sections showing details of resist-coated wafers.

VLSI because of the larger ratio of line height to line width.

Regardless of resist type, a thin coating and a flat resist surface are needed for high resolution and good

line-width control. Multilevel resist systems are being investigated to fulfill this need.* In a typical double-level resist system, a thick organic resist is first spun. This resist conforms to the wafer surface and is planar at the top. Then an intermediate layer such as SiO_2 is deposited. Now another thin coating of resist is made. This resist is capable of achieving high resolution and good line-width control. This is shown in Fig. 4.

After the resist coating is complete, an oven bake, called prebake, is carried out to drive out the solvents, increase sensitivity, and improve resistance to mechanical handling. This finishes the procedure of resist casting.

Resist Exposure—After prebake, the resist is ready for exposure. A wide variety of exposure systems are currently in production use or under development (Fig. 5).

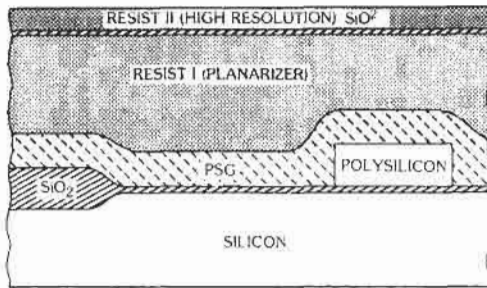


Fig. 4. Schematic cross section of a wafer coated with double-level resist system.

* Reference 10.

The contact printer refers to the scheme in which the mask is in intimate contact with the resist-covered wafer when exposure is made. Contact printers were the predominant lithographic exposure tool in production in the 1970s. They can be used reliably down to 3–4 μm geometries. In a projection printer, the mask or reticle is placed between the exposure source and the wafer to be exposed and does not make contact with the wafer. The mask image is reduced by a factor N and focused on the wafer. The scheme is suitable for high resolution. In a direct-write system, the need for a mask is completely eliminated. The resist is exposed by a finely focused particle beam that can be electrostatically scanned and blinked to form the required exposure pattern. This is inherently a slower procedure, and throughput therefore is adversely affected. Such systems, however, are uniquely suitable for mask generation.

A key feature of an exposure system is its capability to align a given level to an existing level on the wafer. This is done with the help of specially designed geometrical patterns (fiducials), which are placed on each level. When these patterns on the mask are made to overlap the corresponding structure on the wafer within reasonable bounds, a good alignment is said to have been made. The detailed features of alignment mechanisms differ, and the superiority of a particular system is partly determined by its capability to achieve alignment within tighter bounds.

Resist Development—After exposure, the resist is ready for development. Development refers to selective removal of resist as determined by the mask. It is accomplished either by immersing the wafer in a bath of developer or by spraying fine jets of developer on the wafer surface. Control of both the development time and temperature is important for reproducible

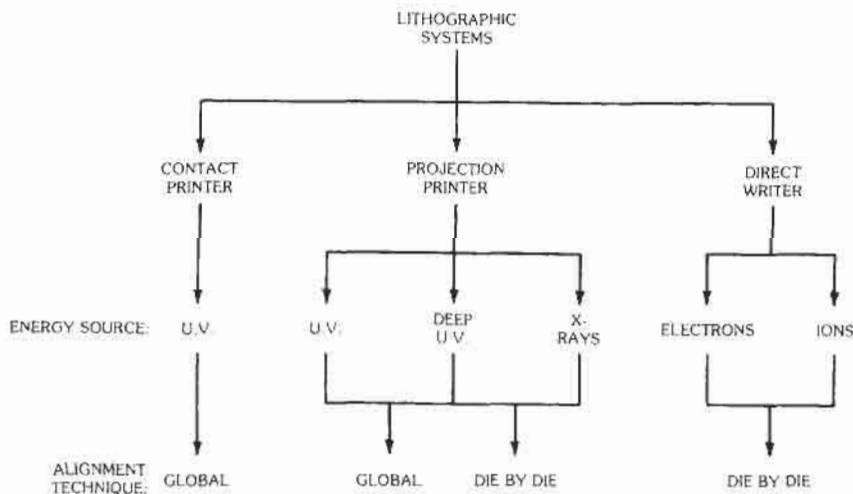


Fig. 5. Current lithographic systems.

line-width control. Fig. 6 shows schematically a positive-resist film after exposure and development. After development, the wafer is rinsed to remove all unwanted resist residues. Then the wafer is oven baked to make the resist more resistant to subsequent chemical etching (postbake). In multilevel-resist systems, only the top high-resolution resist is developed. With this resist used as a mask, the intermediate layer and thick planarizing resist must be dry etched to reveal the underlying films, as shown in Fig. 7.

Etching

All components of an integrated circuit are made up of a selective arrangement of variously doped semiconductors, insulators, and conductors. This requires the capability to etch certain layers selectively without damaging others. The important features that determine the preferential choice of an etchant are:

- (A) Etch rate
- (B) Etch rate selectivity
- (C) Directional selectivity

There are two prevalent categories of etchants in production use.

Wet Etching—Wet etching refers to the use of liquid chemicals. During the early days of IC development, this was the only kind of etchant in use. For many films in use, chemical etchants are available that provide acceptably controlled etch rates and high etch-rate selectivity to other layers. However, because of the very nature of the etching process, these etchants suffer from poor directional etch selectivity. Fig. 8 schematically depicts the line-width loss suffered from lateral etching with wet etches. This precludes their

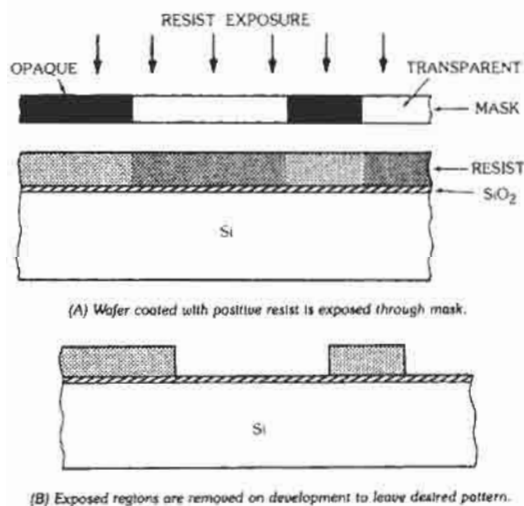


Fig. 6. Patterning with resist.

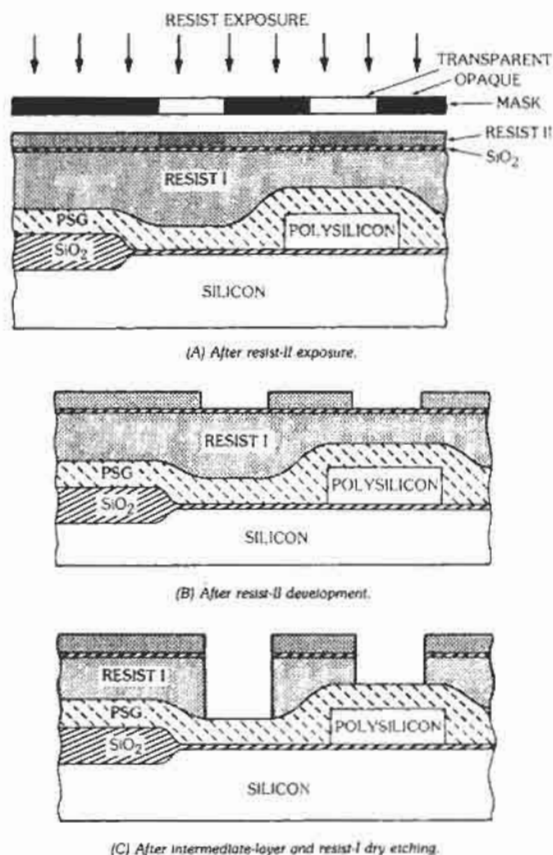


Fig. 7. Double-level resist systems.

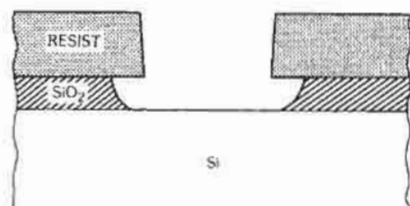


Fig. 8. Schematic cross section illustrating line-width loss from isotropic wet etch processes.

use in etching the narrow lines needed in VLSI. Nevertheless, for noncritical applications, such etchants find widespread use in etching resists, polycrystalline semiconductors, insulating layers such as oxides and nitrides, and metals such as aluminum. Table 5 lists etchants used in etching some common films.

Dry Etching—Dry etching makes use of reactive plasmas to carry out etching in the gas phase.* The key advantages of this methodology are highly directional

TABLE 5. ETCHANTS USED WITH COMMON FILMS

Film	Etch	Composition	Temp. (°C)
Silicon	Planar etch	2 ml HF, 15 ml HNO ₃ , 5 ml CH ₃ COOH	25
Polysilicon	Iodine etch	50 ml HF, 100 ml HNO ₃ , 110 ml CH ₃ COOH, 0.3 g I ₂	25
Silicon dioxide	Bell 2 etch	54 ml H ₂ O, 36 ml NH ₄ F, 10 ml HF	25
Silicon nitride	Hot phosphoric etch	H ₃ PO ₄	165
Aluminum	Phosphoric etch	55 ml H ₃ PO ₄ , 11 ml CH ₃ OOH, 4.5 ml HNO ₃ , 2 ml H ₂ O	25
Titanium	—	90 ml H ₂ O, 10 ml HF	25
Tantalum	—	20 ml HNO ₃ , 10 ml HF, 10 ml H ₂ O	25
Molybdenum	Dalton etch	92 g K ₃ [Fe(CN) ₆], 20 g KOH, 300 ml H ₂ O	25
Resists	J100	*	100

* Unknown, manufactured by Indust-Ri-Chem Lab., Richardson, Tex.

etch anisotropy and the facility to penetrate small resist openings; these make it suitable for etching the small geometries encountered in VLSI. By suitably adjusting the constituents and proportions of etchant species, a good etch-rate selectivity to other layers is achieved. Plasma etching of resists, SiO₂, Si₃N₄, Al, Poly-Si, and metal silicides is in extensive use in the industry today. Fig. 9 shows schematically small geometries etched through a resist mask. A list of gases commonly used for plasma etching is given in Table 6.*

Damascene Processing—The continued increase in speed and size of ICs has resulted in a change from the traditional aluminum metallization and SiO₂ interlevel metal insulation to copper metallization and lower dielectric constant insulators (called “low-k” insulators). The copper provides lower resistivity for the wires that interconnect individual transistors. The lower dielectric constant insulators provide reduced capacitance on these wires. Both of these factors decrease the

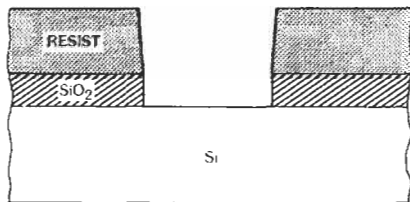


Fig. 9. Because of high anisotropy in etch rate, almost vertical walls are obtainable with plasma etching without any line-width loss.

TABLE 6. GASES USED FOR PLASMA ETCHING

Film	Gases
Silicon dioxide	SiF ₄
	CF ₄
	C ₃ F ₈
Silicon	CF ₄ , O ₂
	CCl ₄ , HCl
Silicon nitride	CF ₄
Tantalum, titanium, tungsten, molybdenum, vanadium	CF ₄
Chrome, chrome oxide	CCl ₄
Aluminum	CCl ₄ , BCl ₃
Resists	Ar, O ₂

RC delay of signals on the lines and therefore increase the speed capabilities of the circuits. The shift from aluminum/SiO₂ to Cu-low k has resulted in an alternative to the conventional thin film deposition, resist pattern application, and etching process sequence. This new process sequence is called damascene processing. In this process sequence the dielectric is deposited first, and the areas where copper conductors are desired are removed by the conventional resist application and etch process. This leaves the dielectric with trenches in a pattern desired for the copper metallization. The thick layer of copper metal is deposited to fill the etched areas as well as depositing over the surface of the dielectric. The entire surface of the wafer is then polished by a chem-mechanical polishing process (called CMP) to remove the surface layer of copper, leaving the desired conducting wires only in the trench regions, separated by the low-k dielectric insulator. There are

* Reference 11.

* Reference 7.

many details and variants of damascene processes that may be found in the references.

Substrate Doping Procedures

All semiconductor devices rely on selective doping of various areas, either n-type or p-type, to a required concentration and depth. There are three predominant techniques of introducing dopants in a controlled manner.

Predeposition and Drive-In—The geometries to be doped are lithographically defined on the substrate in a mask layer such as oxide or nitride that is impervious to the given dopant. Then the wafers are exposed to the dopant source, which may be a gas, solid, or liquid, in a well-controlled furnace at high temperature. An inert gas such as nitrogen or argon is used as a carrier for the dopant species.

The doping density in the substrate is a function of predeposition temperature and time and is given by the equation

$$N(x, t_p) = N_s \operatorname{erfc} \left[x / 2(D_p t_p)^{1/2} \right] \quad (\text{Eq. 1})$$

where

N is the doping density at a distance x below the surface,

t_p is the predeposition time,

D_p is the diffusion constant,

N_s is the surface concentration.

The surface concentration is usually equal to the solid solubility. The solid solubility and diffusion constant of common impurities in silicon are shown in Figs. 10 and 11.*

The total number of dopants introduced in the substrate per unit area is obtained, by integrating Equation 1 over the depth, to be the dose

$$Q_p = 2N_s (D_p t_p / \pi)^{1/2} \quad (\text{Eq. 2})$$

After the predeposition cycle, wafers are generally loaded in a different furnace at a higher temperature to drive in the impurities to obtain a required junction depth. A passivating layer such as an oxide is initially either grown or deposited on the wafers to prevent the escape of impurities. A redistribution of impurities takes place, and the doping density as a function of drive-in time is given by

$$N(x, t_d) = \left[Q_p / (D_d t_d)^{1/2} \right] \exp(-x^2 / 4D_d t_d) \quad (\text{Eq. 3})$$

where

t_d is the drive-in time,

D_d is the diffusion constant at the drive-in temperature.

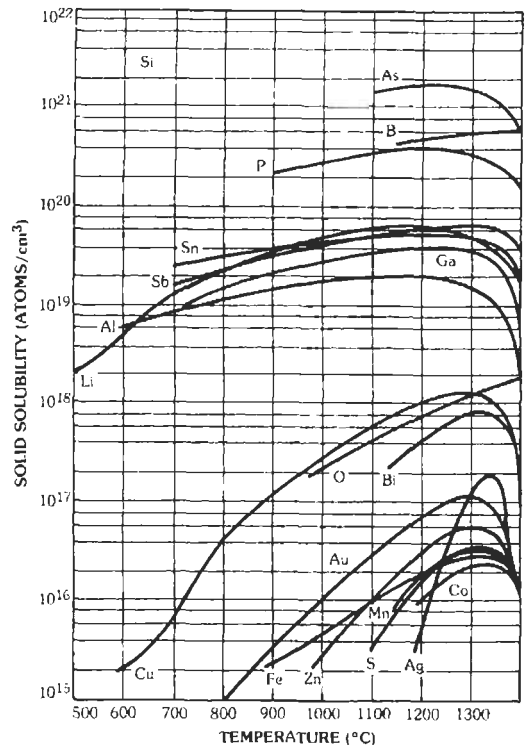


Fig. 10. Solid solubility of elements in silicon. (From F. A. Trumbore, "Solid Solubilities of Impurity Elements in Germanium and Silicon," *The Bell System Technical Journal*, Vol. 39, Jan. 1960, pp. 205–233. © 1960, AT&T; reprinted with permission.)

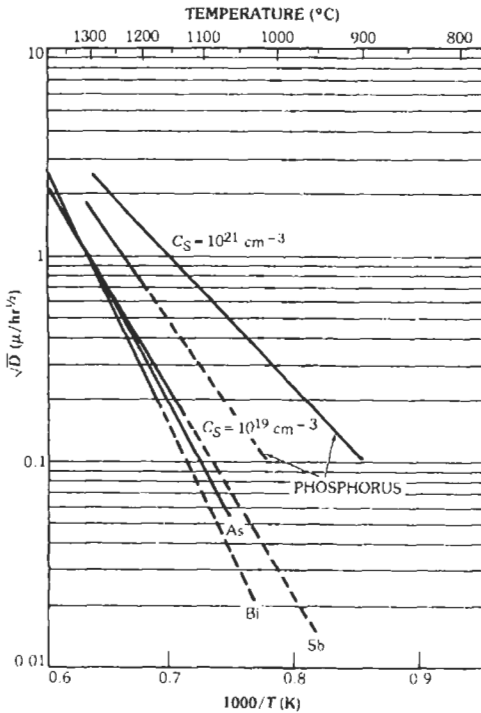
Equation 3 can be used to calculate the resulting junction depth if the background concentration is known.

Diffusion from Doped Oxides—Because of the solubility of the dopant at the surface, the standard two-step diffusion process is not suitable for obtaining shallow junctions of low doping concentration, as required in a number of applications. One alternative is to use doped glass as a diffusion source.† Doped glass containing the required doping can be either chemically deposited at low temperature and low pressure or spun-coated on the wafer surface. The impurities are then driven in at a higher temperature. By limiting the dopant concentration in the doped oxide, shallow junctions with low doping density can be realized.

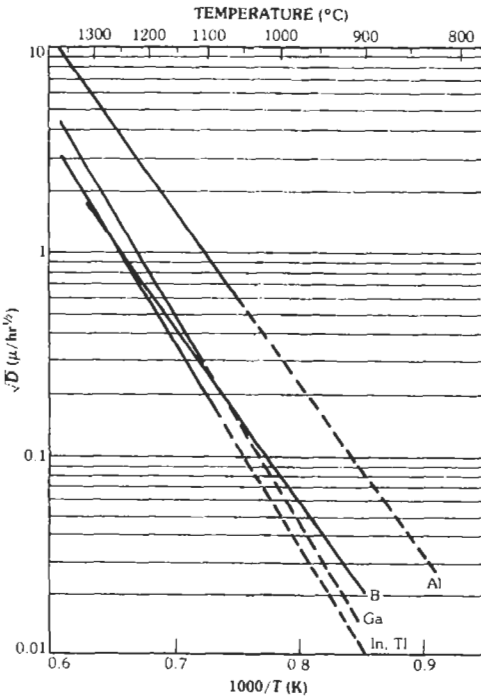
Doping by Ion Implantation—There are certain limitations associated with the diffusion processes that can be overcome by ion implantation. Ion implantation is a technique of extracting dopant species from a source, separating the required ions from other spurious particles, accelerating them to the required energy,

* References 12 and 13.

† Reference 14.



(A) Donor Impurities.



(B) Acceptor Impurities.

Fig. 11. Diffusion constant in silicon. (From A. S. Grove, *Physics and Technology of Semiconductor Devices*, New York: John Wiley & Sons, Inc., 1967.)

and embedding them in the substrate. The implantation energy determines the impurity concentration. To a good approximation, the implanted ions settle with a gaussian distribution given by Reference 15:

$$N(x) = \frac{Q}{\sqrt{2\pi}\Delta R_p} \exp\left[-\frac{1}{2}\left(\frac{R_p - x}{\Delta R_p}\right)^2\right] \quad (\text{Eq. 4})$$

where

- $N(x)$ is the dopant concentration at a distance x below the surface,
- Q is the ion dose,
- R_p is the ion range,
- ΔR_p is the range straggle.

The values for R_p and ΔR_p as a function of implant energy are available in published tables. Table 7 gives values for usual dopants in silicon.*

With this technique, it is possible to exercise a tight control on both the depth and the concentration of dopants, which is of prime importance in VLSI. By choosing multiple implant energies and doses, a variety of intended profiles can be approximated. The ion implantation is followed by a thermal anneal, which activates the dopants and anneals out the crystal damage produced by the process.

An IC Process Example—NMOS Process

The information presented above can best be consolidated by illustrating an IC fabrication process. For this purpose, we have chosen to discuss NMOS technology, whose typical basic process flow can be illustrated with the help of Fig. 12.

The starting material is p-type silicon in wafers that are carefully cleaned and thermally oxidized in steam to obtain an SiO_2 thickness of 300 Å. An LPCVD Si_3N_4 is now deposited to a thickness of 2000 Å. Next, the isolation regions are defined by using a moat mask as shown in Fig. 12A. Sequentially, dry plasma etching of nitride and oxide is carried out, with resist protecting the areas not to be etched. Resist is removed by J100 wet etch, and a low-energy boron implant is made that raises the p-type doping density in the opened windows but does not penetrate through the nitride and oxide stack, as illustrated in Fig. 12B. Slices are then cleaned and thermally oxidized in steam to oxidize the exposed regions to obtain an SiO_2 thickness of 7000 Å; nitride prohibits oxidation of other areas, as shown in Fig. 12C. Nitride is now wet etched in hot phosphoric acid, and an unmasked boron implant is made to adjust the threshold voltage of enhancement-mode devices, as shown in Fig. 12D.

The depletion-mode-device region is defined by a second mask in order to adjust the threshold voltage of depletion-mode devices. An arsenic implant dopes the exposed silicon areas n-type near the surface. This is

* Reference 16.

TABLE 7. VALUES OF R_p AND ΔR_p

Implant Energy (keV)	Dopants in Silicon					
	Arsenic		Boron		Phosphorus	
	R_p (Å)	ΔR_p (Å)	R_p (Å)	ΔR_p (Å)	R_p (Å)	ΔR_p (Å)
10	97	36	333	171	139	69
20	159	59	662	283	253	119
30	215	80	987	371	368	166
40	269	99	1302	443	486	212
50	322	118	1608	504	607	256
60	374	136	1903	556	730	298
70	426	154	2188	601	885	340
80	478	172	2465	641	981	380
90	530	189	2733	677	1109	418
100	582	207	2994	710	1238	456
110	634	224	3248	739	1367	492
120	686	241	3496	766	1497	528
130	739	258	3737	790	1627	562
140	791	275	3774	813	1757	595
150	845	292	4205	834	1888	628
160	898	308	4432	854	2019	659
170	952	325	4654	872	2149	689
180	1005	341	4872	890	2279	719
190	1060	358	5086	906	2409	747
200	1114	374	5297	921	2539	775

shown in Fig. 12E. The slices are now carefully cleaned, a high-quality 500 Å thermal oxide is grown in dry oxygen to serve as the gate insulator, and a 5000 Å LPCVD polycrystalline silicon layer is deposited. This polysilicon layer is doped heavily n-type (n+) by phosphorus diffusion; the wafer cross section at this step is shown in Fig. 12F. The polysilicon gate regions are patterned with a third mask and dry etched, as shown in Fig. 12G. The resist is wet etched, and a heavy arsenic implant is made that converts regions not covered by polysilicon or thick oxide to n+ type, as depicted in Fig. 12H.

A 7000 Å phosphosilicate glass (PSG) is deposited on cleaned wafers and patterned with a fourth mask, and the doped glass is dry etched to open contact areas, as shown in Fig. 12I. Resist is removed by wet etching, and the wafers are cleaned and subjected to a temperature greater than 900 °C, which causes the doped glass to flow so that all edges are smooth. Now a 10 000 Å aluminum film is evaporated or sputter deposited (Fig. 12J). A fifth mask patterns the aluminum interconnect, aluminum is dry etched, and resist is removed by wet etch (Fig. 12K). After careful clean-up, the aluminum is sintered at 450 °C in H₂ to pro-

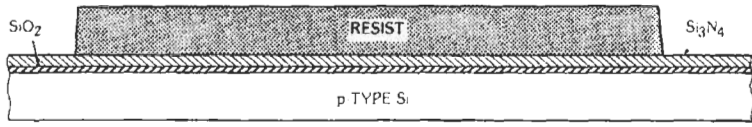
vide good ohmic contacts. A 3000 Å layer of plasma nitride is now deposited (Fig. 12L). A sixth mask defines the bonding-pad regions, where the coated nitride is plasma etched to allow access to the bonding pads. The wafers are cleaned, and this finishes the fabrication process.

The circuits are tested for functionality and then diced. Good dies are bonded in appropriate packages, retested, and shipped to customers.

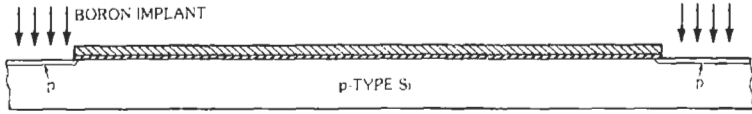
Yield Statistics

Only a fraction of a large number of chips on a silicon wafer are completely functional. Defects in the masks, dust particles on wafer surfaces, nonideality of the basic silicon material, and short or open circuits in the wiring all cause some of the circuits to be nonfunctional. With present design techniques, any single defect of sufficient size will kill an entire chip.

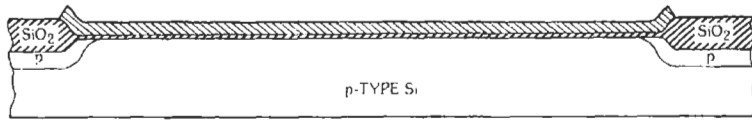
The simplest model for the yield, or the fraction of chips that do not have defects, assumes a random defect distribution across the wafer. If there are D fatal defects per unit area and the area of an individual chip



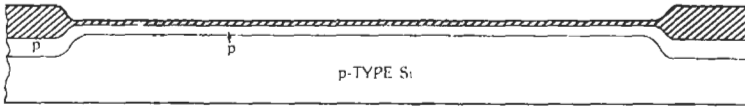
(A) After field pattern.



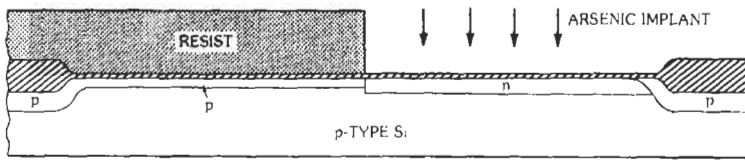
(B) After field implant.



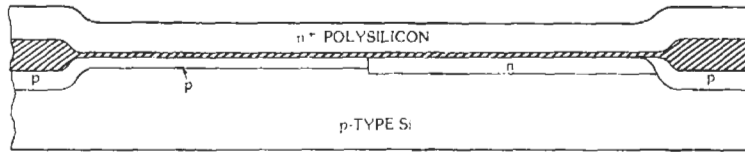
(C) After field oxidation.



(D) After enhancement threshold implant.

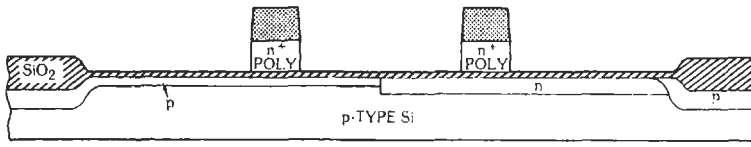


(E) After pattern and depletion implants.

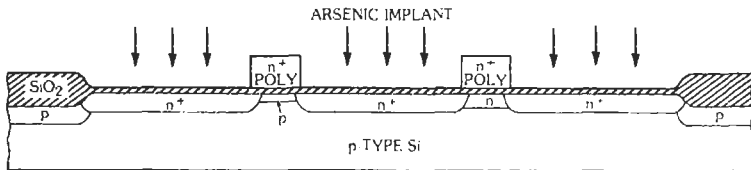


(F) After gate oxidation, poly deposition, and doping.

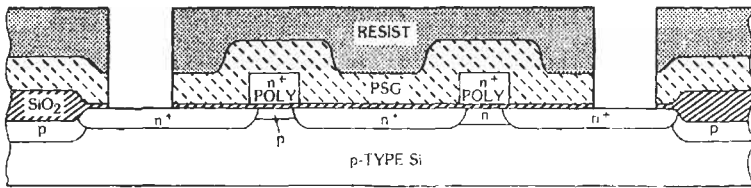
Fig. 12. Schematic cross sections illustrating the purpose and



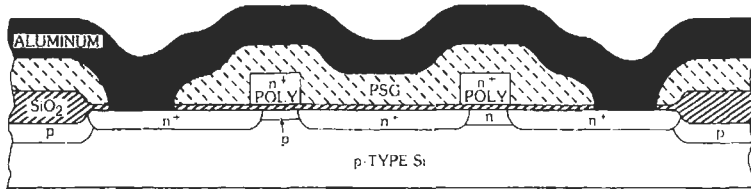
(G) After gate pattern and etch.



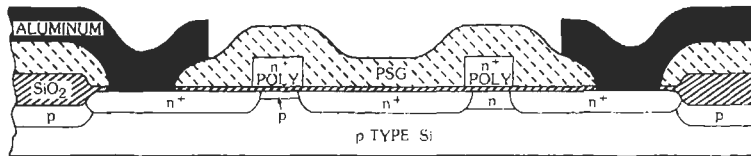
(H) After source-drain implant.



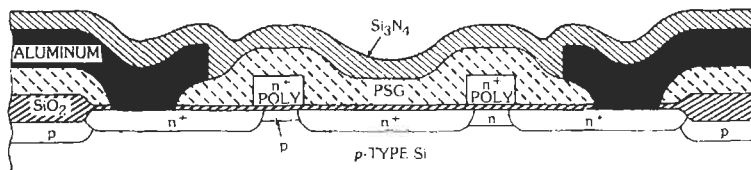
(I) After PSG deposition, contact pattern, and etch.



(J) After aluminum deposition.



(K) After aluminum etch.



(L) After nitride overcoat.

consequence of various steps involved in an NMOS process.

is A , then the probability that a chip has n flaws is, in the simplest case, given by the Poisson distribution $P_n(DA)$. The probability of a good chip is:

$$P_0(DA) = e^{-DA} \quad (\text{Eq. 5})$$

While this equation is not rigorously applicable to fabrication processes, it is a good approximate model for estimating the yield of various design alternatives.

IC DESIGN*

The design of integrated circuits requires the synthesis and analysis of a large number of active elements. The classical form of analysis can be extended to integrated circuits, whereas the focus of synthesis goes beyond the classical notions inasmuch as the geometric layout of the circuit topography for integrated circuits is a major fraction of the circuit synthesis. The electrical aspects of circuit design derive from the considerations that are laid out in Chapter 18, with certain constraints based on the scale of integration. The design of the circuit topography is unique to integrated circuits and has evolved as a major discipline.

Modeling and Simulation

The design of integrated circuits requires the electrical analysis of circuits that contain a large number of elements. For LSI or VLSI circuits, this involves the simulation of up to tens of millions of circuit elements if the total circuit response is to be examined. The evaluation of this class of circuits is extremely computer-intensive. In most cases, it is not practical to simulate the circuit with classical time- or frequency-domain analysis at the transistor level. The current practice for simulation and analysis of such classes of circuits is to adopt a hierarchical procedure with different levels of abstraction at each level. The design of an LSI logic circuit would involve a typical modeling hierarchy such as that in Chart 1.

The top of the hierarchy is an abstract definition of the architecture of the circuit, which is used to provide a guideline for the various ways of accomplishing the objective of the circuit. It trades off, for example, the use of pipeline processing versus serial processing. The behavioral level of simulation actually involves

CHART 1. DESIGN HIERARCHY FOR LSI CIRCUITS

Architecture simulation
Behavioral simulation
Functional simulation
Logic simulation
Transistor simulation

* Reference 9.

the definition of the major blocks of the circuit and their interaction, with the details of the overall data or control flow to accomplish the objective of the chip being examined. The functional level actually describes the overall logical response of the major blocks, relating the logical inputs and outputs with no details of internal realization of logic in the block. The logic-simulation level details the realization of each block at the gate level, provides logic minimization, and in some cases introduces the notion of relative timing. The transistor-level simulation considers the transient response of the circuit, including the detailed simulation of all elements of the circuit. The key to the usefulness of this hierarchical simulation is the ability to mix the different levels of abstraction to examine the performance of the entire circuit with focus on one block at a time. Such mixed-mode simulators are being evolved, and common hardware-description languages that operate on a unified data base that is accessed by any level of the hierarchy are now available.

Topological Realization of Circuits

The design of integrated circuits differs from the design of board-level circuits in the importance of the actual physical realization of the active elements and interconnections. Typical LSI circuits have many thousands of active elements and interconnections that must be topologically related to each other. This requires the interaction of the device physics, process technology constraints, and topological constraints. To obtain a practical solution to this complex interactive design environment, it is generally accepted practice to describe the process constraints in the form of a design rule package. This design rule package is a simple description of the lateral spatial relationship of the various active-element forming geometries as well as the wiring. A simple example of design rules is shown in Fig. 13. The original concept of the design rules was governed by the constraints of the technology only. However, the design data base for the geometrical description of the circuit for VLSI circuits is very large, and some constraints that limit the size of this data base have been introduced into the design rules.† This is simply a quantization of the minimum spatial distance describing any technology constraint. This kind of quantization is equivalent to the definition of the finest grid on which a geometry must fall.

There are two major approaches to the realization of the geometrical data base for an IC. The first is the classical approach in which the complete drawing of each individual section is introduced into the data base in detail. This is a tedious and inefficient way of generating data. The preferred approach to data entry is through a symbolic description of the layout in which the definitions of the often-used elements are

† Reference 9.

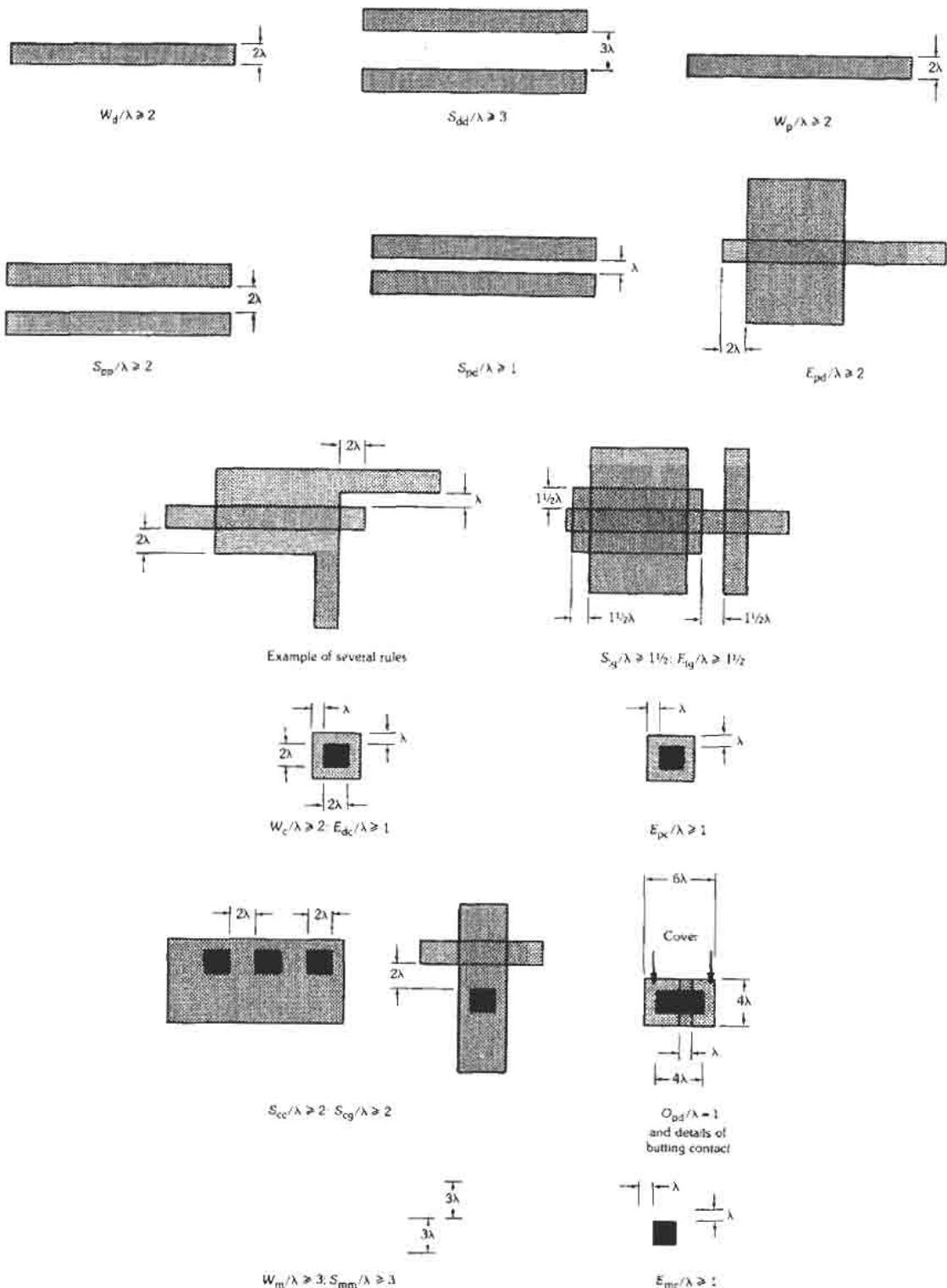


Fig. 13. An example of a design rule set. (From C. Mead and L. Conway, *Introduction to VLSI Systems*. Reading, MA: Addison-Wesley Publishing Co., Inc., 1980.)

built into the system and the designer calls up such elements and strings them together on a graphic entry system. These symbols are a modified description of the circuit elements much like the drawing of a classical circuit schematic (Fig. 14). The design rules are then used to generate automatically the full two-dimensional geometrical description from the symbolic description. The symbolic layout system is very efficient and suffers only from a 10–20% reduction of the packing density of circuit elements compared to a detailed customized input of the complete geometries. In an attempt to reduce this gap, software compaction of the design data base has been developed. The connection of geometries on a layout data base is another aspect of integrated-circuit design that has been traditionally inefficient. Automatic wire-routing software developed for multilayered printed circuit boards has

been adapted for this purpose and is being used. The use of artificial-intelligence concepts in the realization of more efficient routing algorithms is under investigation.

Once the design data base has been created, the design task is to verify the validity of the data base geometrically and electrically. The geometric verification is a checking of all geometries for design-rule violations. This is achieved routinely by all layout systems through software analysis of the data base. The electrical validity of the geometrical data base involves the extraction of the electrical circuit schematic and the comparison of the electrical schematic to the various levels of hierarchy of the initial design. This upward feedback through the hierarchy is of crucial importance to the design cycle and requires the design and geometric data bases to be compatible. It is currently common practice to extract electrical schematics and parasitic resistances and capacitances from the layout data base and compare them to the design data base. The comparison occurs at the transistor level.

Various integrated design systems with the above functions are currently available.

PACKAGING INTEGRATED CIRCUITS

The standards for packaging integrated circuits and for allocation of functions to pins for the purposes of interchangeability are set by international committees EIA/JEDEC (Electronic Industries Association) and IEC (International Electrotechnical Commission). There are two major kinds of specifications for each type of package: mechanical and thermal. A variety of package types exist, the most popular of which will be described.*

Plastic Packages

Plastic dual-in-line packages consist of a circuit mounted on a 16-, 18-, 20-, 24-, or 28-pin lead frame and encapsulated within an electrically nonconductive plastic compound. The compound will withstand soldering temperature with no deformation, and circuit performance characteristics remain stable when the device is operated in high-humidity conditions. These are the lowest-cost packages for integrated circuits. An example of a 16-pin plastic package is shown in Fig. 15.

Ceramic Packages

Another type of dual-in-line package is hermetically sealed and consists of a ceramic base, a ceramic cap, and a 16-, 18-, 20-, 24-, 28-, or 48-lead frame. Hermetic sealing is accomplished with glass. These devices are divided into two categories: side-braze and frit seal. An

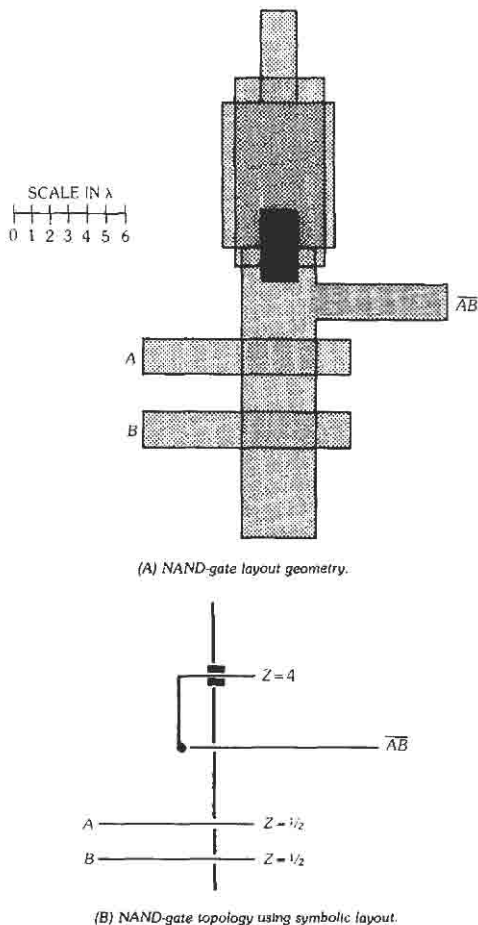
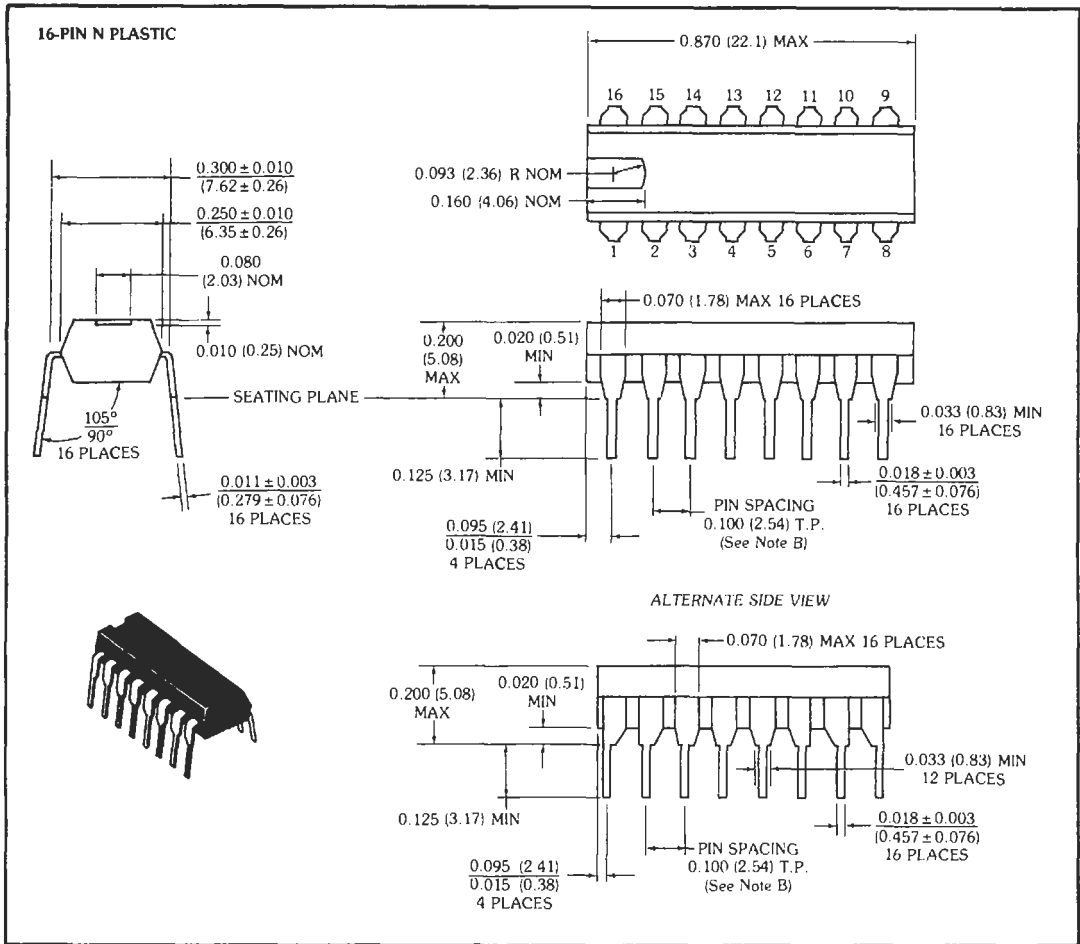


Fig. 14. Symbolic layout and its topological equivalent. (From C. Mead and L. Conway, *Introduction to VLSI Systems*. Reading, MA: Addison-Wesley Publishing Co., Inc., 1980.)

* Details of packages can be found in integrated-circuit data books.



NOTES: A. All dimensions are shown in inches (and parenthetically in millimeters for reference only). Inch dimensions govern.
 B. Each pin centerline is located within 0.010 (0.26) of its true longitudinal position.

Fig. 15. Dual-in-line plastic packaging for integrated circuits. (From *TTL Data Book*, 2nd ed. Dallas, TX: Texas Instruments, Inc., 1976; p. 4-6.)

example of the mechanical dimensions of a 16-pin dual-in-line ceramic package is shown in Fig. 16.

example of a 36-pad leadless ceramic chip carrier is shown in Fig. 18.

Ceramic Flat Package

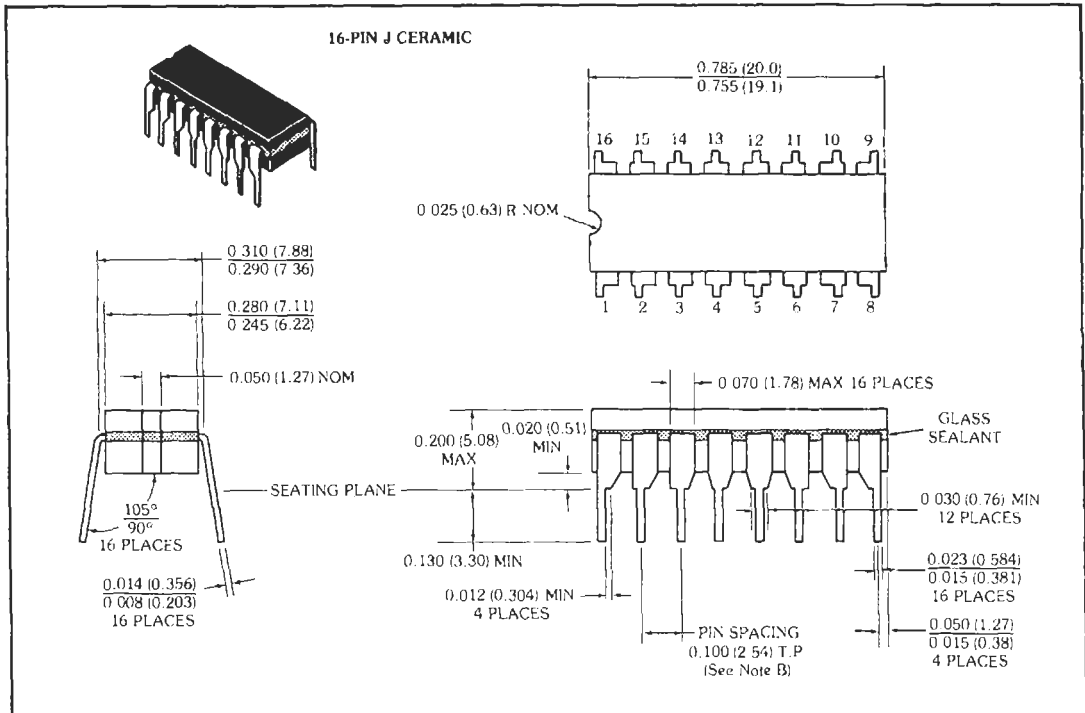
Hermetically sealed ceramic flat packages consist of an electrically nonconductive ceramic base and cap and a 16- or 24-pin lead frame. An example of the mechanical dimensions of such a package is shown in Fig. 17.

Leadless Ceramic Chip Carrier

Leadless ceramic chip-carrier packages are emerging as the most popular for LSI and VLSI chips because of their board-level density advantage. An

Flip-Chip Package

Flip-chip packages attach an IC to a chip carrier by placing solder balls on the interconnect pads and bonding the chip face down ("flipped") on the surface of a chip carrier or circuit board. The advantage of the flip-chip packaging approach is that it can be done with pads over the entire surface of the chip, all of the bonds can be done at one time, and the cost can be reduced by preparing the chips for final packaging in wafer form. This last item allows the use of the usual low-cost IC mass production techniques to replace some serial, more costly, packaging operations. Flip-chip packaging



NOTES A. All dimensions are shown in inches (and parenthetically in millimeters for reference only) Inch dimensions govern
 B. Each pin centerline is located within 0.010 (0.26) of its true longitudinal position

Fig. 16. Dual-in-line ceramic packaging for integrated circuits. (From *TTL Data Book*, 2nd ed. Dallas, TX: Texas Instruments, Inc., 1976; p. 4-4.)

is just being introduced into manufacturing, but it promises to rapidly advance in applications.

Thermal Resistance

The most important system specification of an IC package is its ability to conduct heat, since the IC dissipates power in the package. The thermal properties of a package are specified as a thermal resistance. Junction-to-ambient thermal-resistance values of dual-in-line packaging systems are shown in Fig. 19. This figure of merit is used to relate the internal chip temperature to the ambient system temperature.

DIGITAL INTEGRATED CIRCUITS

Digital logic functions can be realized in many different configurations. The traditional logic realization used resistors or diodes in conjunction with a transistor. This approach has been replaced in the current integrated circuits by the use of active transistors mostly, since these can be realized most efficiently in the technology. The logic realization is

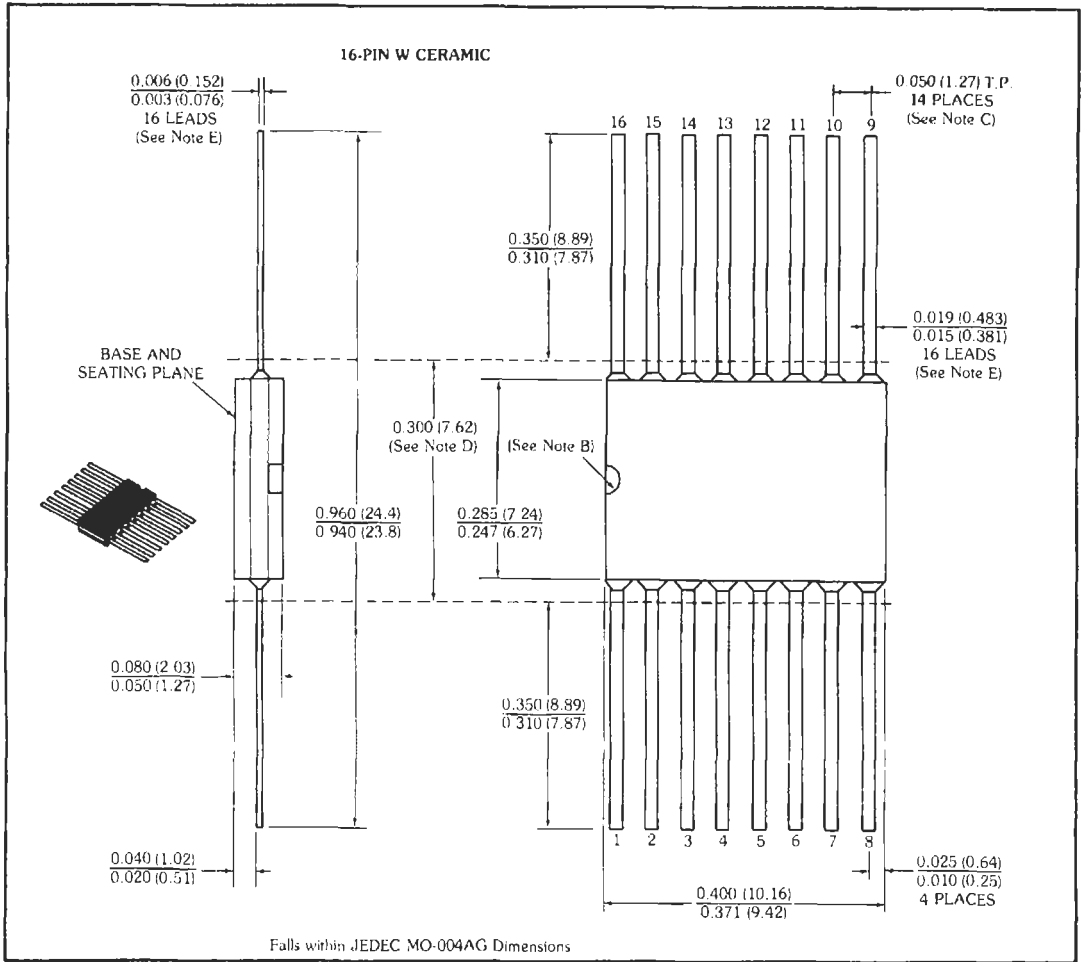
dependent on the property of the transistor used and will be categorized differently for bipolar and field-effect transistors.

The implementation of logic-integrated circuits requires the definition of the voltage, current, and timing standards. These have been defined in Table 8. A very large variety of digital logic circuits spanning SSI, MSI, LSI, and VLSI is commercially available. The various manufacturers' data books and applications books should be consulted for detailed information.

Bipolar Logic Families

Bipolar logic functions are commercially available as SSI and MSI circuits, and they are being developed for LSI and VLSI circuits. The different families described here are graded from the most popular SSI and MSI standard logic to the more advanced families used for LSI and VLSI.

Transistor-Transistor Logic (TTL)*—The most widely used logic family for SSI and MSI logic is transistor-transistor logic. A broad spectrum of TTL circuits are available that allow logic designers to optimize all portions of a system cost effectively.



- NOTES: A. All dimensions are shown in inches (and parenthetically in millimeters for reference only). Inch dimensions govern.
 B. Index point is provided on cap for terminal identification only.
 C. Leads are within 0.005 (0.13) radius of true position (T.P.) at maximum material condition.
 D. This dimension determines a zone within which all body and lead irregularities lie.
 E. Not applicable for solder-dipped leads.
 F. When solder-dipped leads are specified, dipped area extends from lead tip to within 0.050 (1.27) of package body.

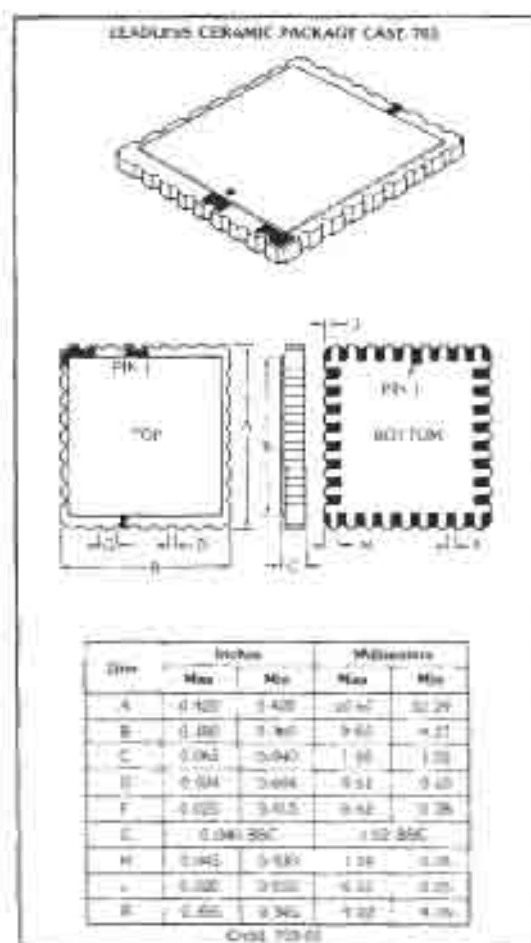
Fig. 17. Ceramic flat packaging for integrated circuits. (From *TTL Data Book*, 2nd ed. Dallas, TX: Texas Instruments, Inc., 1976; p. 4-10.)

The basic schematic of generic TTL logic is shown in Fig. 20. This arrangement utilizes the base-emitter diodes and collector-base diode of the multiemitter transistor for logic and the output transistor for drive. This is very easily realized in planar bipolar technology. Although the basic TTL circuit consists of two transistors, additional components are necessary to increase circuit speed and fan-out capability (Fig. 21). Two major classes of TTL logic exist. These are called

54XX and 74XX series TTL. The 74XX series is a limited-temperature (0 °C to +70 °C) IC, and the 54XX series is a full-military-temperature (-55 °C to +125 °C) IC. A classification by low power and high speed in addition to the standard series exists.

An improvement in the TTL performance is achieved by the addition of a Schottky barrier diode that permits the transistor to operate in the nonsaturated mode; this reduces storage time and improves speed. Initially, the Schottky TTL was developed for high speed. Further evolution of the Schottky technology led to the low-power Schottky TTL family in

* Reference 17.



NOTE: See also proposed mil-STD-883C, 208, 209, 210, 211, 212, 213, 214, 215, 216, 217, 218, 219, 220, 221, 222, 223, 224, 225, 226, 227, 228, 229, 230, 231, 232, 233, 234, 235, 236, 237, 238, 239, 240, 241, 242, 243, 244, 245, 246, 247, 248, 249, 250, 251, 252, 253, 254, 255, 256, 257, 258, 259, 260, 261, 262, 263, 264, 265, 266, 267, 268, 269, 270, 271, 272, 273, 274, 275, 276, 277, 278, 279, 280, 281, 282, 283, 284, 285, 286, 287, 288, 289, 290, 291, 292, 293, 294, 295, 296, 297, 298, 299, 300, 301, 302, 303, 304, 305, 306, 307, 308, 309, 310, 311, 312, 313, 314, 315, 316, 317, 318, 319, 320, 321, 322, 323, 324, 325, 326, 327, 328, 329, 330, 331, 332, 333, 334, 335, 336, 337, 338, 339, 340, 341, 342, 343, 344, 345, 346, 347, 348, 349, 350, 351, 352, 353, 354, 355, 356, 357, 358, 359, 360, 361, 362, 363, 364, 365, 366, 367, 368, 369, 370, 371, 372, 373, 374, 375, 376, 377, 378, 379, 380, 381, 382, 383, 384, 385, 386, 387, 388, 389, 390, 391, 392, 393, 394, 395, 396, 397, 398, 399, 400, 401, 402, 403, 404, 405, 406, 407, 408, 409, 410, 411, 412, 413, 414, 415, 416, 417, 418, 419, 420, 421, 422, 423, 424, 425, 426, 427, 428, 429, 430, 431, 432, 433, 434, 435, 436, 437, 438, 439, 440, 441, 442, 443, 444, 445, 446, 447, 448, 449, 450, 451, 452, 453, 454, 455, 456, 457, 458, 459, 460, 461, 462, 463, 464, 465, 466, 467, 468, 469, 470, 471, 472, 473, 474, 475, 476, 477, 478, 479, 480, 481, 482, 483, 484, 485, 486, 487, 488, 489, 490, 491, 492, 493, 494, 495, 496, 497, 498, 499, 500, 501, 502, 503, 504, 505, 506, 507, 508, 509, 510, 511, 512, 513, 514, 515, 516, 517, 518, 519, 520, 521, 522, 523, 524, 525, 526, 527, 528, 529, 530, 531, 532, 533, 534, 535, 536, 537, 538, 539, 540, 541, 542, 543, 544, 545, 546, 547, 548, 549, 550, 551, 552, 553, 554, 555, 556, 557, 558, 559, 560, 561, 562, 563, 564, 565, 566, 567, 568, 569, 570, 571, 572, 573, 574, 575, 576, 577, 578, 579, 580, 581, 582, 583, 584, 585, 586, 587, 588, 589, 590, 591, 592, 593, 594, 595, 596, 597, 598, 599, 600, 601, 602, 603, 604, 605, 606, 607, 608, 609, 610, 611, 612, 613, 614, 615, 616, 617, 618, 619, 620, 621, 622, 623, 624, 625, 626, 627, 628, 629, 630, 631, 632, 633, 634, 635, 636, 637, 638, 639, 640, 641, 642, 643, 644, 645, 646, 647, 648, 649, 650, 651, 652, 653, 654, 655, 656, 657, 658, 659, 660, 661, 662, 663, 664, 665, 666, 667, 668, 669, 670, 671, 672, 673, 674, 675, 676, 677, 678, 679, 680, 681, 682, 683, 684, 685, 686, 687, 688, 689, 690, 691, 692, 693, 694, 695, 696, 697, 698, 699, 700, 701, 702, 703, 704, 705, 706, 707, 708, 709, 710, 711, 712, 713, 714, 715, 716, 717, 718, 719, 720, 721, 722, 723, 724, 725, 726, 727, 728, 729, 730, 731, 732, 733, 734, 735, 736, 737, 738, 739, 740, 741, 742, 743, 744, 745, 746, 747, 748, 749, 750, 751, 752, 753, 754, 755, 756, 757, 758, 759, 760, 761, 762, 763, 764, 765, 766, 767, 768, 769, 770, 771, 772, 773, 774, 775, 776, 777, 778, 779, 780, 781, 782, 783, 784, 785, 786, 787, 788, 789, 790, 791, 792, 793, 794, 795, 796, 797, 798, 799, 800, 801, 802, 803, 804, 805, 806, 807, 808, 809, 810, 811, 812, 813, 814, 815, 816, 817, 818, 819, 820, 821, 822, 823, 824, 825, 826, 827, 828, 829, 830, 831, 832, 833, 834, 835, 836, 837, 838, 839, 840, 841, 842, 843, 844, 845, 846, 847, 848, 849, 850, 851, 852, 853, 854, 855, 856, 857, 858, 859, 860, 861, 862, 863, 864, 865, 866, 867, 868, 869, 870, 871, 872, 873, 874, 875, 876, 877, 878, 879, 880, 881, 882, 883, 884, 885, 886, 887, 888, 889, 890, 891, 892, 893, 894, 895, 896, 897, 898, 899, 900, 901, 902, 903, 904, 905, 906, 907, 908, 909, 910, 911, 912, 913, 914, 915, 916, 917, 918, 919, 920, 921, 922, 923, 924, 925, 926, 927, 928, 929, 930, 931, 932, 933, 934, 935, 936, 937, 938, 939, 940, 941, 942, 943, 944, 945, 946, 947, 948, 949, 950, 951, 952, 953, 954, 955, 956, 957, 958, 959, 960, 961, 962, 963, 964, 965, 966, 967, 968, 969, 970, 971, 972, 973, 974, 975, 976, 977, 978, 979, 980, 981, 982, 983, 984, 985, 986, 987, 988, 989, 990, 991, 992, 993, 994, 995, 996, 997, 998, 999, 1000.

Fig. 18. Leadless ceramic chip-carrier package. (From Motorola CMOS Data, 1978, p. 2-14.)

which Schottky diodes replace the multimeric transistor as the input element. The Schottky TTL circuits are shown in Fig. 22.

Table 9 defines the voltage and current standards of various TTL logic families. Table 10 shows the performance characteristics of different TTL families.

Emitter-Coupled Logic (ECL)—The highest-speed, lowest-noise bipolar logic family is emitter-coupled logic. It is also known as current-mode logic (CML), since the circuit works on a current-steering principle. It is designed as a nonsaturating form of logic, which eliminates transistor storage time as a speed-limiting characteristic. The typical ECL gate shown in Fig. 23 comprises a differential-amplifier input, an internal bias reference, and an emitter-follower output to restore dc levels. High-fan-out operation is possible because of the high input impedance of

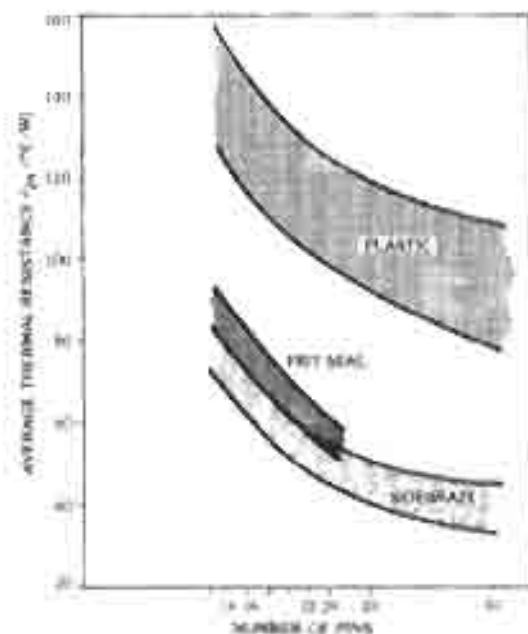


Fig. 19. Typical thermal resistance for various packages. (From Motorola CMOS Data, 1978, p. 2-1.)

TABLE 8. DEFINITIONS OF LOGIC STANDARD TERMINOLOGY

Parameter	Definition
V_{CC}	Most positive power supply voltage for a circuit
V_{DD}	Bias reference supply voltage
V_{EE}	Most negative power supply voltage for a circuit
V_{OH}	High (1) level output voltage
V_{OL}	Low (0) level output voltage
V_{IH}	Input clamp voltage
V_{IL}	Low (0) level input voltage
V_{IH}	High (1) level input voltage
V_{T1}	Negative-going threshold voltage
V_{T2}	Positive-going threshold voltage
$V_{OZ}(\text{off})$	Off-state output voltage
$V_{OZ}(\text{on})$	On-state output voltage
I_{IH}	High (1) level input current
I_{OH}	High (1) level output current
I_{IL}	Low (0) level input current
I_{OL}	Low (0) level output current
$I_{O}(\text{off})$	Off-state output current
$I_{O}(\text{on})$	On-state output current
I_{OC}	Short-circuit output current
I_{CC}	Supply current

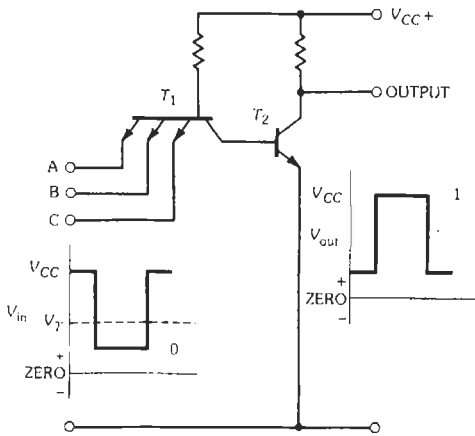
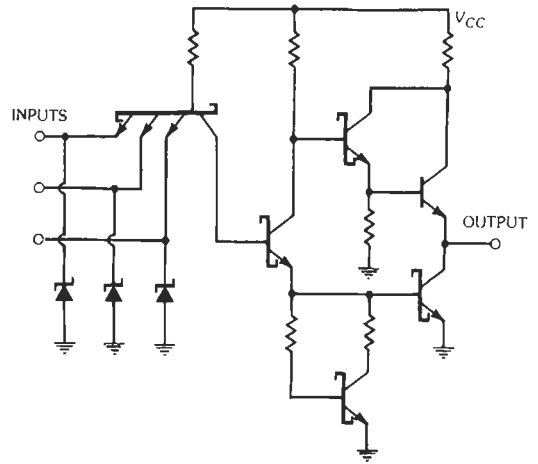
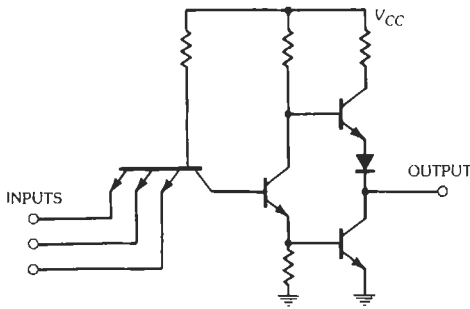


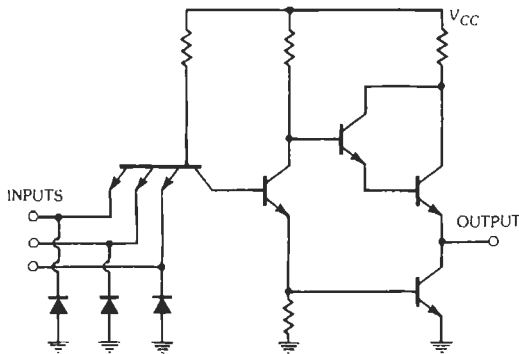
Fig. 20. The basic TTL circuit. (From *Designing with TTL Integrated Circuits*. Dallas, TX: Texas Instruments, Inc., 1971; Fig. 1-13, p. 12.)



(A) Standard.

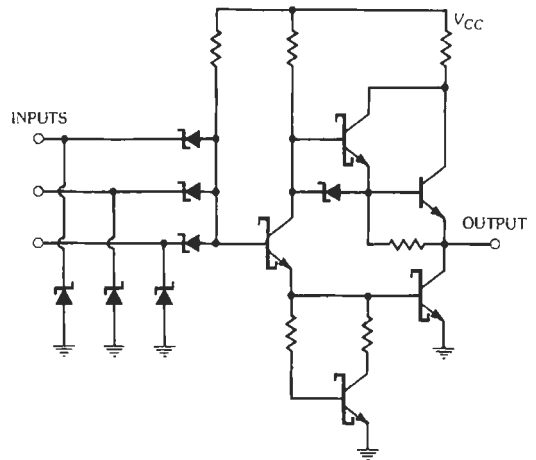


(A) Standard, or low-power, circuit.



(B) High-speed TTL.

Fig. 21. Diagrams of typical implementations of TTL circuits. (From *Designing with TTL Integrated Circuits*. Dallas, TX: Texas Instruments, Inc., 1971; Fig. 2-1, p. 18. Fig. 2-3, p. 19.)



(B) Low-power.

Fig. 22. Schottky TTL circuits. (From G. D. Kraft and W. N. Toy, *Mini/Micro Computer Hardware Design*. Englewood Cliffs, NJ: Prentice-Hall, Inc., 1979, p. 71.)

the differential amplifier and the low output impedance of the emitter follower. Power-supply noise is virtually eliminated by the nearly constant current drain of the differential amplifier even during transition time. The logic-voltage standards for ECL are shown in Table 11. Note that CML circuits tend to adopt TTL voltage standards. Typical gate speed for ECL logic ICs is 2 ns at 100 mW power dissipation for commercially available circuits.

Integrated Injection Logic (I²L)*—Integrated injection logic (I²L), also known as merged transistor logic (MTL), is a bipolar logic comparable in density

* Reference 18.

TABLE 9. TTL VOLTAGE AND CURRENT STANDARDS*

Parameter	TTL	Schottky TTL(LS)	Units
V_{IH}	2 (min)	2 (min)	V
V_{IL}	0.8 (max)	0.8 (max)	V
V_{OH}	2.4 (min)	2.5-2.7†	V
V_{OL}	0.4 (max)	0.4-0.5†	V
I_{IH}	40 (max)	20 (max)	μ A
I_{IL}	-1.6 (max)	-0.36 (max)	mA
I_{OH}	-400 (min)	-400 (min)	μ A
I_{OL}	16 (min)	4-8†	mA

* $V_{CC} = 5\text{ V} \pm 10\%$.

† Different specifications for the 54 and 74 series.

TABLE 10. 54/74 TTL TYPICAL PERFORMANCE CHARACTERISTICS

	Low-Power 54/74L	Standard 54/74	High-Speed 54/74H	Schottky 54/74S	Low-Power Schottky 54/74LS
Power/Gate	1 mW	10 mW	22 mW	19 mW	2 mW
Delay/Gate	33 ns	10 ns	6 ns	3 ns	7 ns

and power dissipation with dynamic MOS logic. It is an LSI- or VLSI-caliber bipolar logic family. The basic I^2L gate is shown in Fig. 24. The basic cell is made up of a multiemitter NPN transistor operated in the inverted mode and a lateral PNP transistor whose base and collector are common to the collector and base of the first transistor. For logic operation, the PNP transistor is used both as an injector of base current for the NPN transistor in the same cell and as a current-source load for an NPN transistor in an adjacent cell. The operation of the basic cell as an inverter can be understood as follows. When V_{in} is high, all the current from the PNP transistor is used up by NPN transistor T_1 operating in the saturated mode, and no current is left for the base of NPN transistor T_2 . Hence V_{OL} is low and V_{O2} is high. When V_{in} is low, the opposite situation applies. Typical gate delays of 5 ns per gate at 0.2 mW have been achieved at 2-micron geometries. Many advances based on self-aligned structures are being made that are enhancing the performance of this logic family.

Schottky Transistor Logic (STL)*—Schottky transistor logic, also called Schottky coupled transistor logic (SCTL) and complementary constant current logic (C^3L), is a generic diode transistor logic (DTL) implemented with Schottky logic diodes and a single NPN switch transistor. Fig. 25 shows a three-output STL gate. For relatively high-gain transistors ($h_{\beta} > 10$)

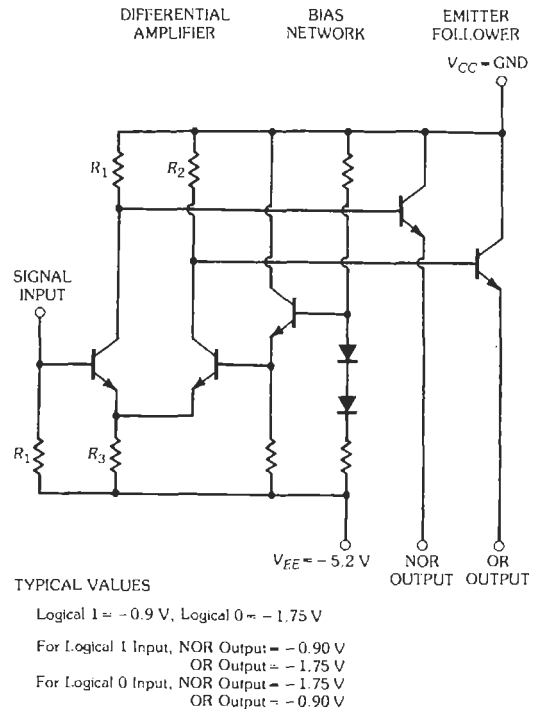


Fig. 23. Diagram of a typical ECL gate. (From *Sigmetics ECL Data Book*.)

* Reference 19.

TABLE 11. LOGIC STANDARDS FOR 10 000 SERIES ECL*

Parameter	10 100	10 500
	10 200	10 600
V_{IL} (min)	-1.85 V	-1.85 V
V_{IH} (max)	-0.81 V	-0.72 V
V_{OL} (min)	-1.85 V	-1.85 V
V_{OL} (max)	-1.65 V	-1.62 V
V_{OH} (min)	-0.96 V	-0.93 V
V_{OH} (max)	-0.98 V	-0.95 V
V_{BB} (typical)	-1.29 V	-1.29 V

* 25 °C, $V_{EE} = 5.2$ V, $V_{CC} =$ ground

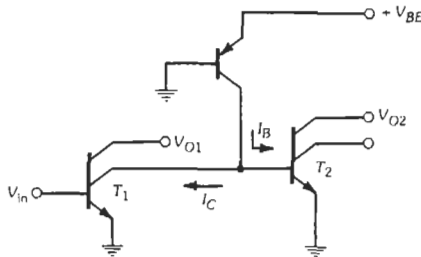


Fig. 24. The integrated-injection-logic, I^2L , gate. (From K. Hart and A. Slob, "Integrated Injection Logic," *IEEE Trans. Solid State Circuits*, SC-7 (1972), p. 346.)

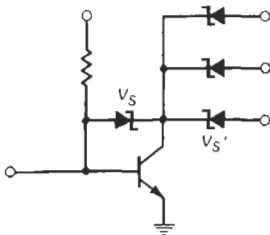


Fig. 25. Three-output STL gate. (From K. Hart and A. Slob, "Integrated Injection Logic," *IEEE Trans. Solid State Circuits*, SC-7 (1972), p. 346.)

and fan-out limited to one per logic diode, the signal-voltage swing is shown to be the difference in the Schottky-diode forward voltages. The STL characteristics are compatible with the supply voltages, power, density, and speed requirements of VLSI circuits. Typical on-chip gate speeds of less than 2 ns at 10 μ W of power are achievable with 2-micron technology.

MOS Logic Families*

There are two classes of MOS logic families, based on the operation of the logic: (A) static, or ratioed,

logic and (B) dynamic, or ratioless, logic. The realization of these two classes of logic is further described by technology type: p-channel (PMOS), n-channel (NMOS), and complementary (CMOS). The unique feature of a MOS technology is the symmetric nature of the MOS-FET. This has been utilized in logic circuits by the use of a pass transistor as shown in Fig. 26. The signal flow can be interrupted by the use of such a series switch when the switch is turned off. In the on state, the switch allows signal flow in either direction.

Static, or Ratioed, Logic—The generalized schematic of a static, or ratioed, logic gate is shown in Fig. 27A. It consists of a load device and a driver device or several driver devices that can be switched to perform a logic function. The ideal characteristic for a load device is a constant-current source. The typical device realization of the load in a PMOS or NMOS technology is an active depletion-type transistor with the gate and source connected as shown in Fig. 27B. The typical load line achieved by this technique is shown in Fig. 28. Note that the nonideal behavior is due to the change in the depletion-mode-device threshold voltage with source-to-substrate bias, commonly known as the body effect. This static circuit has a logic high level determined by the power-supply voltage, since the driver device is switched off in this state and only draws subthreshold current. The logic low level is determined by the ratio of the load current and the on current of the driver device. This ratio is a function of the threshold and aspect ratio of the load and the driver device defined as the beta ratio of the inverter. This is the most important design parameter for ratioed static logic.

$$\text{Beta ratio} = \frac{(\text{Width/Length})_{\text{driver}}}{(\text{Width/Length})_{\text{load}}}$$

The load current is constantly turned on in this circuit for the depletion load realization as in PMOS or NMOS. This causes a static power dissipation in one logic state. The static logic has the advantage of being totally asynchronous and requires no clock pulses. The power dissipation is drastically reduced by the use of CMOS technology (Fig. 27C). The load is a p-channel

* Reference 9.

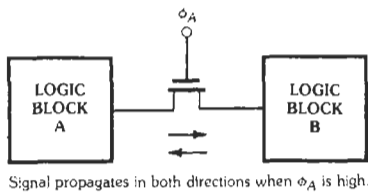


Fig. 26. Pass transistor in MOS logic.

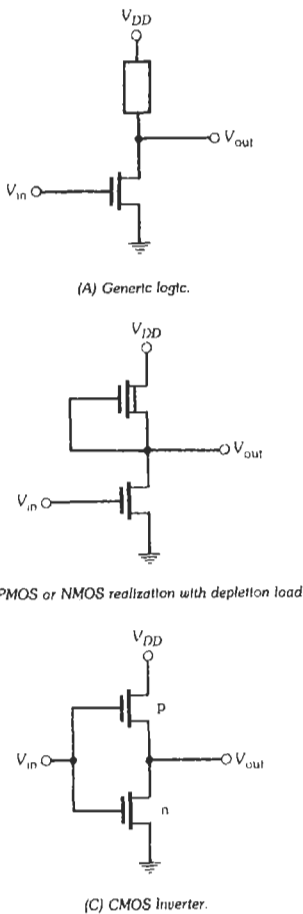


Fig. 27. Static, or ratioed, logic.

transistor, and the driver is an n-channel transistor. The CMOS inverter or logic gate draws power only during state transition. Note also that the CMOS static inverter has an inherently large ratio since one device is in the subthreshold regime in each state.

Dynamic, or Ratioless, Logic—The use of dynamic, or clocked, logic for MOS circuits is very attractive for low-power synchronous applications. Examples of dynamic logic gates realized in NMOS and CMOS technology are shown in Fig. 29. The use

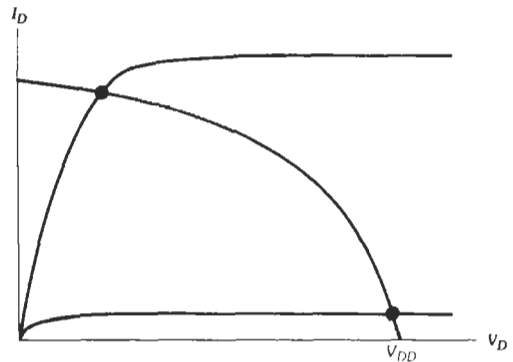


Fig. 28. Typical load line achieved with depletion load.

of two-phase clocks in synchronous systems is standard practice. The timing of these clocks is shown in Fig. 29A. The ϕ_A clock turns on the active load and precharges the output node to the supply voltage independent of the state of the input. The ϕ_B clock discharges the output only if the input logic combination is appropriate. This time-multiplexed logic is possible if the leakage rate through the clocked transistors does not discharge the output stage within a clock period. This scheme requires the generation of a two-phase clock for NMOS and PMOS circuits. However, for the CMOS circuit (Fig. 29B), single-clock operation is possible. The advantages of dynamic logic are the low power required and synchronous operation. It is assumed that the clock propagation speed on-chip is faster than the clock period. This requirement tends to limit the overall gate speed for a given technology.

Microprocessors and Minicomputers*

The implementation of logic in integrated circuits in the decade of the 1960s relied on providing building blocks that integrated unique functions. A major breakthrough in logic implementation was provided in the early 1970s by devices that controlled logic operations with a stored program. These devices made complete computing systems possible with a minimum number of ICs. A simplified block diagram of a microcomputer is illustrated in Fig. 30. The complete computing system consists of input/output (I/O), a central processing unit (CPU), and memory sections for the program and data. In early microcomputers these elements were all separate, and in most cases each element consisted of several different ICs. The function of each of these sections is as follows.

I/O section—The lines at the left of the I/O section shown in Fig. 30 connect the computer to the input and output devices, and receive instructions from the CPU.

* References 20 and 21.

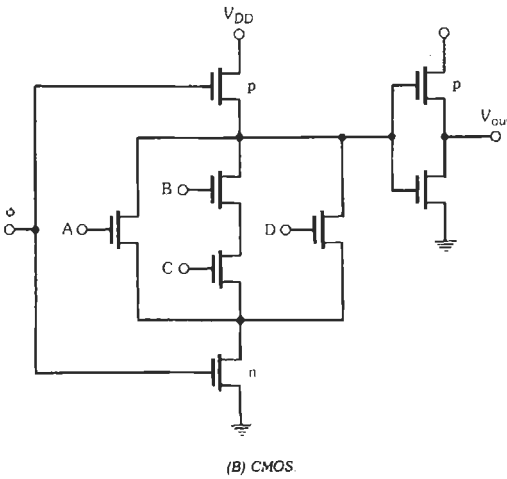
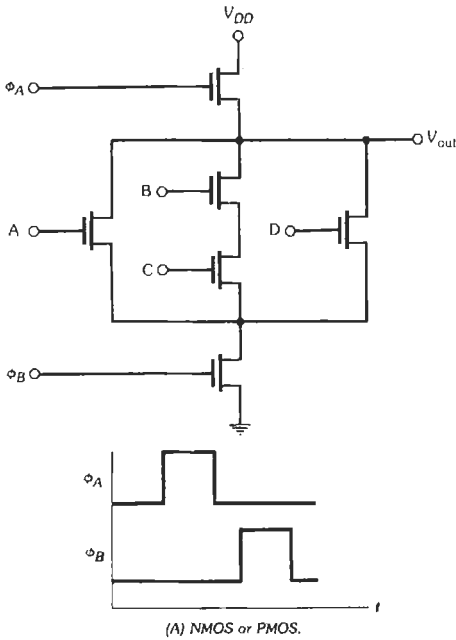


Fig. 29. Dynamic logic.

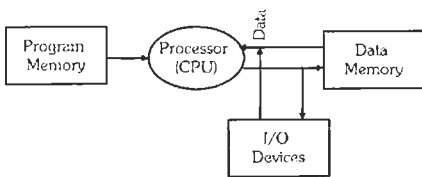


Fig. 30. Simplified block diagram of a computer system. (From H. Corporaal, *Microprocessors Architectures: From VLIW to TTA*. West Sussex, England: John Wiley & Sons, Inc., 1998.)

A simplified diagram of a typical I/O section is shown in Fig. 31. Input and output multiplexers (MPX or MUX), also known as data selectors, perform selection of the I/O devices. Output information is stored in the output buffers. The I/O register provides temporary storage during the transmission of information between the CPU and I/O sections.

CPU—The heart of this type of computer is the central processing unit, usually called the CPU. The CPU is the part of the computer that controls all the other parts. Designs vary widely but, in general, the CPU consists of the control unit, the arithmetic and logic unit (ALU) and memory (registers, cache, RAM and ROM) as well as various temporary buffers and other logic. The control unit fetches instructions from memory and decodes them to produce signals that control the other parts of the computer. This may cause it to transfer data between memory and ALU or to activate peripherals to perform input or output. A parallel computer has several CPUs that may share other resources such as memory and peripherals. The term “processor” is sometimes used in place of “CPU.” Fig. 32 shows a

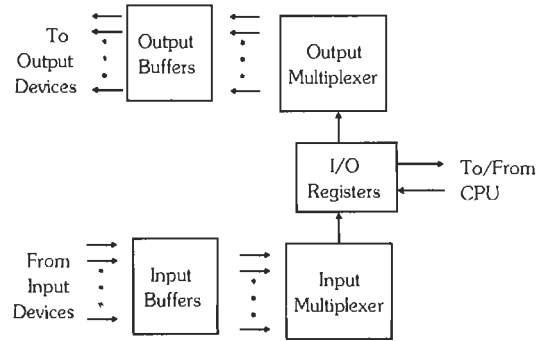


Fig. 31. Typical I/O section. (From R. Henle and B. Kuvshinoff, *Desktop Computers*. New York: Oxford University Press, Inc., 1992.)

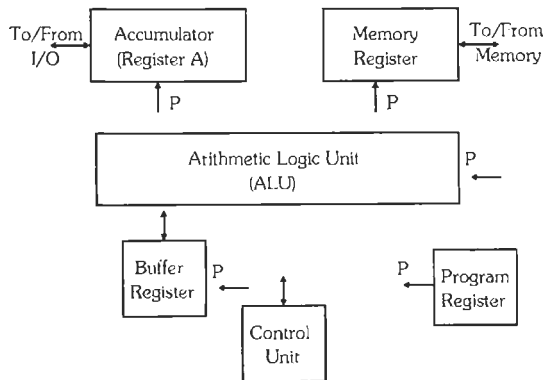


Fig. 32. The CPU. (From R. Henle and B. Kuvshinoff, *Desktop Computers*. New York: Oxford University Press, Inc., 1992.)

block diagram of the components of a very simple CPU. The connection between the arithmetic logic unit (ALU) and the accumulator and other registers determines the number of bits that the ALU can process in parallel. This is called the word length of the CPU. Word lengths of 32 and 64 bits are not uncommon. The ALU performs arithmetic operations, such as addition and subtraction, and logic operations, such as detection of equality. The CPU uses several registers as data registers, working registers, and scratch pad memory (A and B); these help in storing intermediate results and communicating with the I/O section and memory sections (M). A specific register (P) acts as the program counter and determines the operation sequence for the system. The program counter gives specific sequence instructions as it steps up during execution of the program. The control unit provides direction for the computer operation.

Memory—There are two sections for memory in a typical computer, one for data and one for instructions for what to do with the data, called the program. The I/O section and CPU contain several temporary storage registers for data in the form of digital information. The majority of the data, however, is stored in the Data memory. The details of memory functions will be discussed in the section on memory integrated circuits. The operations that transform the data have to be specified in a form that is directly “understandable” by the CPU. This language is called machine language, and a program written in this language is stored in the Program memory. When data is entered into the CPU, the set of operations defined by the program is performed on the data, and then the data is either returned to the Data memory, transferred to the I/O section, or both.

Available Microprocessors—A microprocessor is a computer whose entire CPU is contained on one (or a small number of) integrated circuits. The important characteristics of a microprocessor are the widths of its internal and external address bus and data bus (and instruction), its clock rate, its instruction set, and its available memory.

The first commercial microprocessor was the Intel 4004 that appeared in 1971. This was the CPU member of a set of four integrated circuits called the MCS-4, which was originally designed for use in a calculator

but was marketed as a “programmable controller for logic replacement.” The 4004 is referred to as a 4-bit microprocessor since it processed only 4 bits of data at a time (the word length). This very short word size is due mainly to the limitations imposed by the maximum integrated circuit density then achievable.

As integrated circuit densities increased with the rapid development of integrated circuit manufacturing technology, the power and performance of the microprocessors also increased. This is reflected in the increase in the CPU word size from the original 4 bits in 1971 to 64 and even 128 bits by the end of the century. The smaller microprocessors have relatively simple instruction sets, e.g., no floating point instructions, but they are nevertheless suitable as controllers for a very wide range of applications such as car engines and microwave ovens.

The Intel 4004 was followed with over ten generations of Intel microprocessors, with current devices being called Pentium III and Pentium IV class devices. Many other manufacturers of microprocessors have also become prominent. The larger, more recent microprocessor families have gradually acquired most of the features of large computers. As the microprocessor industry has matured, several families of microprocessors have evolved into de facto industrial standards with multiple manufacturers and numerous “support” chips including RAM, ROM, I/O controllers, etc. A single chip microprocessor may include other components such as memory (RAM, ROM, PROM), memory management, caches, floating-point unit, input/output ports, and timers.

In the early 1990s the Semiconductor Industry Association (SIA) began developing a roadmap for integrated circuit features and characteristics. This roadmap, now called the International Technology Roadmap for Semiconductors (ITRS), is updated each two years, and identifies the technological challenges and needs facing the semiconductor industry over the ensuing 15 years. This roadmap has been instrumental in identifying the desired performance of integrated circuits, particularly memories and microprocessors, and providing a set of challenges for the research community. A sample of the 2000 ITRS update for expected microprocessor performance for the next several years is shown in Table 12.

TABLE 12. MICROPROCESSOR CHARACTERISTICS*

Year of Production	1999	2001	2003	2005	2008	2011	2014
On-chip clock (MHz)	1,250	1,767	2,490	3,500	6,000	10,000	13,500
Technology (microns)	0.18	0.13	0.11	0.08	0.06	0.04	0.03
Chip size (cm ²)	1.70	1.70	2.14	2.35	4.68	5.36	6.15
Mtransistors/cm ²	7.0	14.0	25.6	46.7	115	284	701

* Ref. 2000 *International Roadmap for Semiconductors*.

A sampling of available microprocessor characteristics as of the year 2000 is shown in Table 13. These features compare well with the projections in Table 12.

It is clear from Tables 12 and 13 that the microprocessors of today will be significantly overshadowed by those of the future. In many respects, these microprocessors make the one-chip microcomputer a landmark development in computer technology because they reduce the computer to a small, inexpensive, and easily replaceable design component.

MEMORY INTEGRATED CIRCUITS*

The major advantages of semiconductor memories are that they can utilize the most advanced technology and they are required in very large volumes in all systems so that they provide the economies of scale in the learning process for the maturing of the technology in time. The large volume of memory integrated circuits thus acts as a catalyst in the timely development of the technology it uses and thereby allows increases in yield and decreases in unit costs, which in turn propels an increase in volume. The progress in the density of memories has been phenomenal because of this synergism between the system pull for more memory at low cost and the technology push via geometry scaling that allows improved bit density, improved performance, and decreased cost. Fig. 33 shows the actual and projected progression of memory integration vs. feature size over a 30-year time frame. There seem to be no fundamental physical roadblocks to prevent devices from being produced at least down to the 25nm feature

* Reference 22.

DRAM Production Level

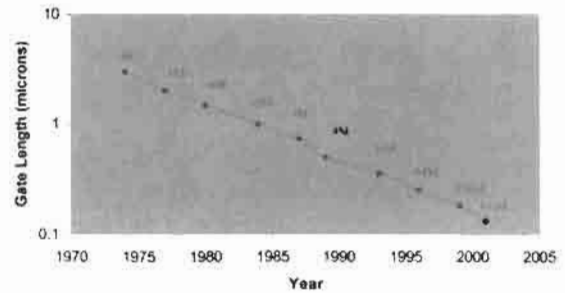


Fig. 33. DRAM density vs gate length and year, (after Iwai et al., *Microelectronics Engineering*, Vol. 28, 1995, p. 147.)

size. The slight slowing of the growth is an indication of the substantial technological hurdles that must be overcome for very small device geometries.

Semiconductor memories are classified according to their function as read-only memories (ROM) or random access memories (RAM). The use of serial access memories has been attempted, but has not been accepted due to both cost and performance limitations. Memories are also classified as static and dynamic. The static memories are able to hold data as long as power is applied (and sometimes even without power), while dynamic memories need to have each memory element periodically refreshed to retain data. The dynamic memories generally have smaller cell sizes, so are preferable for very high-density applications. The highest density memories are dynamic random access memories, called DRAM. Table 14 illustrates

TABLE 13. AVAILABLE MICROPROCESSOR FEATURES

	Celeron	Pentium III Xeon	K6-III	Athlon (K7)	Power PC 750CXe
Manufacturer	Intel	Intel	AMD	AMD	IBM
Transistors	7,500,000	28,100,000	21,300,000	22,000,000	N/A
CPU Speed (MHz)	800	1,000	450	1,200	700
Technology (microns)	0.25	0.18	0.18	0.18	0.18

TABLE 14. DRAM CHARACTERISTICS*

Year of Production	1999	2001	2003	2005	2008	2011	2014
DRAM size	256M	512M	1G	2G	6G	16G	48G
Technology (microns)	0.18	0.15	0.12	0.10	0.06	0.04	0.03
Chip size (cm ²)	1.32	1.45	1.59	1.74	2.05	1.91	2.68
Gbits/cm ²	0.2	0.42	0.68	1.46	2.97	8.99	18.1

* Ref. 2000 *International Roadmap for Semiconductors*.

the ITRS projections for these memory features over the next several years.

Externally, memory chips typically require the application of a binary address input to locate the information. A set of control or mode-setting signals is needed to tell the memory to read or write, and the memory circuit provides the desired information after a time interval that is called the access time for the memory. Memories are organized to provide either one bit of data at a time or one byte of data at a time. The choice of the data word is a system partitioning function. In addition, other system options such as power-down modes are common. Many different memories are available that fall into the categories described here. The various manufacturers' data books and applications books should be consulted for detailed information.

Read-Only Memories (ROM)

Read-only memories have information programmed into them during manufacture. They act as tables of data that can be accessed by the system at any time. This class of memories is extremely useful in the storage of programs, operating systems, fixed utilities for a system, etc.

The read-only memory function is realized by the use of single transistors as memory cells connected in an X - Y matrix as shown in Fig. 34. The gate of each transistor is turned on by an X address signal, and the current through the transistor is monitored through the Y address line. The programming of the memory results in the presence or absence of the transistor in a specific location. This may be done at various points in the fabrication process. The presence of current through the transistor is detected as a 1, and the absence of current is detected as a 0. Note that the memory cell in a ROM is a single transistor. Further, only a read function is required. For this reason, ROMs are the densest semiconductor memories.

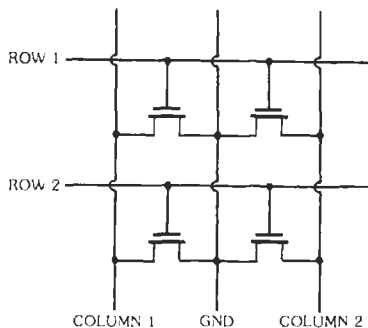


Fig. 34. A 2×2 ROM memory-cell array.

Nonvolatile Read-Only Memories

Nonvolatile read-only memories may be classified as "read-mostly memories." They are programmable after manufacturing. They are referred to as nonvolatile because they retain data even when no power is applied. Three subclasses of nonvolatile ROMs exist. The first is a programmable read-only memory (PROM), which can be electrically programmed. The second is an erasable and programmable read-only memory (EPROM), which can be programmed by electrical signals but erased only by exposure to ultraviolet light. The third is an electrically erasable and programmable read-only memory (EEPROM).

Generally, PROMs are constructed with bipolar technology. A simple 2×2 array of bipolar PROM cells is shown in Fig. 35. Programming is accomplished by applying a sequence of 20–30 mA current pulses to blow a fusible link.

Both EPROMs and EEPROMs are generally constructed by the use of MOS technology. The physical mechanism by which these memory cells store data is a storage of charge in the gate insulator of a MOSFET. This has been accomplished by two categories of transistor structures, shown in Fig. 36. The first structure incorporates multiple dielectric layers (Fig. 36A). The most common version of the multiple-insulator structure is the MNOS (metal-nitride-oxide-silicon) transistor. The oxide layer in contact with the silicon surface is very thin (2–3 nm). The charge is injected from the substrate into the nitride layer, by way of hot electron

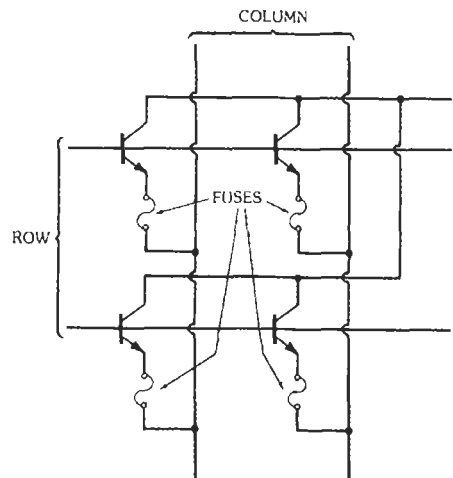


Fig. 35. A 2×2 array of bipolar PROM cells. (From G. D. Kraft and W. N. Toy, *Mini/Micro Computer Hardware Design*, Englewood Cliffs, NJ: Prentice-Hall, Inc., 1979; p. 97.)

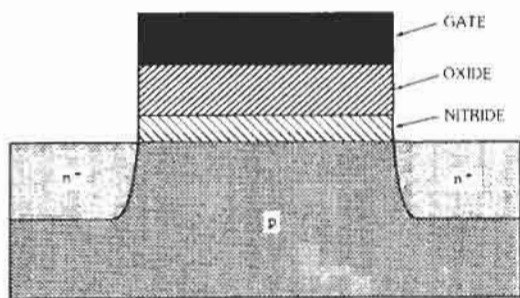
tunneling through the thin oxide. The charge is trapped in the nitride layer, which generally has a large trap density. The band gap of silicon dioxide is larger than that of nitride, and thus the charge is unable to flow back into the silicon. This presence of charge in the insulator alters the threshold of the transistor, thus changing the transconductance and causing a differential in the current through the transistor in the presence or absence of charge. This is read much the same way as the ROM cell.

The second structure that achieves the same purpose is the floating-gate transistor (Fig. 36B). The charge is injected into the floating gate from the silicon by the application of the proper electric field. The charge on the floating gate provides the change in threshold voltage. The floating-gate structure has been used in both EPROMs and EEPROMs, since the floating gate acts as an equipotential surface and can be biased with a reversed electric field to result in an erase function. The control of the floating-gate process is easier than the control of the MNOS process, and most of the non-volatile memories are fabricated this way. There are various subclasses of the floating-gate structure that are based on the biasing scheme or the injection mechanism used for the program and the erase modes.

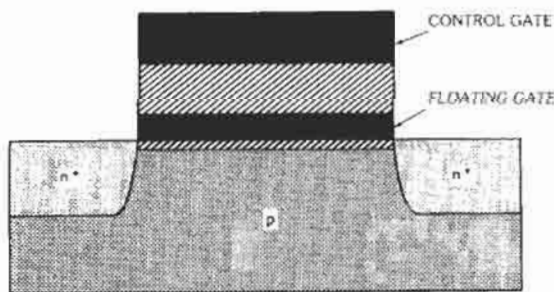
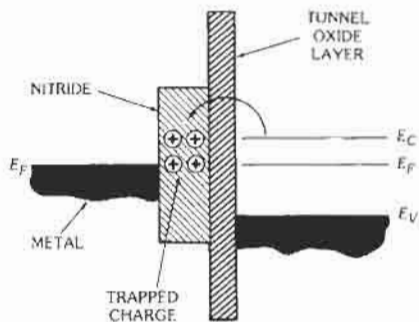
Random-Access Memories (RAM)

A random-access memory is one in which data may be written in or read out at very high speed. The bit or set of bits into which data may be written is chosen randomly in an *X-Y* matrix. There are two types of random-access memories, static and dynamic. Static memories retain data indefinitely while power is applied to the chip; dynamic memories retain data for a short period of time, after which they must be refreshed. The system overhead associated with the refreshing of the memories is offset by the fact that dynamic memories are four times denser than static memories at the same design rules.

Static RAM—Static RAMs are realized from a memory cell with a bistable circuit. A simple flip-flop is commonly used for the memory cell, as shown in Fig. 37. Note that this can be constructed with either bipolar or MOS technology. The bipolar versions have traditionally been faster than the MOS versions. However, the bit density of MOS sRAMs has been larger than for bipolar versions. The important elements of the cell are the load (M_1, M_2) and the cross-coupled driver



(A) NMOS.



(B) Floating gate.

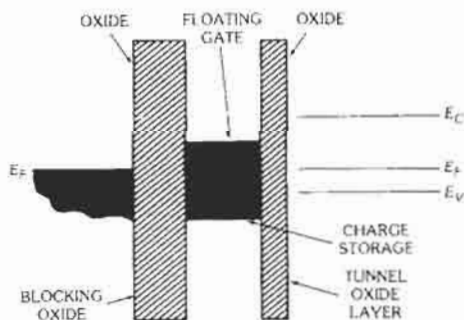


Fig. 36. Transistors for EPROM and EEPROM memories.

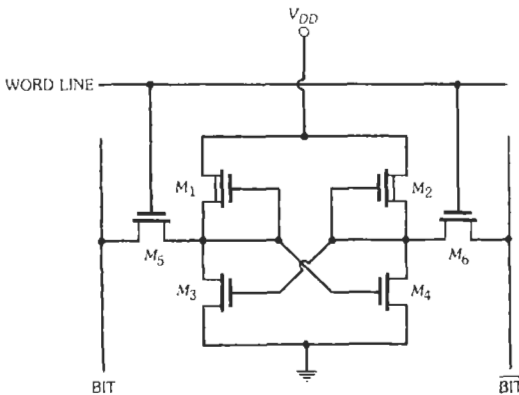


Fig. 37. A static RAM cell.

devices (M_3 , M_4) that form the bistable latch. The access transistors (M_5 , M_6) are used to address the cell. The major considerations in the design of an sRAM cell are power, speed, and data hold. The power in the cell is proportional to the load current. The load element has been realized as a depletion MOSFET in 4K-bit memories. However, it has been replaced by high-impedance resistors realized by undoped polysilicon, or by subthreshold current in a short-channel MOSFET, for higher-density memories. The sRAM cell with lowest power is realized from complementary MOS technology. These memories can hold data with extremely low power (50–1000 nW) and are being used with batteries as low-cost alternatives to nonvolatile memories in some systems. Because of the large number of elements and connections to the memory cell, static RAMs tend to be lowest in bit density per chip among semiconductor memories. However, the differential cell provides a large signal-to-noise ratio and very high-speed operation. Bipolar memories are all static and have demonstrated access times below 10 ns.

Dynamic RAM—Dynamic RAMs are constructed commercially with MOS technology only, and they represent the largest-volume, most widely used semiconductor memory. The memory is physically realized by the storage of charge on an MOS capacitor that is accessed by a MOSFET, as shown schematically in Fig. 38. The data bit is stored as a charge packet on the MOS capacitor. Since the MOS capacitor is not an ideal element and has nonzero leakage, it is discharged through these leakage mechanisms as a function of time. The design of dynamic RAMs is thus critically dependent on the reduction of leakage in the MOS capacitor. This requires the realization of very long carrier lifetimes in the semiconductor and very high crystalline quality. The signal-to-noise ratio depends on the ratio of the storage capacitor to the capacitance of the digit or bit line on which the charge is dumped when the cell is accessed. As the bit density per chip grows, this ratio is degraded, since the cell area is

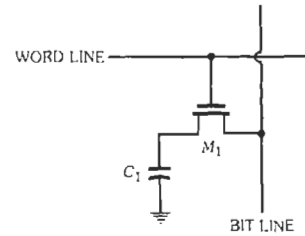


Fig. 38. The one-transistor/capacitor dRAM cell.

decreased and more cells are added on the bit line. It is thus important to store as high a voltage as possible on the cell capacitor and increase its capacitance per unit area as much as possible. The sensing of this small signal is one of the most difficult circuit design problems, and rather elaborate dynamic circuit techniques have been developed for the purpose.

New Developments

The growth of all the above-described memories following the scaling of MOSFET technology is on an evolutionary path. The state-of-the-art memories in the development stage include up to 2 gigabits of ROM and 2 gigabits of dynamic RAM. The memory market for all sectors except static RAM has been dominated by MOS technology and will continue to be so. However, the bipolar memories have had a strong edge in speed in the past. One-micron MOS static RAMs are speed-competitive with bipolar memories. However, the bipolar memories are significantly more radiation hard and will tend to have major applications in hostile environments.

LINEAR INTEGRATED CIRCUITS*

Unlike digital integrated circuits, which respond to and produce two-state logic signals, linear ICs give an output signal that can be made to vary linearly with respect to a varying input signal. Since linear ICs can be used in a variety of applications and new designs appear regularly, standardization of circuit elements as in the case of digital ICs is impractical. Several important families of linear ICs have evolved. They include the device categories described below. The various manufacturers' data books and applications books should be consulted for detailed information.

Differential Amplifier

The basic differential amplifier consists of two identical input transistors connected to respond to the difference between two input signals while simultaneously

* References 23, 24, 25, 26, and 27.

blocking the identical part of these signals. This common-mode rejection is useful in noisy environments. The amplifier may consist of a single stage or two stages and usually has a low-impedance output stage. The output is also differential. These amplifiers are typified by a wide bandwidth (dc to VHF), moderate gain (less than 1000 times) and moderate common-mode rejection ratio (CMRR) (50 dB). Internal feedback is provided for linearity control, and external feedback is difficult to apply without severely degrading the bandwidth. Applications include linear amplification, mixers, product detectors, amplitude modulators, frequency multipliers, voltage-controlled oscillators, and disk or tape memory read amplifiers.

Operational Voltage Amplifier

The basic operational amplifier, or op amp, is a dual-input differential amplifier followed by one or more direct-coupled gain stages and a low-impedance output stage. The typical op amp exhibits very high voltage gain at dc (80 to 120 dB) and a unity-gain bandwidth of 1 to 20 MHz. The op amp may be internally compensated, i.e., 100% negative feedback can be applied externally without oscillation problems, or pins may be provided to tailor the compensation for the application. Pins may also be provided to adjust the input offset voltage. Op amps are commonly made with a bipolar process but may have JFET or MOSFET input stages. Complete CMOS designs are also available.

Any op amp can be connected as shown in Fig. 39 to provide a linear voltage gain equal to R_f/R_{in} , provided the op amp is ideal and properly compensated. The ideal op amp would have the following characteristics:

1. A negligibly small differential input voltage can produce any desired output.
2. The input currents are negligibly small.
3. The output impedance is zero.

Application of these rules to negative-feedback circuits such as in Fig. 39 immediately leads to the ideal transfer function. No existing op amp has these ideal characteristics, but many come close to meeting one or more of them. A variety of op amps are available with

emphasis on perfecting different properties such as low noise, low input-voltage offset, low input current, high slew rate, single-power-supply operation, and large power-supply-range operation.

The op amp was originally designed to perform such mathematical operations as integration, differentiation, comparison, summation, multiplication, and others. Though op amps are still used in these applications, complete ICs performing these functions are available. The use of various types of external feedback networks permits a vast array of additional applications. Accordingly, the op amp is the most versatile of linear ICs, and its circuit applications include analog-to-digital conversion, peak detection, averaging, function generation, oscillator and pulse circuits, biological function amplification, voltage-controlled oscillator circuits, automatic gain control, voltage comparators, active filters, precision diodes, sample-and-hold circuits, and numerous others. Many of these complete circuits in turn are available in IC form.

Fig. 40 summarizes some of the elementary op-amp circuit configurations and gives ideal formulas that describe the circuit performance. The op amp is used in a number of modes besides the basic operational one, such as a high-impedance buffer to minimize loading in RC filters or to compress the nonlinearity of other devices by the op-amp open-loop gain for voltage comparison.

Other Amplifiers

The operational transconductance amplifier (OTA) provides an output current in response to a voltage input. The OTA bias current controls the transconductance, and considerable control over the performance of an OTA is possible. The operating characteristics make the OTA usable in a wide variety of circuits such as multiplexers, sample-and-hold circuits, gain controls, modulators, multipliers, comparators, multistable circuits, and (with the addition of an output stage) high-gain op amps.

The current-differencing amplifier, or "Norton" amplifier, provides a voltage output proportional to the difference of two input currents. It is especially suited for single-power-supply operation without loss of common-mode range. Again, a large variety of circuit functions can be obtained.

Voltage Reference

Although classified under linear ICs, the voltage reference is a highly nonlinear device and provides a constant output voltage almost independent of input voltage, load current, temperature, and time. Output voltages range typically from 1.22 volts to 10 volts with a precision of 0.05% to 5%. Temperature coefficients range from 100 parts per million down to 0.5 ppm. Output current capabilities are less than 10 mA. Both fixed-output-voltage and programmable-output-voltage devices are available.

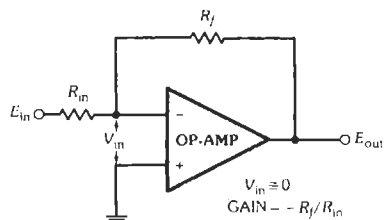
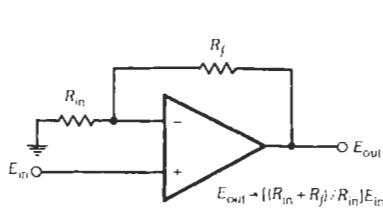
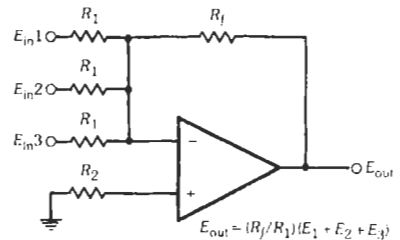


Fig. 39. Inverting dc amplifier.



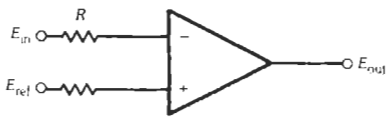
(A) Dc amplifier (noninverting).

$$E_{out} = (R_f/R_{in} + 1)E_{in}$$

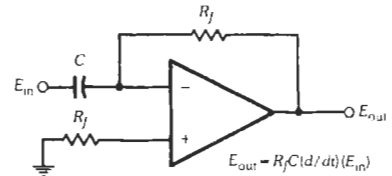


(B) Summing amplifier.

$$E_{out} = (R_f/R_1)(E_1 + E_2 + E_3)$$

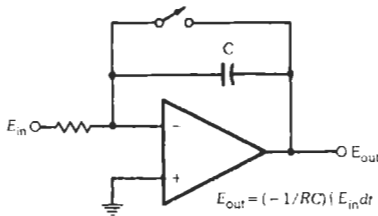


(C) Analog-to-digital converter.



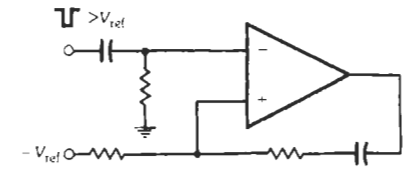
(D) Differentiator.

$$E_{out} = R_f C (d/dt)(E_{in})$$



(E) Integrator.

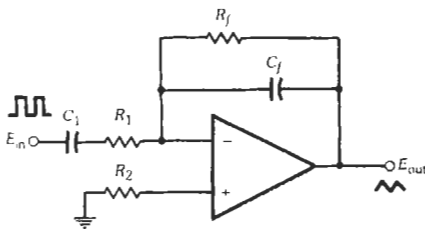
$$E_{out} = (-1/RC) \int E_{in} dt$$



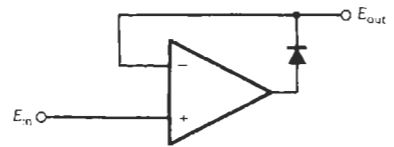
(F) Monostable multivibrator.

$$V > V_{ref}$$

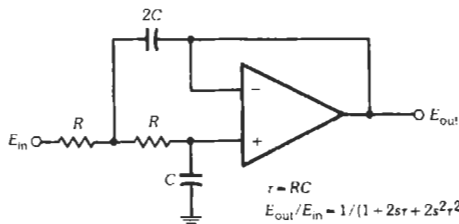
$$-V_{ref}$$



(G) Sweep generator.



(H) Precision diode.



(I) Low-pass filter.

$$\tau = RC$$

$$E_{out}/E_{in} = 1/(1 + 2s\tau + 2s^2\tau^2)$$

Fig. 40. Simple ideal op-amp circuits.

Reference diodes also come in IC form and simulate the characteristic of a zener diode, but with a much sharper breakdown characteristic.

Some applications are amplifier biasing for temperature independence, constant-current-source circuits, level detectors, and low-voltage regulators.

Voltage Regulators

Bipolar voltage-regulator ICs incorporate a voltage reference and a sense amplifier and maintain the output voltage at a value almost independent of load, input voltage and ripple, and temperature. Both fixed- and programmable-output-voltage units are made. Voltages up to 50 volts can be regulated. Current capability is usually on the order of tens of milliamperes, and the IC is intended as the driver for large-pass transistors that can accommodate the higher currents and operating temperatures that occur in power supplies. External components for current limiting, noise reduction, and compensation are desirable additions in the design of a complete power supply.

Comparators

The voltage comparator is a differential amplifier design with a small delay time between the application of a differential input signal and the output transition. The output swing is made compatible with TTL inputs. Emphasis is placed on low-input bias current and its offset and low-input offset voltage, and the usual amplifier specifications are of lesser importance. Common applications are high-speed analog-to-digital converters, fast zero-crossing detectors, tape- and disk-file read channels, and differential line receivers.

Special-Purpose Linear Integrated Circuits

Since any circuit that can be assembled from discrete semiconductor components can almost always be duplicated or simulated with monolithic IC technology, a great many special-purpose ICs exist. These include custom designs for use in various kinds of commercial products as well as a wide range of off-the-shelf numbered units. The commercial units include such devices as music synthesizers, phase-locked loops, tone decoders and encoders, function generators, programmable filters, FM and video demodulators, IF detectors, subcarrier regenerators, AGC, FM stereo demultiplexers, AM-receiver functional blocks, timers, audio and power amplifiers, voltage-to-frequency converters, analog-to-digital and digital-to-analog converters, instrumentation amplifiers, and sample-and-hold circuits.

Though many of the ICs are dedicated to a specific application or function, the addition of a few external components can often result in a wide range of additional applications. The list of special-purpose linear

ICs grows constantly, and the design engineer is well advised to consult current manufacturers' data and application literature.

Miscellaneous Linear ICs

Other useful building blocks are timers, transistor-diode arrays, digitally controlled analog pass-gates, optical isolators, bar or dot LED-display drivers, and analog shift registers.

TRENDS IN INTEGRATED CIRCUITS

Scaling and Miniaturization

The explosive growth of integrated circuits has been fueled by the ability to scale the minimum lithographic dimensions of an integrated circuit. This results in a threefold advantage:

1. Increase in density of circuit elements per chip
2. Increase in circuit performance due to increased device gain and reduced load capacitance
3. Decrease in cost per function, which provides economic incentive

The most significant improvement in integration complexity has been achieved in MOS circuits. An exponential growth of the number of circuit elements in time has been noted. This has been made possible by the simple scaling laws that relate the scaling of vertical and lateral dimensions to the scaling of doping and voltages. Table 15 shows three sets of scaling laws that have been used to various degrees. The first assumes that the electric field in the device must be held constant. The second is based on complying with currently set voltage standards, and the third proposes a change in the voltage standards only when the electric field is high enough to cause problems as the result of secondary effects. Another scaling theory has been proposed on the basis of the off-state switching behavior of MOS devices. Scaling of bipolar devices does not proceed along scaling laws because the base width that is the critical dimension is much smaller than all other dimensions. Scaling of bipolar technology is aimed at reduction of parasitic device capacitance.

Scaling of ICs has proceeded unabated for the past two generations, despite several predictions along the way that the "end of scaling is imminent." These predictions proposed unbreakable roadblocks along the way at minimum IC feature sizes of 1.0 micron, 0.5 micron, 0.25 micron, and smaller. In each case, ingenuity and hard work have overcome critical roadblocks. At the present time, manufacturing is occurring at 0.15-micron feature sizes, and 0.1-micron feature sizes are fully anticipated in manufacturing in 2005. The "end of scaling" is still being predicted, but at feature sizes in the .02 micron range.

Integration complexity in ICs has also benefited greatly from improvements in process equipment and

TABLE 15. DEFINITION OF SCALING LAWS

Scaling Law	Constant Field	Constant Voltage	Quasi-Constant Voltage
Dimensions (λ)	λ	λ	λ
Gate Oxide (λ_O)	λ	$\sqrt{\lambda}$	λ
Doping (λ_N)	λ	λ	λ
Voltage (λ_V)	λ	1	$\sqrt{\lambda}$

Full-scale ($\lambda = 1$) values of $L = 3 \mu\text{m}$, $t_{ox} = 500 \text{ \AA}$, $N_A = 2.5 \times 10^{15} \text{ cm}^{-3}$ are used.

process understanding over the past decades. This has allowed much larger silicon area to be used for a single IC, while still attaining acceptable yield and therefore acceptable cost. This trend to bigger chips is continuing even with reduced feature size to provide the continuing increases in functionality demanded by expanding applications. The decision between continued reduced feature size and increasing chip area is expected to be totally cost based.

Image-Sensing ICs

The application of solid-state devices to high-resolution imaging began in earnest with the invention of the charge-coupled device (CCD) at Bell Laboratories in 1969. Many approaches to the solid-state imaging problem have been undertaken. These approaches may generally be divided into three distinct areas: (1) memory arrays, (2) charge injection devices (CIDs), and (3) charge-coupled devices (CCDs). Memory arrays, as the name implies, are simply solid-state memory devices that have been packaged in a format that allows an image to be placed on the area that is the usual memory array. The primary advantage of the memory-array approach to solid-state imaging is the fact that these devices are X - Y addressed as well as being process compatible with other memory components. The primary disadvantage of this approach to solid-state imaging is the fact that memory-cell architectures as normally designed for digital applications have relatively high-capacitance sense nodes, and therefore they have relatively low sensitivity and high noise, which precludes these devices from being used for most high-performance applications.

Charge injection devices employ construction similar to memory arrays in that they are X - Y addressed, but the charge detection scheme is different. The device in its modern form consists of an array of charge-coupled pairs of capacitors as illustrated in Fig. 41. The readout is achieved by applying a capacitively coupled pulse to the row while a certain column is selected by lowering its potential.

The most widely explored architecture for solid-state imaging is the CCD. These devices store charge in a potential well and transfer the charge from one

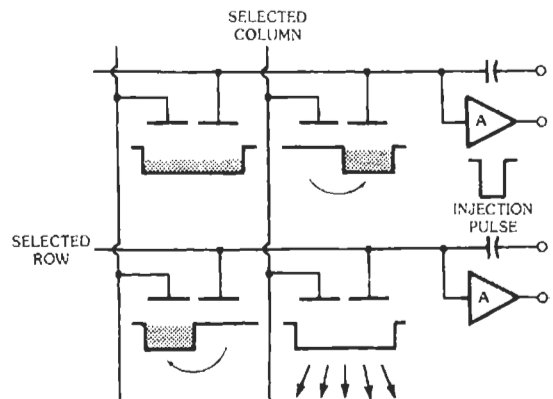


Fig. 41. Basic operation of a charge injection device.

electrode area to an adjacent electrode area by application of voltages with appropriate time and phase relationships, as illustrated in Fig. 42. The initial CCD devices employed three separately clocked phases to store and transfer charge in an MOS structure. The CCD concept has been extended, however, to include four-phase, three-phase, two-phase, and virtual-phase (one-phase) architectures, with significant advantages for devices with fewer gates. The construction of these different structures is illustrated in Fig. 43.

Solid-state imaging devices provide the combined advantages of high performance, low power, high reliability, and ease of use. These features are available in a wide variety of imagers fabricated as memory arrays, CIDs, and CCDs, and in a wide spectrum of formats and performance ranges. These features will combine to make solid-state imaging the mainstay of advanced imaging systems for future applications.

Speech-Synthesizer IC*

Speech is generated with these chips by the excitation of a time-varying digital filter. The excitation and

* Reference 28.

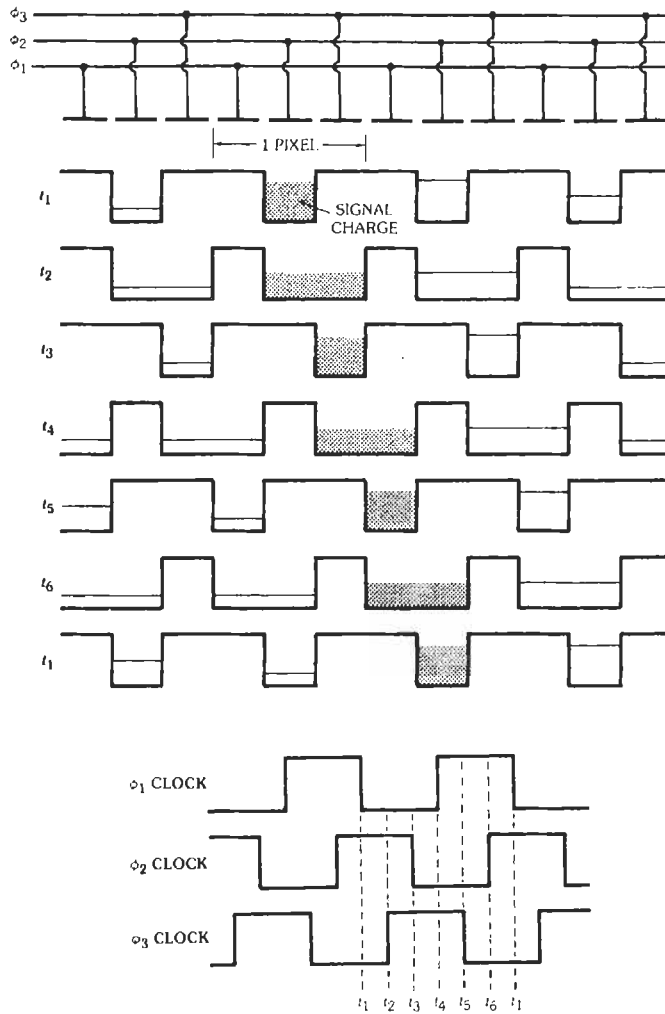


Fig. 42. Basic operation of a charge-coupled device.

filter parameters are stored in ROM, EPROM, RAM, or disk memory or are generated by a program. Overall control is provided by a microprocessor. The process is entirely digital up to the analog-to-digital converter, which provides the analog signal to drive a speaker.

An alternative to the variable-filter approach is to store digitized speech and sample the data at an 8-kHz rate, which leads to a digital data rate of about 100 kHz. Logarithmically compressed amplitude data could be used, analogous to digital telephone systems, which results in a rate of 64 kilobits/second with very good quality. The time-varying filter techniques provide acceptable speech quality but at a much lower digital input data rate, down to an average rate of 1200 bits/second for a ten-pole filter derived from a linear prediction model of speech.

The low data rates are made possible because of the redundancy in speech and by using a simplified simu-

lator of the human speech-generating system. The vocal tract is simulated by a dozen or so connected pipes of differing diameter, and the excitation is represented by a pulse stream at the vocal-chord rate for voiced sound or by a random noise source for the unvoiced parts of speech. The reflection coefficients at the junctions of the pipes can be obtained from a linear prediction analysis of the speech waveform, and the names associated with the ICs, such as LPC or PARCOR, refer to the particular analysis used to obtain the coefficients for the equivalent electrical digital filter.

System in a Package (SIP) and System on a Chip (SOC)

As ICs continue to become pervasive in consumer applications, the need for smaller size and lower cost

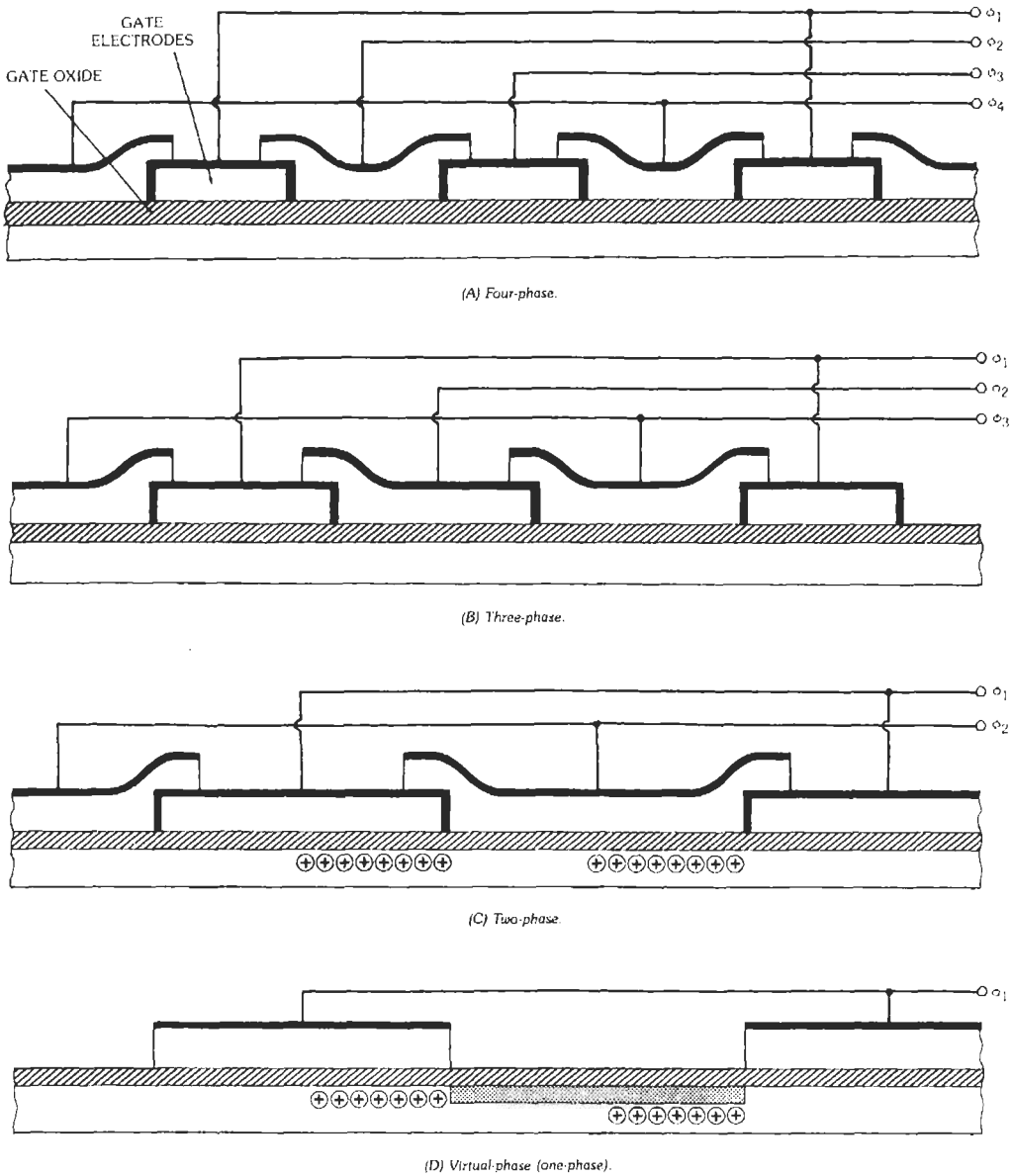


Fig. 43. Alternate structures for CCDs.

final products drives the need for fewer and fewer individual components to be assembled. This trend naturally leads first to a System in a Package (SIP) and then to a System on a Chip (SOC). In the SIP a single package will host a set of different computation chips (microprocessor, controllers, DSPs, I/O unit, RF modules, etc.), memory blocks, and interface elements such as capacitors, resistor, and inductors to perform a total system function. It is reasonable to predict that these components will be designed by different design teams in different companies to lower the costs and

meet the stringent time-to-market requirements. To freely mix and match these different components will demand that each of them be "insensitive" to its environment at every level of abstraction.

In the SOC the same demands that were described above for a SIP will be placed on the individual chip. These demands will result in the individual chip having many more complex elements and using very complex processes in its manufacture. This will be a major challenge to the IC design and processing teams in the next few years.

Micro-Electro-Mechanical Systems (MEMS)

Micro-Electro-Mechanical Systems (MEMS) are integrated mechanical elements, sensors, actuators, and electronics that are fabricated using microfabrication technology similar to that used for integrated circuits. In most cases MEMS exist on a common silicon substrate; however, MEMS can be combined as discrete microcomponents in a single package. While MEMS electronics are fabricated using IC process materials and process sequences, many times the micromechanical components, sensors, and actuators are fabricated using compatible "micromachining" processes that deposit new layers or selectively etch away parts of the silicon wafer to form the mechanical and electromechanical devices. MEMS devices that are manufactured using batch fabrication techniques like those used for ICs provide the outstanding levels of functionality, reliability, and sophistication that have been demonstrated on ICs, with the coupled advantage of relatively low cost.

MEMS have emerged as significant contenders to satisfy the needs of several major existing and emerging markets. Probably the most familiar application to the IC industry is in the MEMS accelerometers that are quickly replacing conventional accelerometers for air-bag deployment systems in automobiles. The conventional approach uses several bulky accelerometers made of discrete components mounted in the front of the car with separate electronics near the air bag. MEMS have made it possible to integrate the accelerometer and electronics onto a single silicon chip at a cost that is a factor of 10 below the mechanical counterpart. A second application that is receiving significant publicity is the TI DMD device that is used in projection displays. The high resolution and brightness provided for this application, combined with the low cost and ease of transport and modification of digital data, promises to revolutionize the movie and commercial display markets. A third MEMS application that has a major market potential is in optical switching. The recent emphasis on expanded optical communications to meet the requirements of wireless communications and the Internet is driving a wave of need for high-density, low-power, low-cost optical switches. This need is being met by the micromechanical optical switch, which is expected to serve a major market in the future. In addition to the applications noted above, MEMS are appearing in a multitude of special applications in such diverse fields as health care, metrology systems, defense, and space science applications.

APPLICATION SPECIFIC INTEGRATED CIRCUITS*

One of the fastest growing semiconductor market sectors is application specific integrated circuits (com-

monly referred to as ASICs).† The increasing use of ASICs continues to be market driven by the rapid advance of technology, which has shortened the useful life of most electronic products as improved and less expensive products make existing models quickly obsolete.‡ Furthermore, because most of the profit from electronic systems is earned early in the product life cycle, the penalty for being late to market is severe. The time required to design, manufacture prototype systems, and perform design revisions with ASICs is much less than with full-custom integrated circuits. ASICs are cost-effective solutions for all but the very highest volume products in which the smaller chip size of full custom design and layout can offset the advantages of ASICs if the product life is sufficiently long. The next paragraphs describe the major classes of ASICs, and then anticipated future trends in ASICs are discussed. The use of ASICs rather than "off-the-shelf" standard logic parts is also influenced by ASIC benefits such as reduced cost and size, and improved reliability resulting from higher levels of integration so that fewer system components are required. Systems based on proprietary ASICs are also much more difficult to replicate by competitors.

Gate Arrays

Gate arrays are the largest single ASIC market segment.§ They offer a good balance of cost, capability, time to initial silicon, and performance. A gate array is a chip from a master slice in which the active devices are regular arrays of digital logic "base cells" that can be interconnected into numerous types of simple and complex logic gates or sequential elements. The designer selects elements from a set of functional circuit blocks called the "macro library." Master slices are semiconductor wafers that have been fabricated through much of the IC process so that all active devices are present, though not yet connected by the metal layers. An inventory of wafers is built to provide manufacturing economies of scale; then the custom interconnect layers are patterned as orders are received.

There are two primary gate-array architectures: channelled and sea-of-gates. Channelled architectures allow the use of selected rows and/or columns of base cells to use those regions as channels for the horizontal and/or vertical routing of signals and power supplies (Fig. 44). Triple-level metallization process technology has enabled the sea-of-gates architecture to begin replacing channelled gate arrays for higher-density applications that require more gates. With three interconnect layers, routing over utilized base cells is possible without setting aside channels. "Gate utilization" is a term that describes the percentage of

† References 29, 30, and 31.

‡ References 32 and 33.

§ Reference 29.

* This section was contributed by Mark G. Harward.

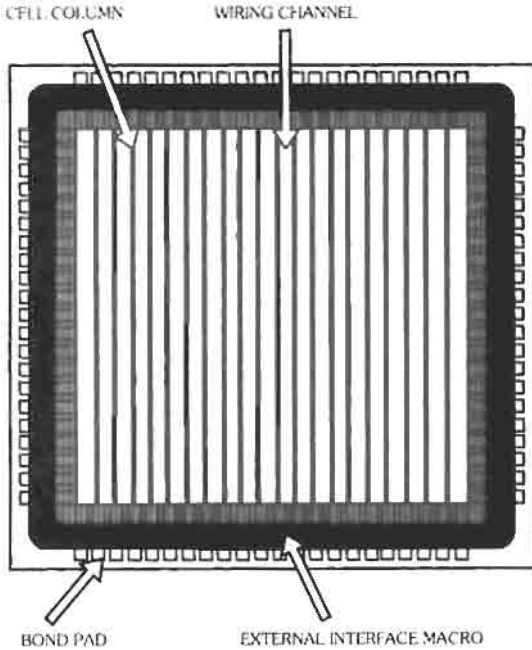


Fig. 44. Gate array—channelled architecture.

the available gates that can be utilized. Channelled gate arrays typically have approximately a 50% best-case gate utilization, including the cost of base-cell region lost to the channels. Sea-of-gates gate arrays with three levels of metal can achieve as high as 80% utilization to allow a more complex design to be done in a smaller chip size.

The gate-array manufacturer spends a great deal of engineering effort to devise the best possible base-cell design and layout. The base cell is required to efficiently implement as many different logic functions as possible by itself and in combination with other base cells (Fig. 45). As designs become more complex, there has been a need to incorporate more and more

memory onto many gate-array designs. There are two means for this; the most common today is to use "memory compilers" provided by the vendor. Compilers use the base-cell arrays to implement memory and the support circuitry (address decoders, input/output, read/write paths) as efficiently as possible. The other technique is actually a hybrid combination of standard-cell and gate-array methodologies in which a complex standard-cell block of memory is embedded in the gate-array region that has the base cells removed. The embedded-memory method provides the advantage of custom memory design density but requires a full manufacturing process cycle and is more expensive. Interestingly, the base-cell design is greatly affected if embedded memory is the primary means of incorporating RAM, as the base cell then does not have to efficiently implement compiled memory cells.

A significant fraction of gate-array designs are input-output-pin limited rather than gate limited. This means that many users find their selection of gate array to be based on the number of I/Os rather than the utilizable gates of the chip. Recently, gate-array vendors have introduced families of chips to address this need that have a higher ratio of pads to gates, achieved either by using tighter pad pitch or by placing two rows of pads at the chip perimeter.*

The gate-array design flow is similar to that of standard cells, which are described below. The keys to correct designs are the same: thorough and correct description of the chip requirements, correct capture of the design, verification of the design with simulation of logic functionality and timing delays, generation of a set of test vectors sufficient to verify adherence to specification of the completed silicon, and generation of the test program.

Standard Cells—Designers select circuit blocks from a predefined library of standard cells that are custom designed and characterized for maximum density and performance. Automated placement and routing tools are used to define the circuit layout. There are no

* Reference 33.

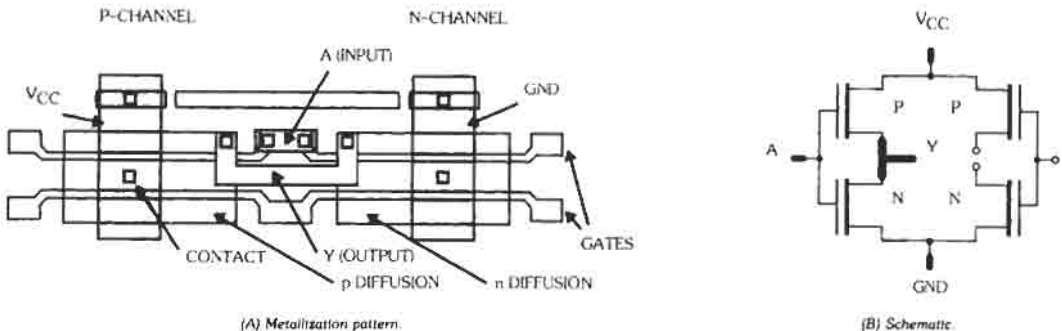


Fig. 45. CMOS gate-array base cell (configured as a 1 × drive inverter).

master slices prefabricated, as all levels must be generated to make the standard cells. Thus, of ASICs, standard cells cost the most and take the longest time to get initial silicon. However, the density, complexity, and performance of standard cells are typically greater than those possible with other ASICs because of the high complexity of the library elements compared to the gate-array macro library and nonoptimal architectures of PLDs.*

An example of a standard-cell chip arrangement is shown in Fig. 46. For maximum area efficiency, careful floorplanning is done manually or automatically by the placement and routing CAD tools. The idea is to group similar circuit blocks and select their relative locations to reduce the overhead area required for interconnections between the blocks of circuitry. The flow of data from input to output is also considered in order to minimize the wire lengths and the resultant parasitic capacitances that degrade circuit performance.

The library of standard cells typically includes "hardwired" cells (cells laid out and internally interconnected by the vendor) and "soft macros" (read as software-generated macros), which are circuit blocks of higher complexity that the computer-aided-design (CAD) tools generate (compile) to match the designer's specification. These compiled macros include complex circuit blocks such as single and dual I/O port static random access memory (sRAM), first-

in-first-out (FIFO) data buffers, shift registers, and many of the standard TTL family of MSI logic blocks implemented in CMOS.

An example of a typical standard-cell design flow is depicted in Fig. 47. There are many tasks required to assure properly functioning circuits. The foundation of a good design is always laid with a thorough design/product specification based on the finalized system requirements. The biggest relative weakness of standard cells compared to other ASIC solutions is the time required to receive prototype silicon from the vendor because all fabrication levels must be processed. Thus, it is most critical for standard-cell designs to be correct on first pass. Interestingly, multiple passes are more frequently the result of changes in system requirements or specifications than of design errors. However, design of functionally correct circuits that meet all specifications on the first pass is best achieved by rigid adherence to the vendor's design flow. This includes much simulation of logic functionality and timing to verify the design before the design is submitted for fabrication.

Programmable Logic Devices (PLDs)—Programmable logic devices are digital integrated circuits completely fabricated including the interconnection levels; the end-user configures the PLD using the on-chip programming circuitry and programming elements to specify the logic-circuit functions and the programmable interconnection points. Several types of programming elements are used by the many PLD manufacturers, including: fuses, EPROM, EEPROM, static RAM, and antifuses. Non-PLD ASIC families are collectively termed "mask-programmed" to differentiate the configuration methods. The history of programmable logic device architectures is nicely described in Roger Alford's *Programmable Logic Designer's Guide*.† This reference also gives a more detailed description of each PLD architecture than it is possible to present here. There are several classes of PLDs that must be understood by the designer considering a PLD solution. These include PLA (programmable logic array), PAL™ (programmable array logic), and FPGA (field programmable gate array).

User-programmable PLAs were invented in 1975 by Napoleone Cavlan at Signetics. This combined a programmable fuse element with National Semiconductor's mask-programmed logic-array architecture of a programmable AND array combined with a programmable OR array, which allowed digital designers to form Boolean sum of product (SOP) logic terms. An example of the PLA architecture is shown in Fig. 48. The PLA architecture allowed the flexibility of programming both the AND and OR arrays, but this resulted in large die area (i.e., high cost), which required a wide-body DIP package, and in relatively long propagation delays of 50 ns. The other factor that

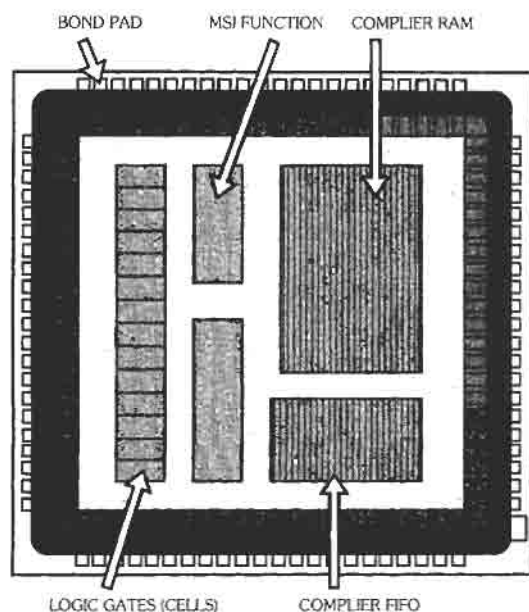


Fig. 46. Example of standard-cell chip arrangement.

* Reference 37.

† Reference 36.

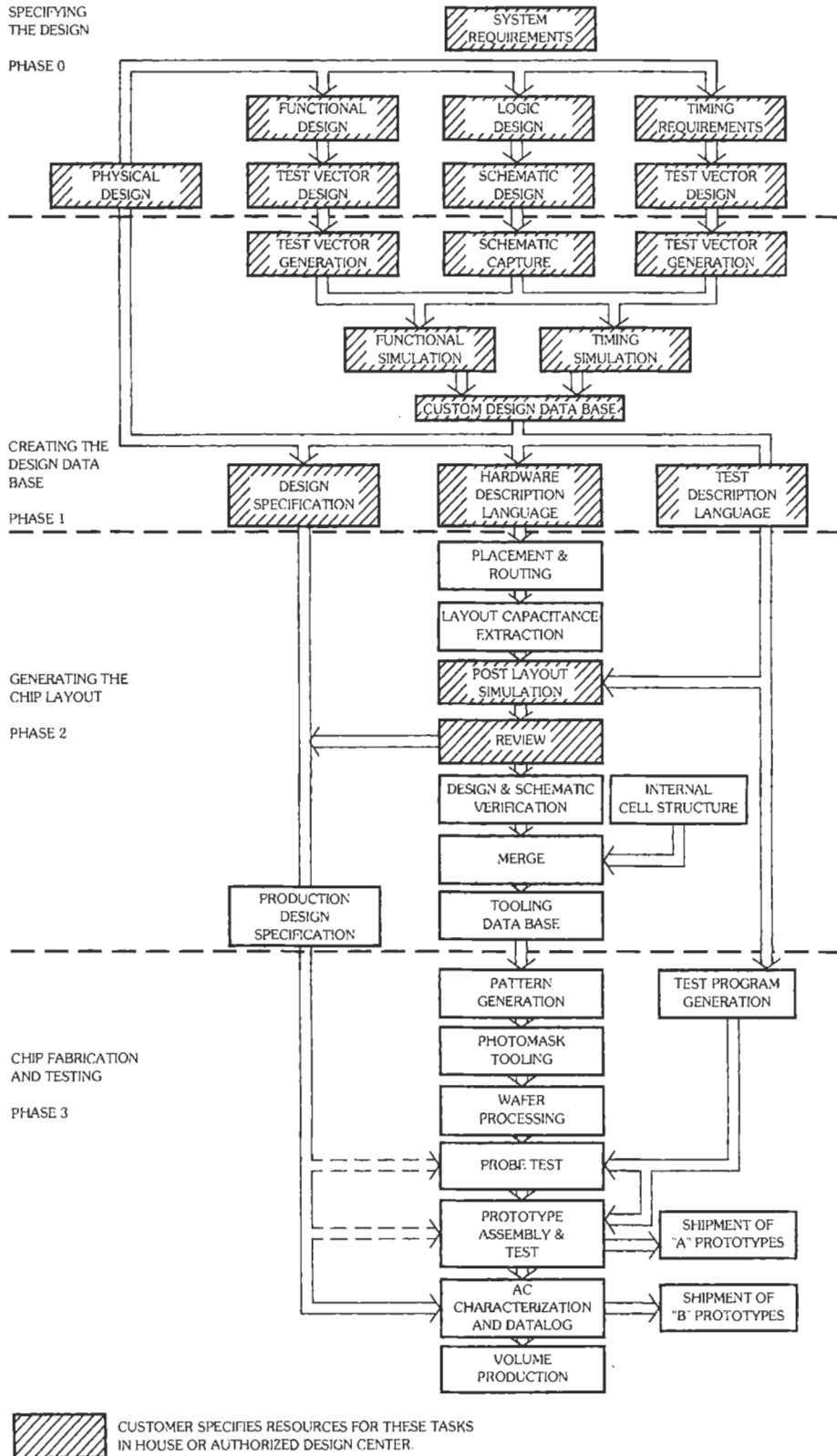


Fig. 47. Typical standard-cell design flow.

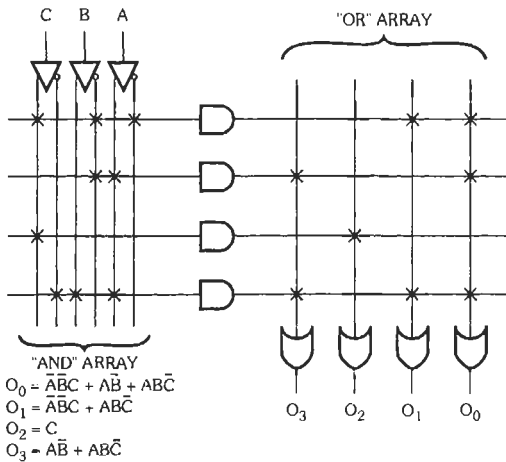


Fig. 48. PLA™ architecture.

retarded market acceptance of the PLA was the lack of software tools to ease the PLA user's design entry task.

John Birkner invented the PAL™ at Monolithic Memories Inc. (MMI) after much experience with the problems of digital system design using discrete components and study of the PLA weaknesses. The basic PAL architecture still had an AND array followed by an OR array with a fuse programming element, but with the fundamental difference of reduced flexibility by allowing programming only in the AND array. Birkner decided this trade-off was worthwhile based on a study of his digital system design work, which showed that a limited number of product (AND) terms input into each OR gate was acceptable if over eight (Fig. 49). Thus, significant die size reduction was achieved. Advanced Micro Devices, Inc. improved on the PAL with its 22V10 architecture, which introduced the concept of programmable output macrocells (Fig. 50). To better address the needs of designers to perform multiple levels of logic, Signetics unveiled the folded architecture in 1985; this allowed the outputs of the single array of NAND gates to be "folded back" to the inputs of the NANDS. Several variations of "folded" PAL and FPGA architectures have since been marketed.

The next major PLD architecture innovation was the field programmable gate array. There are several very different architectures that fit into the FPGA category. Several of the FPGA architectures resemble arrays of PALs with programmable interconnect resources provided to wire between the individual PALs. Altera's architecture shown in Figs. 51, 52, and 53 depicts this well. Altera employs an ultraviolet-light-erasable EPROM programming element. Xilinx introduced the logic cell array (Fig. 54) in 1985 using an sRAM cell programming element. Actel introduced the next major FPGA architecture in 1988 based on its antifuse programming element* (Figs. 55 and 56). This architec-

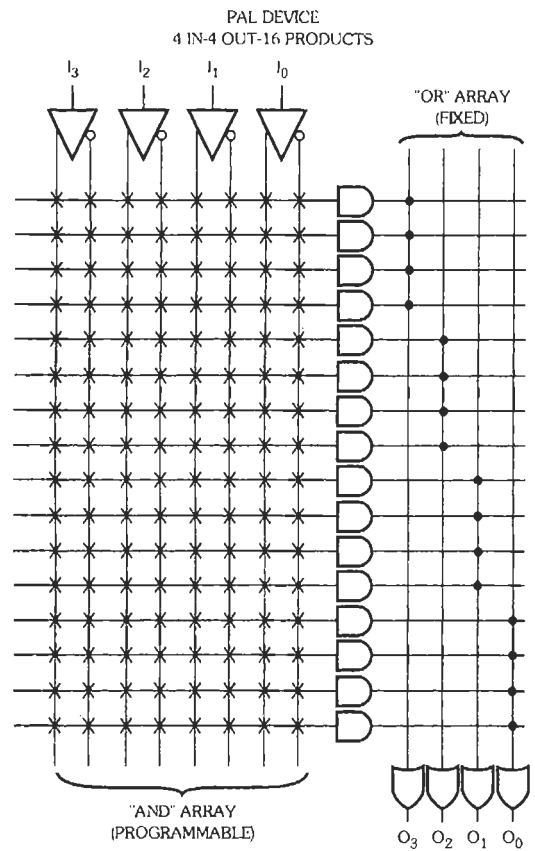


Fig. 49. PAL™ architecture. (Copyright © 1986 Advanced Micro Devices, Inc. Reprinted with permission of copyright owner. All rights reserved.)

ture most resembles a gate array because of the smaller granularity of its logic module and highly segmented routing architecture. New FPGA architectures based on new antifuse technologies have been announced by start-ups QuickLogic† and CrossPoint.‡

Analog ASICs—Analog and mixed digital/analog ASICs are emerging as an important new family of integrated circuits as CAD tools and design methodologies have been improved and adapted from the digital IC domain.‡ This is significant, as analog integrated circuit design has historically required much hand-crafting, which virtually precluded the availability of application specific analog ICs. The main classes of

* Reference 37.

† References 38 and 39.

‡ References 40 and 41.

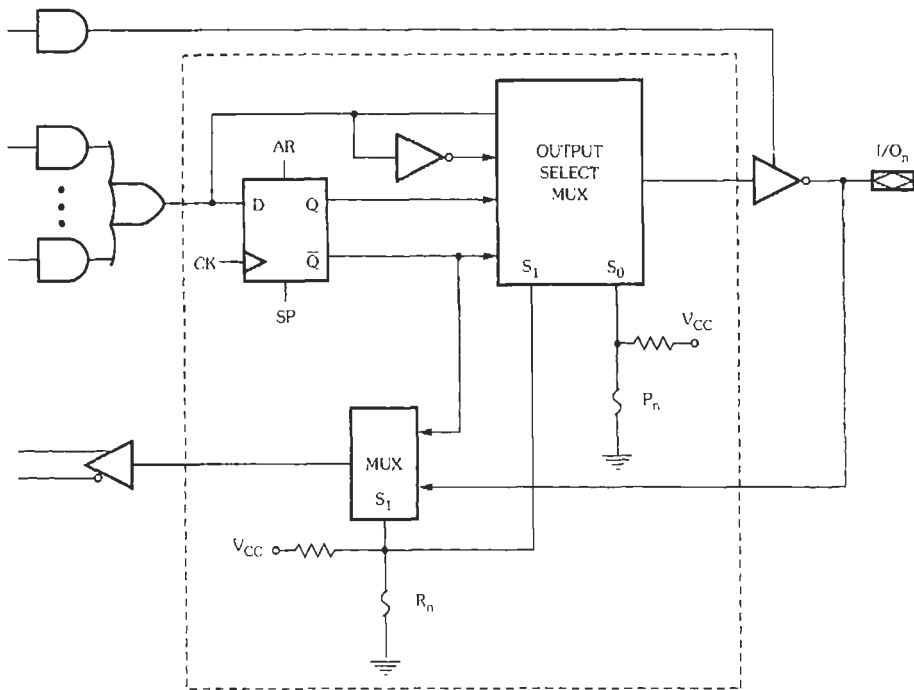


Fig. 50. PAL™ programmable input/output macrocell. (Copyright © 1987 Advanced Micro Devices, Inc. Reprinted with permission of copyright owner. All rights reserved.)

analog ASICs fall into these classes, each of which is also valid for merged digital/analog:*

1. Analog arrays: A chip from a master slice in which analog or a combination of analog and digital circuit elements are fabricated. The interconnection layers are defined by the end-user designer as in gate arrays.
2. Analog standard cells: Fully predesigned analog circuit blocks that reside in a library from which the user can select the needed items. This is the same concept as digital standard cells discussed previously.
3. Analog parameterized cells: Partially predesigned analog circuit blocks that allow the user not only to select suitable circuit blocks from the standard-cell library but also to set key parameters such as gain via mask-programmable connections of circuit elements.
4. Analog programmable chips: Similar to programmable logic devices in that they are completely fabricated and configured by the user. The programmable elements are configured by the end user to set the functional parameters of blocks of predesigned analog circuitry and to define the requisite interconnectivity.†

* Reference 42.

Future Trends in ASICs

Application Specific Memory—Application specific memory is based on a random access memory (RAM) chip that utilizes ASIC techniques to allow configuration of such attributes as input/output bit width (1, 4, 8, ...). Normal commodity static and dynamic memory chips require very large amounts of engineering to achieve the necessary performance while minimizing chip area to reduce cost due to the very high volumes of generic memories. Memory suppliers are beginning to develop and introduce products that leverage the engineering deployed on generic RAMs while adding the capability of customizing some aspects for applications such as video RAM and image processing.‡ It is interesting to note that one approach to application specific memory is to embed an array of logic gates, the converse of embedding RAM inside gate arrays as previously described.

Multiple-Chip Modules (MCM)—Multiple-chip modules may become an ASIC growth area for the most complex and volume/performance-sensitive systems. These devices have multiple chips attached to a common substrate and are then interconnected

† Reference 43.

‡ Reference 44.

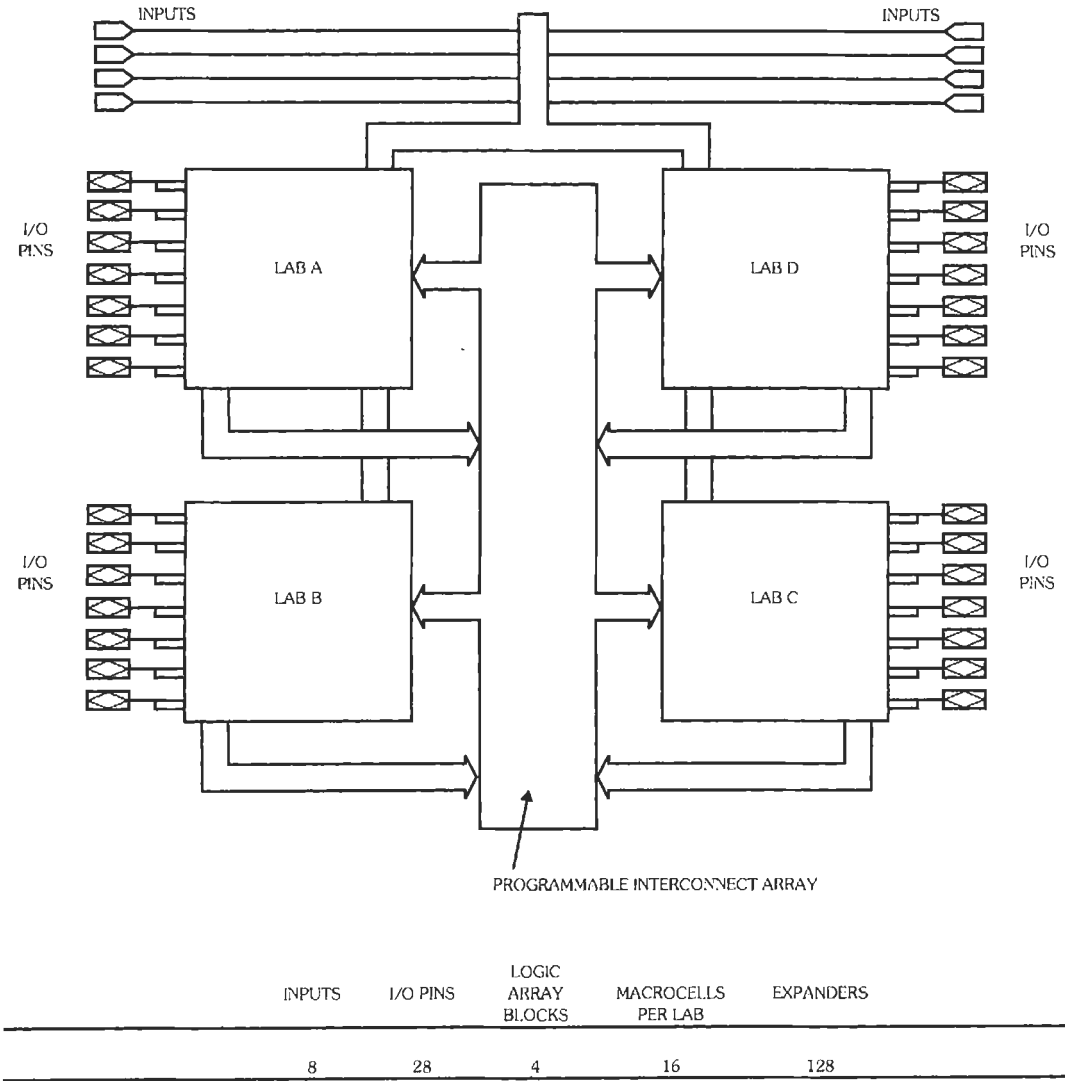


Fig. 51. Altera FPGA architecture (EPM5064).

together with advanced packaging technologies such as tape-automated bonding (TAB). Multiple-chip modules may become an ASIC-like solution as ever more complex chips become available for commodity use as in today's standard cell library.*

Field Programmable Gate Arrays—Field programmable gate arrays are expected to continue to be one of the fastest growing and most dynamic areas of ASIC for the near future as new architectures† and technologies are introduced to advance FPGA density

and performance closer to those of gate arrays while reducing the relative cost.‡ Accelerated displacement of the gate-array market by FPGAs has occurred as FPGA cost reductions proceed faster than those for gate arrays.

Merged Analog and Digital ASICs—The next volatile area of growth may be merged analog and digital ASICs, both mask-programmed and programmably configurable versions. Many high-volume consumer products such as personal stereos and

* Reference 45.

† References 46 and 47.

‡ References 29 and 48.

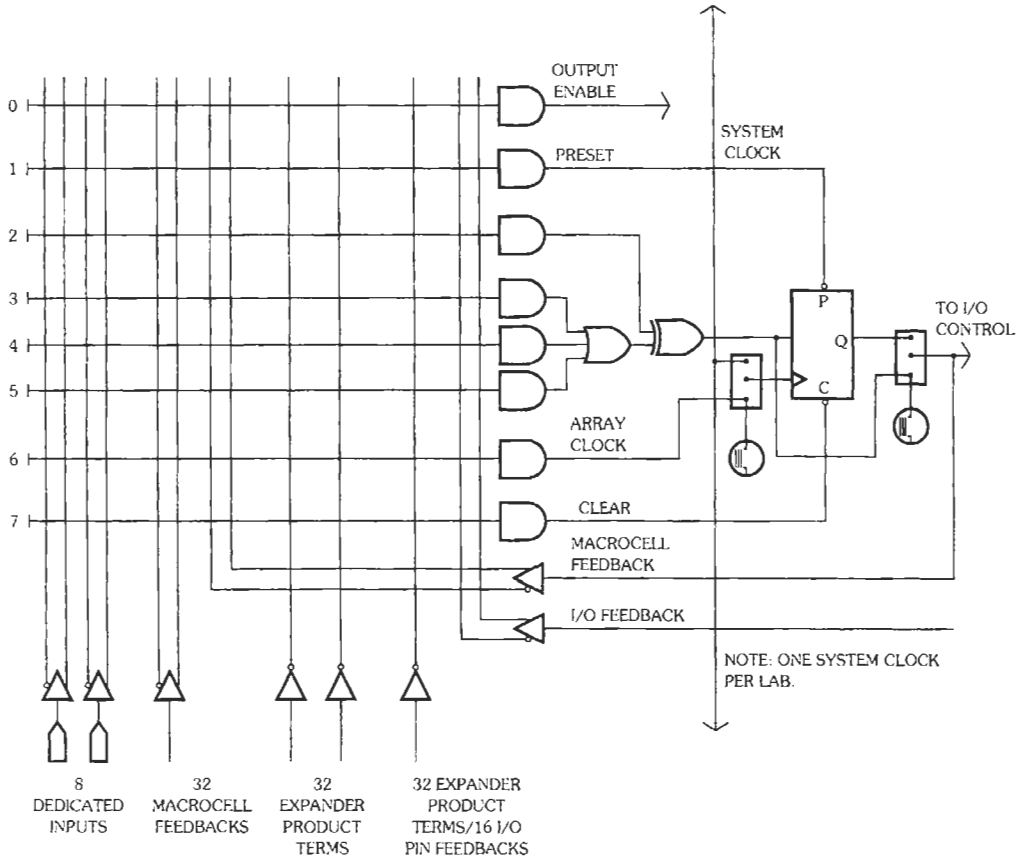


Fig. 52. Altera FPGA input/output macrocell.

televisions benefit from merged linear/digital ASIC product advances.*

Silicon Compilers—Silicon compilers provide automatic design and layout from high-level circuit descriptions. This has been a promising full-custom design methodology for some time that has not yet proved to be a popular design solution. The great success of other ASIC solutions may have retarded the demand for development of silicon compilers to date. Silicon compilation may become more widespread as IC complexity increases and the IC design task shifts to the system designer who will demand CAD tools that accept high-level design specifications and early access to the latest IC technology. The rapidly increasing computation power of engineering workstations should also aid commercialization of silicon compilers.

SINGLE-WAFER INTEGRATED DEVICE PROCESSING†

This section presents an overview of various single-wafer fabrication techniques for integrated processing of microelectronic devices. Numerous processing modules, sensors, and associated fabrication processes have been developed for advanced semiconductor device manufacturing.‡ Specific developments and results will be presented in the areas of rapid thermal processing (RTP), dry and vapor-phase surface cleaning, epitaxy, plasma processing, and in situ sensors.

Single-Wafer Semiconductor Manufacturing

The continuing demands for higher chip packing densities and enhanced system-level performance are the main technology and device scaling drivers. The

* References 40, 41, and 42.

† This section was contributed by Mehrdad M. Moslehi.

‡ References 49 and 50.

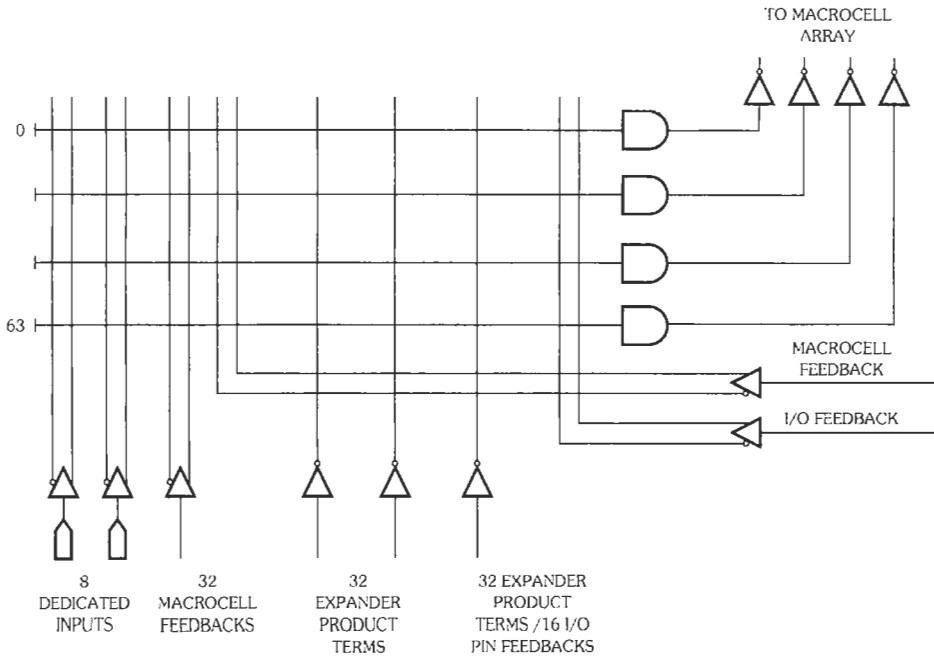


Fig. 53. Altera FPGA expander product term array to “fold back” signals between macrocells.

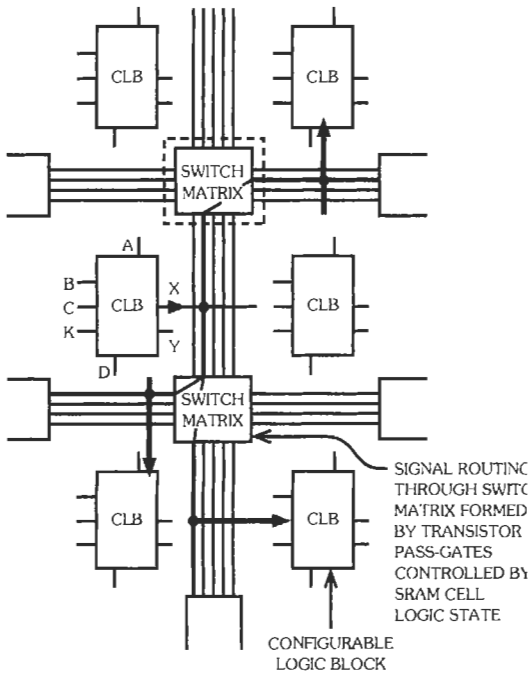


Fig. 54. Xilinx FPGA architecture. (Courtesy Xilinx, Inc.)

high-performance processor and dense memory chips now require device integration levels on the order of 100–1000 million transistors based on advanced sub-micron (e.g., 0.13–0.18 μm) CMOS and BiCMOS technologies. This level of integration density places stringent demands on the processing equipment, device fabrication processes, and manufacturing process control. State-of-the-art semiconductor manufacturing factories usually employ batch equipment for high-throughput wafer processing. The high manufacturing throughputs (e.g., thousands to over 20,000 wafers/month) and stringent manufacturing process control requirements are achieved at the cost of reduced flexibility as well as rapidly escalating cost of new generations of advanced processing equipment. Semiconductor technology scaling also affects wafer size. Currently, the state-of-the-art chip factories are set up for 200-mm wafer processing. The wafer size is expected to grow to 300 mm by 2005.

Single-wafer processing (SWP) is an alternative to batch equipment for various device fabrication processes. As the wafer size increases beyond 200 mm, advanced SWP equipment will be preferred for many applications. The factors in favor of SWP include compatibility with multichamber cluster equipment for vacuum-integrated processing, improved fabrication cycle time, and enhanced process repeatability. An argument against SWP has been its inferior processing throughput compared with the batch equipment.

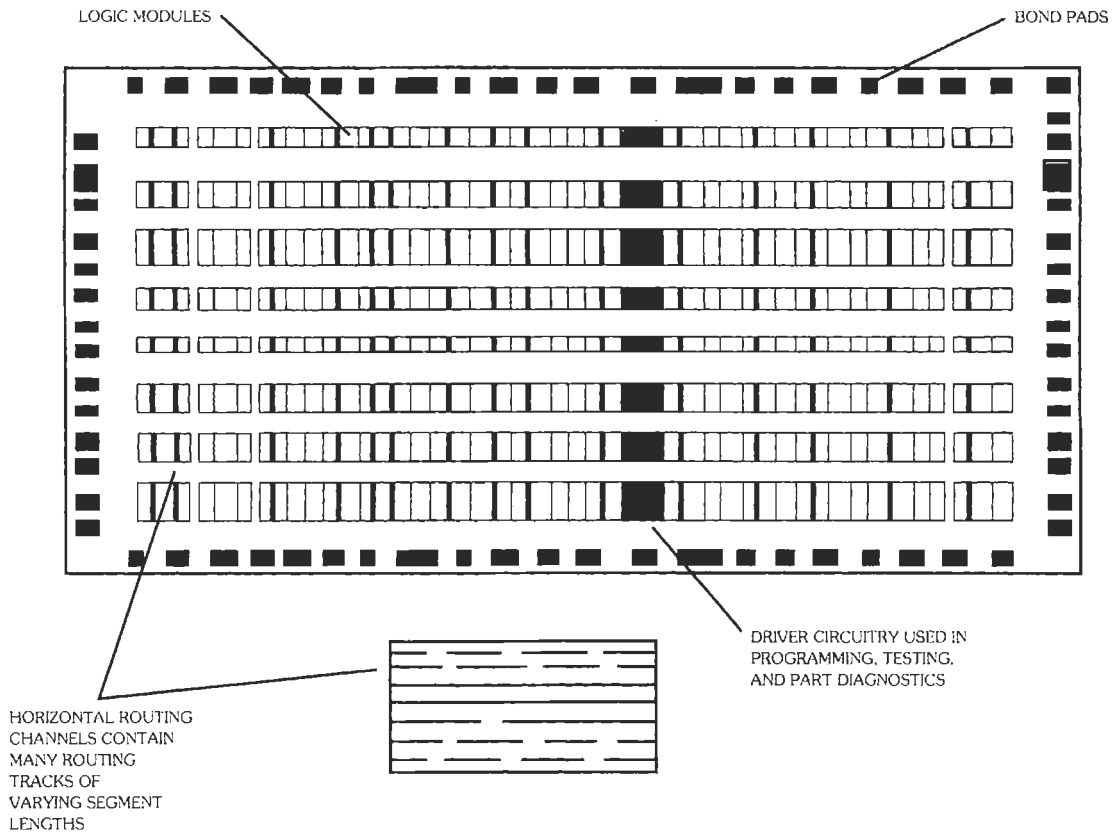


Fig. 55. Actel FPGA antifuse based logic module architecture.

However, many applications can afford somewhat lower throughput in return for capabilities for integrated processing and improved manufacturing process control.

As a result, an alternative semiconductor-device manufacturing methodology and associated techniques have been developed to meet several objectives, including: (1) lower overall factory cost due to the use of single-wafer processors, (2) sensors for process control/diagnosis based on computer-integrated manufacturing (CIM) methods for enhanced factory utilization and flexibility, (3) rapid manufacturing cycle time due to SWP and real-time process control techniques, and (4) enhanced flexibility for manufacturing numerous products based on different technologies. The following subsections will describe some selected elements of the single-wafer semiconductor minifactories and their impact on advanced submicron technology integration.

Integrated Device Manufacturing

The combination of in situ monitoring sensors, SWP, vacuum-integrated cluster tools, and computer-

integrated manufacturing (CIM) can provide a suitable fabrication environment for flexible and high-yield semiconductor device manufacturing. The use of in situ sensors for monitoring of equipment, process, and wafer parameters results in increased equipment/process up-time, reduced process and device parameter spread, improved cluster tool reliability and functionality, and reduced overall device manufacturing cycle time. This subsection will present an overview of the features and impact of sensors for semiconductor device manufacturing applications. Specific examples will be presented for the use of critical sensors in conjunction with cluster tools for advanced CMOS device processing.

Sensors: Motivations and Impacts—Semiconductor factories usually employ well-established statistical process control (SPC) techniques to minimize the process parameter deviations and to increase the device fabrication yield. The conventional fabrication environments rely on controlling a limited set of critical equipment and process parameters; however, most of the significant wafer, process, and equipment parameters of interest are not monitored in real time or in situ. These critical parameters may include conductive

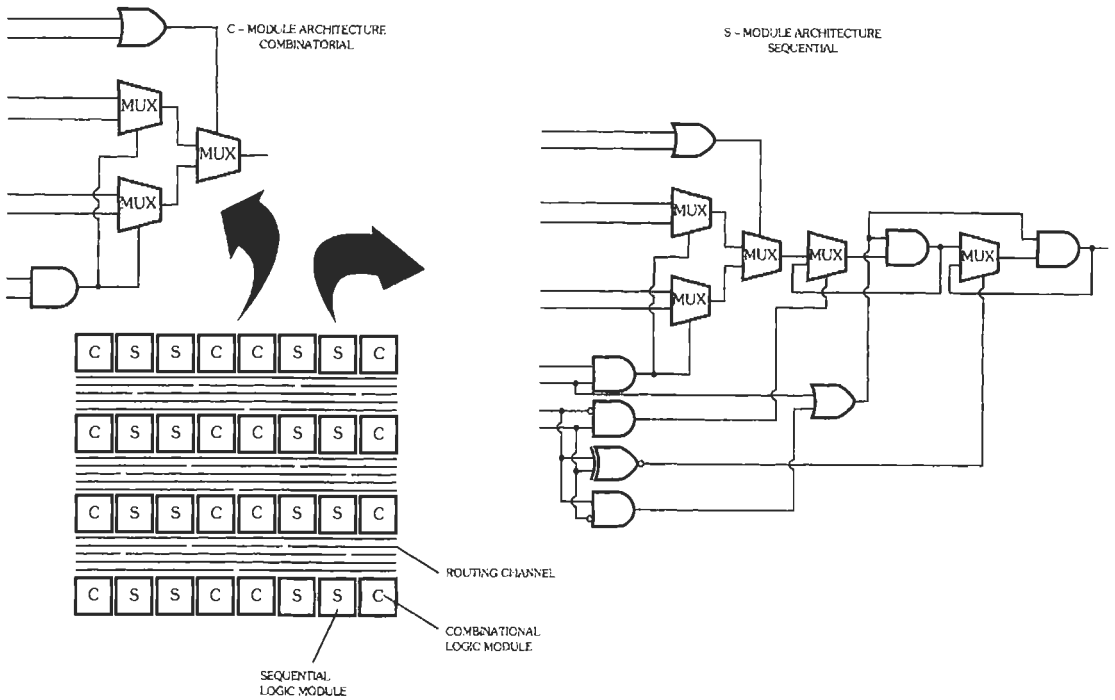


Fig. 56. Actel FPGA logic modules and examples of array configurations.

layer sheet resistance values, film thicknesses, process uniformity, equipment cleanliness, etc. The existing commercial processing equipment usually lacks suitable in situ sensors for real-time equipment diagnosis and prognosis. In many cases, the processing reactors have not been designed for modular implementation of new sensors. As a result, the detrimental effects of any unexpected process parameter variations and equipment-related drifts are not detected in real time. The relatively long time delays between detection of process or equipment problems and implementation of suitable corrective measures increase the overall chip manufacturing cost.

The conventional SPC techniques operate based on extensive ex situ or off-line measurements on pilot and device wafers. These off-line measurements provide useful statistical data related to the status of equipment and various wafer lots in process. However, the off-line measurements are usually time consuming and labor intensive and degrade the device manufacturing cycle time. In general, the following trends and driving forces affect the semiconductor technology evolution in microelectronics manufacturing: (1) increased chip fabrication yield, (2) reduced device manufacturing cycle time, (3) real-time process control to maintain process and wafer parameters within specified bounds, (4) improved semiconductor device reliability and performance, and (5) reduced manufacturing cost per chip.

The practical limitations and high cost of the conventional SPC methodology in conjunction with the above-mentioned driving forces as well as the rapidly rising costs of new semiconductor processing equipment and factories have promoted the evolution of flexible semiconductor minifactories based on in situ sensors and vacuum-integrated cluster tools. Cluster systems consisting of single-wafer processing modules and in situ sensors allow sequential integrated processing in the CIM-supervised factories. In situ sensors improve the manufacturing cycle time by eliminating the need for the off-line measurements and by enhancing the process and equipment reliability. Moreover, the in situ sensors employed for equipment and process diagnosis/prognosis will increase the equipment/process up-time and improve the cluster equipment reliability and functionality. Real-time in situ process and wafer monitoring sensors are also very effective in reducing the process and device parameter spread, resulting in manufacturing yield improvements.

Sensor Hierarchy for Semiconductor Processing—Fig. 57 shows the sensor hierarchy for various semiconductor device processing applications. The process control sensors can be divided into the invasive and noninvasive groups. The invasive sensors are intrusive to the process environment and/or wafer surface within processing equipment and, as a result, are not suitable for in situ process monitoring applications. For

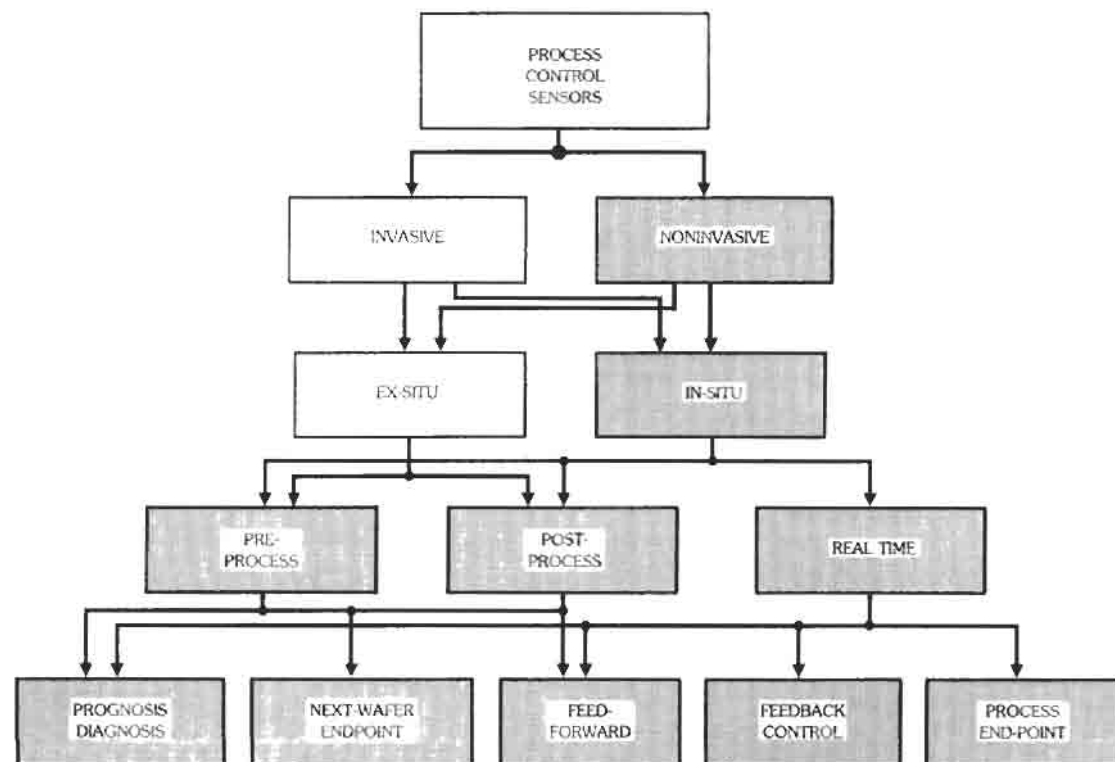


Fig. 57. Hierarchy of process monitoring and control sensors for semiconductor manufacturing.

example, most contact-type sensors that require direct mechanical contact to the wafer surface (e.g., thermocouples for wafer temperature measurements) are invasive. On the other hand, the noninvasive in situ sensors are the most suitable type for device fabrication process control needs. These sensors usually employ noninvasive probe beams or characteristic detected beams of energy with well-defined physical signatures (e.g., optical, acoustical, electronic, and electromagnetic) to monitor the desired parameters based on a known physical phenomenon. The ex situ sensors are essentially for off-line preprocess and postprocess monitoring applications; however, the in situ sensors are designed for real-time as well as on-line pre- and postprocess monitoring and process control.

As indicated in Fig. 57, the real-time in situ sensors can be employed for process end-point detection, real-time feedback control, and process/equipment diagnosis and prognosis. For instance, an in situ real-time sheet resistance sensor can be used within a metal CVD process module to provide feedback process control (for instance, to control the deposition rate) and to achieve precise process end-point detection. The same sensor may also provide data to feed-forward links and to expert systems for equipment and process

prognosis/diagnosis purposes. The in situ preprocess and postprocess sensors are mounted within a load-lock chamber or a dedicated sensor (metrology) module interconnected to the vacuum-integrated cluster equipment. These sensors are useful for next-wafer (delayed) process end-point control, feed-forward links, and equipment/process prognosis/diagnosis. The specific choice of various types of in situ sensors depends on various factors such as the application and cost. For instance, depending on the process application, the task of in situ dielectric thickness measurement may either be performed in real time (for real-time rate control and/or end-point detection) or following the process by a postprocess sensor (for delayed end-point control and feed-forward links).

The in situ sensors are intended to monitor various parameter sets within three domains. These domains are associated with the equipment, process, and wafer parameters. Depending on a specific application, the critical equipment parameters to be monitored may include gas showerhead temperature (e.g., in PECVD TEOS oxide), background contamination (e.g., in epitaxial Si growth and PVD processes), quartz window temperature (in lamp-heated RTP systems), chamber wall deposits, and energy source status (e.g., status of

lamps and the total lamp power in RTP). Additional equipment parameters are the gas flow rate and pressure settings in a processing reactor. However, the process parameters are directly governed by the equipment parameters (both controlled and uncontrolled parameters). Depending on a specific application, the process parameters of interest to monitor may be wafer temperature, process gas flow field, plasma density, and plasma electron energy. In general, a semiconductor wafer exposed to a process contains the most important parameters that have to be controlled. This is due to the fact that the wafer parameters have a direct impact on device performance and yield. For example, in a CVD or PVD (sputter deposition) process used for deposition of dielectric or conductive films, the important wafer-level parameters include film thickness, refractive index, stress, conductive-layer sheet resistance, and thickness uniformity.

Among the above-mentioned three parameter domains, the equipment parameters are the easiest to monitor. The equipment parameters can be directly adjusted and controlled, whereas the process and wafer parameters are dependent variables. Monitoring and control of the desired process and wafer parameters require specialized noninvasive in situ sensors and advanced control techniques. Proper adjustments of process and wafer parameters are usually performed

by varying the equipment parameters. The process parameters are functions of the equipment variables via equipment models, and the wafer parameters are determined by the process parameters via process models.

Sensors for Direct and Indirect Feedback Process Control—Fig. 58 shows the direct and indirect methods of parameter measurement and device manufacturing process control. The conventional methods of manufacturing process control employ equipment-level sensors to monitor the equipment parameters. These measurements (e.g., lamp power, mass-flow controller settings, etc.) provide the necessary real-time signals for an indirect feedback control loop that maintains the critical equipment parameters at or near the desired predetermined set-point values. The use of closed-loop control based on measured equipment parameters usually provides a significant improvement (in terms of process parameter spread) over the open-loop mode of operation; however, this type of process control relies on two levels of equipment and process models in order to relate the equipment parameters to the wafer parameters that are the ultimate parameters of interest for manufacturing process control. These models or the process-vs.-equipment and wafer-vs.-process relationships can be affected by various

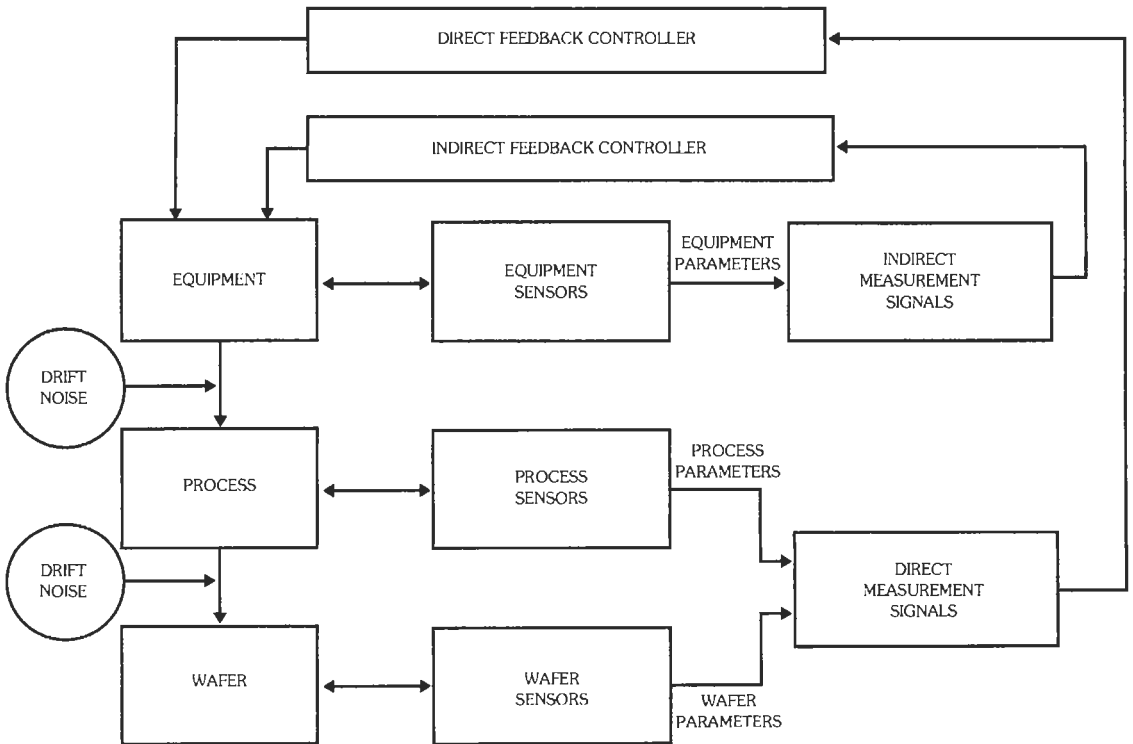


Fig. 58. Schematic illustration of direct and indirect process control techniques.

sources of drift and noise resulting from the uncontrolled variables. For instance, the optical/vacuum quartz window temperature in an RTP system behaves as an uncontrolled variable that can cause process drift. Although the uncontrolled equipment parameters are not controlled directly by an operator, most of them should be monitored to further refine the equipment model for improved process control. For example, real-time monitoring of the optical/vacuum quartz window temperature in an RTP system can provide a useful signal for improved closed-loop wafer temperature control and for equipment diagnosis.

As illustrated in Fig. 58, the direct process control techniques employ suitable in situ sensors to directly measure the important process and/or wafer parameters. These direct measurement signals are fed to a controller, which continuously adjusts the equipment parameters in real time to minimize the differences between the measured signals and the process and/or wafer parameter set points. It is evident that the direct control methods should result in a much more robust process control and smaller process and/or wafer parameter spread. This is because the direct feedback control loops contain the drift/noise sources within the overall control loop. When in situ noninvasive sensors

are used for real-time monitoring of the desired wafer parameters, the process control feedback loop will minimize the error signal representing the difference between the measured and set-point wafer parameter values. As a result, the equipment-related drift effects will be essentially eliminated.

As described earlier, the indirect equipment-related process control signals are related to the direct process/wafer-related process control parameters via the process and equipment models. As an example, Fig. 59 illustrates some representative equipment and process/wafer parameters in a typical plasma deposition or etch process. The plasma process and wafer parameters (direct process control signals such as plasma density, ion flux, wafer temperature, and processing rate) are related to the plasma equipment parameters (indirect process control signals such as RF/microwave power, dc bias voltage, and susceptor temperature) via the specific plasma equipment model. Of course, the functional relationships between the wafer and process parameters are established by the physics of plasma and the plasma process model. In general, a good knowledge of the equipment and process models is necessary to ensure optimal and fast real-time process control based on either direct or indirect feedback control techniques.

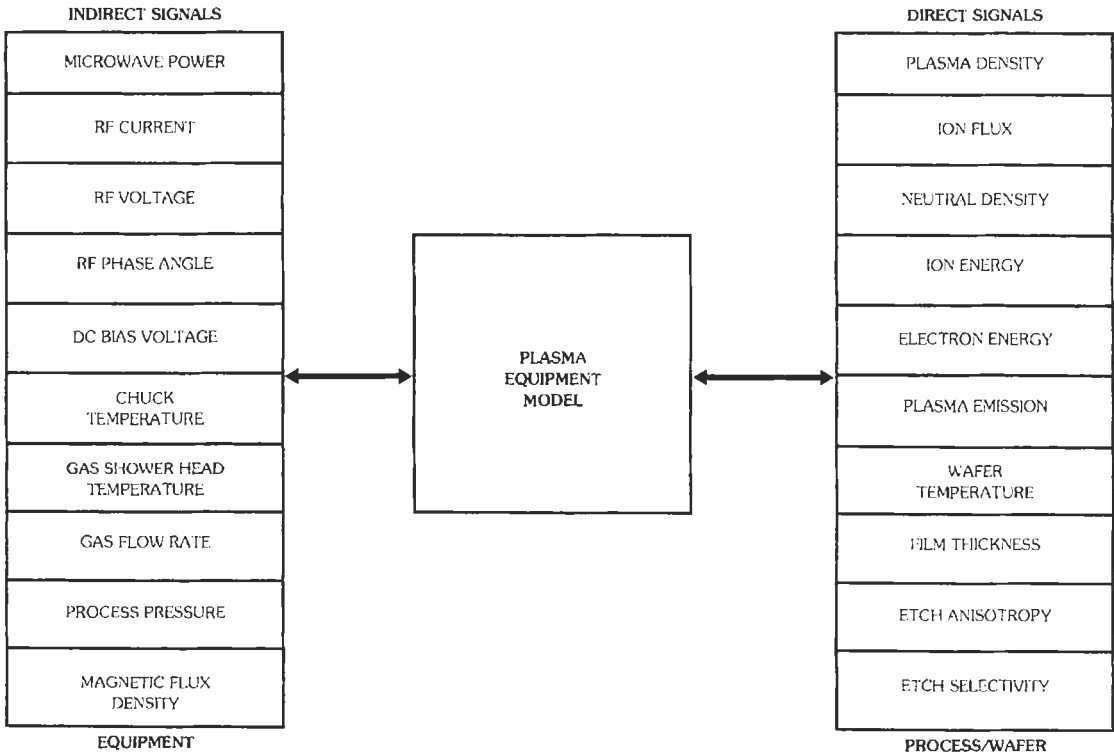


Fig. 59. Typical wafer/process and equipment parameters in a plasma process.

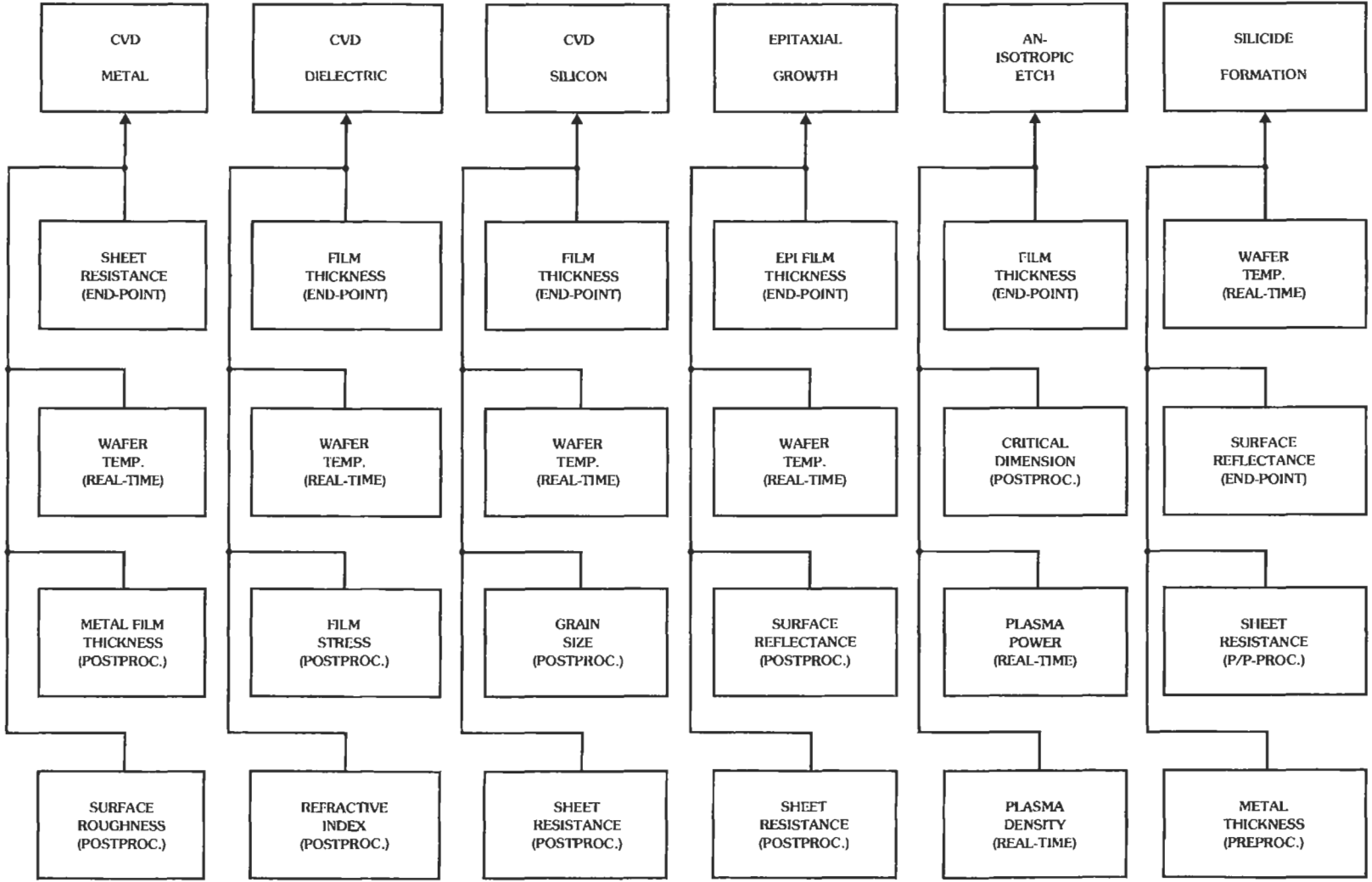
Critical Sensors for Device Fabrication Process Control—The critical in situ process control sensors required for a given fabrication process are determined based on various considerations such as the impact of the specified process on the overall device performance and fabrication yield. The process requirements for a given process indicate if any real-time and/or pre- and postprocess in situ sensors must be utilized to meet the target specifications for the process and wafer parameters. Fig. 60 illustrates various single-wafer processes and their associated “critical” in situ process control sensors. For example, consider an RTP module used for CMOS gate dielectric formation by rapid thermal oxidation and nitridation (RTO/RTN) processes. Thickness, growth temperature, and uniformity of gate dielectrics are the main critical wafer and process parameters that can influence the overall device performance and fabrication yield. As a result, one important sensor is a real-time multipoint temperature sensor for closed-loop temperature control to improve process repeatability. Another very useful sensor is an in situ ellipsometer for postprocess measurement of gate dielectric thickness and refractive index. This sensor will help to further reduce the overall gate dielectric thickness spread via delayed (next-wafer) feedback. An additional postprocess in situ sensor may be used for monitoring the gate dielectric thickness uniformity. The in situ postprocess uniformity sensor will assist in meeting the process uniformity specification and can provide delayed feedback process control information for process uniformity control, assuming the RTP module employs a multizone illuminator for dynamic and real-time uniformity control. If the temperature sensor monitors multiple points on wafer for dynamic uniformity control during the gate dielectric formation process, the postprocess uniformity sensor will be rather noncritical.

Various blocks shown in Fig. 60 illustrate some selected critical in situ sensors for CVD processes (metals, dielectrics, amorphous and polycrystalline silicon), epitaxial growth, plasma etch (anisotropic and isotropic), resist processing, ion implant processing, dry surface cleaning, glass reflow, and silicide formation. As another example, consider the critical sensor needs for metal CVD processes such as tungsten CVD using rapid thermal CVD (RTCVD) modules. Again, a temperature sensor is needed for real-time wafer temperature control. Precise control of wafer temperature and its uniformity in the CVD tungsten temperature range (300–550 °C) are critical requirements because of their impact on the deposition rate and process uniformity. Another very useful sensor is a noninvasive sheet resistance sensor for real-time monitoring of tungsten sheet resistance during the RTCVD process. This sensor will provide the process end-point information and will help to reduce the process parameter spread. Besides these real-time sensors, the process can benefit from two postprocess sensors for in situ monitoring of the metal film thickness and/or its surface roughness.

Intelligent Cluster-Tool-Based Device Manufacturing—The earlier discussions have indicated that the use of low-cost in situ sensors in conjunction with single-wafer vacuum-integrated cluster equipment can provide a fast-cycle-time processing medium to manufacture advanced semiconductor devices. This type of factory relies on intelligent real-time process control based on in situ sensors to minimize the device parameter spread and to enhance the overall equipment reliability and process up-time.

Fig. 61 shows the simplified block diagram of a CMOS process flow and its associated cluster equipment for a fully cluster-based fabrication sequence. Each cluster equipment contains three or four process modules. It is expected that increasing the number of process modules on a cluster equipment much beyond four to six may result in reduced equipment reliability and increased cost. The process flow starts with the CMOS tank (or well) formation module and is followed by the isolation, pregate, CMOS gate, source/drain, silicide, planarized interlevel dielectric, contact, metal, and passivation modules (assuming one metal level). As shown, the CMOS tank module can be completed using two cluster tools: one for microlithography marker pattern formation, and another one to complete the CMOS well formation. The device isolation module employs three cluster tools: one to complete the oxidation mask stack deposition, another cluster equipment for completion of the patterned oxidation mask, and a third module for field oxidation and oxidation mask removal. Other process modules employ suitable cluster tools, as indicated. This CMOS process flow employs a total of about 16 cluster tools (each with three or four individual process chambers) besides some additional equipment requirements such as the wet cleaning stations, ion implanter, and microlithography stepper. Appropriate sensors are included in each cluster equipment to allow real-time process control and device manufacturing based on intelligent equipment and factory control.

Figs. 62 and 63 show schematic diagrams of two different cluster tools with some real-time and pre- and postprocess in situ sensors for the two CMOS flow modules shown in Fig. 61. Fig. 62 illustrates a three-module cluster tool for CMOS gate structure formation. The first step in the process sequence uses the pre-clean module (such as vapor-phase clean) to remove surface contaminants including native oxide. Then, the RTP module is used to form the gate dielectric by RTO. The subsequent step is silicon gate deposition by LPCVD. All three modules require real-time temperature sensors (critical for RTP and LPCVD/RTCVD but less critical for the preclean module) for wafer temperature measurements and process uniformity control. The preclean module may also employ an ultraviolet-light intensity sensor for photochemical cleaning process monitoring and control. Moreover, the RTP module uses a real-time thickness sensor for process end-point detection. An alternative option would be to use a postprocess in situ thickness sensor. Similarly, the



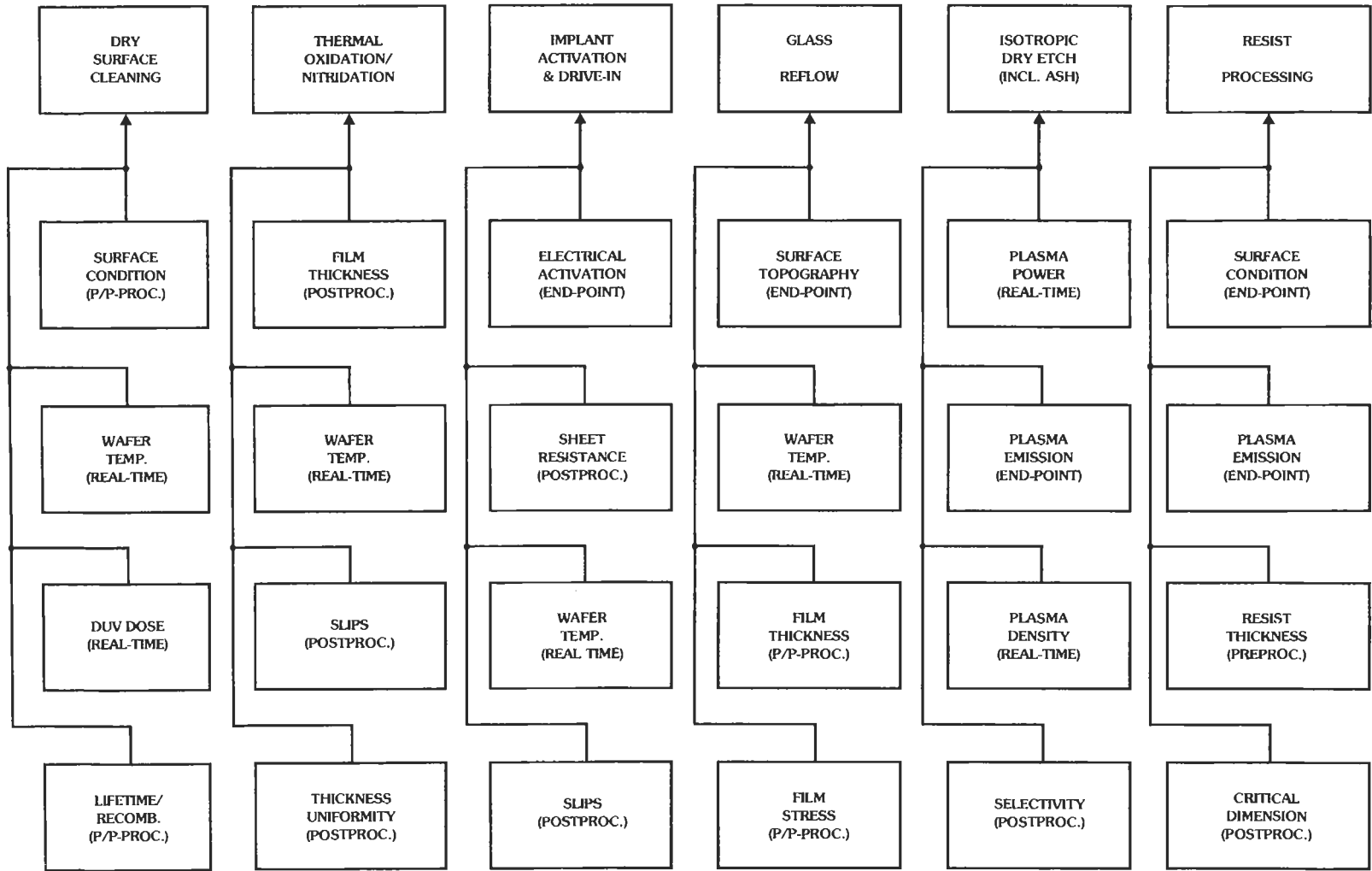


Fig. 60. Useful in situ sensors for process control in various semiconductor device fabrication processes.

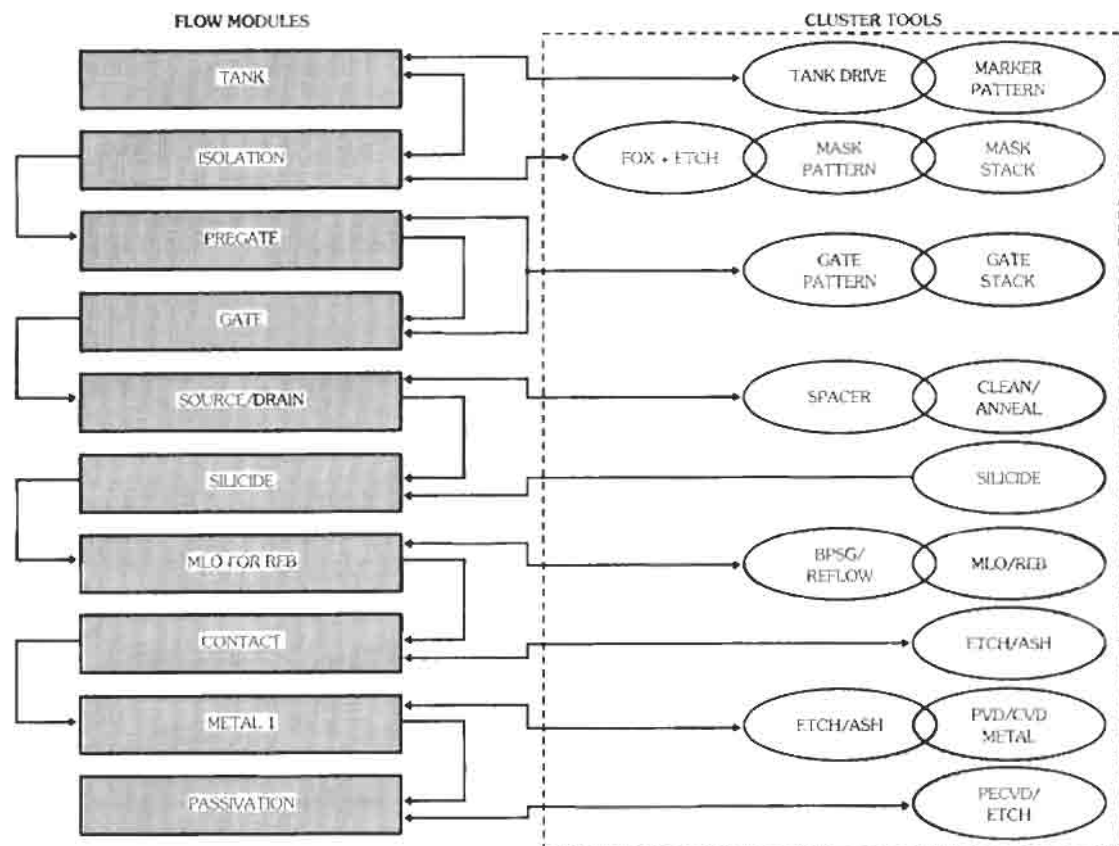


Fig. 61. Schematic example of a CMOS process flow and its associated cluster equipment.

LPCVD silicon process module also uses a silicon thickness sensor for process end-point detection. The load-lock (or metrology) chamber contains a slip sensor for postprocess monitoring of possible RTO-induced slip dislocations. A sheet resistance sensor may be used for postprocess monitoring of LPCVD in situ-doped silicon sheet resistance. Another load-lock sensor can be used to monitor the effectiveness of the surface preclean process.

Fig. 63 illustrates the example of a four-module cluster tool for self-aligned silicide formation. The process sequence employs preclean for native oxide removal, refractory metal sputtering, RTP silicide formation, plasma etch for removal of unreacted refractory metal and metal nitride, vapor-phase clean, and finally another RTP cycle for silicide annealing. All the critical real-time and pre- and postprocess sensors are also shown.

RTP for Integrated Device Processing

Conventional and state-of-the-art semiconductor factories employ batch furnaces for many thermal fab-

rication steps. Various applications include gate and field oxidations, low-pressure chemical-vapor deposition (LPCVD) (e.g., for polysilicon and silicon nitride), and thermal anneals (e.g., for junction formation and glass reflow). A significant number of other fabrication processes, however, are already based on SWP. Examples include plasma etch and deposition, ion implantation, and microlithography. The current trend toward replacement of batch processes with SWP is evident from the increasing market share for cluster tool equipment with SWP modules. RTP is an obvious choice to replace batch furnaces for thermal processing in a semiconductor minifactory environment. Advanced RTP tools are based on modular designs compatible with vacuum-integrated cluster equipment. High-performance multizone illuminators with cylindrical symmetry have been applied to various processes including RTO, CVD of tungsten and amorphous/polycrystalline silicon, silicide formation, sinter, epitaxy, as well as high-temperature rapid thermal annealing (RTA).^{*} Concurrent use of multizone illuminators and multipoint temperature sensors provides a capability for real-time wafer temperature control and process uniformity optimization.[†] Moreover,

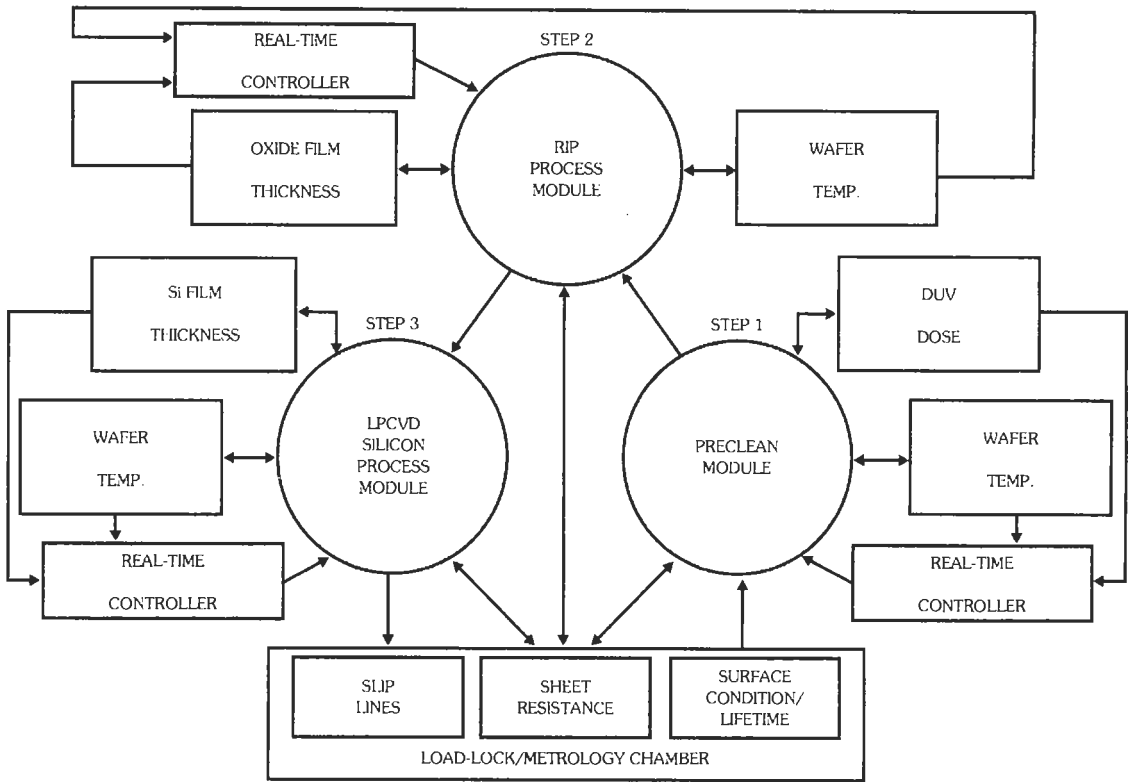


Fig. 62. Schematic illustration of a cluster equipment for CMOS gate structure formation.

numerous in situ process monitoring sensors have been developed to meet the RTP monitoring and control requirements. Some process applications of advanced load-locked single-wafer RTP techniques and associated sensors will be described.

Multizone RTP Illuminators—The known RTP systems mostly employ single-zone tungsten-halogen or arc lamp modules for wafer heating. The illuminator is usually designed to provide relatively uniform steady-state wafer heating over a narrow range of wafer temperatures. The conventional RTP illuminators provide real-time lamp power control capability for only a single zone. This has been known to impose major limitations in terms of dynamic wafer temperature uniformity and slip dislocation control.‡ For a given optical flux distribution on a semiconductor wafer, temperature uniformity can be affected by various additional parameters such as gas flow rates and process pressure. Moreover, the transient heat-up and cool-down temperature uniformity of a wafer can be

quite different from the steady-state uniformity.‡ These observations imply that the conventional single-zone lamp RTP systems may not be able to meet the stringent wafer temperature and process uniformity requirements over an extended domain of process parameters. Some developments in the area of RTP equipment have addressed these requirements. One approach has been to use a crossed lamp configuration, which employs two parallel crossed banks of linear tungsten-halogen lamps above and below the wafer.§ This technique is superior to the conventional single-zone RTP designs; however, owing to lack of cylindrical symmetry, only a limited real-time uniformity control can be achieved.§ This limitation is due to the circular shape of semiconductor wafers and linear/crossed configuration of the lamps.

Another approach has been to employ multizone illuminators with overall cylindrical symmetry. As an example, Fig. 64 shows a schematic illustration of Texas Instruments' Advanced/Automated Vacuum Processor (to be called TI-AVP) using a three-zone illuminator. The TI-AVP three-zone illuminator employs a

* Reference 50.
 † Reference 51.
 ‡ Reference 52.

§ Reference 53.

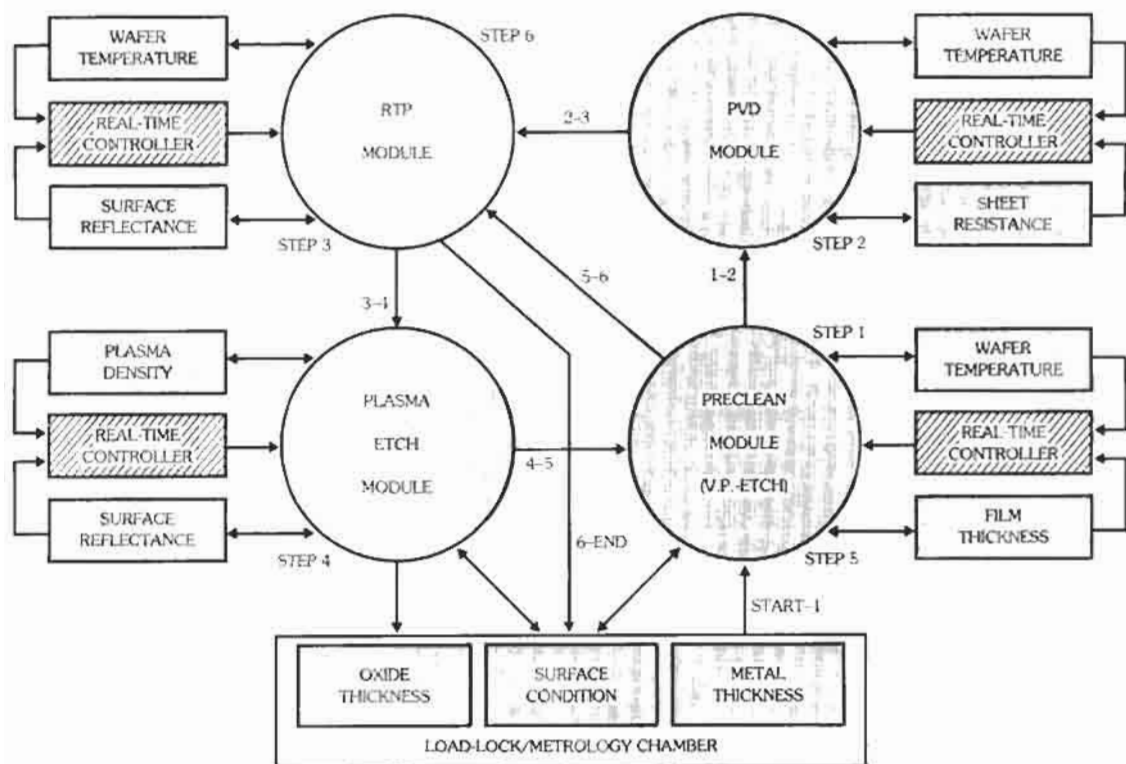


Fig. 63. Schematic picture of a cluster equipment for self-aligned silicide formation.

single center lamp, a first ring of 12 lamps, and an outer ring of 24 lamps. The combination of the point-source lamps in each ring simulates a continuous optical flux ring. All the zones are concentric with the central axis of the reactor process chamber and the center of the semiconductor wafer. The multizone lamp module heats the semiconductor wafer from the rear. Each lamp zone can be controlled independently in real time via its dedicated remote and computer-controlled power supply. The illuminator design provides multiple hollow light-pipe holes through the illuminator housing and associated reflector assembly for multipoint real-time wafer temperature measurements and control. This feature is important because various RTP fabrication processes are performed over an extended process parameter domain (e.g., 0.25–760 torr pressure for LPCVD and RTO processes). Recent modeling results have shown that the process pressure and gas flow rates can have significant effects on wafer temperature uniformity, particularly near the wafer edge.* In the RTP examples presented here, wafer processing and heating is performed face-down. Fig. 65 shows a representative open-loop response of four radially distributed temperature sensors over a 150-mm wafer using a three-zone RTP illuminator (1 torr/250 sccm [standard cubic centimeters per minute] flow

of nitrogen). With optimum zone power ratios, this illuminator provides temperature nonuniformity of less than $\pm 1.5^\circ\text{C}$.

Multipoint Temperature Sensors—RTP applications require accurate and repeatable real-time control of wafer temperature. The most commonly used non-contact temperature measurement technique is optical pyrometry. The advantages of optical pyrometry for RTP include its noninvasive nature and relatively fast measurement speed. Both of these requirements are critical for proper control of transient heat-up and cool-down temperature ramps and accurate and repeatable temperature control. The fundamental disadvantage with conventional pyrometry, however, is its dependence on the wafer surface roughness and emissivity. Surface roughness is typically wafer dependent, while emissivity depends on a range of parameters, including temperature and chamber reflectivity, as well as wafer-doping characteristics and surface optical condition.† Further, the direct reflection of heating lamp light from the wafer surface must be eliminated from the pyrometer measurement. Similarly, the heat-

* Reference 54.

† References 55–57.

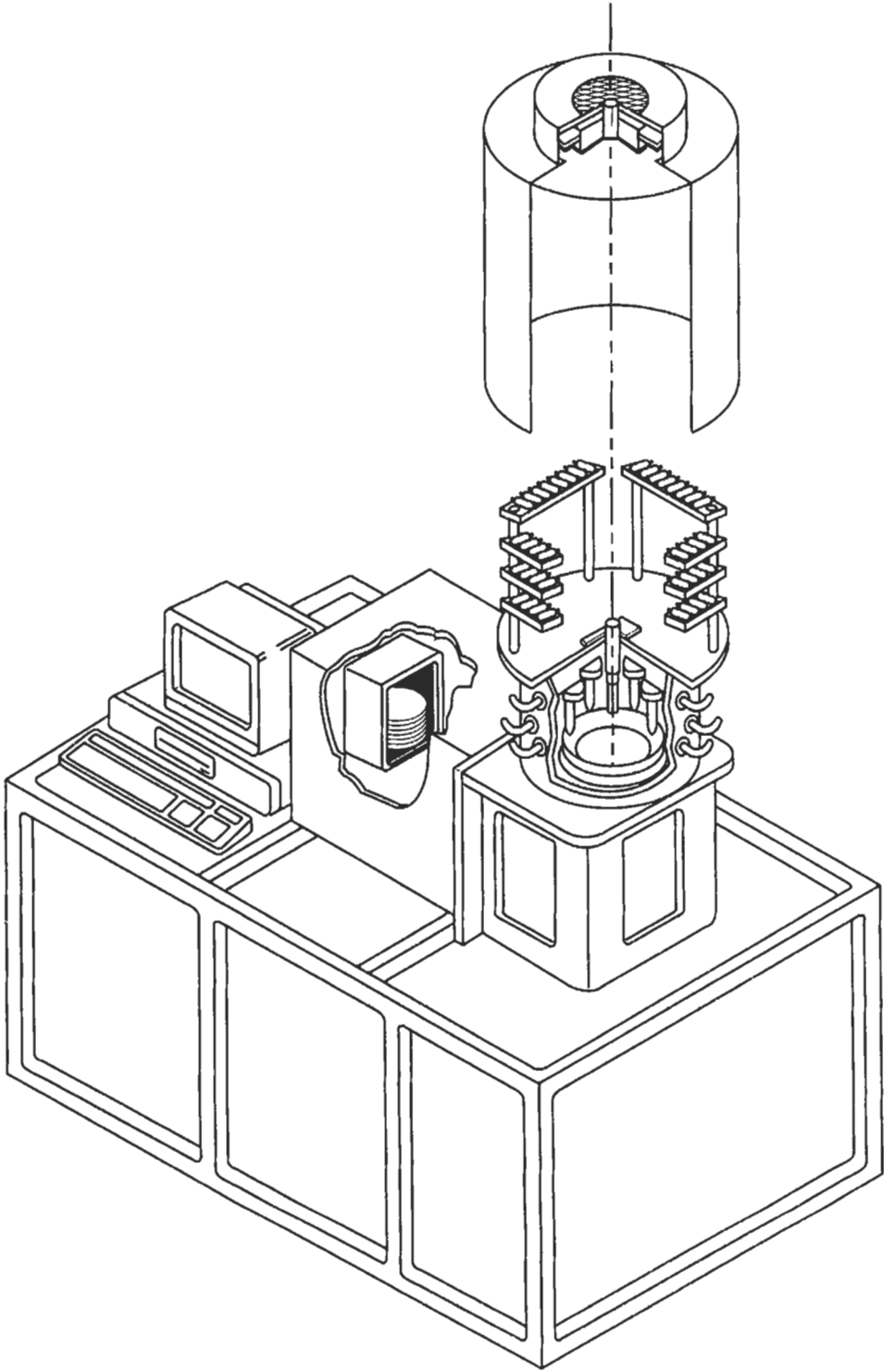


Fig. 64. Schematic illustration of a Texas Instruments single-wafer (TI-AVP) reactor with a multizone lamp heating module.

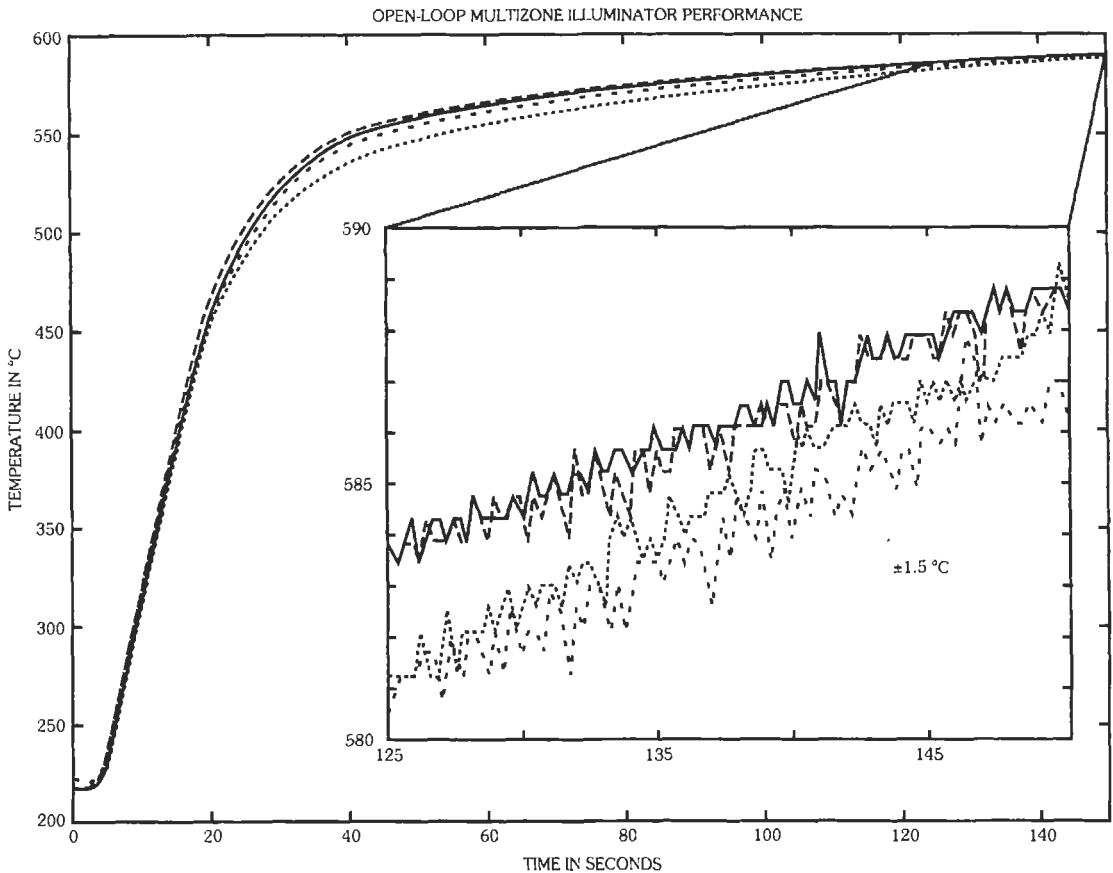


Fig. 65. Open-loop temperature-vs-time characteristics of a 150-mm wafer heated using a three-zone RTP illuminator with optimum zone power ratios.

ing of any optical chamber windows through which the wafer radiance is monitored, and any ensuing changes in the window index of refraction and transmission characteristics, must be corrected for in the measured radiance. Without proper compensation, these factors can result in large temperature-measurement errors. Frequent calibrations against thermocouple instrumented wafers have allowed somewhat improved temperature control for a specific process/reactor but are not satisfactory in general. Preprocess measurement of wafer reflectance and backside roughness may be used to correct for wafer-to-wafer changes in surface optical conditions. The use of preprocess wafer reflectance measurements has been reported for improved pyrometry measurements.* Real-time wafer emissivity measurement has been developed to allow correction for emissivity changes. Aside from the above issues affecting pyrometry measurement accuracy and repeatability, wafer temperature uniformity is a key requirement to prevent slip dislocations in high-tem-

perature RTP.† Multipoint fiber-optic pyrometry allows real-time monitoring and control of wafer temperature using multizone lamp heat sources.‡

Fig. 66 is a schematic diagram of a three-zone lamp-heated RTCVD reactor with multipoint temperature measurement and control. The wafer is processed face-down to minimize particulate deposition on the device side. A three-zone tungsten-halogen lamp source provides radiant heating to the back side of the wafer. A quartz window provides a vacuum seal and allows for optical access to the wafer. Multiple pyrometers are arranged at various radial positions, each looking at the back side of the wafer through a $3.3\text{-}\mu\text{m}$ interference filter and a suitable fiber light pipe. The light pipe is water cooled to prevent excessive heating caused by reflection and radiation from the hot wafer. The pyrometer/light pipe assembly detail is shown in Fig. 67. The light pipe material used for $3.3\text{-}\mu\text{m}$ transmission is either a sapphire rod or a fluoride fiber bundle.

* Reference 58.

† Reference 52.

‡ Reference 51.

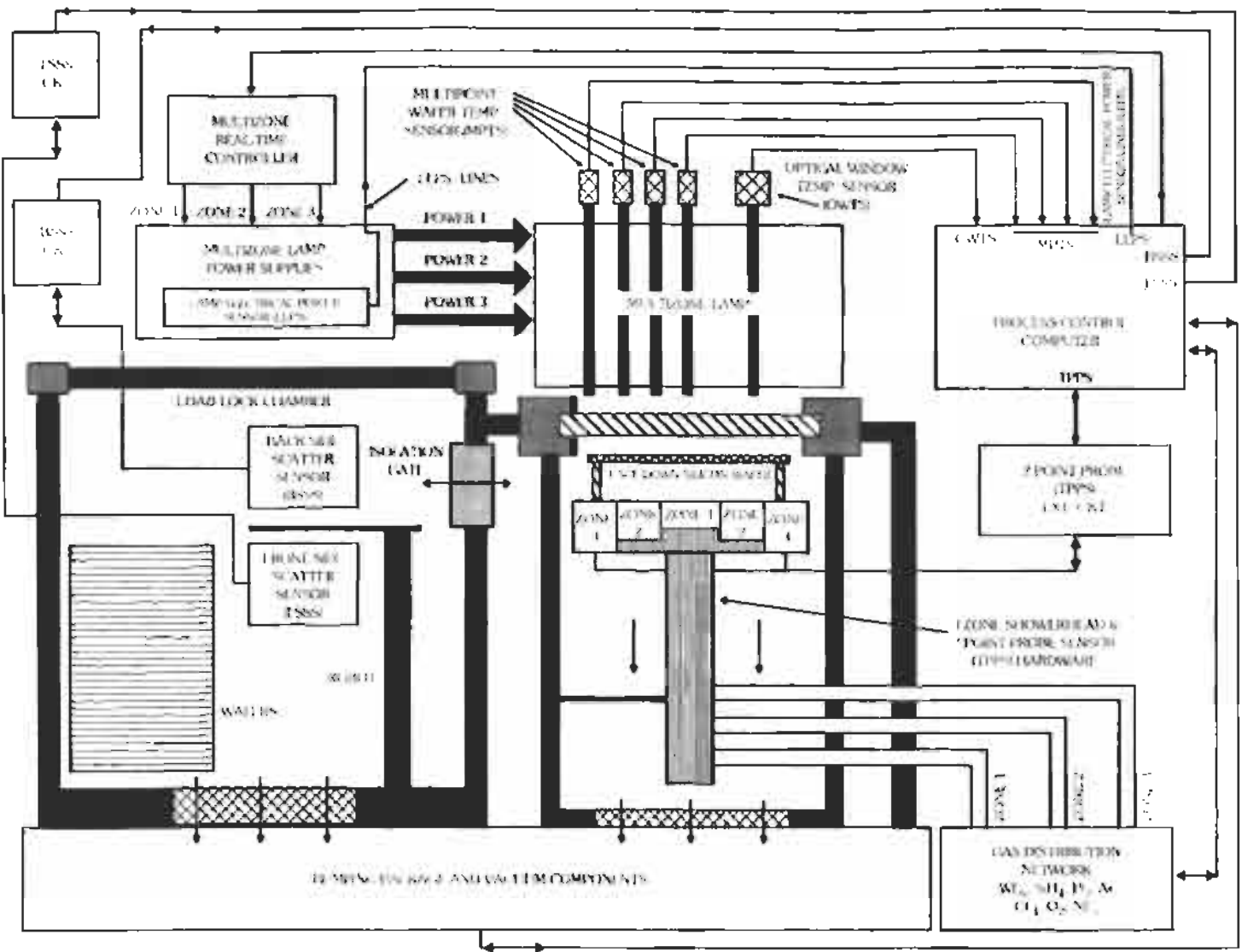


Fig. 20. Schematic diagram of an RT-CVD reactor used for impurity deposition. Various sensors are shown for temperature monitoring and CVD-W process control.

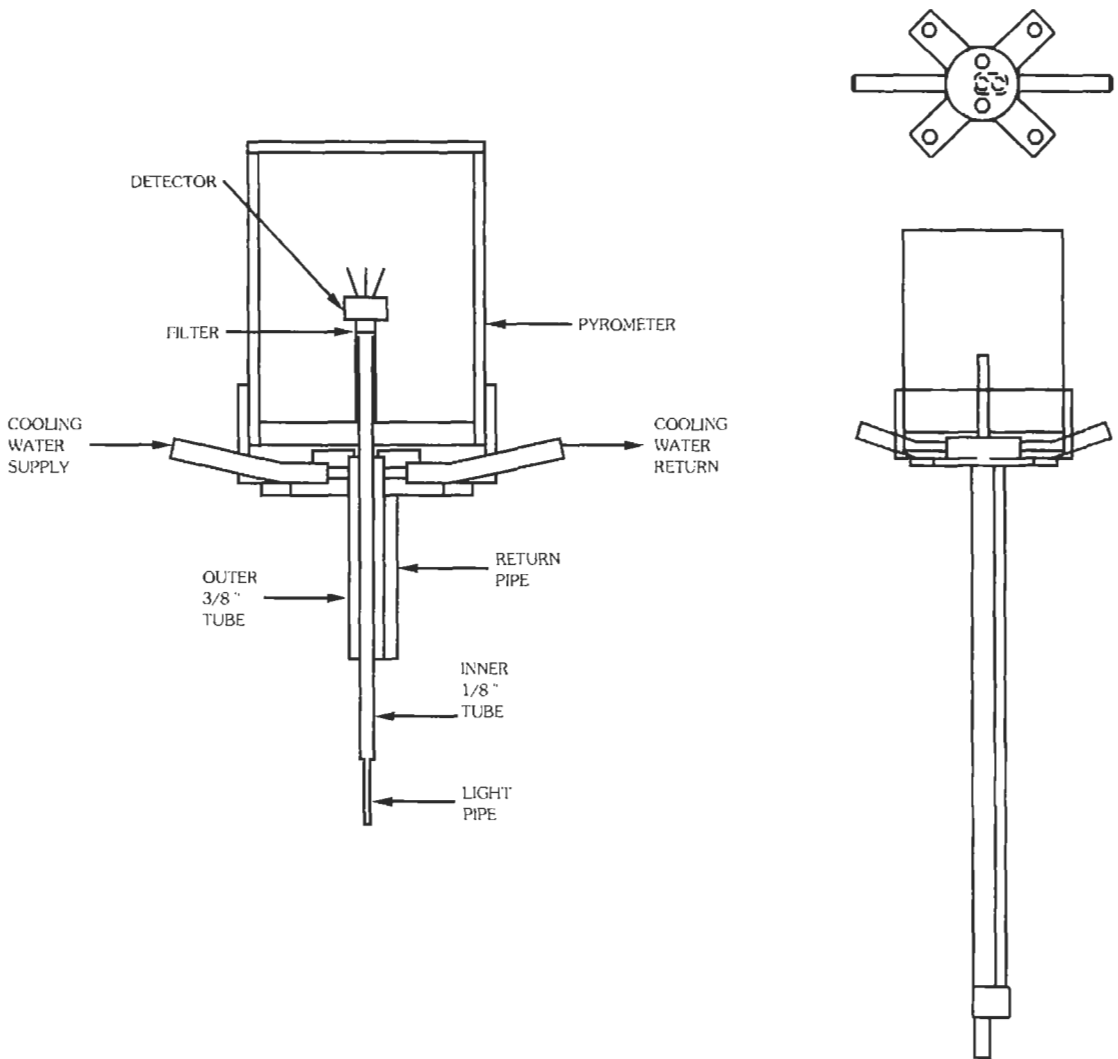


Fig. 67. Detail of a water-cooled light pipe, filter, and pyrometer assembly for multipoint temperature measurement.

The temperature of the water-cooled light-pipe tip, facing the wafer, is maintained at less than 80 °C. The temperature readout from each pyrometer is relayed to the process equipment with the process computer providing a temperature set-point to the controller. Any corrections for surface roughness, emissivity, lamp interference, or window temperature are then used to correct the pyrometer temperature reading.

The real-time emissivity measurement system is shown in Fig. 68.* This technique involves the use of a chopped CO laser beam, at 5.4 μm , to probe the wafer back-side surface during processing. This choice of wavelength eliminates the measurement of reflected lamp light, since the quartz jackets of the tungsten-

halogen lamp bulbs provide a necessary cutoff of direct filament light at above 3.5 μm . The incident laser beam power is measured using a ZnSe beam splitter and an infrared (IR) detector. The light transmitted through the beam splitter is then focused onto a chalcogenide fiber and transmitted to the probe tip at the wafer end. Light that leaves the fiber tip and is transmitted through the chamber window arrives at the wafer. The chamber window is typically fused quartz, which has a small thermal expansion coefficient. The window material where the 5.4- μm beam is incident, however, needs to allow maximum transmission at this wavelength. A thin (1 mm) sapphire section is ideal, given the inertness of sapphire. The light arriving at the wafer surface is partly reflected, absorbed, or transmitted. Since the transmitted portion is essentially

* Reference 51.

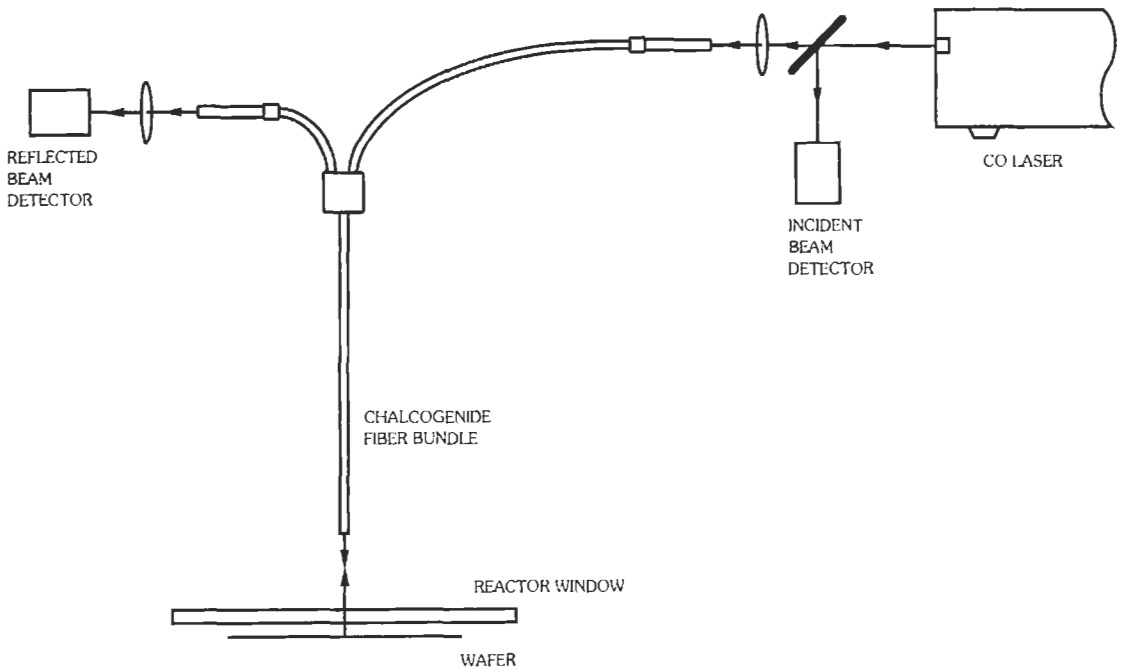


Fig. 68. Schematic of a 5.4- μm laser based emissivity and temperature measurement system, using chalcogenide IR fiber.

eliminated above 600–700°C, it can be neglected in most RTP applications. Additional optical hardware can be arranged to measure the transmitted portion if the operating temperature is sufficiently low. However, the transmission measurements are not required at all if the semiconductor substrates are heavily doped (e.g., p⁺/p⁺ epi material used for CMOS/BiCMOS technologies). Therefore, with the transmission neglected, the reflected light is collected via a fiber bundle arrangement at the probe tip and relayed back to a suitable detector. Spectral wafer emissivity can be extracted based on measurements of the reflected and the incident beam power levels. The reflected beam detector is also used to measure the light emitted by the heated wafer at the 5.4- μm wavelength, such that radiance and emissivity are measured simultaneously, allowing precise calculation of wafer temperature.

It is essential to note that the surface roughness of the back side of the wafer affects the measurement of the reflected light, since the scattering profile of the reflected light generally depends on the surface roughness. The reflection from a typical wafer back side at 5.4 μm contains a strong specular component, with a smaller but significant diffuse component. The diffuse component is not measured by the reflection bundle; consequently, a preprocess surface roughness measurement is used to provide a ratio of expected diffuse-to-specular reflection.

Fiber-Optic Scatter Sensor—The optical reflection and transmission properties of semiconductor

wafers have important effects on various process variables such as wafer temperature, dynamics of wafer heating and cooling, and fabrication process uniformity in RTP. For instance, knowledge of these properties is required to evaluate substrate emissivity and temperature. Reflection and transmission parameters depend not only on the intrinsic properties of the surface films but also on the roughness present on the wafer surface. Surface roughness may be due to the wafer back side or a polycrystalline film present on the wafer (front side or back side). Reflection and transmission of wafers with a smooth surface will be essentially specular or coherent. In wafers with a rough surface film (or a rough back side), only a fraction of the reflected signal will be specular; the remaining fraction will be scattered within a space cone angle around the specular direction. The latter will increase as the surface roughness increases (or the light wavelength decreases). Surface roughness can also have an effect on the substrate emissivity and should be taken into consideration in any real-time pyrometry-based temperature sensor. In situ preprocess/postprocess measurement of the surface scattering characteristics and reflection properties of a wafer can be used to monitor its surface roughness and spectral emissivity. An important application of this technique is that measurement of these properties can be used in the real-time determination of the emissivity and temperature of a silicon substrate at different steps during a process flow. Another important application relates to measurements of CVD films on silicon wafers. For

instance, it has been shown that surface roughness of LPCVD W films is directly proportional to the film thickness.* Thus, an in situ noninvasive sensor that can evaluate surface roughness and spectral reflectance of semiconductor wafers can be used as a tool to monitor thickness and quality of LPCVD W (and other polycrystalline films). This can be done both for process control and for process prognosis/diagnosis purposes.

An in situ sensor to determine surface roughness, reflectance, and spectral emissivity of silicon wafers with various surface films has been developed.† The sensor has been designed for implementation in the vacuum load-lock chambers of SWP reactors. Operation of this in situ sensor is based on the relation between the surface roughness and the specular and scattered reflection properties of semiconductor wafers for an incident monochromatic electromagnetic wave. Thus, measurement of the total reflectivity as well as the scattering and specular parameters can be correlated to surface roughness and spectral emissivity of silicon substrates with various surface films. These parameters depend on the type and thickness of films on both the front side and back side of these substrates.

Metal Sheet Resistance Sensor—Wafer-to-wafer repeatability of CVD processes including CVD metal is of primary concern in ULSI manufacturing. This subsection describes a sensor device that can perform real-time monitoring of the CVD metal process with a high degree of accuracy and reliability.‡ A variation of the well-known four-point probe technique commonly used to measure metal-film sheet resistances *ex situ* has been employed to measure the mean sheet resistance of blanket metal films across the entire wafer during metal CVD. These measurements are used to determine the deposition process end point, thus improving wafer-to-wafer repeatability compared to standard timed process end points. This method differs from four-point probe measurements in that the driving current is passed between the same probe points across which the voltage drop is measured. As the current is passing through the voltage measurement probes, there is an inherent contact potential between the probe and the measurement surface. This technique is viable only if this contact potential is repeatable enough to be accounted for in the conversion algorithm. Extensive experimentation indicates that this is indeed the case.

A linear equation relates the mean (room temperature) sheet resistance to the voltage drop across the blanket metal film for a given current:

$$V_p = A(T) \times R_{sh} + B(T) \quad (\text{Eq. 6})$$

where

V_p is the measured probe voltage,

R_{sh} is the mean room-temperature sheet resistance,

$A(T)$ and $B(T)$ are temperature-dependent coefficients that are linear functions of temperature.

This equation underscores the importance of accurate temperature measurement for taking full advantage of this technique. Excellent agreement with experimental results has been achieved for CVD-W films. This agreement continues to hold even at elevated temperatures. An example of this is given in Fig. 69, which plots the probe voltage at 450 °C versus the mean room-temperature sheet resistance.

One possible problem relates to the effect of contacting metal pins with the wafer edge surface during the deposition process. There is a shadowing effect if the 0.8-mm diameter probe is kept in contact with the wafer edge during the entire process. The probe acts as a heat sink and shadow mask at that location, inhibiting tungsten nucleation. To overcome this effect, the probe can be retracted from contacting with the wafer when measurements are not being taken, especially at the beginning of the process when nucleation is taking place. Measurements may then be taken at various points during the process, indicating when the desired process end point is near. At this time, the probes are finally extended to measure the exact CVD process end point, with negligible disturbance to the wafer.

The initial experiments utilizing this technique for CVD process end-pointing have resulted in improved wafer-to-wafer repeatability by a factor of four over the conventional timed end-pointing technique. It is also effective for process prognosis and diagnosis, providing real-time information about the deposition rate.

RTP for Integrated CMOS Processing—The RTCVD processes described here prevent back-side depositions by wafer back-side and window purge. The RTCVD process modules (LPCVD amorphous and polycrystalline silicon, LPCVD tungsten, LPCVD oxide, LPCVD nitride) require in situ chamber cleaning for improved wafer-to-wafer process repeatability and reduced particle counts. Various RTP and RTCVD processes (e.g., dry and wet RTOs, LPCVD amorphous silicon, LPCVD polysilicon, source/drain and gate anneal, and CMOS well formation) have been developed that use wafer rotation for enhanced process uniformity.§ In situ sensors can be used to monitor wafer emissivity, film thicknesses, and conductive layer sheet resistance values.

An important requirement in RTP-based process integration is good temperature and process uniformity control throughout the entire integrated process flow. Conventional integrated processes based on the batch equipment (or a combination of batch and SWP) modify the wafer back-side structure repeatedly throughout

* Reference 59.

† Reference 51.

‡ Reference 51.

§ Reference 50.

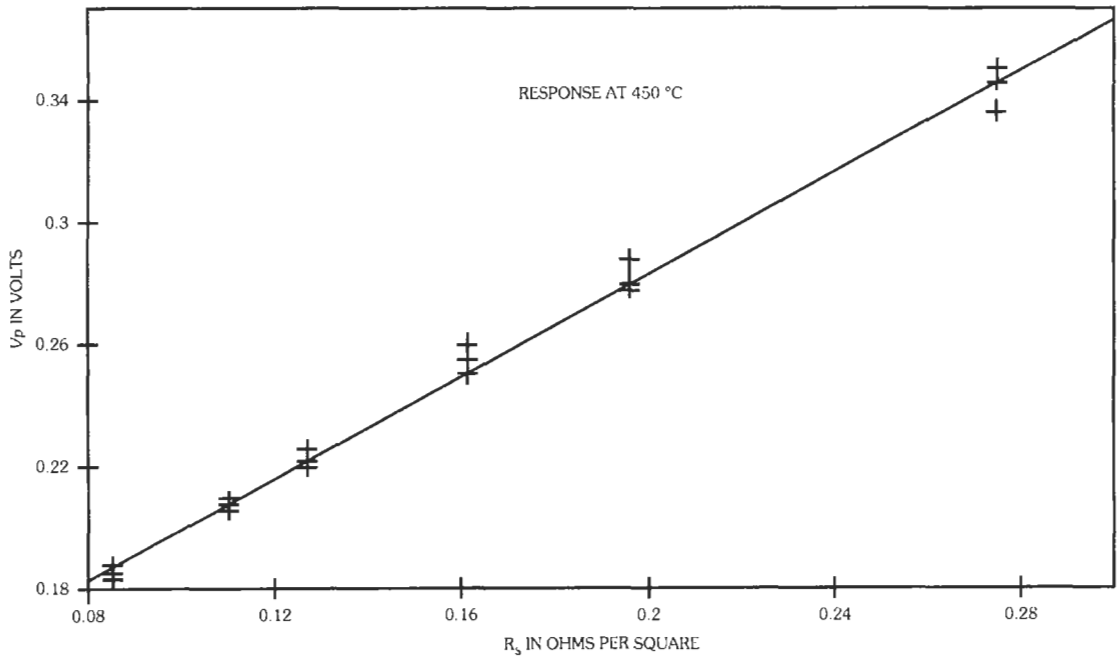


Fig. 69. Two-point probe voltage at 450 °C vs the mean room-temperature tungsten sheet resistance.

the process flow. For instance, at various process steps the wafer back sides may be bare or may have a combination of silicon dioxide, silicon nitride, and/or polysilicon layers. Moreover, the thicknesses of these layers can vary as the process flow proceeds owing to many wafer cleaning and etch processes (wet and plasma). The integrated process flow can also cause back-side layer thickness nonuniformities. These variations translate into large-scale wafer back-side emissivity changes, which can cause significant RTP temperature and uniformity control problems. The process-control difficulties can be somewhat alleviated by using multipoint pyrometry with real-time emissivity compensation. Even with emissivity compensation, the wafer back-side nonuniformities can degrade RTP uniformity as the result of distorted or modified optical energy absorption and radiative losses.

To minimize the above-mentioned complications, a sub-0.50 μm CMOS process integration has been performed that uses starting silicon wafers with customized back-side seal layers.* The $p^- \text{epi}/p^+$ wafers are processed to receive a uniform back-side stack of 1300 Å silicon nitride on 1000 Å silicon dioxide. This back-side seal serves several useful purposes. (1) It prevents dopant outdiffusion from the heavily doped substrate; (2) it prevents back-side emissivity variations during oxidation and oxide deglaze processes (due to the oxidation and etch resistance of silicon nitride); and (3) it

ensures uniform and stable back-side emissivity throughout the entire integrated RTP-based process flow. It should be emphasized that the one-sided depositions in RTCVD processes help to maintain a relatively constant and uniform back-side emissivity on the wafers in process. The above-mentioned back-seal process also eliminates any process nonuniformities caused by nonuniform wafer emissivities. Small wafer-to-wafer back-side emissivity variations are easily tracked and compensated via suitable in situ sensors. All the RTP-based processes described here employ reflective showerheads facing the front side of the wafer. As a result, the front-side pattern effects on process uniformity are essentially eliminated (as the result of the optical black-body cavity produced between the wafer front side and the reflective showerhead).

Numerous RTP-based fabrication processes have been developed for an SWP-based sub-0.5 μm twin-well double-level metal CMOS technology.† Fig. 70 shows a schematic cross-sectional view of this technology. Two technology versions have been investigated: one has buried-channel PMOS and surface-channel NMOS transistors without salicide and with silicided contacts; the other has surface-channel NMOS and PMOS transistors with salicide. This device cross section shows multiple RTP-fabricated device regions including: (1) gate electrode (LPCVD amorphous and polycrystalline silicon), (2) gate

* Reference 50.

† Reference 50.

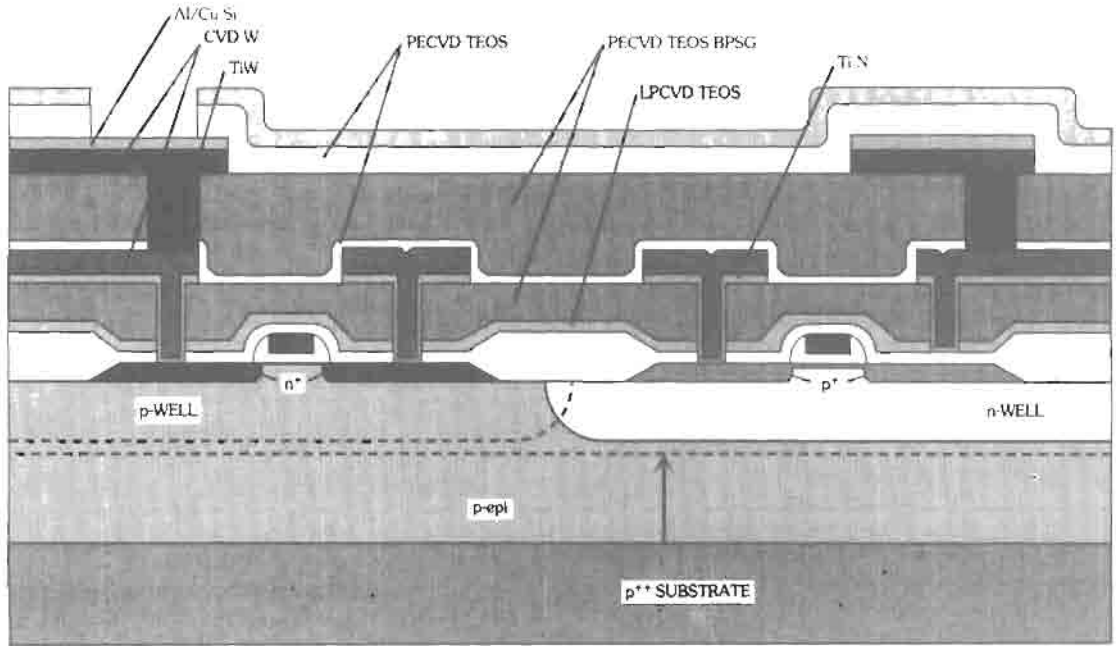


Fig. 70. Schematic cross-sectional diagram of a sub-0.50 μm RTP-based CMOS technology.

dielectric (dry RTO), (3) n and p wells (wet RTO for initial oxide and RTP oxynitridation-enhanced diffusion for well formation), (4) source/drain junctions (S/D RTA), (5) silicided contacts (RTP silicide react), (6) poly-buffer-LOCOS or PBL isolation (LPCVD silicon nitride on LPCVD polysilicon on dry RTO oxide oxidation mask), (7) oxide spacers (LPCVD silicon dioxide), (8) RTP forming gas anneal, and (9) double-level metal system (includes LPCVD tungsten). The total number of RTP steps (depositions, oxidations, and anneals) in the integrated CMOS flow exceeds 15. A single-wafer RTCVD-based epitaxial growth process was also developed for blanket boron-doped epitaxial silicon deposition to fabricate starting epitaxial wafers.*

RTP-Based Process Development Overview—

Table 16 shows a list of various RTP-based processes developed for an integrated sub-0.50 μm CMOS flow† along with their specific applications and process parameter domains.

The epitaxial growth process was developed using SiH_2Cl_2 (DCS)/ H_2 chemistry at low pressures (5–15 torr). A typical in situ-doped process employs an initial low-temperature (650–750 °C) in situ oxide removal process step (GeH_4/H_2) followed by epitaxial silicon deposition ($\text{DCS}/\text{H}_2/\text{B}_2\text{H}_6$) at 1000 °C for blanket p-epi on p⁺ substrate.

The RTO processes are used for oxide growth in the thickness range of 75 to 250 Å. Gate oxide (80 Å) and PBL initial oxide (90 Å) are grown by using dry RTO. It was determined that the dry and wet RTO processes should not be performed in the same process chamber because of possible cross-contamination and process memory effects of steam. As a result, dedicated reactors were used for these two processes. The wet RTO process employs an external pyrogenic steam generator where H_2 and O_2 flow rates can be controlled. The wet RTO process is used to grow the initial 250 Å oxide for CMOS well formation. The well implants are performed through this oxide. Moreover, this oxide provides the oxynitridation-enhanced diffusion (ONED) effect for well formation using ammonia anneal. Wet RTO is also used to grow a sacrificial (dummy) gate oxide after the isolation process module and before the dry gate oxidation step. To maintain reasonable process throughputs, both dry and wet RTO processes are performed near the atmospheric pressure (650 torr). Much lower pressures are not acceptable because the oxidation rates are significantly reduced and the resulting RTO times become excessively long.

Source/drain rapid thermal anneal (S/D RTA) is used for activation and formation of shallow p⁺n and n⁺p source/drain junctions. This process is performed at 950–1000 °C in argon ambient. Although the process can be performed at lower pressures (1 torr), high pressures (650 torr) are preferred in order to ensure sufficient cooling of the showerhead plate facing the wafer. This choice degrades the process throughput in

* Reference 49.

† Reference 50.

TABLE 16. RTP-BASED PROCESSES DEVELOPED FOR A SUB-0.5 μm CMOS TECHNOLOGY

RTP-Based Processes	Applications	RTP Parameter Domain
Epitaxy/in situ clean	epi material (p^+/p^+)	700–1000 °C, low-pressure DCS/ H_2 , GeH_4/H_2
Dry RTO	Gate oxide, PBL oxide	1000–1050 °C, 650 torr oxygen ambient
Wet RTO	Thick oxides (ONED tank)	900–950 °C, 650 torr pyrogenic steam
Source/drain RTA	S/D activation, gate doping	900–1000 °C, 80–650 torr argon ambient
RTP tank formation	CMOS n & p well formation	1050–1100 °C, 650 torr ammonia ambient
LPCVD polysilicon	CMOS gate formation	650–700 °C, low pressure SiH_4/Ar , $\text{Si}_2\text{H}_6/\text{Ar}$, (SiH_4/H_2)
LPCVD amorphous Si	CMOS gate formation	500–590 °C, low pressure SiH_4/Ar , $\text{Si}_2\text{H}_6/\text{Ar}$, (SiH_4/H_2)
LPCVD tungsten	Multilevel metal	300–550 °C, low pressure $\text{SiH}_4/\text{H}_2/\text{WF}_6$
TiN/TiSix RTA react	Salicide, silicided contacts	600–750 °C, low pressure N_2 or NH_3 ambient
RTP sinter (FGA)	Forming gas anneal	450–500 °C, low pressure FGA or N_2 ambient
LPCVD nitride	PBL nitride deposition	800–850 °C, low pressure DCS/ NH_3 or SiH_4/NH_3
LPCVD oxide	Oxide spacers, undoped oxide	700–750 °C, low pressure TEOS/ O_2

standard mode of processing, mainly because of the time intervals required to ramp up and ramp down the chamber pressure between vacuum (for wafer transport) and process pressure (for wafer processing).

CMOS n and p well formation is accomplished by RTA in an ammonia ambient following the n and p well implants. This process is the highest “Dr” RTP fabrication step used in this integrated CMOS flow. It is usually done at 1100 °C for 5 minutes in 650 torr of ammonia in the presence of a 250 Å oxide grown using a wet RTO process. The ONED effect reduces the time/temperature needed to form the well profiles with the desired junction depths. It is essential that the RTP/ONED process is performed at higher pressures (e.g., 650 torr) to induce the maximum ONED effect for reduced RTP time.

LPCVD polysilicon and α -Si processes are used for gate electrode formation based on a split-deposition gate process. Both of these processes are performed at a low (15 torr) pressure using a mixture of SiH_4/Ar . Argon is preferred over hydrogen as a carrier gas. This is because the deposition rates are somewhat higher in SiH_4/Ar than in SiH_4/H_2 (hydrogen causes a chemical retardation effect since it is a byproduct of the deposition process). Moreover, hydrogen may cause some gate oxide reliability degradation at higher temperatures (above 650 °C). The polysilicon

and amorphous silicon deposition processes are performed at 650 °C and 560 °C, respectively. The lower α -Si deposition temperature ensures continuous, thin (≈ 200 Å), crystallite-free layers of α -Si directly over the thin gate oxide for its protection during a subsequent photoresist processing and patterned ion implantation step. Deposition processes based on an alternative disilane/argon chemistry have also been developed. The disilane-based LPCVD-Si processes provide deposition rates comparable to those of the silane-based processes using smaller Si source gas flows.

Tungsten LPCVD is employed at low pressures (0.5–8 torr) for use in a double-level metal system. The first metal (CVD-W) level is deposited directly over an RTP-reacted TiN layer using a two-step process (silane-reduced WF_6 followed by a hydrogen-reduced WF_6). The second metal layer employs a layered structure including an underlayer of CVD-W, also deposited by RTCVD.

The TiN/TiSi₂ react process has been the most widely used application for RTP and the first RTP used for commercial product manufacturing. The RTP react process described here is performed at 650–750 °C in a low-pressure (1 torr) nitrogen ambient. This process can also be done at higher N_2 pressures (e.g., 650 torr) and/or in an ammonia ambient. It is used for formation

of silicided contacts to the source/drain junctions as well as formation of TiN contact barrier and nucleation layer for the metal-I CVD-W layer.

The RTP sinter process is performed in a forming gas ambient at 1–650 torr and at 450–475 °C. The higher process pressures (e.g., 650 torr) enhance the amount of hydrogen flux to the gate dielectric interface. This process step is used to reduce the gate oxide surface-state density and restore the desired subthreshold and threshold voltage characteristics. The RTP sinter process is used following the last damage-producing plasma etch process in the integrated flow (e.g., after metal 2 etch and/or after the passivation-overlayer etch).

The nitride LPCVD process can be performed in SiH_4/NH_3 or DCS/NH_3 ambients. This process is usually carried out at 800–850 °C and at low pressures (1–5 torr). The oxide LPCVD process is performed at 700–750 °C with TEOS/O_2 used for dielectric spacer formation.

The above-mentioned RTP-based processes can be conducted in one of two modes of operation: (1) standard pressure cycling (SPC) mode, and (2) reduced pressure cycling (RPC) mode. The SPC mode operates as follows: (1) transfer the first wafer from the vacuum load-lock to the vacuum process chamber, (2) start process gas flows and stabilize pressure, (3) start wafer processing by activating the lamp energy source, (4) stop the process by inactivating the lamp energy source, (5) stop the gas flows and pump down to vacuum, (6) transfer the processed wafer from the process chamber to the vacuum load-lock chamber, and (7) proceed with processing the next wafer in the cassette. As a result, the process throughput in the SPC mode is limited by the time intervals needed to ramp the chamber pressure between vacuum and process pressure. This is particularly true for the higher-pressure processes such as RTO. The RPC mode of operation can increase the effective process throughput by eliminating these pressure cycling segments from wafer to wafer. For instance, in a dry RTO process, the load-lock chamber and process chamber are first pumped down to vacuum after loading the wafer cassette into the load-lock chamber. Before or after transferring the first wafer to the process chamber, oxygen flow is initiated into both the load-lock and process chambers, and their pressures are stabilized at the desired process pressure. Wafer processing and transport are done at constant pressure without any pressure cycling. This mode of operation will increase the effective process throughput. The RPC mode is compatible with the dry oxidation and anneal processes that employ oxygen, argon, ammonia, nitrogen, and forming gas (e.g., dry RTO, S/D and gate RTA, RTP/ONED tank formation, silicide). However, the SPC mode is the preferred choice with all the RTCVD processes because of safety and reactor cleanliness/reliability considerations.

The RTP reactor count and process partitioning should be done based on the considerations to mini-

mize the number of RTP modules without compromising the quality of fabrication processes and without cross-contamination issues. As a result, the RTP modules used in an integrated CMOS flow and their dedicated processes may be chosen as follows:*

- LPCVD-Si module for LPCVD α -Si, polysilicon, and in situ germane cleaning
- Dry RTO module
- Wet RTO module
- S/D and gate RTA and ONED tank module
- Metal anneal module for silicide react and sinter
- LPCVD tungsten module
- LPCVD silicon nitride module
- LPCVD silicon dioxide module
- Epitaxial growth module (with germane cleaning)

These are a total of nine RTP modules used for the integrated sub-0.5 μm CMOS flow. Besides these modules, one other thermal process module is also employed. A single-wafer high-pressure oxidation system (SWP-HIPOX) has been used to grow thick field oxides (rates of 1000–3000 Å/min) and to perform BPSG reflow. This can be considered a tenth RTP module even though the SWP-HIPOX module employs resistive wafer heating.

In this minifactory, wafer transport between critical sequential steps (e.g., from gate RTO to LPCVD-Si gate electrode) is performed in a controlled vacuum environment to minimize native oxide growth and particles.

Selected RTP Results—As discussed earlier, the wafer back-side emissivity has significant effects on RTP uniformity and control. Fig. 71 shows the thickness of gate oxide grown at 1000 °C (for a fixed time) with closed-loop pyrometry control without in situ emissivity compensation versus thickness of back-side oxide layer (back-seal layer with 1300 Å nitride on top of oxide with varying thicknesses). These data indicate well over 10% variation in gate oxide thickness when the back-side oxide thickness varies between 800 Å and 1400 Å. The use of a stable back-seal structure in conjunction with in situ emissivity measurement and compensation is an effective solution to this problem.

Compared to dry RTO, wet RTO is capable of growing oxides under reduced temperature/time conditions. As an example, Fig. 72 shows the oxide thickness versus wet oxidation time at an oxidation temperature of 950 °C and a total pressure of about 600 torr. These oxidations were performed using a nonstoichiometric hydrogen-to-oxygen ratio of 1. Therefore, the oxidation ambient consisted of a mixture of pyrogenic steam and excess oxygen. Under these conditions, a 5-minute process could grow about 250 Å of oxide. A 950 °C/5 min dry RTO process would grow about 80 Å.

The wet RTO growth rate is also affected by the hydrogen-to-oxygen flow ratio. Fig. 73 shows the

* Reference 50.

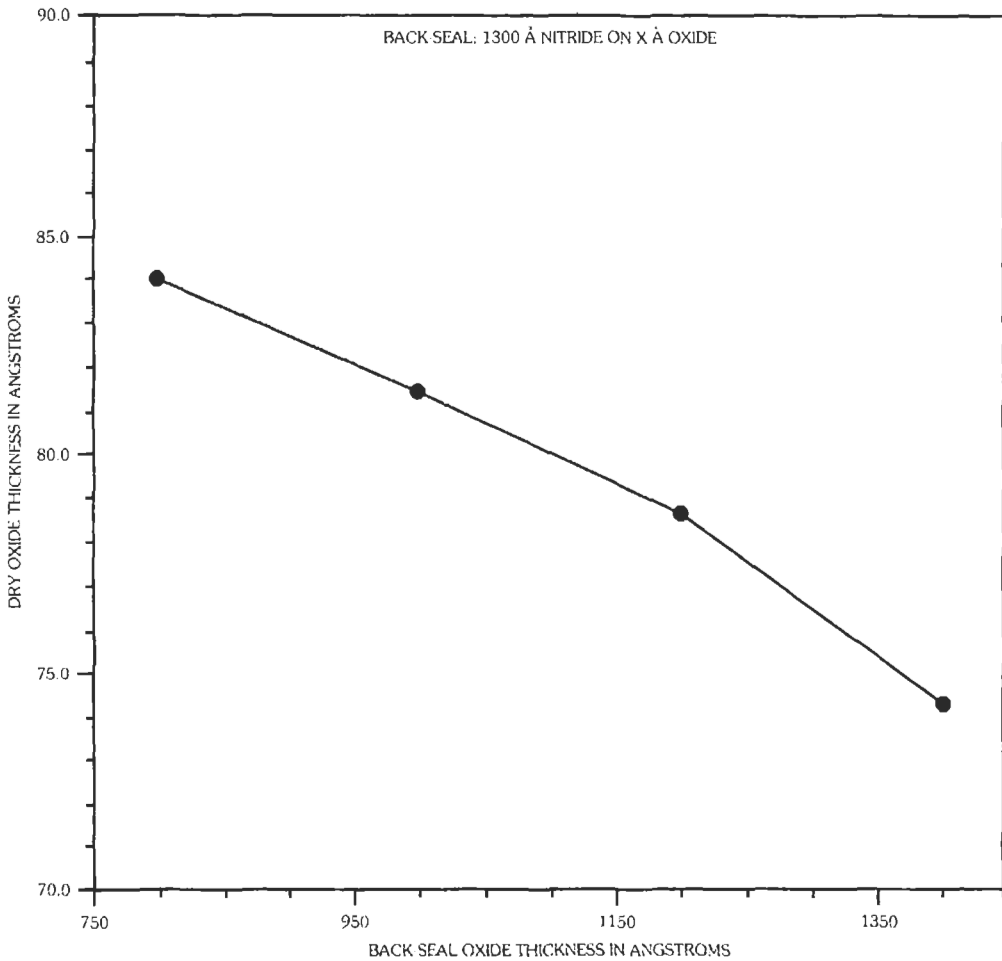


Fig. 71. Effect of variation in back-side film thicknesses on gate oxide thickness.

oxide thickness data versus flow ratio for 950 °C/5 min oxidations at 600 torr. A flow ratio of zero corresponds to steam-free dry oxidation, and a flow ratio of 2 corresponds to stoichiometric steam generation in the torch. The oxide thickness corresponding to a flow ratio of zero is over 100 Å. This is thicker than the thickness expected for a fully dry RTO process under similar temperature and time conditions. The difference can be attributed to the residual water vapor in the wet RTO process chamber, which causes some growth rate enhancement. This is clearly a good reason for not using one chamber for both dry and wet RTO processes. A stoichiometric flow ratio of 2 offers the largest thickness because it has the highest steam partial pressure (600 torr without any excess oxygen). For process repeatability reasons, a flow ratio of 1 has been observed to be superior to the stoichiometric ratio of 2. The local minimum corresponding to a flow ratio of 1.4 may be attributed to steam condensation on cold

surfaces, resulting in reduced steam partial pressure in the process chamber.

Figs. 74 and 75 show examples of metal-oxide-semiconductor (MOS) breakdown distribution data for capacitors with various gate electrode formation process splits (80 Å gate oxides). The silicon gates were formed in two separate deposition steps (split-deposition process): a first layer of thin (≈ 200 Å) amorphous silicon and a second layer of polysilicon (≈ 2800 Å). The data in Fig. 74 are for split-deposition gate MOS capacitors with various gate electrode processing conditions, with and without an in situ germane-based native oxide removal process. In this figure, "PR" represents photoresist processing (coat and strip) over the initial amorphous Si layer, and "HF" represents wet HF native oxide removal (in dilute HF), before the subsequent polysilicon RTCVD (AVP poly) at various temperatures. The photoresist process was used to simulate photoresist-induced contamination over the gate α -Si layer.

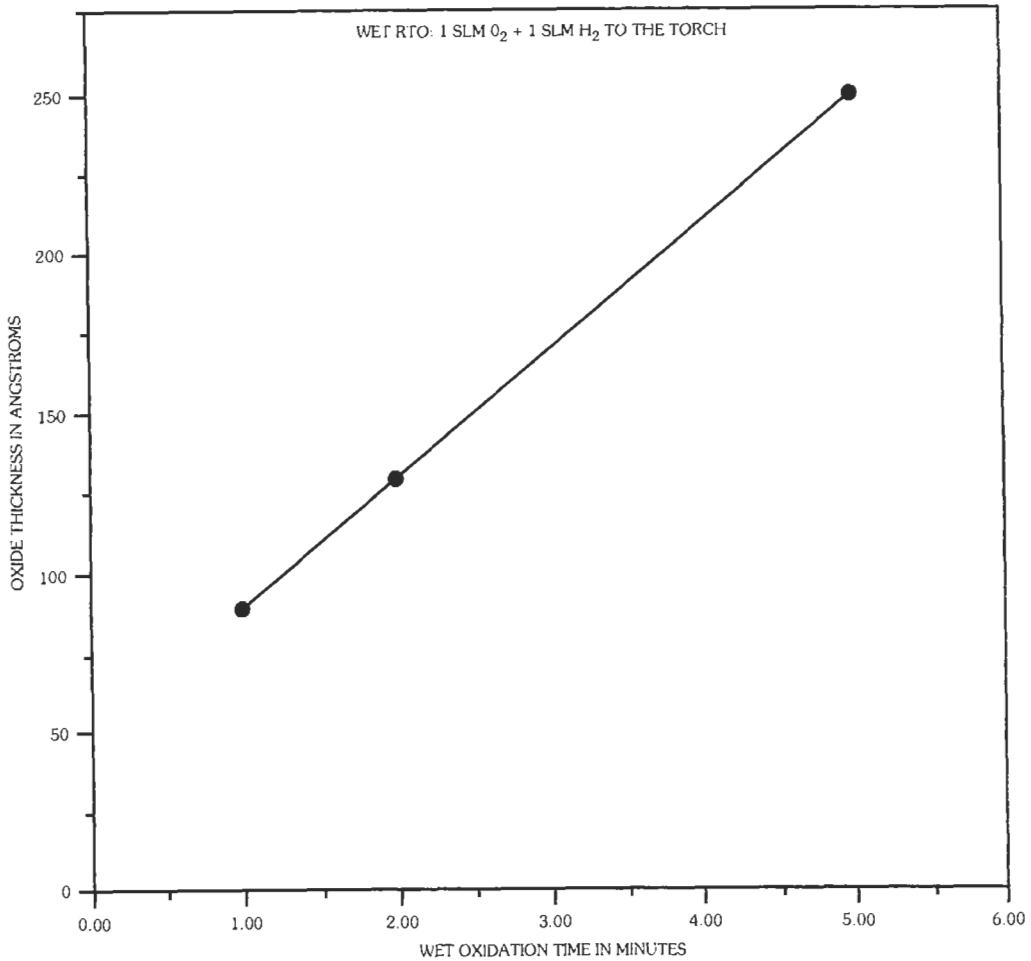


Fig. 72. Oxide thickness vs wet RTO time at 950 °C and 600 torr. The oxidations were performed using 1 slm (standard litre per minute) of O₂ and 1 slm of H₂ to the pyrogenic steam generator.

These data indicate that the devices with in situ germane cleaning (prior to polysilicon RTCVD) have breakdown distribution comparable to that for other devices without germane cleaning. The use of germane cleaning ensures MOS process reproducibility and repeatable gate patterning by removing the oxide from the original α -Si/polysilicon interface. The devices with gate polysilicon deposited at various temperatures (590, 620, and 650 °C) also show comparable breakdown characteristics. The low-field breakdowns in Fig. 74 could be mostly attributed to particles. The breakdown distribution data in Fig. 75 are for split-deposition gate devices with various methods of gate formation. These include RTCVD (silane versus disilane chemistry) and furnace (silane chemistry) α -Si processes, and RTCVD (silane versus disilane chemistry) polysilicon processes at various temperatures. The breakdown distributions are nearly comparable and indicate acceptable characteristics for silicon RTCVD.

RTP has also been employed to form silicided shallow source/drain junctions with low leakage current densities (as low as 5 nA/cm²).^{*} The buried-channel PMOS version of the integrated sub-0.50 μ m CMOS flow employs nonsalicide source/drain junctions with silicided contacts whereas the surface-channel PMOS/NMOS version consists of silicided structures. Figs. 76 and 77 show the sheet resistance values of RTP-reacted titanium layers over doped (p⁺ and n⁺) silicon and silicon dioxide, respectively. The RTP react processes were performed in 1 torr nitrogen using various one-step and two-step processes. For an initial titanium layer thickness of 400 Å, the two-step processes result in TiN/TiSi_x sheet resistance values of less than 4 Ω /sq. on n⁺p (As/P doped) and p⁺n junctions. For 800 Å starting titanium on oxide, the RTP react process results in a relatively oxygen-free TiN layer with

* Reference 50.

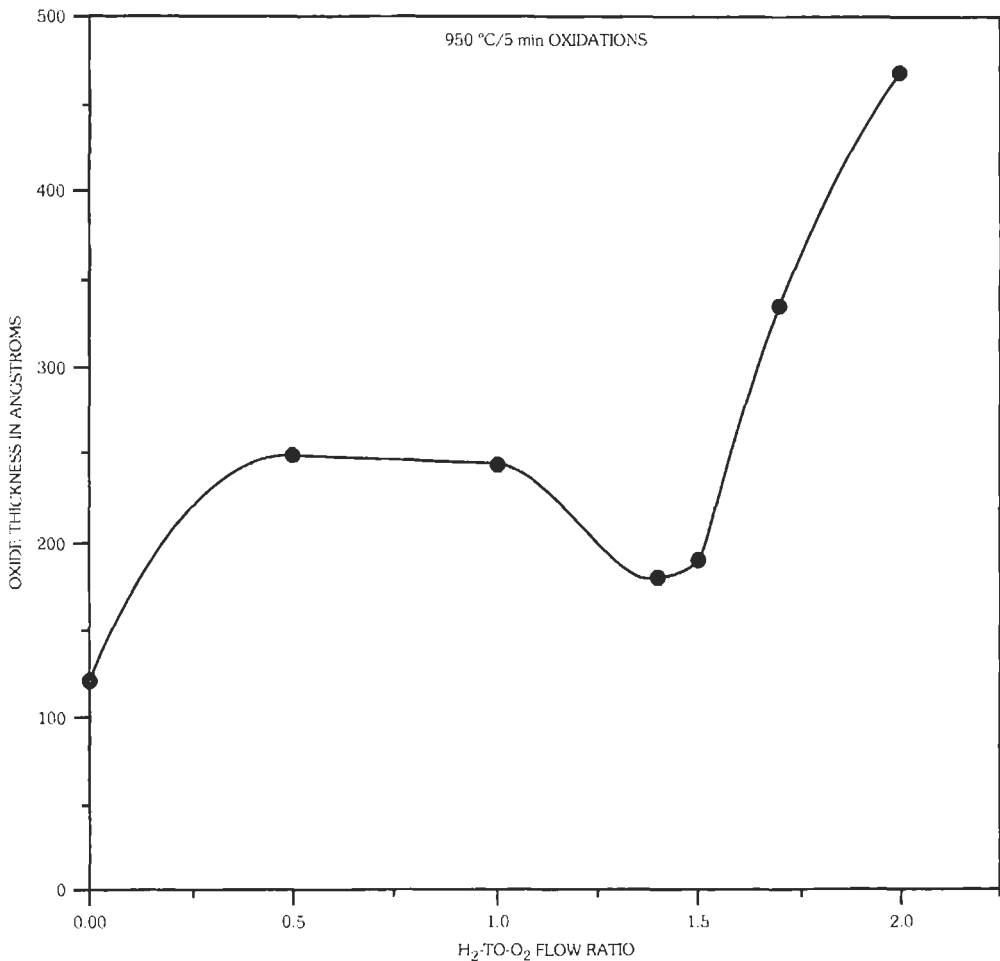


Fig. 73. Effect of hydrogen-to-oxygen flow ratio on wet oxide thickness for 950 °C/5 min oxidations at 600 torr total pressure.

sheet resistance values on the order of 20 Ω/sq. The first RTCVD tungsten metal layer is deposited over the RTP-reacted layer (TiN over the interlevel oxide dielectric and TiN/TiSi₂ in the contact holes).

Suitable tungsten RTCVD processes have been developed for a double-level metal system for the sub-0.50 μm CMOS technology. For the first metal level, RTCVD tungsten is used to fill the contact holes and form the metal lines over TiN (using the same deposition step).

Single-Wafer Surface Cleaning

Many steps in the processing of silicon devices are preceded by a “clean” to remove trace surface contaminants or to chemically pretreat the wafer surface. Traditionally, such cleans have been based nearly entirely on chemical solutions,* which selectively attack and remove the surface contaminants. A number of factors

render these “wet” cleans increasingly undesirable. The chemical wastes generated by such processes are hazardous to the environment, and their proper treatment and disposal are expensive. Besides the environmental and economic concerns, stringent requirements on device performance also dictate a trend to integrate the processing steps with the corresponding precleans, so that exposure to uncontrolled ambient can be minimized. Equipment and process incompatibility has prevented such integration involving the wet cleans.

Alternative cleans based on gas-phase/vapor-phase chemistries are more amenable to the desired process integration. Notable examples of such “dry” cleans are (1) the incorporation of an in situ high-temperature hydrogen and HCl bake prior to silicon epitaxy and (2) slight contact sputter etch prior to the deposition of metal interconnections. In the remainder of this

* Reference 60.

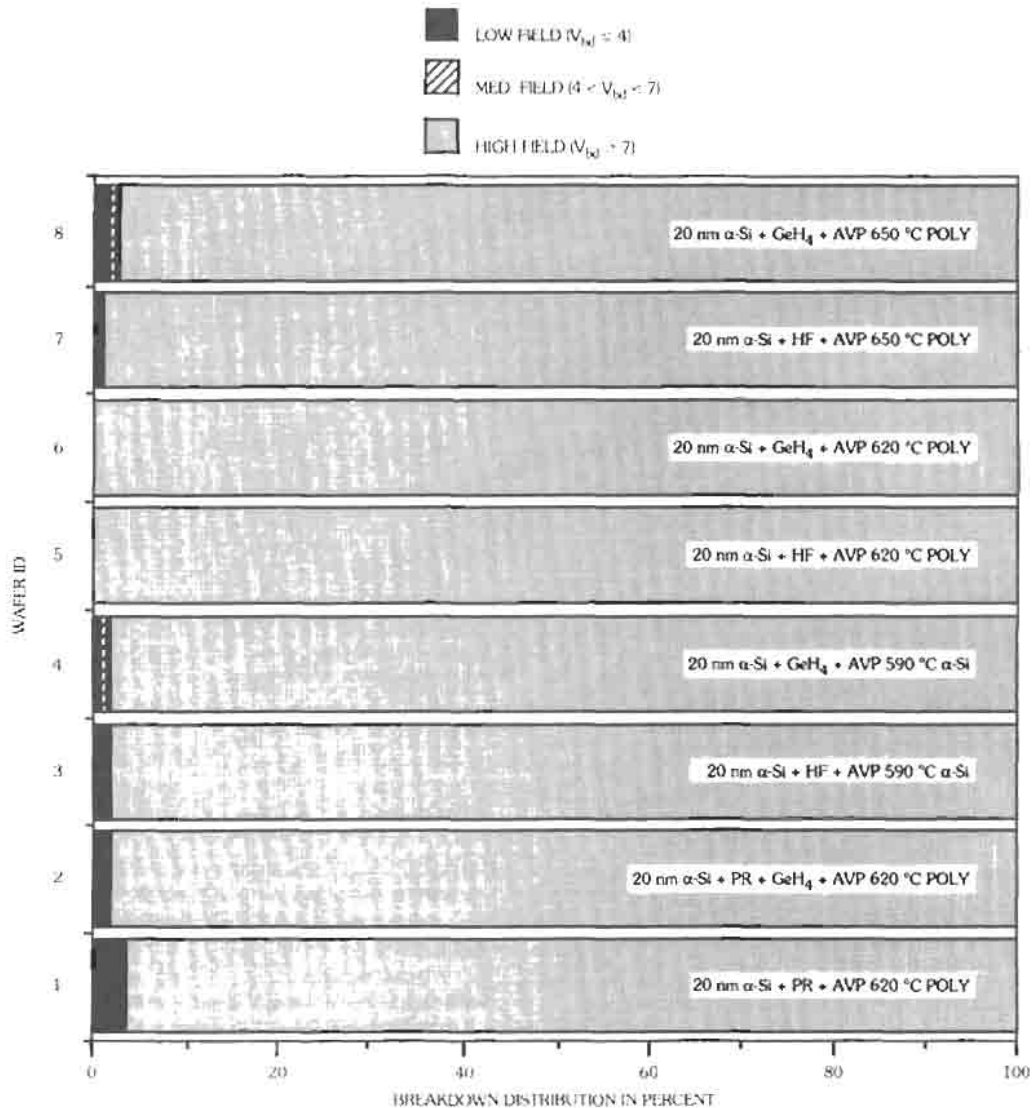


Fig. 74. Breakdown characteristics for split-deposition gate MOS capacitors with 80 Å gate oxides and various polysilicon RTCVD processes. The data are for devices with various surface treatments between the initial amorphous silicon and the subsequent polysilicon RTCVD (germane cleaning, wet HF deglaze, and photoresist coat/strip).

subsection, the possible alternative dry cleans will be discussed along the traditional classification of surface contaminants: native/chemical oxides of silicon, trace organic, and trace metallic contaminants.

Surface Cleaning in IC Manufacturing—The microstructure, material purity, and electrical properties of various interfaces between different layers (semiconductor/dielectric, semiconductor/semiconductor, metal/semiconductor, and others) strongly influence the overall reliability and performance characteristics of ULSI circuits. The critical device

interfaces include polysilicon/single-crystal Si (in bipolar transistors), Si/SiO₂ (in MOS and bipolar transistors), silicide/single-crystal Si, and others. As a result, high-performance device fabrication requires effective preprocess and postprocess cleaning of the semiconductor surfaces to remove thin native oxide layers and other contaminants such as residual metallic and organic impurities.

Fig. 78 shows various possible sources of surface contamination and damage in device processing (represented as "causes") and their corresponding effects (shown as "effects") on the semiconductor surface.

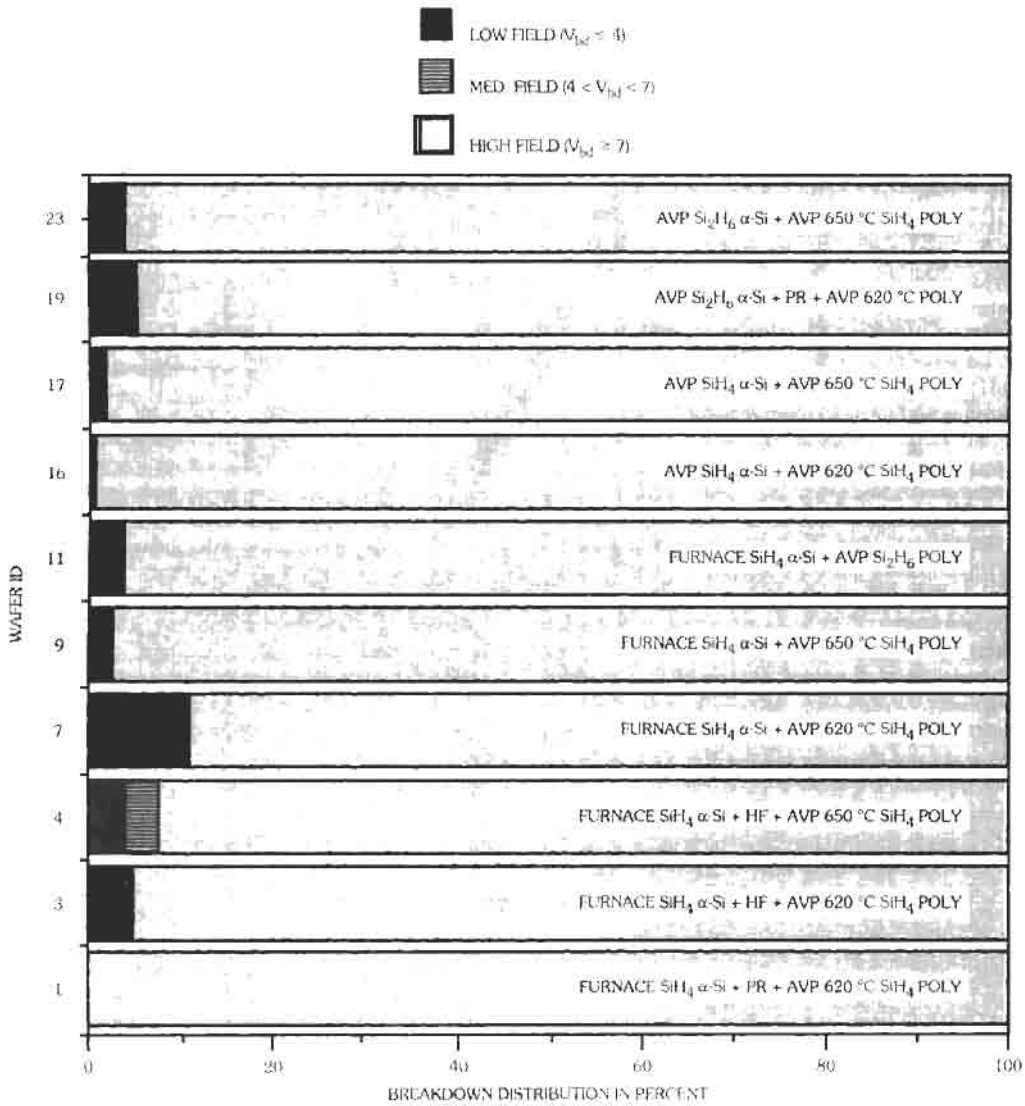


Fig. 75. Breakdown characteristics for split-deposition gate MOS capacitors with 80 Å gate oxides and various amorphous silicon and polysilicon deposition processes.

The dark bullets show definitive cause-effect correlations, and the light bullets represent likely cause-effect links. Wet chemical treatments (e.g., NH₄OH:H₂O₂), atmospheric ambient exposure, low-purity process gases, and poor process equipment base pressure can form thin (up to ≈30 Å) native oxide layers on silicon surfaces.* The precise thicknesses of these native oxide layers depend upon the specific source; for instance, oxides resulting from ambient exposure are

usually less than 12 Å thick,† whereas the chemically grown native oxides resulting from exposure to hot NH₄OH:H₂O₂ may be as thick as 17 Å or more.‡ Native oxide layers are usually detrimental to device performance and increase device parameter spread owing to process repeatability problems.§ A dilute HF treatment is usually used as the last pre-furnace cleaning

† References 63 and 64.

‡ Reference 62.

§ Reference 65.

* References 61 and 62.

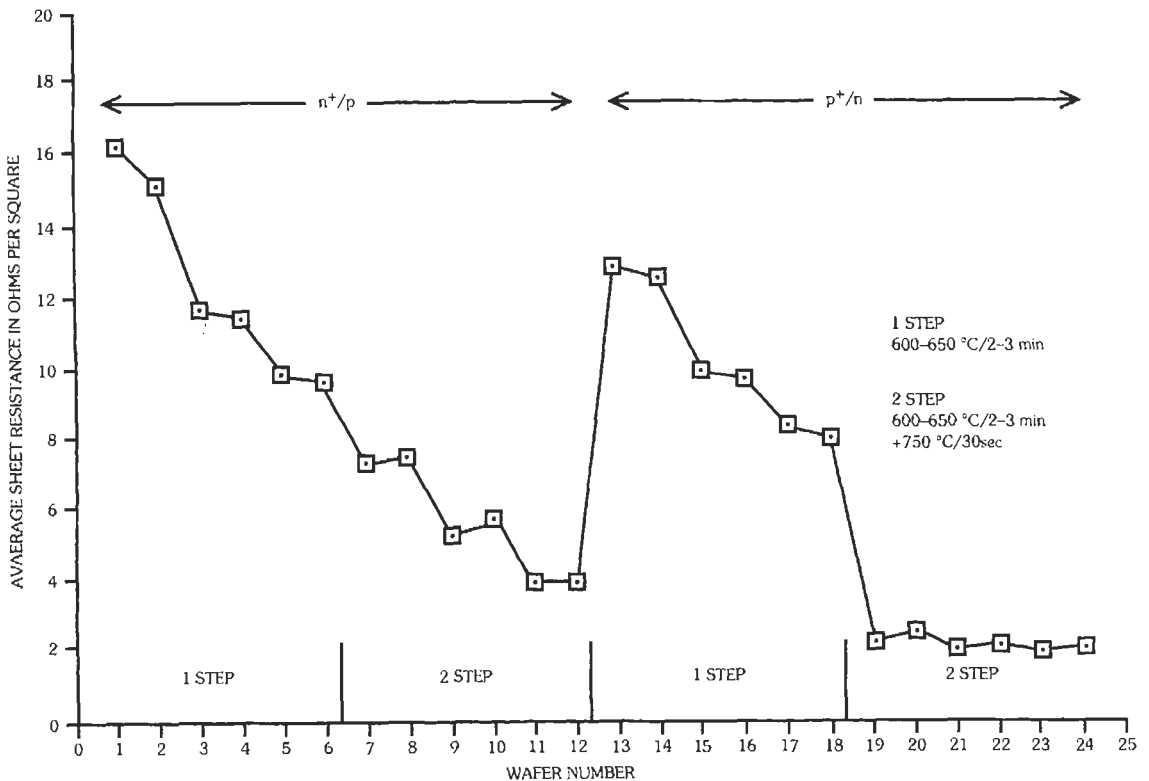


Fig. 76. Sheet resistance values of RTP-reacted TiN/TiSi₂ layers over n⁺p and p⁺n source/drain junctions (400 Å initial titanium). Various 1-step and 2-step RTP react processes were used.

step to remove any residual native oxide layer.* Poor processing equipment base pressure and vacuum integrity can also result in oxygen incorporation and/or formation of surface native oxide layers.† Fig. 78 also summarizes various possible sources of other contaminants in silicon device processing. For example, metallic contaminants may be incorporated into a semiconductor substrate by wet chemicals, reactive-ion etching, photoresist processing, and ion implantation. The current semiconductor device manufacturing environments essentially rely on wet cleaning techniques for removal or reduction of various types of contaminants including the metallic impurities (e.g., Ni, Cr, Fe, Cu, Zn). The wet chemical cleaning processes used prior to the gate oxidation and other high-temperature furnace processing steps are known to have significant effects on the device manufacturing yield and reliability characteristics.‡

Fig. 79 presents a list of different processing techniques for removal of native oxide layers from semi-

conductor substrate. The common procedure for removing native oxides (and other contaminants) prior to epitaxial silicon (or polysilicon contact) deposition involves an ex situ aqueous cleaning (such as the so-called RCA cleaning and a final HF dip), followed by an in situ high-temperature (around 900–1200 °C) hydrogen prebake.§ However, some UHV-CVD epitaxial growth and silicon deposition techniques essentially depend on the effectiveness of an ex situ cleaning and native oxide removal process and maintaining an atomically clean silicon surface in the absence of any form of in situ cleaning process." A suitable wet cleaning step is expected to remove most of the native oxide and other trace contaminants, while a high-temperature H₂ prebake process removes residual native oxide left on the semiconductor surface caused by ambient exposure and/or chemical cleaning. As shown in Fig. 23, other native oxide removal processes have been developed that operate at lower substrate temperatures compared with the high-temperature hydrogen bake process. The germane-assisted native

* References 66 and 67.

† Reference 68.

‡ References 62 and 69.

§ References 70–72.

" References 73 and 74.

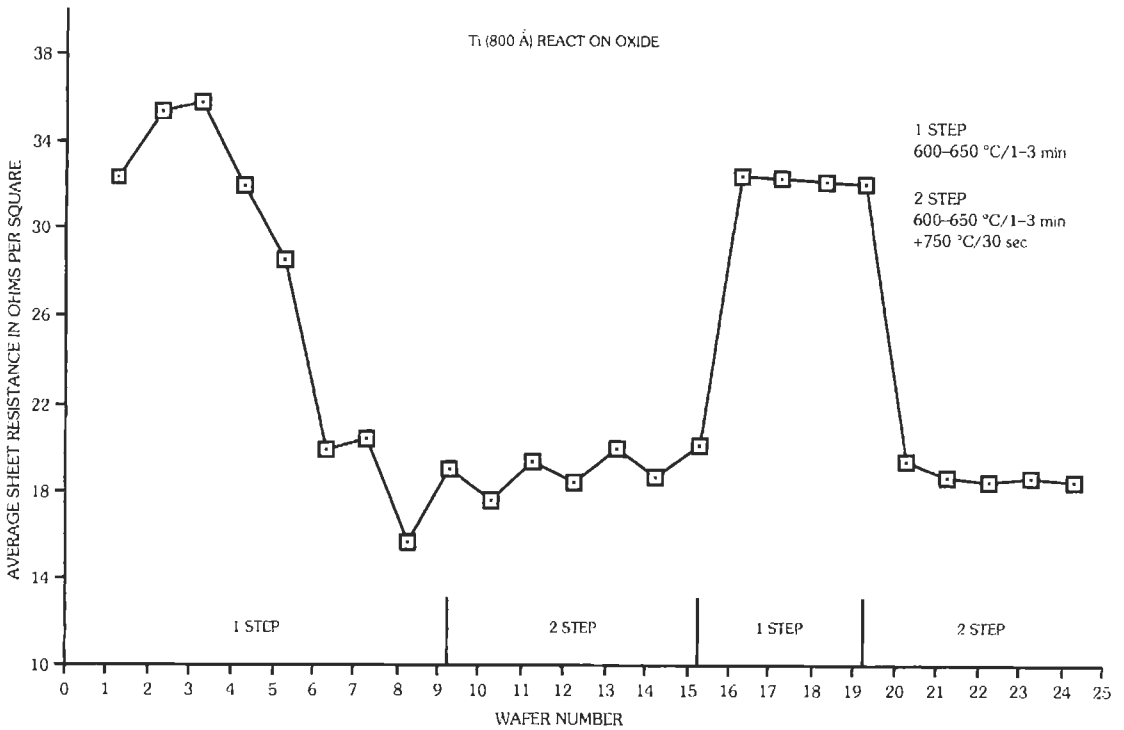


Fig. 77. Sheet resistance values of RTP-reacted TiN layers over silicon dioxide (800 Å initial titanium). Various 1-step and 2-step RTP react processes were used.

oxide removal process is in the low-temperature chemical dry cleaning category.

Fig. 80 shows the effectiveness of various surface cleaning techniques (as specified by the process energy sources and/or cleaning chemistries) for removal of surface contaminants (e.g., native oxides) and other undesirable surface effects (e.g., damage, unwanted surface passivation bonds such as Si-F). The dark bullets represent strong cleaning action for removal of the specified contaminant or undesirable surface effect using the specified processing technique. Moreover, the light bullets demonstrate possible cleaning effect for the specified linked pair of process technique and surface contaminant. All the shown techniques can be quite effective for removal of native oxide layers. Thermal cleaning techniques such as the germane process may also be useful for removal of some metallic and organic contaminants besides native oxide layers. Another example is the ECR plasma surface cleaning technique which has proven effective for removal of native oxide and organic contaminants.*

Low-Temperature Germane-Based Surface Cleaning—A low-temperature (650–750 °C) dry surface cleaning process is described for in situ removal

of native oxide layers and other surface contaminants.† The cleaning process chemistry consists of a mixture of germane (GeH₄) and hydrogen (H₂) gases with very low germane-to-hydrogen flow ratio (10–20 ppm). The germane-based surface cleaning process can be integrated with various thin film growth and deposition processes such as epitaxial growth, gate dielectric formation, and polycrystalline or amorphous silicon deposition for MOS and bipolar device fabrication.

The stacking fault density measurements on epitaxial Si films grown following the in situ germane-assisted native oxide removal have indicated that the two most significant cleaning process parameters are the substrate temperature and GeH₄:H₂ gas flow ratio. Substrate temperatures much below 600 °C result in insufficient native oxide removal and, as a result, higher stacking fault densities in the grown epi layers (for the range of small GeH₄:H₂ flow ratios). Moreover, temperatures in excess of 850–900 °C result in increased stacking fault density in the grown epi layers. The latter is probably due to nucleation of germanium on the silicon surface during the in situ cleaning

* References 75 and 76.

† Reference 49.

CAUSES	EFFECTS						
	NATIVE OXIDE	ORGANIC CONT.	METALLIC CONT.	FLUORO- CARBON	SURFACE DAMAGE	UNWANTED SURFACE BONDS	PARTICLES
WET CHEMICALS	●		●			●	●
RIE		○	●	●	●	●	●
POOR EQUIPMENT BASE PRESSURE	●	○					
LOW PURITY PROCESS GASES	●	○	●			●	●
PHOTORESIST PROCESSING		●	●				○
ATMOSPHERIC AMBIENT EXPOSURE	●	●	○				●
WAFER HANDLING		○	●				●
VACUUM PUMPING AND CHAMBER VENT		●	○				●
UNCONFINED PLASMA PROCESSES			●		○		●
DIRTY EQUIPMENT		●	●				●
ION IMPLANTATION		○	●		●		

Fig. 78. Sources of surface contaminants and damage in semiconductor device processing.

process. The germane-assisted cleaning processes work well in the temperature range of 650–750 °C.*

For a given substrate temperature, process pressure, and total gas flow rate, there is a critical GeH_4/H_2 flow ratio beyond which germanium can nucleate on the surface during the cleaning process. Germanium nucleation can be prevented by reducing the germane-to-hydrogen flow ratio below the critical values. As the total process pressure increases, the germane-to-hydrogen flow ratio may have to be further reduced to prevent germanium deposition or nucleation during the in situ cleaning cycle. The experimental results have indicated that GeH_4/H_2 flow ratios in the range of 10–20 ppm (e.g., 12.5 ppm) work well for effective native oxide removal processes at 5 torr and in the substrate temperature range of 650–750 °C. The germane cleaning process can selectively remove native oxide layers without etching thermal oxides. This selectivity makes

the germane-based cleaning processes superior to the fluorine-based native oxide removal techniques.

It is expected that the GeH_4/H_2 processes clean the substrate surface by direct reaction of germane with the "loose" native oxide layer. The possible oxide reduction reactions are the ones that produce volatile germanium oxide species (GeO). The presence of a large partial pressure of hydrogen balances the cleaning process and prevents germanium nucleation during the in situ cleaning process.

Native and Chemical Oxide Cleans by Vapor-Phase HF—The vapor-phase equivalent of the aqueous HF has been extensively investigated.† Among the advantages of this approach are low chemical consumption, low particle addition, and improvements in electrical performance as measured by a number of parameters, such as silicide sheet resistance and thin oxide stress endurance. More than 3000 wafers processed in a four-month period using only 1 liter of

* Reference 77.

† References 70–80.

- UHV HEATING:
 - Works for ultrathin native oxides,
 - Requires medium/high T's ($T > 750$ °C).
- REACTION WITH ADSORBED LAYERS OF Si OR Ge:
 - Works for ultrathin native oxides.
- HYDROGEN BAKE:
 - Works for ultrathin native oxides; requires high T's ($T > 950$ °C).
- SPUTTERING (e.g., ARGON IONS):
 - Medium/high T's required to eliminate sputter-induced surface damage.
- SPUTTER/CHEMICAL CLEANING (e.g., Ar/H₂ PLASMA)
- HF/H₂O OXIDE REMOVAL:
 - May have to eliminate the residual Si-F bonds.
- ION-ASSISTED CHEMICAL CLEANING (e.g., ECR PLASMA)
- OTHER LOW-T CHEMICAL DRY CLEANING PROCESSES

Fig. 79. Various types of native oxide removal processes.

azeotropic HF solution resulted in excellent particle performance (Fig. 81). Gate oxides grown after vapor-phase HF cleaning have shown several times improvement in oxide endurance under constant-current stress compared to aqueous HF pregate cleaning.*

Comparison of Various Native Oxide Removal Techniques—Table 17 summarizes the overall features of various native oxide removal processes. Among these processes, the conventional hydrogen prebake process requires the highest substrate temperature.† Its native oxide removal selectivity with respect to thermal oxide and silicon is good; however, thermal oxide undercut on patterned substrates may be a problem. The hydrogen prebake technique is essentially a simple process and compatible with the conventional epitaxial growth reactors and processes. Sputter etch processes using RF or ECR plasmas (Ar

or Ar/H₂ ambients) can also be used for native oxide removal.‡ These processes require medium substrate temperatures, on the order of 750–850 °C, in order to eliminate sputter-induced surface damage. The sputter etch techniques offer poor etch selectivity with respect to thermal oxide and silicon and are not usually compatible with the conventional thermal CVD equipment. For in situ cleaning applications, the use of sputter etch native oxide removal modules may be acceptable in vacuum-integrated cluster tools for some applications.

According to Table 17, the chemically enhanced ion-assisted ECR or RF plasma native oxide removal processes are effective for low-temperature ($T \leq 500$ °C) removal of native oxide layers without leaving any residual surface damage.§ Hydrogen or a mixture of hydrogen and silane have been used for chemically enhanced plasma removal of native oxide layers. These processes require special plasma equipment

* Reference 49.

† References 81–84.

‡ References 76 and 85–87.

§ References 75 and 88–92.

PROCESS TECHNIQUES	CONTAMINANTS						
	METALLIC	ORGANIC	NATIVE OXIDE	FLUORO CARBONS	DAMAGED SURFACE	PARTICLES	UNWANTED SURFACE BONDS
PHOTOCHEMICAL	●	●	●	?	○	?	●
VAPOR PHASE PROCESS (HF/H ₂ O, ITC)	?	?	●	?	⊗	⊗	⊗
REMOTE PLASMA	○	●	●	?	?	⊗	○
SPUTTER AND SPUTTER + CHEM ETCH	⊗	⊗	●	○	⊗	⊗	?
THERMAL CLEANING AND MULTIPROCESSING (SWIM, GERMANE)	○	○	●	⊗	⊗	⊗	●
ECR PLASMA	?	●	●	○	○	?	○
HYBRID METHOD 1 (ECR + DUV)	●	●	●	○	○	?	●
HYBRID METHOD 2 (PHOTOCHEM + RMT PLASMA)	●	●	●	○	○	?	●

Fig. 80. Effectiveness of various surface cleaning techniques for removal of native oxide layers and other surface contaminants in semiconductor device fabrication.

modules and are incompatible with standard epitaxial growth or deposition equipment. Remote plasma-enhanced processes such as those with remote hydrogen plasma can remove thin native oxide layers at relatively low temperatures (150–500 °C).^{*} These processes require special plasma modules and are rather ineffective for removal of thicker (e.g., over 10–15 Å) native oxides. They may also generate particles if a quartz discharge tube is used.

As shown in Table 17, the low-temperature thermal native oxide removal processes based on the fluorine-containing gases (e.g., NF₃ or ClF₃) require special corrosion-resistant equipment and offer rather poor native oxide etch selectivity with respect to thermal oxide and silicon.[‡] Another potential problem with these processes is the possibility of formation of Si–F surface passivation bonds and/or residual chlorine remaining on the surface. Special care is needed to prevent surface etch pits. Substrate heating in a UHV

environment such as in an MBE chamber can remove thin native oxide layers at medium temperatures with relatively good etch selectivity with respect to thermal oxide and silicon;[§] however, this technique requires expensive UHV hardware and is not compatible with the conventional silicon processing equipment. In situ deposition of a submonolayer of germanium or silicon on the substrate prior to a UHV heating cycle can help to further reduce the substrate temperature required for removal of thin native oxide layers.[§]

Vapor-phase removal of native oxide layers is performed based on the HF/H₂O chemistry at room temperature.^{||} One advantage of the HF-based vapor phase processes is the capability for providing very good native oxide etch selectivity with respect to other material layers including thermal oxides; however, the HF/H₂O-based vapor-phase processes cannot be performed in the conventional processing equipment and

* Reference 93.

‡ References 94 and 95.

§ References 96 and 97.

|| Reference 98.

||| References 99–101.

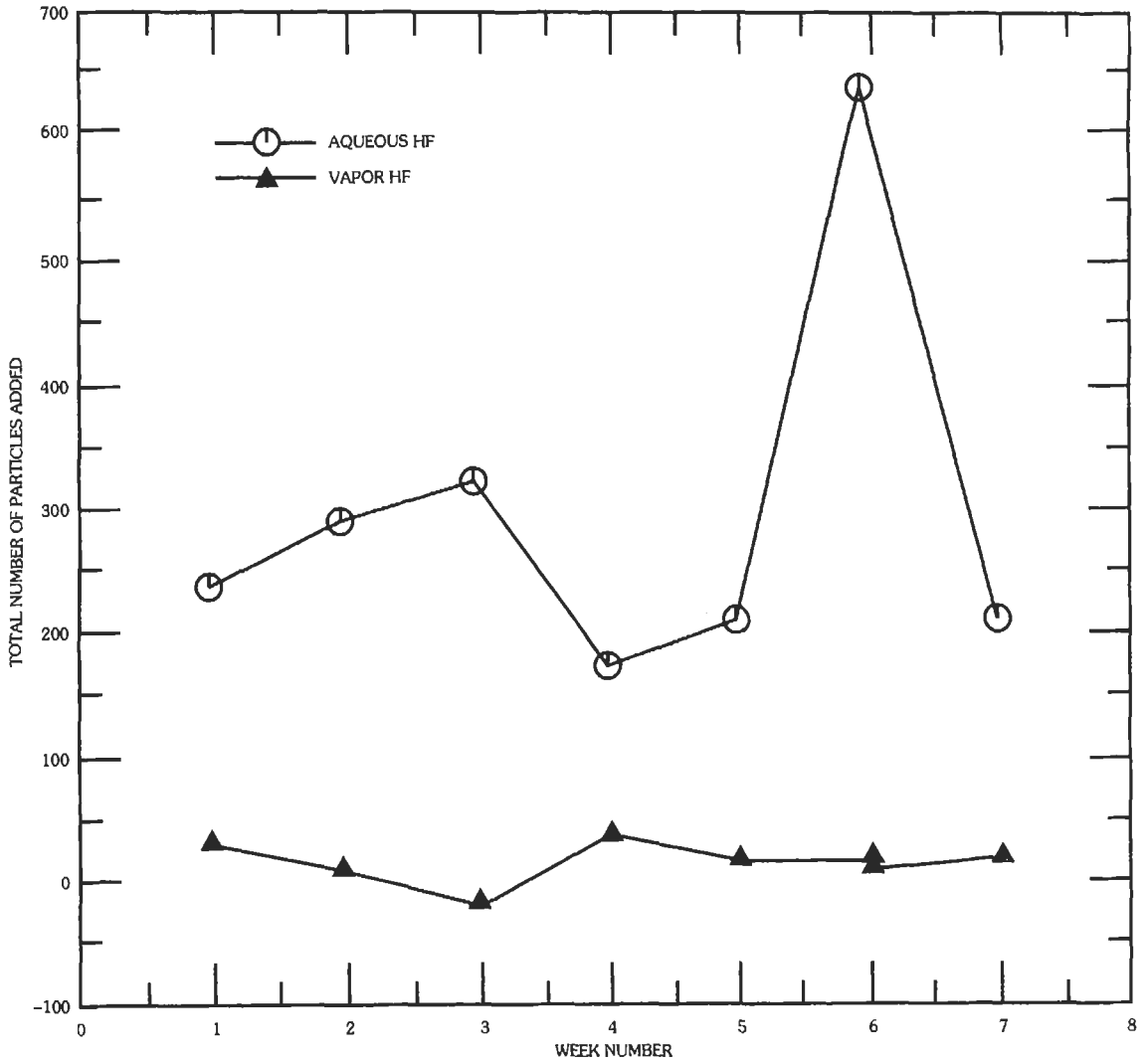


Fig. 81. Particle performance of aqueous and vapor-phase HF native oxide etch.

require a dedicated special process chamber. As shown in Table 17, photochemical cleaning refers to the ultraviolet-enhanced processes for native oxide removal at medium temperatures.* These processes are excellent candidates for low-temperature removal of metallic and organic contaminants as well. They require special photochemical equipment and are incompatible with the standard thermal processing equipment. Finally, as indicated in Table 17, the germane-based native oxide removal process is compatible with the conventional silicon processing equipment. This process is effective

for native oxide removal at low to medium substrate temperatures and provides good native oxide etch selectivity with respect to thermal oxide layers and silicon.

Organic Cleans—The removal of trace organic contaminants using ozone and atomic oxygen generated by irradiating an oxygen containing ambient with deep UV† is now firmly established as a viable substitute for the traditional wet organic clean. The photoactivation of oxygen and the subsequent generation of

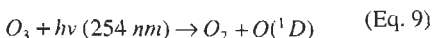
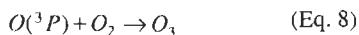
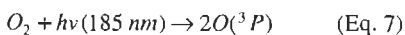
* References 102 and 103.

† Reference 104.

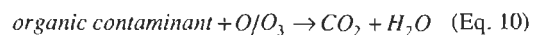
TABLE 17. FEATURES OF VARIOUS NATIVE OXIDE REMOVAL PROCESSES

Oxide Removal Process	Characteristics			
	Process Temperature	Selectivity to Thermal Oxide and Silicon	Equipment Technology	Comments
Hydrogen Prebake	High ($T > 900\text{ }^{\circ}\text{C}$)	Good	Simple	Oxide undercut problem Not suitable for prepoly clean Works for ultrathin oxides
Sputter Etch (ECR, RF) (Ar, Ar+H ₂)	Medium ($T > 750\text{ }^{\circ}\text{C}$)	Poor	Special module (plasma)	Residual surface damage Incompatible with std. equip. Narrow process window
Ion-Assisted Chemical (ECR: H ₂ , H ₂ +SiH ₄)	Low ($T = 400\text{--}500\text{ }^{\circ}\text{C}$)	Poor-Fair	Special module (plasma)	Superior to microwave H ₂ Incompatible with std. equip. Can integrate with ECR dep.
Remote Plasma Clean (Microwave H ₂)	Low ($T = 150\text{--}500\text{ }^{\circ}\text{C}$)	Good	Special module (plasma)	Works for ultrathin oxides Incompatible with std. equip. Must prevent particles
NF ₃ Oxide Removal	Low ($T = 580\text{ }^{\circ}\text{C}$)	Poor	Special module corrosion-res.	Corrosive chemistry Special chamber passivation Si-F surface bonds
ClF ₃ Oxide Removal	Very Low ($T = 25\text{--}60\text{ }^{\circ}\text{C}$)	Poor	Special module corrosive-res.	Corrosive chemistry Special chamber passivation Surface Cl and F and pitting
UHV Heating	Medium ($T > 750\text{ }^{\circ}\text{C}$)	Good	Special module (UHV chamber)	Works for ultrathin oxides Expensive UHV equipment Must prevent surface pitting
Reaction With Adsorbed Layers of Si or Ge	Low-Medium ($T = 625\text{ }^{\circ}\text{C}$)	Good	Special module (UHV chamber)	Works for ultrathin oxides Superior to UHV-only bake Expensive UVH-MBE equipment
HF/H ₂ O Oxide Removal	Very Low (Room temp.)	Good	Special module (vapor-phase)	Special process chamber Can control selectivity Must prevent surface residues
Photochemical Clean	Medium ($T = 730\text{ }^{\circ}\text{C}$)	Fair-Good	Special module (UV photochem.)	Incompatible with std. equip. Better for organic/metallic Must prevent DUV damage
Germane-Assisted Surface Cleaning	Low-Medium ($T = 650\text{--}750\text{ }^{\circ}\text{C}$)	Good	Simple	Compatible with std. equip. Simple and effective Requires scram room set-up

ozone and excited oxygen atoms are summarized in the following equations:



Both the ozone and the atomic oxygen are powerful oxidizing agents that react with trace organic contaminants via the following reaction:



Instead of using deep UV activation, oxygen radicals can also be generated in a plasma ambient. Wafers can be placed downstream from the plasma generator

to minimize possible ion-induced damage. This well-proven technique forms the basis of a number of commercial photoresist ash reactors.

Metallic Cleans—To keep the contaminants at or accessibly close to the surface of the wafer and because most metallic elements are fast diffusers at moderately elevated temperatures, the chemistries chosen for trace metallic removal must be able to form volatile reaction products with the metallic contaminants without requiring too much thermal activation.

Various halogen-based techniques have been investigated in this regard. Some success has been achieved by photochemically activated chlorine chemistry.* Since most chlorides of metals are not volatile at low temperatures, it is believed that the cleaning action is achieved by the slight etching of the silicon substrate and that the metallic contaminants are “lifted off” the surface of the wafer together with the volatile chlorides of silicon. As in the case of trace organic clean, chloride chemistries based on plasma activation have also been investigated and have achieved some degree of success. Alternative nonhalogen-based chemistries have also been investigated by using a reactor with sources for both the deep UV irradiation and the remote microwave plasma activation.† Gases can be injected into the reactor with or without the UV or the plasma excitation. The chemistries are chosen so that the reaction products are sufficiently volatile at relatively low temperatures and so that the etching of exposed silicon can be reduced to a minimum. The cleaning efficiencies of these chemistries have been verified by using trace amounts of radioactive isotopes of some of the most commonly encountered contaminant elements (Fe, Cu, etc.) plated onto bare silicon surfaces.

Single-Wafer Epitaxy

Single-wafer RTCVD technology has been used for fabrication of in situ-doped homoepitaxial Si and heteroepitaxial Si/Ge_xSi_{1-x} multilayer structures. Epitaxial Si films have been grown in a mixture of DCS, H₂, and the dopant gas at 5 torr.‡ Fig. 82 shows the thickness of undoped epitaxial films grown on heavily doped p⁺ (100) wafers versus the DCS flow rate in 200-second runs at ≈1000 °C. These data cover a DCS:H₂ flow ratio range of 1–6%. Under these process conditions, the deposition takes place in the mass-transport-limited regime and, as a result, the growth rate and uniformity are sensitive to the DCS flux in the process chamber. A wide range of controlled growth rates up to more than 1 μm/min can be achieved. Thermal switching of the process temperature provides a useful capability for growth of multilayer epitaxial structures with fairly abrupt interfaces.§ Fig. 83 plots the secondary ion mass

spectroscopy (SIMS) boron concentration depth profile for a seven-layer step-boron-doped epitaxial structure grown following a low-temperature GeH₄ + H₂ surface cleaning (total of eight process steps integrated in situ). This structure was grown at ≈1000 °C on a p⁺ substrate and consists of: (1) an undoped epi layer, (2 through 5) four successively higher step-boron-doped layers, (6) an undoped epi layer, and (7) a heavily doped layer at the surface. The SIMS profiles were obtained with an oxygen primary beam under conditions that provided poor depth resolution. The layer-to-layer transition width for the top two layers is 280 Å/decade of dopant concentration. The actual transition width is smaller than the observed SIMS value because of the SIMS broadening effects. The transition widths are affected by the growth temperature and the thermal treatments during the subsequent growth cycles. Growth of the epitaxial layers at lower temperatures with remote microwave plasma excitation should significantly reduce the transition widths in multilayer epitaxial structures. Growth of the same structure shown in Fig. 83 without thermal switching (temperature kept at 1000 °C even during gas switching) results in much wider transition widths because of the long process gas stabilization times.

Fig. 84 plots the B and Ge SIMS concentration depth profiles in a seven-layer structure consisting of a strained p⁺-Ge_xSi_{1-x} layer. This structure was grown after a GeH₄-assisted in situ surface cleaning. The first five layers are i/p⁺/i/p⁺/i silicon regions (i = intrinsic), each grown at ≈1000 °C. The B₂H₆ flow rate during the growth of the first two p⁺ layers was ramped in time, which is one reason for the gradual B concentration transition between the adjacent i/p⁺ layers. The sixth and seventh layers are strained p⁺-Ge_xSi_{1-x} and n⁺ Si, respectively. Because of the overall thermal treatments and SIMS broadening, the first intrinsic layer near the p⁺ substrate interface cannot be detected very clearly. The SIMS Ge depth profile indicates a Ge fraction of x = 6% in the Ge_xSi_{1-x} layer. The GeH₄:DCS flow ratio during the growth of this Ge_xSi_{1-x} layer was 2.5 sccm:100 sccm and the growth was performed at 650–700 °C. The net B₂H₆ flow rate was about 0.0003 sccm. The B concentration depth profile indicates a Ge_xSi_{1-x} doping concentration on the order of 1 × 10¹⁸ cm⁻³. Based on the SIMS data, the B-to-Ge concentration ratio in the grown Ge_xSi_{1-x} film is nearly equal to the ratio of the number of B and Ge atoms available in the gas phase (or twice the B₂H₆ flow rate divided by the GeH₄ flow rate).

Single-Wafer Plasma Processing

Plasma processes are routinely used in semiconductor fabrication for etching of submicron geometries, deposition and growth of thin films, and in situ organic and metallic contamination removal. Gas composition, pressure, and flow rates define the overall chemistry,

* Reference 105.

† Reference 49.

‡ Reference 49.

§ Reference 84.

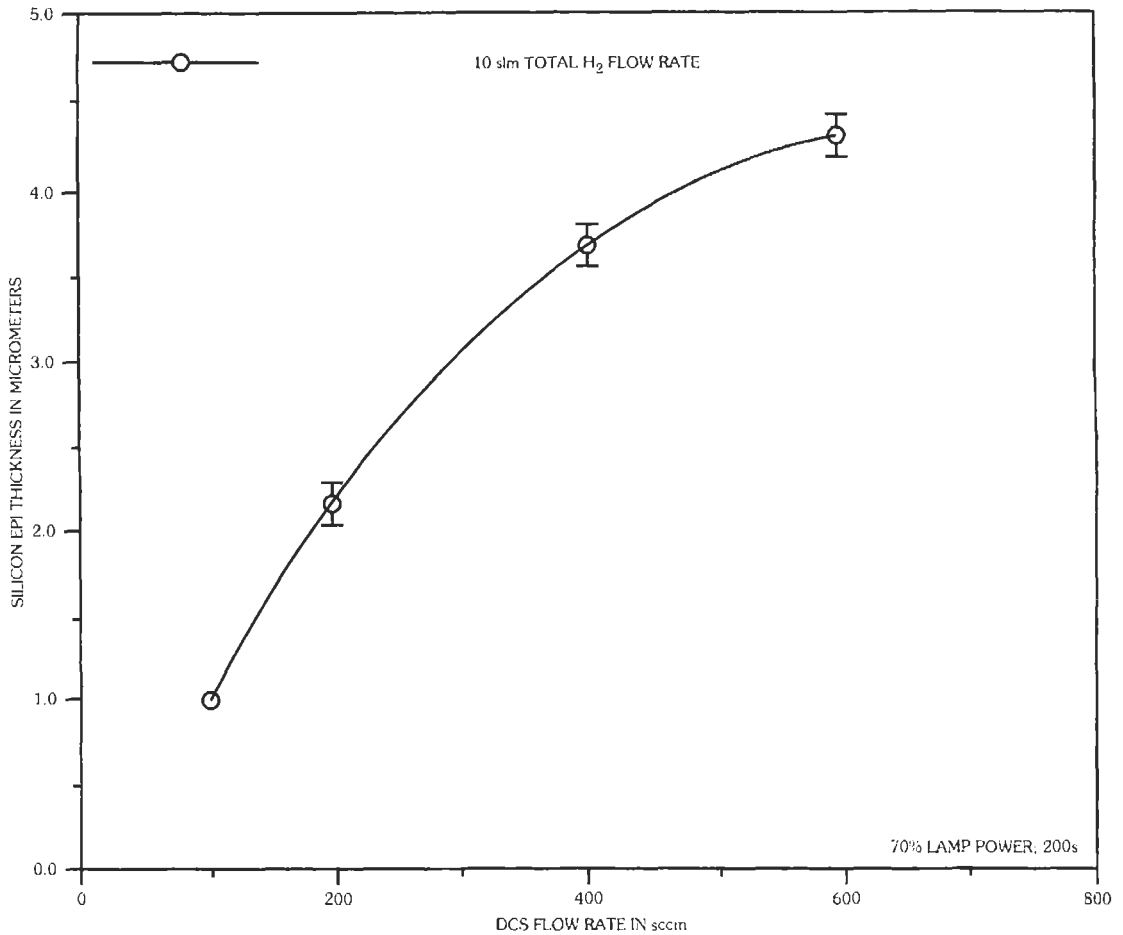


Fig. 82. Dependence of the epitaxial silicon layer thickness on SiH₂Cl₂ flow rate in 200-second growth runs performed at ≈1000 °C.

but the mechanisms and kinetics are strongly influenced by the plasma parameters. Many chemical reactions in the plasma are initiated by plasma electrons, while surface reactions are often ion-activated. Some of the important plasma parameters are the electron and ion concentrations, the mean electron energy, and the energy and flux of ions incident on the surface. Processes that provide high etch selectivity, good anisotropy, good uniformity, low ion-induced damage, and low particulate contamination and that meet throughput requirements are essential for the etching of submicron features.

Capacitively Coupled Discharges—Capacitively coupled RF plasma sources generate large volume plasmas of moderate densities. Although these systems operate over a wide range of pressures (50 mtorr–1 torr), few of the plasma parameters can be independently controlled. For a given RF power, gas pressure,

and gas composition, all the plasma parameters with the exception of the ion energy are fixed. These attributes make it very difficult to engineer processes that simultaneously attain high selectivity, good anisotropy, and low ion-induced damage. The potential for improved uniformity, reduced microloading effects, reduced contamination, and improved anisotropy are driving discharge operation to lower pressures. Deterioration in plasma confinement limits the minimum operating pressure to 30 mtorr. These factors have accelerated the use and development of alternative plasma sources.

Two alternative configurations are derived from standard capacitively coupled RF discharges: the hollow cathode discharge and the magnetron-enhanced etch system. A hollow cathode discharge is formed when two electrodes parallel to each other are powered by the same RF source. Secondary electrons released by primary ion bombardment of either electrode

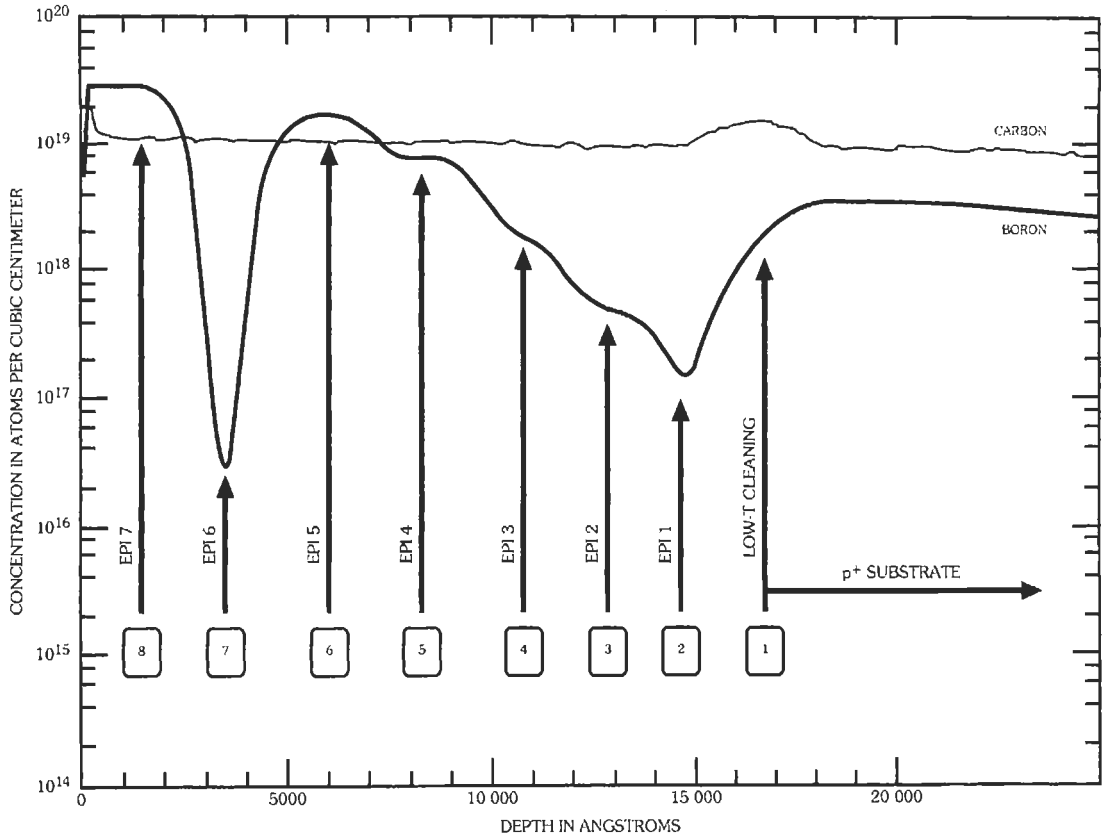


Fig. 83. The SIMS boron concentration depth profile for a seven-layer epitaxial structure grown at $\approx 1000^\circ\text{C}$ with multiple rapid switching of the wafer temperature following a low-temperature GeH_4+H_2 pregrowth cleaning cycle.

bounce back and forth between the equipotential electrodes, enhancing the ionization.* This translates into an enhanced plasma density, a large radical production rate, and lower sheath voltages.† Impressive etch and deposition results from hollow cathode systems have been reported. For example, SiO_2 etch rates in excess of $1\ \mu\text{m}/\text{min}$ with selectivity to Si in excess of 50:1 and uniformities of 1% have been reported.‡ In addition, hollow cathode systems are ideal for processes involving simultaneous etch and deposition, either for etch chemistries that rely on passivation through polymer on side walls to achieve anisotropy and selectivity or in the deposition of planarized films.§ Sputtering and deposition on the counter electrode and plasma instability have prevented this discharge from gaining wide acceptance, although methods for alleviating some of these problems have been reported recently.§

* Reference 106.
 † Reference 107.
 ‡ Reference 108.
 § Reference 109.

The magnetron-enhanced RIE system employs magnetic fields either parallel or perpendicular to the surface of the powered electrode, to boost the ionization efficiency and hence the plasma density. Since the magnetic field reduces the electron mobility, discharge operation can be extended to lower pressures (< 30 mtorr), plasma confinement is improved, and the sheath voltage decreases. The magnetron-enhanced system is described more fully below. The hollow cathode and magnetron systems provide good selectivity (because of the large radical flux and low ion energies), while the high ion flux ensures that anisotropy is not sacrificed. The low ion energies also imply reduced ion-induced damage. Throughput is also enhanced.

Some of the limitations of capacitive coupling cannot be overcome; sheath voltages increase with power. The obvious solution is either to use inductive coupling (inductively coupled RF discharges) or to use higher frequencies (microwave discharges).

Inductively Coupled Discharges—Two generic types of inductively coupled RF discharges have evolved. The first is based purely on induction coupling.

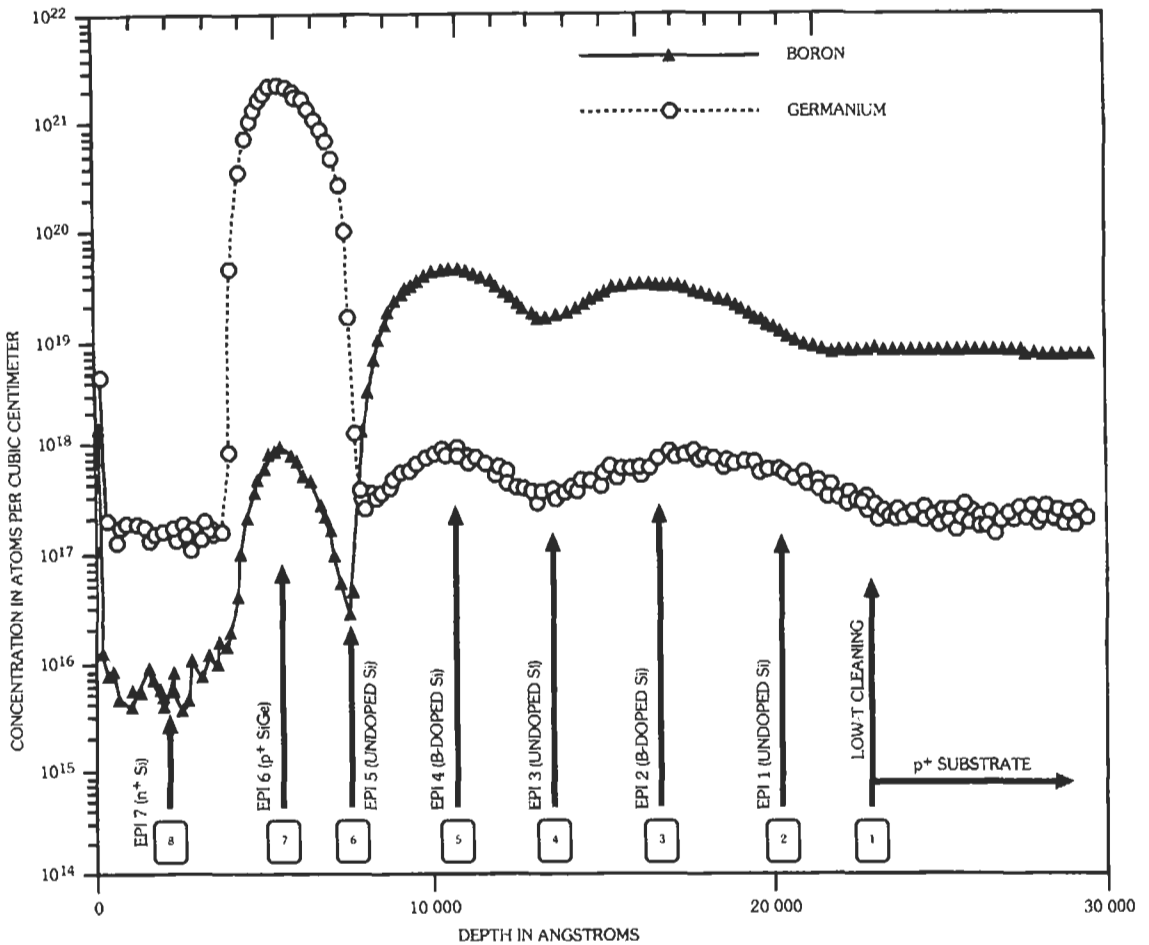


Fig. 84. The SIMS concentration depth profiles for boron and germanium in a seven-layer heteroepitaxial structure with a strained $\text{Ge}_x\text{Si}_{1-x}$ layer grown at $\approx 675^\circ\text{C}$. The surface layer is n^+ Si, and the structure was grown on a p^+ substrate.

An RF-powered antenna coil placed around the periphery of the chamber induces eddy currents that sustain the plasma.* With magnetic plasma confinement, plasma densities in argon exceed $5 \times 10^{11} \text{ cm}^{-3}$ with low ion energies ($< 50 \text{ eV}$), over a wide pressure range (0.5–50 mtorr).† If a uniform magnetic field region exists in the vicinity of the coil, energy coupling becomes extremely efficient since helicon or whistler waves are excited in the plasma. With resonant coupling, the ionization efficiency peaks, and a very dense plasma is generated.‡ Extremely high plasma densities exceeding 10^{12} cm^{-3} have been reported with low intrinsic ion energies ($< 30 \text{ eV}$).§

Preliminary process results are impressive for both these systems. Etch rates for most materials including Si, SiO_2 , and Si_3N_4 exceed $1 \mu\text{m}/\text{min}$. The ion energy is independent of RF power and can be increased from its intrinsic value by RF biasing the substrate. The high radical generation rates and low ion energies result in high selectivity with no ion-induced damage, while the high ion flux promotes anisotropy.¶ The generation region becomes increasingly concentrated owing to skin effects as the plasma density increases. The wafer is located remotely from the source for good uniformity ($< 5\%$, 1σ).** Also, large RF voltages across the antenna coil may sustain subsidiary capacitively

* References 110–113.

† Reference 113.

‡ Reference 114.

§ Reference 115.

" References 113 and 115.

Reference 111.

** References 113 and 115.

coupled discharges that cause sputter damage of the chamber walls and possible contamination.

Magnetron-Enhanced Discharges—Magnetron effects induced by external magnetic fields extend discharge operation to lower pressures. Electron mobility along a direction perpendicular to the magnetic field is reduced by the factor^{*}

$$1/\left(1 + \omega_b^2 / \nu_m^2\right) \quad (\text{Eq. 11})$$

where

ω_b is the electron cyclotron frequency,
 ν_m is the electron momentum transfer collision frequency.

The reduced electron mobility produces better plasma confinement and a reduction in sheath voltages. In addition, the ionization efficiency of secondary electrons emitted from the powered electrodes is enhanced. Both these effects lead to an increase in the plasma density. A factor of 10 enhancement in plasma density is possible with moderate fields (≈ 250 G).[†] Because of the significant dissociation, the plasma chemistry can differ between RIE and magnetron-enhanced systems.[‡]

Various configurations have been investigated. Magnetron systems have traditionally been used for sputtering applications, where discharge operation at submillitorr pressures is vital. The most common configurations include the planar magnetron discharge, with the magnetic field perpendicular to the electric field, and cylindrical magnetron discharges. The magnetic field is generated by permanent magnets or electromagnets located behind the powered electrode. Typically the magnetic field is oriented so that the $E \times B$ electron drift forms a closed loop, thereby trapping the electrons in the vicinity of the electrode.[§] In the conventional planar magnetron, a disk magnet surrounded by a ring magnet produces radial magnetic fields. The magnetic field parallel to the wafer is most intense at midradius. Plasma nonuniformity is a big concern in these systems. Because the plasma is most intense in regions where the magnetic field component parallel to the wafer is a maximum, it causes a race track etch characterized by a fast etch at midradius.[§] While this configuration results in $E \times B$ drift currents that form a closed loop, the severe nonuniformity has prompted the use of alternative magnetic field patterns. One of the options is to keep the magnetic field lines parallel to each other and to the electrode. To counter the effect of the asymmetric $E \times B$ drift, the magnets or wafer can be rotated or scanned.[#]

An adjustable field planar magnetron discharge has been reported.^{**} Rare-earth permanent magnets located behind the wafer induce a magnetic field parallel to the wafer surface. Also, the field lines are parallel to each other in the region of the wafer. The magnet lengths are scaled so as to achieve a uniform magnetic field over the wafer. Movable soft-iron slugs can be used to tailor the magnetic field distribution further. The slugs are located between the magnets and behind the powered electrode. The magnetic field at the wafer is ≈ 500 G. The magnetron RIE system is otherwise similar to a parallel plate system. In addition to powering the electrode, a cylindrical ring around the counter electrode can also be powered with a low-frequency or high-frequency RF power supply. Grounding or RF biasing the ring can be used to optimize uniformity and vary the energy of ions incident on the substrate. In addition, a contoured counterelectrode further improves the etch uniformity. For ion-assisted processes there is an optimum field strength. For very strong fields, the ion energy is too low, and the etch is hindered.^{††} However, the plasma density increases with increasing field strength. The cylindrical ring permits operation with strong fields without the attendant reduction in etch rate. By manipulating the slugs and the ring, the location of fastest etch can be moved progressively from the wafer center to the wafer edge. This flexibility may be exploited to compensate for nonuniformities resulting from previous process steps.

This magnetron RIE system has been specifically optimized for DESIRE, a dry develop lithography process. The basic steps involved in DESIRE are as follows. Exposed photoresist is exposed to a silylating agent such as HMDS (hexa methyl disilane). This causes silicon to be incorporated in the resist. The depth of silicon incorporation is a function of exposure dose. When this resist is exposed to an O_2 plasma, the silicon in the photoresist is converted to SiO_x . Ion bombardment continuously sputters the SiO_x as it is formed. Ultimately the sputter etch depth in the unexposed (or marginally exposed) regions exceeds the depth of silicon incorporation, and no further SiO formation occurs. The etch now proceeds rapidly as the photoresist is consumed by the O_2 plasma. In the strongly exposed regions, the SiO_x that is continually formed forms the etch mask that hinders the etch.

The nature of this process places very stringent demands on the capabilities of the etch process. The ion energy has to be sufficiently low for reasonable etch selectivity between SiO_x and the photoresist. At the same time, a very high etch selectivity exaggerates etch nonuniformities and promotes residue formation. Selectivities between 5:1 and 8:1 are ideal. Since the silylation depth does not exhibit any strong thresholding behavior with respect to exposure, the budget for etch nonuniformities is very small. The requirements

* Reference 116.

† Reference 117.

‡ Reference 118.

§ Reference 119.

** Reference 120.

Reference 121.

** Reference 49.

†† Reference 117.

become especially tight for contact patterns. To engineer an acceptable process, the resist etch and SiO sputter etch uniformity have to be within 4% (1σ). Finally, vertical etch profiles with no undercut require high anisotropy. Undercut can be eliminated by reducing the O₂ flow and cooling the wafer to below -35°C. These requirements of good selectivity and high anisotropy uniformity preclude the use of an RIE etch system. The magnetron-enhanced RIE system with adjustable magnetic fields and ion energies appears to be a good choice. Good results have been reported for 0.4 μm lines and spaces processed using such a system.* Nominal process conditions were 115 W, 3 mtorr, 40 sccm O₂, and a DC bias of -85 V. A similar system may be used for implementing other etch processes that are also tightly constrained.

Microwave Discharges—Microwave discharges have almost become synonymous with electron cyclotron resonance (ECR) discharges. However, cavity based microwave discharges still maintain their niche for afterglow plasmas. A remote microwave plasma is a copious source of radicals and is generally used in conjunction with a capacitively coupled plasma to enhance the etch or deposition rates.

ECR discharges are a source of large volume dense plasmas. At the resonance condition that occurs for a magnetic field of 875 G for 2.45 GHz microwave excitation, the power coupling becomes resonant. At resonance, the electron cyclotron frequency and the microwave frequencies are synchronized, so that electrons are continuously accelerated by the electric field. ECR is most effective for pressures between 10⁻² and 10 mtorr, when an electron can undergo several orbits before suffering a scattering collision. Typical plasma densities exceed 10¹¹ cm⁻³, with ion current densities to the wafer in excess of 10 mA/cm².† At the same time, the intrinsic ion energy is fairly low (< 30 eV), which can be increased by RF biasing the wafer. Even in the absence of cyclotron resonance, magnetized plasmas can still couple microwave power. However, the mean electron energy and consequently the ionization efficiency is much lower. For a given input power, the higher the ionization efficiency, the higher is the plasma density. Thus the size and coupling efficiency of the ECR region can dramatically affect the plasma density for a given input power.‡ Still, magnetized microwave plasmas generate high plasma densities compared to capacitively coupled discharges.

ECR plasmas offer tremendous process flexibility. High-quality etching of a range of materials including Si, SiO₂, Al, resist, and III-V materials** is possi-

ble. Similarly, deposition processes can also be accomplished at comparatively low temperatures.††

Various types of ECR sources have been developed.‡‡ Perhaps the most common source uses a set of electromagnets around the discharge chamber that generate axial magnetic fields.§§ A divergent or mirror magnetic field configuration can be attained with appropriate solenoid and pole-piece configurations. The axial magnetic field is tailored so that it is a maximum at the microwave input end and decreases toward the wafer. The microwave coupler launches the right-hand circularly polarized wave, which couples very efficiently into the plasma, generating a dense plasma.¶¶ The ECR region is fairly large since the axial magnetic field is quite uniform in the radial direction and decreases away from the microwave coupler. The axial magnetic field also enhances the ion current density at the wafer but may degrade uniformity and directionality of ions.## Remote wafer location is required for good uniformity. However, placing the wafer too far away reduces the radical flux and affects throughput. A third electromagnet in the process chamber can be used to mitigate some of these trade-offs.***

Another option is to use a set of permanent magnets that generate a multipolar cusp field around the periphery of the discharge chamber.††† The ECR region is then a thin shell around the inside of the plasma formation chamber. Because the plasma formation regions are much closer to the wall, and the size of the ECR region is limited, the attainable plasma densities appear to be lower. Since the plasma diffuses freely from the chamber toward the wafer, good uniformity can be attained much closer to the source.††††

The last option is to combine the plasma formation chamber with the process chamber as in the distributed ECR design (DECER).‡‡‡ Permanent magnets establish a multipolar cusp field around the periphery of the process chamber. Multiple microwave antennas are placed around the periphery of the chamber. The ECR region once again is a thin shell around the inside of the process chamber. The plasma diffuses radially toward the wafer, which is located centrally in a field-free region. However, the plasma densities reported for such a source are a factor of five lower than those obtained with other sources. The discharge behaves more like a magnetized microwave plasma, with resonant power coupling occurring in a very small region.

* Reference 49.

† Reference 122.

‡ Reference 123.

§ Reference 124.

¶ Reference 125.

Reference 126.

** Reference 127.

†† References 126, 128, and 129.

‡‡ Reference 122.

§§ Reference 130.

¶¶ Reference 123.

Reference 131.

*** Reference 132.

††† Reference 133.

‡‡‡ Reference 134.

A star-shaped ECR source has been reported* that combines the advantages of the various configurations. It consists of a six-chambered source in a star configuration. Each chamber is surrounded by permanent magnets that generate cusp fields so that a broad spherical resonance region is formed around the center of each formation chamber. The cusp fields also confine the plasma. The six-chambered source opens into a multipolar bucket chamber with the wafer located at the far end. Each chamber is powered by a separate microwave antenna. The six antennas are connected to a six-way microwave power divider. The multipolar bucket provides a very high transfer efficiency (Fig. 85). The ion and electron concentrations were measured by a Langmuir probe. Most of the measurements have been made along the center of the system, where the magnetic field is zero. The plasma density as measured by a Langmuir probe decays exponentially as expected. The plasma density decays by a factor of 10 over 8 inches at 1 mtorr (the source is located 10 inches from the wafer). As the pressure is increased, the plasma density at the source does not change

appreciably, but the plasma is better confined to the source (Fig. 86). Consequently, the plasma density at the wafer decreases. Because the diameter of the multipolar bucket is 8 inches, good uniformity is obtained over a 6-inch wafer merely 6 inches from the source.

This chamber design provides dense uniform plasmas at the wafer without the use of axial magnetic fields. In addition, the design can easily be scaled for larger wafers by moving the individual chambers radially outward. Finally, since the ECR region does not intersect the walls of the chamber, as it does in some of the other designs, sputtering contamination is reduced.

Cryogenic Magnetron Etch—The drive toward higher device integration density presents stringent requirements on the anisotropy of etch profiles and on the precise control of the widths of etched lines and spaces. This is to be achieved while maintaining the necessary selectivity and etch rate specifications. One process technology that is expected to achieve the above requirements is RF-assisted remote microwave plasma etch. It is well recognized that this etch technology allows the achievement of fast etch rates, high selectivity, as well as low damage. However, it is also recognized that this technique suffers from unacceptable etch

* Reference 49.

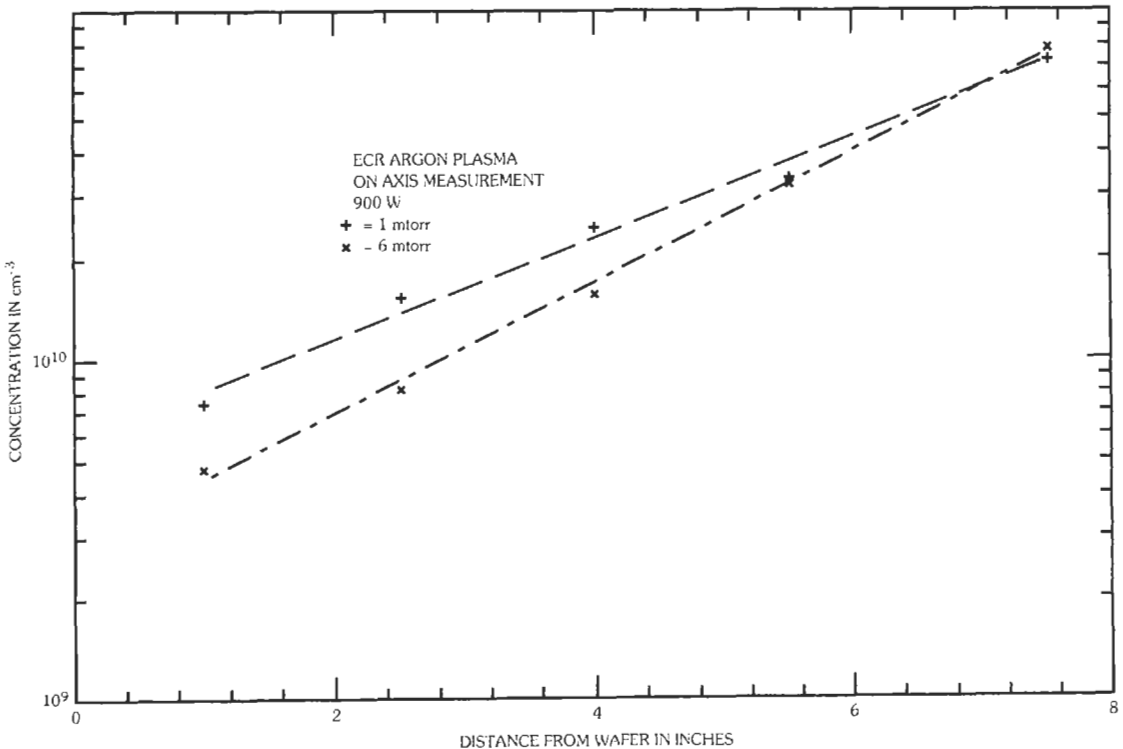


Fig. 85. Measured electron and ion concentrations vs distance from the wafer in a TI custom-designed ECR plasma system. The curves show a high transfer efficiency in the multipolar magnetic bucket separating the process chamber from the distributed ECR plasma source chambers.

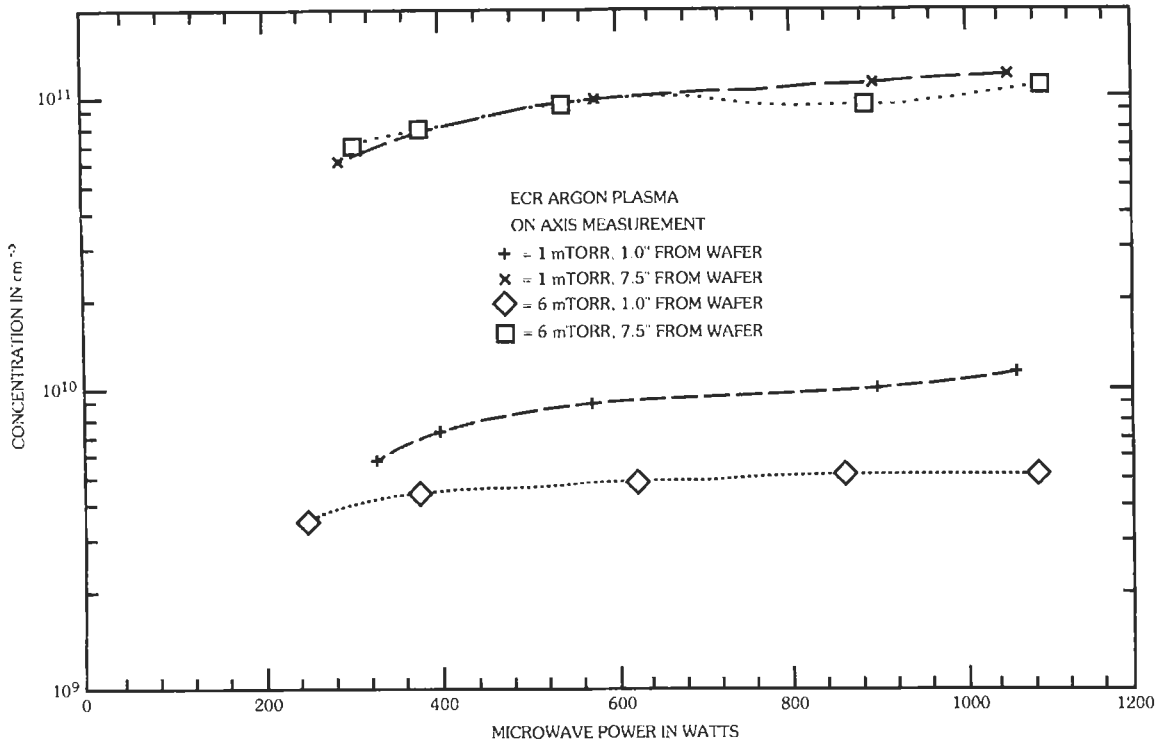


Fig. 86. Measured electron and ion concentrations in an ECR argon plasma vs incident microwave power for two different pressures and two different locations.

profile anisotropy, particularly at room temperature, as evidenced by the observed isotropic etching of the sidewalls. Tachi et al.* demonstrated the sensitivity of the isotropic etch component to wafer temperature. Specifically, by cooling the wafer down to cryogenic temperatures (-100 to -150 °C) the isotropic sidewall etch was completely eliminated.

The achievement of these temperatures is not possible with conventional liquid-coolant-based (e.g., Freon) refrigeration equipment. To meet this temperature requirement, as well as to maintain a cooling load up to 1 kW at the wafer, a liquid nitrogen (LN₂) based cryogenic cooling system has been reported.† The system uses high-pressure helium gas as a coolant, which is pumped in a closed loop between a low-temperature reservoir—an LN₂ dewar—and the wafer chuck. Temperature control is achieved by a bypass valve that allows the necessary fraction of the helium flow to bypass the LN₂ dewar, and then mix with the colder helium leaving the dewar on the way to the chuck, to achieve a required temperature at the wafer. The schematic of the cryogenic cooling system is shown in Fig.

* Reference 135.

† Reference 49.

87. A helium compressor provides the required pumping action in the system. Cold He returning from the chuck must be brought up to ambient temperature before entering the compressor. The heat exchanger does this effectively by precooling the He stream leaving the compressor to the LN₂ bath. An in-line electrical heating element upstream of the compressor is used to maintain a safe inlet He temperature, in case the heat exchanger output falls below the required value. After being precooled in the heat exchanger, the gas leaving the compressor is driven through a copper pipe loop immersed in liquid nitrogen to provide further required cooling. The two valves shown above the LN₂ dewar provide the temperature control capability. Maximum cooling capacity and lowest temperatures are achieved by operating with the bypass valve fully closed and the main valve fully open. Higher temperatures are achieved by allowing a fraction of the flow to go through the bypass path. The LN₂ level in the dewar is maintained by an automatic LN₂ fill system. The cold He leaving the dewar is routed to the wafer chuck, where it maintains the required temperature, and then returns to the compressor via the heat exchanger.

The system has been used successfully to provide cooling down to -160 °C, with a sustained heat load of

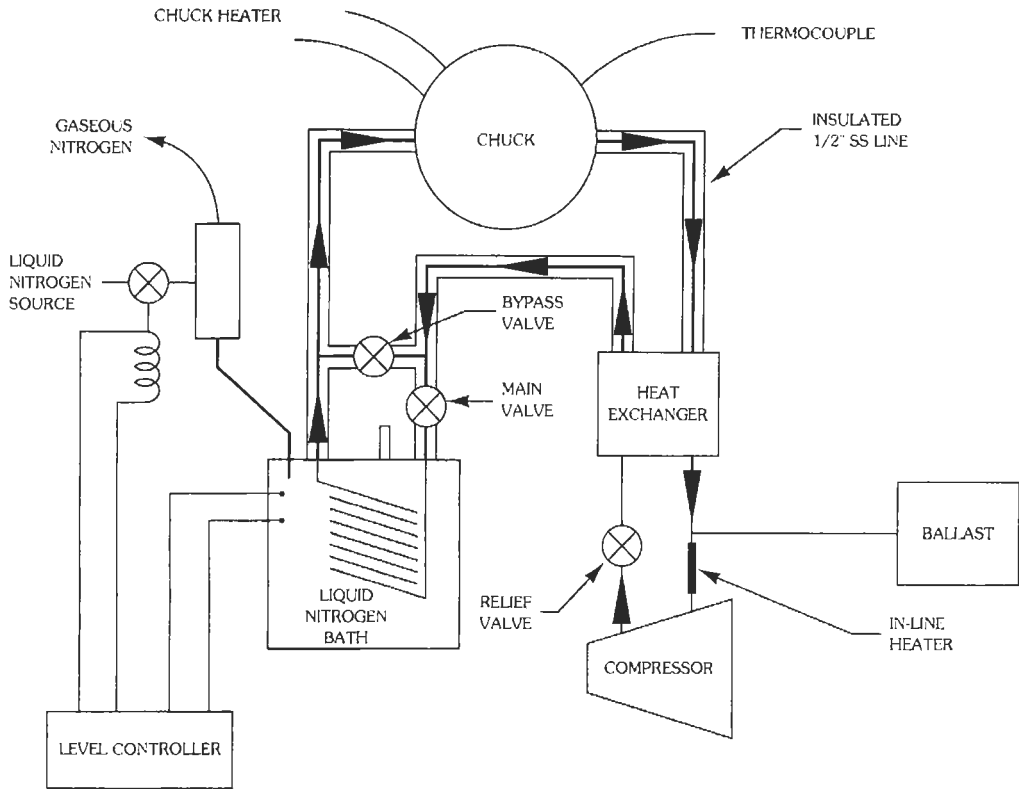


Fig. 87. Schematic diagram of the cryogenic cooling system showing the closed-cycle pumping loop for the high-pressure helium gas coolant. Major system components are indicated.

500 watts. This was done with a loop pressure of around 250 psig and a He mass flow rate of 6 g/s (79 scfm). The heat exchanger effectiveness was 0.75, and the He stream was precooled by 120 °C, down to -110 °C, before arriving at the LN2 dewar. The rate of use of LN2 is no more than 0.4 gpm per 1 kW of cooling load, in this operating range. Initial results with magnetron SF₆ etching of silicon at -150 °C, using this system, confirm the transition to an anisotropic etch regime in this low temperature range.

SWP-Based Process and Device Integration

The greatest potential advantage of SWP for circuit/device design is improved control of structural dimensions and materials quality. Dry etching has shifted to SWP and away from batch processing because each wafer can be etched under the same conditions, thereby decreasing variability between wafers. Furnace operations such as CVD depositions and anneals are now changing to single-wafer RTP for much the same reason. The use of RTP offers some potential advantages for design of MOS devices with design dimensions in the submicron regime where heat treat-

ments must be short to prevent excessive diffusion of dopants.

One major fabrication control problem that seriously impacts device performance is variability in polysilicon gate linewidth. Control margins on linewidth must be held in the range of 10% of the target linewidth. Variability in linewidth across a wafer is a serious problem that can be controlled by improved in situ monitoring of etchant species and temperatures as a function of location on the wafer.

Other serious control problems addressed by SWP with in situ monitoring are contact, via, and metal interconnect etching. It is difficult to know when a contact opening has been sufficiently etched through an interlevel insulator when the opening height-to-width ratio is very large. Another difficulty often encountered is filaments or runners of metal left between patterned and etched metal lines where they run over topographic steps in the insulator; this also can be improved by in situ monitoring. Single-wafer processing can also reduce variability between wafers.

In many cases, quality and parameter control can be improved by grouping several processes into a module to be processed in a cluster tool without manual wafer handling that could introduce particles and

contaminants or cause native oxides to form on silicon. This subsection will be divided into subsubsections discussing groupings of process steps that could be adapted to device fabrication in cluster tools.

SWP techniques can be usefully applied to CMOS, bipolar, or any BiCMOS-like combination. The application to CMOS device processing will be discussed here as an example of the application of SWP. The interrelationship between CMOS device design parameters and SWP with in situ process control will be discussed with reference to Fig. 70. The process modules will be discussed in the order used in fabrication: the discussion will start with the features at the bottom of Fig. 70 and work toward the features at the top of the figure.

Starting Material and CMOS Well Formation—The bottom portion shows the formation of the CMOS n-well and p-well in a lightly doped silicon epitaxial layer deposited on top of a p^+ silicon substrate. In other cases, the CMOS wells may be formed in silicon without an epitaxial layer over p^+ silicon, but this design is more sensitive to CMOS latchup. The use of epitaxial silicon is often avoided for low-cost commodity parts by using n^+/p^+ diffused guard bands around the wells of input/output circuits to prevent CMOS latchup. The cpi thickness is set such that the boron diffusing up from the p^+ substrate just reaches the bottom of the n-well and p-well.

To prevent punchthrough for sub-0.5 μm transistors, the well dopant concentration is greatly increased. All dimensions must be scaled, including the n-well depth into the silicon. By using a 5-min at 1100 °C well drive/anneal in an ammonia ambient (ONED), it is possible to obtain a sufficiently deep n-well to leave an undepleted region of electrons with a sheet resistivity of 1.3 $\text{k}\Omega/\text{square}$ when the n-well is at 3.3 V and the p^+ source/drain and substrate are at zero volts. The short RTP anneal/drive time results in less up-diffusion of boron from the p^+ substrate with the result that: (1) the epitaxial silicon thickness can be decreased to 2.5 μm and (2) the substrate dopant gradient is increased. Both results are useful for preventing CMOS latchup. The phosphorus and boron implants can be both masked by resist patterns and implanted near 300 keV and at a component at lower energy. It is important to note that scaling a design to submicron dimensions requires the use of higher concentration CMOS wells, which then permits the use of implants in the 300 keV regime without long anneal/drive times. For circuits with larger design dimensions, the lower concentration wells would have required the use of MeV implant energies to obtain wells deep enough to leave sufficient undepleted electron concentration in the n-well if short RTP anneal/drive times are used.

It is also important that the p-well concentration reach the boron up-diffusing from the p^+ substrate. This is essential for preventing the formation of an intervening high resistivity layer that would cause premature CMOS latchup. The proper p-well/p-substrate

structure is obtained for the same anneals used for the n-well cases.

Active Area Isolation—The field oxide regions are shown in Fig. 70 isolating the active regions, which include the transistor gates. The field oxide shown is formed by the PBL process, but many other techniques have been used. In the PBL process, an oxide/poly/nitride layer is deposited, patterned, and etched followed by the growth of a thick (0.6 μm) layer of oxide in the isolation regions. Later, the nitride and poly layers remaining after field oxidation must be removed. Several parts of this process are adaptable to implementation in a cluster tool.

Most CMOS fabrication processes use modifications of LOCOS such as PBL, MF3R, or similar but not identical processes. These processes are easily controlled and result in a narrow “bird’s beak” region of gradual increase in field oxide thickness at the edge of the active region. At the time when the poly gate is patterned and etched, the smoothness of the transition makes it easy to etch away undesired polysilicon without leaving any residue.

It is difficult to scale LOCOS isolation to the deep submicron regime for which the isolation regions must be 0.4 μm wide or less. For the deep submicron regime, several shallow trench isolation techniques have been demonstrated. All of these shallow trench techniques share a common problem in that a planarizing etchback must be carefully controlled to prevent topographical steps (up or down) at the edge of the trench. This is an ideal case for the use of SWP with in situ control of the etchback to obtain the exact surface required.

Gate Oxide and Polysilicon Electrode Deposition—One of the most critical steps in CMOS processing is the formation of a high-quality gate oxide and the subsequent deposition of a polysilicon layer that will form the MOS gates. Before these steps, a dummy gate oxide is usually grown and then etched away before growing the gate oxide. The growth of gate oxide and deposition of polysilicon is one of the best candidates for implementation in a cluster tool. If no implants are to be made through the dummy gate oxide, the growth and removal of the dummy gate oxide could also be incorporated into the cluster operation.

In most cases, dopants are added to both the NMOS and PMOS channels to set threshold voltage and increase the punchthrough voltage. The dopants are usually added by implanting through the dummy gate oxide. One of the major design considerations for submicron CMOS is whether to use a surface channel or a so-called buried-channel PMOS device. Both the NMOS and PMOS transistors shown in Fig. 70 are surface-channel. For surface-channel operation, the NMOS gate must be doped n-type and the PMOS gate must be doped p-type. This is usually accomplished by doping the gate only with the implants used to create

the NMOS and PMOS source/drains. For a PMOS "buried-channel" transistor, the polysilicon is doped n-type for both the NMOS and PMOS transistors. For buried-channel, a boron implant must be made into the surface of the silicon. This creates a pn junction immediately underneath the gate oxide. For a properly designed transistor, the boron implant is shallow enough that the p-type region is completely depleted by the built-in voltage of the pn junction. This depleted region makes it easier for the drain voltage to reduce the potential barrier at the source with resultant excess subthreshold leakage for very-short-channel PMOS transistors. For very-short-channel PMOS devices, the surface-channel device is the best choice because the depletion layer is much thinner for surface-channel than for buried-channel.

It is especially important to grow the gate oxide at a low temperature when a buried-channel PMOS device is fabricated by using a boron implant through the dummy gate oxide. If the oxidation temperature is too high, boron diffusion will make it difficult to keep the boron profile shallow. This is especially true for transistors designed for the deep submicron where the pn junction may need to be $0.05\ \mu\text{m}$ deep or less. Surface-channel PMOS design seems better adapted to RTP in dry oxygen ambient at high temperatures. In this case, the boron is implanted into the poly after the growth of the oxide so boron diffusion does not limit the choice of RTP parameters for gate oxidation. Boron diffusion through oxide is enhanced by the presence of hydrogen or fluorine, and this must be considered in the design of the subsequent process steps.

Gate Structure Formation—Following the deposition of polysilicon, the gate is patterned and etched. Later, one or two layers of oxide will be deposited and etched anisotropically to form sidewall spacers at the side of the gate as shown for the NMOS gate and PMOS gate in Fig. 70. As stated earlier, the control of poly linewidth is the most critical parameter in CMOS fabrication. SWP with multiple sensors for etchant species and temperature will be required to obtain sufficient control of linewidth to build transistors with $0.35\ \mu\text{m}$ linewidths and below. It is also extremely important that the etch removes all polysilicon from regions not covered by resist during the etch: sensors are needed to determine when the poly has been removed from all regions, including topography step regions where the poly will be thicker after conformal deposition.

In some cases, the sidewall oxide spacer may be formed by thermal oxidation rather than by deposition. In most cases, implants of light or medium dose may be made before the deposition of one of the sidewall oxide spacers. These so-called LDD or lightly doped drain implants must go under the edge of the gate by lateral straggle at implant or later diffusion. Also of major importance is that the LDD junction be as shallow as possible to decrease the ability of the drain voltage to influence the source potential barrier.

As a final step, a pad oxide will be deposited before the implants of the next module. This pad oxide protects the silicon from contaminants in the photoresist process, randomizes the source/drain implants to minimize channeling, and protects the silicon during anneal of the source/drain implants.

Source/Drain Junctions with or without Silicide—Relatively high-dose arsenic and/or phosphorus implants must be made for the NMOS source/drains and a boron implant for the PMOS source/drains. These implants go through the pad oxide and must be blocked by the field oxide, poly gates, and sidewall oxide spacers adjacent to the gates so that the source/drains are formed only in the active areas not covered by gate. A particular design concern for surface-channel CMOS is that the poly must be thin enough that the source/drain dopants reach the bottom surface of the poly (after anneal). Rapid thermal anneal of the source/drain implants is an obvious choice for SWP. This is particularly important when the junctions must be kept shallow.

Many CMOS processes include the formation of silicide over the active regions and the polysilicon gates. After anneal of the source/drain implants, an additional oxide may be deposited and anisotropically etched to increase the width of the sidewall oxide spacer before silicide processing (see Fig. 70) and leave the polysilicon and active areas free of oxide. In the "salicide" direct react process, titanium metal is deposited; layered titanium nitride/titanium silicide is formed by reacting in a nitrogen ambient; the titanium nitride is etched away leaving silicide over exposed silicon and polysilicon; and finally the remaining silicide is annealed at a slightly higher temperature to obtain the lowest possible sheet resistivity from the TiSi_2 compound. Precise control of titanium thickness, ambient pressures, and react/anneal temperatures is required for this silicide process, making this process an obvious choice implementation with a cluster tool.

It should be noted that the formation of a silicide over the source/drain junctions requires that these junctions be deep enough to prevent shorting by the silicide. Usually when silicide is used, it is not desirable to have very shallow junctions for source/drain; the shallow junction for the CMOS transistors is obtained by the LDD structure formed earlier in the process.

Metallization and Interlevel Isolation—Fig. 70 shows a schematic representation of the formation of two levels of metal interconnect separated by interlevel insulators. The first metal may use tungsten plugs in the vias with copper-doped aluminum as the top interconnect level. The use of other combinations of metals is being investigated. The interlevel insulators will be combinations of deposited undoped and doped oxides. In some processes, a layer of oxide doped with phosphorus and/or boron will be reflowed at a moderate temperature to smooth the insulator surface before

subsequent deposition of a metal layer. Phosphorus-doped oxides also play an important role in preventing contamination to the gate oxide by mobile ions. In other cases, a resist coat or spin-on glass followed by etchback planarization will be used. Etchback planarization processes in general require very accurate measurement of oxide thicknesses and feedback control on etch duration.

One of the most critical problems that can be addressed by SWP with sensors for feedback control is etch of contact and via. For devices designed for deep submicron dimensions, the height-to-width aspect ratio can be in the range of 2–3. This can be a special problem for contacts from metal-1 because the aspect ratio will be different for contacts to active area and to poly over field oxide (see Fig. 70). Vias always connect metal-1 and metal-2, so the height–diameter ratio should be the same for all vias.

If silicide is used to clad the active areas, it is extremely important to limit the temperatures used in the metal/interlevel oxide processes. The problem is that the dopants in the source/drains will diffuse into the silicide. The problem is especially severe for PMOS, in which case it is best to keep these later temperatures 700 °C or lower. In some applications, temperatures as high as 800 °C can be used.

Gate oxide is subject to damage during processing steps after poly gate etch. Dry etching of oxides and metals using high-energy ions can cause radiation damage to the gate oxide. Mechanical stress from subsequent layers can influence the effect of radiation damage on the gate oxide. The effect of damage to the gate oxide is most noticeable for NMOS transistors and can affect the resistance of the NMOS transistor to hot carrier stress. It is well known that hydrogen will weakly bond to the damage centers and also make it easier to create additional damage centers. Hydrogen can also accelerate boron diffusion through gate oxide if the temperatures are too high. In principle, it is desired to minimize the use of hydrogen ambient during later processing, but this is difficult to accomplish. It is important to limit processing temperatures if an ambient with considerable hydrogen partial pressures is used. All of the above parameters should be monitored during SWP.

Protective Overcoat Layer—The final step after etch of the last metal interconnect is to place a protective layer over the metal and remove the protective overcoat over the bond pads.

SILICON-ON-INSULATOR TECHNOLOGY AND DEVICES*

Silicon-on-insulator (SOI) technology is a primary candidate to provide the improved performance, cheaper processing, and unique device structures required for future generations of ICs.† Development

of this technology has been significantly enhanced by the fact that in addition to improved performance and yield, devices fabricated on SOI substrates exhibit significantly increased tolerance to transient radiation and single event upset (SEU).‡ This radiation tolerance has resulted in considerable interest in SOI technology development for devices useful to both space science missions and military electronics.

SOI Material Options

The unique feature of SOI is fabrication of the entire integrated-circuit structure in a thin film of high-quality silicon on the surface of an insulating layer, as illustrated in Fig. 88. This construction provides advantages in both performance and yield for both bipolar and MOS IC structures. Current SOI materials can be considered as follow-ons to the early dielectric isolation technology and SOS technology. The advantages of devices fabricated on these isolated substrates have been demonstrated in both bipolar and CMOS technologies; however both of the prior substrate technologies have significant limitations. The dielectric isolation scheme relies on separation of individual devices by “V” grooves, which are etched into the wafer using an orientation-dependent etch. These “V” grooves require spacings of several mils between devices, which is not compatible with densities required for VLSI circuits. In addition, the thickness tolerances obtainable with dielectric isolation technology do not provide the thin silicon films required for most CMOS device applications.

SOS substrates can be fabricated with thin silicon films and are routinely used for isolated-substrate CMOS device applications. However, the SOS substrates are costly and not readily available in large volumes required for VLSI applications. In addition, until recently the silicon films on SOS were not available in subhalf-micron thicknesses, as anticipated to be required for advanced CMOS devices. A further limitation of SOS technology has been the reluctance (due to concern for cross contamination) of device manufacturers to run SOS in high-volume front ends designed for silicon wafer fabrication.

The current view of SOI materials is a silicon substrate with a “buried oxide” layer in the 200–1000 nm range, covered by a single-crystal silicon film with a thickness from less than 100 nm to greater than 2000 nm. The lower values of buried oxide and silicon film thickness are used for CMOS device applications, while the thicker layers are used for bipolar device applications. This material is illustrated in Fig. 89. This construction is very similar to the dielectric isolation case, without the “V” groove isolation. The advantage of this construction over the conventional dielectric isolated technology is that isolation can now

* This section was contributed by H. H. Hosack.

† References 136–139.

‡ Reference 140.

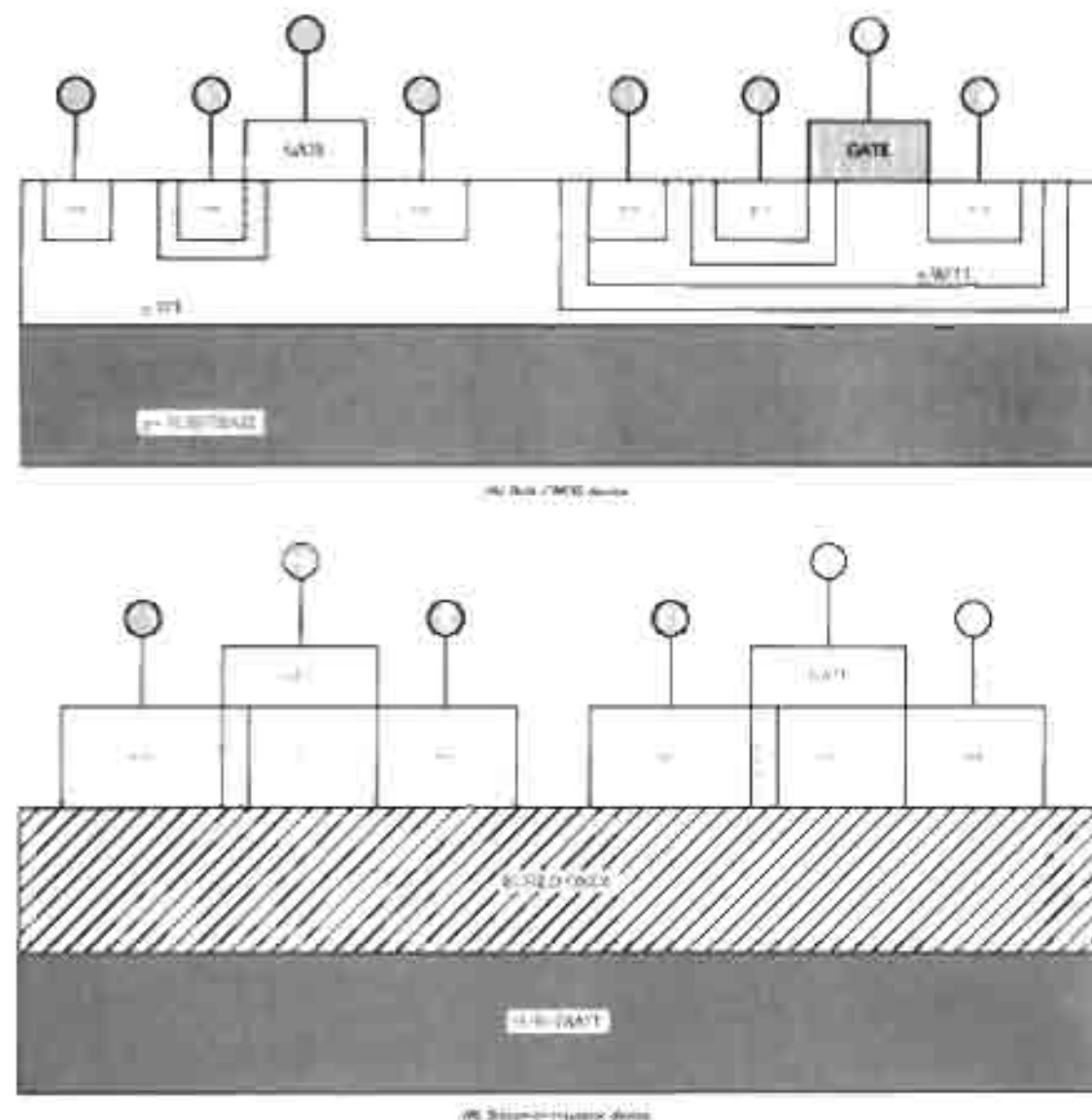


Fig. 88. Silicon-on-insulator device structures.

be provided by simply etching the silicon layer to form "rings" around isolated silicon islands, which can be accomplished at geometries as narrow as allowed by lithographic tools. The advantage of the current SOI materials over SOS is that the all-silicon construction ensures compatibility with large volume (non custom-designed specifically for silicon technologies), as well as allowing potential availability from commercial silicon wafer manufacturers.

Several alternative constructions have been proposed and investigated for producing SOI materials.*

Among these are Separation by Implantation of Oxygen (SIMOX), Bond and Etchback (BESOI), Zone Melting Recrystallization (ZMR), Field Isolation by Porous Oxidation of Silicon (FIPOS), and Epitaxial Lateral Overgrowth (ELO). Each of these techniques has its own specific advantages and disadvantages, but thus far only two, SIMOX and BESOI, have shown the potential for being manufacturable technologies.

* Reference 111.

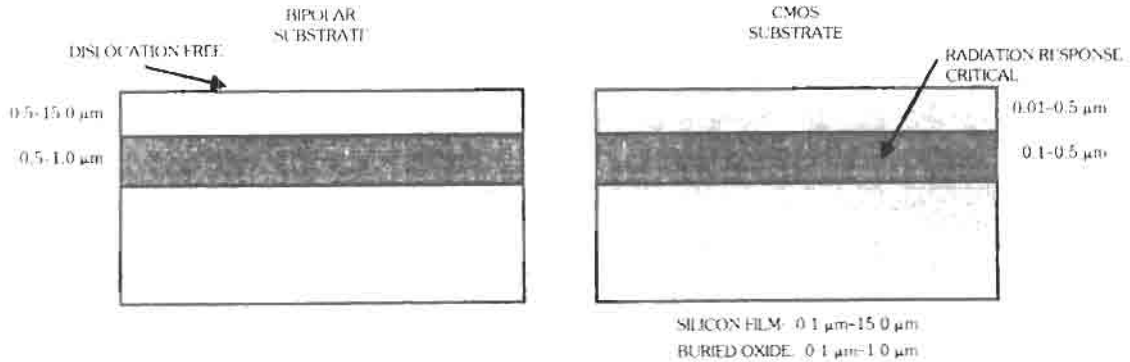


Fig. 89. Comparison of bipolar and CMOS substrates.

compatible with fabrication of generic wafers for both CMOS and bipolar large-area IC fabrication. This conclusion comes primarily from the fact that SIMOX has currently demonstrated the potential to produce both the very thin superficial layers needed for advanced performance CMOS device applications* and BESOI has been demonstrated to produce the thicker, dislocation-free layers needed for bipolar device applications.†

Device Structures Using SOI Materials

Devices built on SOI material have significant advantages to offer in the search for advances in IC technology. These include performance advantages, which add to the goal of increased circuit size and complexity, as well as yield advantages, which add to the ability to fabricate advanced large-area devices at acceptable cost. These advantages are applicable to both commercial and military ICs.

Among the performance advantages of SOI devices, the most significant are reduced susceptibility to "soft errors" from high-energy particles and upset from high-flux transient doses of X rays, elimination of latchup effects, high speed owing to reduced junction capacitance gained from the buried oxide, potential for very low standby power requirements, and reduced threshold lowering in short-channel transistors. In addition, and perhaps more significant for the short-term implementation of SOI in commercial manufacturing, SOI material offers advantages in higher packing density and lower "defect susceptible" area, both of which imply higher yield and lower final IC cost. As an added advantage, it is possible to integrate bipolar circuits into SOI substrates along with CMOS circuits, allowing the use of various device technologies on a

single chip to optimize a broad range of performance parameters. This feature is particularly appropriate when applied to cases where devices may have significantly different voltage-level requirements, as may be needed for true "system-on-a-chip" alternatives—perhaps a key to future low-cost consumer markets and high-density, high-performance military systems.

The performance advantages of SOI ICs noted above are primarily a consequence of the unique construction allowed by the use of SOI substrates. These features can be noted by considering the typical structure of a CMOS/SOI transistor as shown in Fig. 90. Fig. 90A shows a surface view of the device; Fig. 90B shows a cross section through the gate; and Fig. 90C shows a cross section through the current flow channel. This device is constructed in a "mesa" format, where the individual transistors are built on separated single-crystal silicon islands formed on the buried oxide. This construction is uniquely allowed in SOI technology by the use of the insulated substrate to produce dielectric isolation between devices fabricated in the silicon islands. The first performance advantages noted above are reduced susceptibility to the effects of high-energy particles and resistance to high-dose X rays. These features are a direct consequence of the minimum charge collection volume available for an SOI device. This is apparent from Fig. 90, where only the volume bounded by the gate, buried oxide, and junction edges produces susceptibility to charge generation. This compares to the bulk case where there is a much larger depletion region, as well as the possibility of charge collection from the volume below the depleted area.

Since SOI transistors are completely dielectrically isolated, the advantages of no latchup and minimum parasitic capacitance are "built-in" for this construction. The advantage of no latchup is significant, particularly as smaller geometries are considered. Although various techniques are being considered to eliminate latchup in minimum-geometry bulk devices, in many

* Reference 142.

† Reference 143.

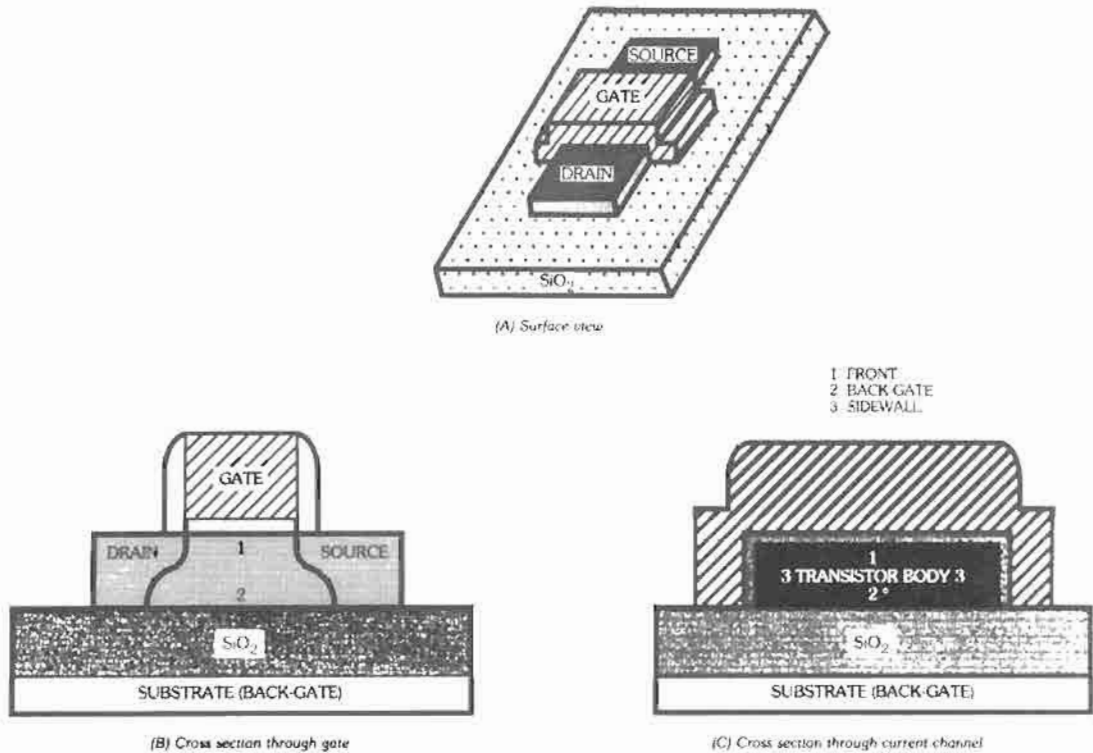


Fig. 90. Schematic of mesa-isolation SIMOX structure showing different channels present.

cases the techniques used impose significant added costs in area or other constraints on both device fabrication and performance.

The magnitude of the advantage to IC performance offered by reduction in capacitance resulting from the insulating substrate in SOI technology varies with device design and construction. Advantages obtained for sRAMs are typically calculated to be in the range of 10–20% reduction in access time. Advantages in propagation delay measured on CMOS dual-mode prescaler circuits have been shown to be more than a factor of 2.* The difference between these advantage factors is a function of the dominant source of capacitance in the circuit.

Perhaps the most significant potential performance advantage for VLSI applications offered by SOI is the possibility of making devices requiring very low standby power. In conventional bulk technology, a significant portion of the standby current can result from the leakage from junction areas that are not critical to transistor performance, such as the bottom of source-drain areas, which are used for contacts. As can be seen in Fig. 90, in the SOI transistor the leakage current associated with an "off" transistor is due only to

the vertical area at the drain/channel junction. This area represents only a small fraction of the total source-drain area that could cause standby current in a bulk transistor of corresponding geometry.

The two fabrication advantages noted above, higher packing density and lower defect-susceptible area, can also be noted from Fig. 90. For the mesa-isolated device, separation between individual circuit elements is accomplished by producing individual silicon mesas. For this case, the minimum isolation width is determined by the minimum level of lithography available, rather than process limitations, such as are often found in LOCOS and trench isolation used for bulk devices.

The advantage of the lower defect-susceptible area noted above can also be seen from Fig. 90. As can be observed, the only "active area" of this device is the small cross section lying directly below the gate of the transistor. All of the remaining areas of the device are only effective in providing contact to the area beneath the gate, and these areas are entirely surrounded by insulation—either the buried oxide on the bottom or the isolation between the mesas. Even if a portion of the mesa is erroneously removed, or if a portion of the contact to the mesa is not opened, the liability is only a reduction in contact area, which in many cases will not

* Reference 144.

significantly impact the device performance. Compare this situation with the bulk silicon case. In this situation, the entire area beneath both of the junctions that comprise the source and drain are critical areas, and any defect in the silicon in this area, or lithography defect that erroneously opens these areas, will cause a fatal flaw in the resulting circuit. This difference in "active area" between the SOI device and the usual bulk technologies can be extremely significant in providing the ability to produce adequate yield on advanced large-area, high-performance structures.

Devices on SOI

The electrical performance of bipolar devices on SOI materials is similar to that of bipolar devices on bulk material, with the advantage of reduced capacitance and reduced susceptibility to latchup and upset by radiation sources. CMOS devices in SOI materials may be fabricated to operate in two different modes.* These are designated as "fully depleted" and "nonfully depleted" devices. From a physical standpoint, the difference between these two configurations is that in the nonfully depleted device the SOI film is doped in such a way that in operation a neutral region exists between the front and back gate depletion regions, thus elimi-

nating coupling between the two parallel MOS devices making up the SOI transistor. In the fully depleted case, there is no neutral region between the two parallel devices, and direct coupling exists between the front gate and the back gate of the SOI device.

Nonfully Depleted Devices

The electrical characteristics of nonfully depleted CMOS/SOI devices are similar to those of bulk devices, with two distinct differences: the appearance of a "kink effect" and enhanced subthreshold slope. Both of these effects are directly related to the fact that the body of the SOI transistor is electrically "floating" between the source and drain junctions. The kink effect is a variation in the I_D versus V_{D_S} characteristics of n-channel MOSFETs caused by impact ionization near the drain. (This effect is not usually as significant for p-channel devices because of the low impact ionization cross section for holes.) This effect is illustrated in Fig. 91; the "kink" is clearly evident. Physically, this effect is one of impact ionization and accumulation of holes in the neutral body region of the device. As the drain bias is increased, impact ionization occurs near the drain junction owing to the high fields. In a bulk device, the holes generated from this mechanism flow into the substrate as substrate current, but in the SOI device this current path is cut off, so the holes accumulate in the neutral body. This raises the

* Reference 145.

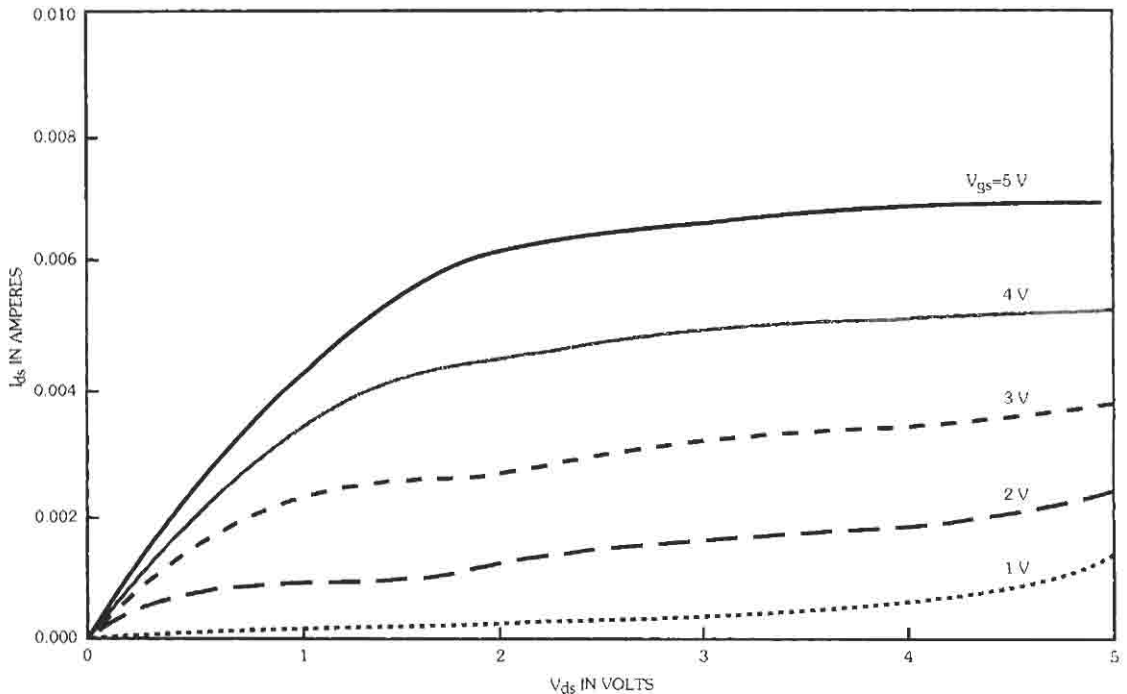


Fig. 91. I_{DS} vs V_{DS} curves for an n-channel MOSFET.

body potential, which in turn raises the drain current through an effective reduction in the threshold voltage and results in the "kink" in the I_{ds} versus V_{ds} curve. The kink effect is significantly affected by the body doping of the device—a higher body doping produces a more pronounced kink effect in SOI devices.

The subthreshold characteristics of MOS transistors on SOI materials can also be affected by the floating body effects. At high V_{ds} , the subthreshold current increases more rapidly with V_{gs} (giving a steeper subthreshold slope) because as V_{gs} increases, the weak inversion current causes ionization near the drain. The holes accumulated in the body from this impact ionization raise the body potential, resulting in a faster increase of the drain current with V_{gs} than would otherwise occur.

In many device structures the kink effect and enhanced subthreshold slope are not significant limitations. However, since these features become more pronounced for smaller geometries, significant efforts on "drain engineering" are used to minimize these effects.

Fully Depleted Devices

When the thickness of the silicon film is thinner than the combined depletion regions of the two parallel MOS transistors making up an SOI device, the V_t is not substantially influenced by the body potential, thus eliminating the kink effect. In addition, other advantages of the fully depleted configuration are higher drive current, reduced hot-carrier effects, and reduced short-channel characteristics.

High drive current in completely depleted SOI transistors is a result of lower channel doping, increased mobility due to lower vertical electric fields, and a more effective use of gate voltage to influence inversion charge. In a fully depleted device, a portion of the surface transistor gate voltage drop occurs across the active device film, and a portion of the drop is across the buried oxide and substrate. In this case, the vertical electric field is lower in the completely depleted device than it is in a comparable bulk device, which results in an increased mobility and increased current. In addition, in the presence of the buried oxide the changes in the gate voltage are more significantly reflected in the inversion charge than is the case for a bulk device. This results in an even greater increase in drive current. Calculations show that a 20% to 50% increase in drive current over comparable bulk devices is available with fully depleted technology.

The inverse subthreshold slope of fully depleted SOI transistors is typically smaller (resulting in faster turn-on) than that for comparable geometry bulk or nonfully depleted SOI devices. In bulk and nonfully depleted SOI devices, a change in the value of V_{gs} is reflected in a change in the voltage drop over the gate oxide and the depletion region in a ratio determined by the capacitive voltage divider between these two device components. In the usual construction of fully depleted devices, however, changes in V_{gs} occur

mostly over the buried oxide, since there is no change in the depletion charge in the silicon film. This results in a much higher percentage of V_{gs} being reflected in changes in the surface potential of the device, and therefore a much smaller inverse subthreshold slope and faster turn-on characteristics. This rapid turn-on of completely depleted SOI devices, coupled with the availability of high current drive, makes this technology ideally suited for high-performance, low-voltage circuit applications.

SOI Circuits

A major drive for the initial development of SOI integrated circuits was applications in radiation hardened CMOS memory devices for defense applications. Prior to SOI, radiation hardened memory circuits were fabricated on silicon on sapphire substrates. These devices performed well in radiation environments, but the silicon on sapphire substrates typically had defect densities that were too high to use in fabrication of VLSI circuits, and they were too brittle to be run in standard silicon fabrication lines. The initial SOI VLSI devices were memories, produced to densities from 64 Kbits to 1 Mbit, fabricated on standard VLSI fabrication facilities with high yield and excellent radiation hardness.* These devices demonstrated the advantage of SOI over the prior silicon on sapphire technology.

The above-mentioned ICs were produced using nonfully depleted SOI technology. Various circuits have been fabricated using fully depleted transistors, and in all cases results superior to comparable bulk ICs have been achieved. A 16K sRAM fabricated as both a completely depleted SOI device and as a bulk device has shown a 40% increase in speed for the completely depleted structure. A high-speed CMOS prescaler circuit was built in completely depleted SIMOX and was found to operate 50% faster than a comparable bulk device. In addition, completely depleted devices with excellent short-channel characteristics at gate lengths below 0.5 μm have recently been discussed.†

SOI based devices received their first major commercial introduction in the year 2000 in the form of advanced microprocessors for IBM server applications. These systems provide a 20% performance advantage because of the use of SOI devices. The SOI device used is designated the AS400 processor. This device operates in the partially depleted SOI mode, and employs design features uniquely applicable to SOI circuits. This effort not only demonstrates that partially depleted SOI circuits can be fabricated in usual IC manufacturing environments, but it also shows examples where SOI unique features can be exploited for a larger performance gain than achievable through a simple bulk-to-SOI re-map. Several

* Reference 146.

† Reference 147.

other microprocessor manufacturers have announced plans to follow the IBM lead.

Three-Dimensional ICs

One of the most advanced applications proposed for SOI is in 3D ICs. In this application, SOI transistors are used as elementary building blocks for 3D ICs, as illustrated schematically in Fig. 92. Transistors in this format would require a succession of SOI process implementations and would allow for extensive "system-on-a-chip" type applications, particularly those requiring massive parallel processing. Several uses of this technology have been proposed, with applications ranging from image-data gathering and processing to neural network applications.*

Conclusion

The anticipated advantages of devices produced on SOI materials have been discussed for several years. These advantages include better performance due to higher drive current and lower capacitance, simpler isolation due to mesa structures, elimination of latchup due to completely oxide-isolated transistors, and extreme tolerance to both transient radiation and SER. Recent progress in fabrication and understanding of SOI materials has not only allowed the experimental verification of many of these advantages but has also resulted in the discovery of significant new effects and possibilities for new device structures. As SOI materials become more readily available in manufacturing quantities, it is anticipated that the advantages of SOI technology will be used extensively in the manufacture

* References 148 and 149.

of high-performance commercial and military electronics.

REFERENCES

1. *Journal of Electrochemical Society*, The Electrochemical Society, Inc., Manchester, NH.
2. *Journal of Vacuum Science and Technology*, published for the American Vacuum Society by the American Institute of Physics, New York.
3. *Semiconductor International*, Cahners Publishing Co., Chicago.
4. Kern, W., and Puotinen, D. "Cleaning solutions based on hydrogen peroxide for use in silicon semiconductor technology." *RCA Review*, June 1970, pp. 187-206.
5. Burger, R. M., and Donovan, R. P. *Fundamentals of Silicon Integrated Device Technology*, Vol. 1. Englewood Cliffs, NJ: Prentice-Hall, Inc., 1967, p. 254.
6. Rosler, R. S. "Low pressure production processes for poly, nitride and oxide." *Solid State Technology*, April 1977, p. 63. Kern, W. and Ban, V. S. "Chemical vapor deposition of inorganic thin films." *Thin Film Processes*, T. L. Vossen and W. Kern, eds. New York: Academic Press, Inc., 1978.
7. Hollahan, J. R., and Rosler, R. S. "Plasma Deposition of inorganic thin films." *Thin Film Processes*, T. L. Vossen and W. Kern, eds. New York: Academic Press, Inc., 1978.
8. Maissel, L. I., and Glang, R., eds. *Handbook of Thin Film Technology*. New York: McGraw-Hill Book Co., 1970, pp. 1-37.
9. Mead, C., and Conway, L. *Introduction to VLSI Systems*. Reading, MA: Addison-Wesley Publishing Co., Inc., 1980.

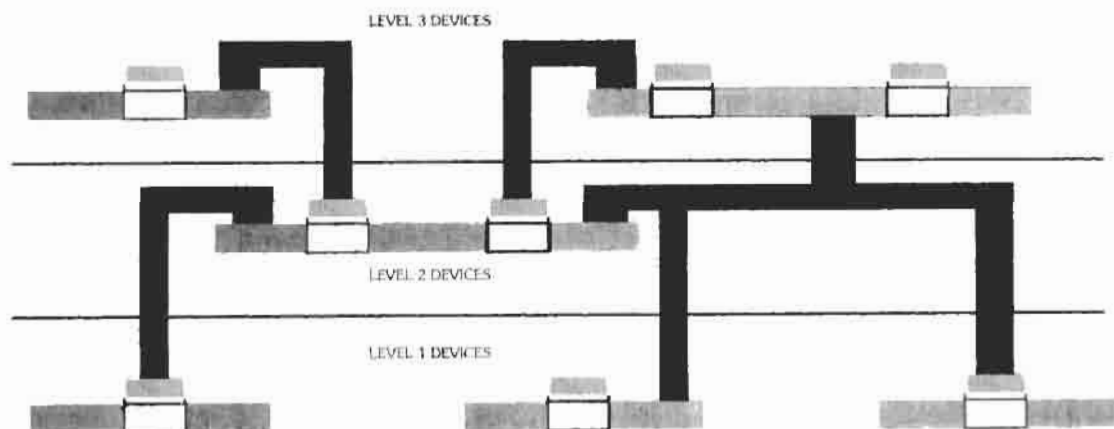


Fig. 92. 3D IC structure.

10. Moran, J. M., and Maydan, D. "High resolution, steep profile resist patterns." *J. Vac. Sci. Technol.*, Vol. 16 (1979), p. 1620.
11. Poulsen, B. "Plasma etching—A review." *J. Vac. Sci. Technol.*, Vol. 14 (1977), p. 266.
12. Trumbore, F. A. "Solid solubilities of impurity elements in germanium and silicon." *Bell Sys. Tech. J.*, Vol. 39 (1960), p. 205.
13. Research Triangle Institute Technical Report, "Integrated silicon device technology." *RTI Review-Reports*, Vol. 4 ASD-TDR-63-316, 1964.
14. Brown, D. M., and Kennicott, P. R. "Glass source B diffusion in Si and SiO₂." *J. Electrochem. Soc.*, 1971, p. 293.
15. Meyer, J. W., Ericksson, L., and Davies, J. A. *Ion Implantation in Semiconductors*. New York: Academic Press, Inc., 1970.
16. Gibbons, T. F., Johnson, W. S., and Mylroie, S. W. *Projected Range Statistics—Semiconductors and Related Materials*. 2nd ed. New York: Halsted Press, 1975.
17. Morris, R. L., and Miller, J. R. *Designing with TTL Integrated Circuits*. Texas Instruments Electronics Series. New York: McGraw-Hill Book Co., 1971.
18. Hart, K., and Slob, A. "Integrated injection logic." *IEEE Trans Solid State Circuits*. SC-7 (1972), p. 346.
19. Sloan, B. J. "STL technology." *IEEE IEDM Tech Digest*, 1979, p. 324.
20. H. Corporaal. *Microprocessors Architectures: From VLIW to TTA*. West Sussex, England: John Wiley & Sons, Inc., 1998.
21. R. Henle and B. Kuvshinoff, *Desktop Computers*. New York: Oxford University Press, Inc., 1992.
22. Elmasry, M. I., ed. *Digital MOS Integrated Circuits*. IEEE Press, John Wiley & Sons, Inc., distributor, 1981.
23. Roberge, J. K. *Operational Amplifiers: Theory and Practice*, New York: John Wiley & Sons, Inc., 1975.
24. Tobey, G. E., Graeme, J. G., and Huelsman, L. P. *Operational Amplifiers: Design and Applications*. New York: McGraw-Hill Book Co., 1971.
25. Graeme, J. G. *Applications of Operational Amplifiers: Third-Generation Techniques*. New York: McGraw-Hill Book Co., 1973.
26. Wong, Y. J., and Ott, W. E. *Function Circuits: Design and Applications*. New York: McGraw-Hill Book Co., 1976.
27. Graeme, J. G. *Designing with Operational Amplifiers: Application Alternatives*. New York: McGraw-Hill Book Co., 1977.
28. Markel, J. D., and Gray, A. H., Jr. *Linear Prediction of Speech*. New York: Springer-Verlag, 1976.
29. *ASIC Outlook 1991*. Scottsdale, AZ: Integrated Circuit Engineering Corporation.
30. "ASICs—Catalyst for growth." *ASIC Technology & News*, Vol. 3, No. 2, June 1991.
31. Haskard, M. R. *An Introduction to Application Specific Integrated Circuits*. Englewood Cliffs, NJ: Prentice Hall, 1989.
32. "Application specific LSIs for specialized and short life products enjoy an expanding market." *J. Electronic Engineering* (Japan; language, English), Vol. 26, No. 265, January 1989, p. 28.
33. "Gate arrays ride ASIC popularity to greater demand." *J. Electronic Engineering* (Japan), Vol. 27, No. 277, January 1990, p. 28.
34. Wilson, R. "Hitachi arrays focus on I/O." *EE Times*, No. 659, September 16, 1991.
35. Bursky, D. "Core CPU and cells ease X86 design." *Electronic Design*, Vol. 39, No. 22, November 21, 1991, p. 157.
36. Alford, R. C. *Programmable Logic Designer's Guide*. Indianapolis, IN: Howard W. Sams & Co., 1989.
37. El-Ayat, K. "A CMOS electrically configurable gate array." *Proceedings of the 1988 IEEE International Solid-State Circuits Conference*, New York, p. 76.
38. Baker, S. "Quicklogic unveils FPGA." *Electronic Engineering Times*, April 22, 1991.
39. Bursky, D. "FPGAs mirror masked gate-array architecture." *Electronic Design*, Vol. 39, No. 22, November 21, 1991, p. 63.
40. Goodenough, F. "Cut product size and cost with mixed signal ASICs." *Electronic Design*, Vol. 39, No. 22, November 21, 1991.
41. Pryce, D. "Analog-digital ICs provide versatility." *Electronic Design*, Vol. 35, No. 5, March 1, 1990.
42. Ismail, M., and Bibyk, S. "CAD latches onto new techniques for analog ICs." *IEEE Circuits and Devices*, Sept. 1991, pp. 11–17.
43. Lee, E., and Gulak, P. "A CMOS field-programmable analog array." *IEEE Journal of Solid-State Circuits*, Vol. 26, No. 12, Dec. 1991, pp. 1860–1867.
44. Andrews, W. "ASIC memories: Bigger, faster, and customized." *Computer Design*, Vol. 22, No. 18, October 1, 1988, pp. 44–62.
45. Licari, F. J., and Enlow, L. R. *Hybrid Microcircuit Technology Handbook*. Westwood, NJ: Noyes Publications, 1998, p. 20.
46. Ahrens, M. "An FPGA family optimized for high densities and reduced routing delay." *Proceedings of the 1990 IEEE Custom Integrated Circuits Conference, Boston*, p. 31.5.
47. Hsieh, H. C. "Third-generation architecture boosts speed and density of field-programmable gate arrays." *Proceedings of the IEEE 1990 Custom Integrated Circuits Conference, Boston*, p. 31.2.
48. Tuck, B. "Denser, faster FPGAs encroach further on masked gate arrays." *Computer Design*, Vol. 29, No. 17, September 1, 1990, p. 30.
49. Moslehi, M., Chapman, R., Wong, M., Paranjpe, A., Najm, H., Kuehne, J., Yeakley, R., and Davis, C. "Single-wafer integrated semiconductor device

- processing." *IEEE Trans. on Electron Devices*, Vol. 39, No. 1, 1992, p. 4.
50. Moslehi, M., Kuehne, J., Velo, L., Yin, D., Yeakley, D., Huang, S., Jucha, B., and Breedijk, T. "RTP for advanced CMOS process integration." *Proc. of SPIE on Rapid Thermal and Integrated Processing*, 1992.
 51. Moslehi, M., Kuehne, J., Yeakley, R., Velo, L., Najm, H., Dostalick, B., Yin, D., and Davis, C. "In-situ fabrication and process control techniques in rapid thermal processing." *MRS Proc. on Rapid Thermal and Integrated Processing*, Vol. 224, 1991, p. 143.
 52. Moslehi, M. M. "Process uniformity and slip dislocation patterns in linearly ramped-temperature transient rapid thermal processing of silicon." *IEEE Trans. on Semiconductor Manufacturing*, Vol. 2, No. 4, 1989, p. 130.
 53. Gyurcsik, R. S., Riley, T. J., and Sorrell, F. Y. *IEEE Trans. on Semiconductor Manufacturing*, Vol. 4, No. 1, 1991, p. 9.
 54. Campbell, S., Knutson, K., Ahn, K., Leighton, J., and Liu, B. "Gas flow patterns and thermal uniformity in rapid thermal processing equipment." *IEDM Tech. Dig.*, 1990, p. 921.
 55. Sorrel, F. *SRC workshop on temperature measurement*, Santa Fe, New Mexico, Feb. 1990.
 56. Nulman, J., Cohen, B., Blonigan, W., Antonio, S., Meinecke, R., and Gat, A. "Pyrometric emissivity measurements and compensation in an RTP chamber." *MRS Symp. Proc.*, Vol. 146, 1989, p. 461.
 57. Sato, T. *Jpn. J. Appl. Phys.*, Vol. 6, 1967, p. 339.
 58. Gelpey, J., and Liao, J. "Integrated temperature control with automatic emissivity compensation for RTP." *SRC workshop on temperature measurement*, Santa Fe, New Mexico, Feb. 1990.
 59. Kamins, T., Bradbury, D., Cass, T., Laderman, S., and Reid, G. "Structure of LPCVD tungsten films for IC applications." *J. Electrochem. Soc.*, Vol. 131, Dec. 1986, p. 2555.
 60. Kern, W. "The evolution of silicon cleaning technology." *Proc. of the First International Symp. on Cleaning Technol. in Semiconductor Device Mfg. (The Electrochem. Soc.)*, Vol. 90-9, 1989, p. 3.
 61. Deal, B., and Kao, D. "The physics and chemistry of thin native oxide films on silicon," *Mat. Res. Soc. Symp.*, 1987.
 62. Moslehi, M., Shatas, S., Saraswat, K., and Mcindl, J. "Interfacial and breakdown characteristics of MOS devices with rapidly grown ultra-thin SiO₂ gate insulators." *IEEE Trans. Electron Devices*, Vol. ED-34, No. 6, 1987, pp. 1407-1410.
 63. Morita, M., Ohmi, T., Hasegawa, E., Kawakami, M., and Suma, K. "Control factor of native oxide growth on silicon in air or in ultra-pure water." *Appl. Phys. Lett.*, Vol. 55, No. 6, 1989, pp. 562-564.
 64. Hirashita, N., Kinoshita, M., Aikawa, I., and Ajioka, T. "Effects of surface hydrogen on the air oxidation at room temperature of HF-treated Si(100) surfaces." *Appl. Phys. Lett.*, Vol. 56, No. 5, 1990, pp. 451-453.
 65. Crabbe, E., Hoyt, J., Moslehi, M., Pease, R., and Gibbons, J. "Reduced process sensitivity of polysilicon emitter contacts for VLSI bipolar transistors." *Symp. on VLSI Technol. Dig. Tech. Papers*, 1987, pp. 33-34.
 66. Burrows, V., Chabal, Y., Higashi, G., Raghavachari, K., and Christman, S. "Infrared spectroscopy of Si(111) surfaces after HF treatment: Hydrogen termination and surface morphology." *Appl. Phys. Lett.*, Vol. 53, No. 11, 1988, pp. 998-1000.
 67. Miyauchi, A., Inoue, Y., and Suzuki, T. "Interface impurities of low-temperature (900 °C) deposited Si epitaxial films prepared by HF treatments." *Appl. Phys. Lett.*, Vol. 57, No. 7, 1990, p. 676.
 68. Friedrich, J., and Neudeck, G. "Limitations in low-temperature silicon epitaxy due to water vapor and oxygen in the growth ambient." *Appl. Phys. Lett.*, Vol. 53, No. 25, 1988, pp. 2543-2545.
 69. Ruzyllo, J., Hoff, A., Frystak, D., and Hossain, S. "Electrical evaluation of wet and dry cleaning procedures for silicon device fabrication." *J. Electrochem. Soc.*, Vol. 136, No. 5, 1989, pp. 1474-1476.
 70. Mieno, F., Nakamura, S., Deguchi, T., Maeda, M., and Inayoshi, K. "Low temperature silicon epitaxy using Si₂H₆." *J. Electrochem. Soc.*, Vol. 134, No. 9, 1987, pp. 2320-2323.
 71. Yamazaki, T., Minakata, H., and Ito, T. "Continuous growth of heavily doped p⁺-n⁺Si epitaxial layer using low-temperature photoepitaxy." *Appl. Phys. Lett.*, Vol. 55, No. 9, 1989, pp. 879-881.
 72. Yamazaki, T., Watanabe, S., and Ito, T. "Heavy boron doping in low-temperature Si photoepitaxy." *J. Electrochem. Soc.*, Vol. 137, No. 1, 1990, pp. 313-318.
 73. Meyerson, B., LeGoues, F., Nguyen, T., and Hareme, D. "Nonequilibrium boron doping effects in low-temperature epitaxial silicon films." *Appl. Phys. Lett.*, Vol. 50, No. 2, 1987, pp. 113-115.
 74. Meyerson, B., Uram, K., and LeGoues, F. "Cooperative growth phenomena in silicon/germanium low-temperature epitaxy." *Appl. Phys. Lett.*, Vol. 53, No. 25, 1988, pp. 2555-2557.
 75. Shibata, T., Kondo, N., and Nanishi, Y. "Si surface cleaning and epitaxial growth of GaAs on Si by electron cyclotron resonance plasma-excited molecular-beam-epitaxy at low temperatures." *J. Electrochem. Soc.*, Vol. 136, No. 11, 1989, pp. 3459-3462.
 76. Salimian, S., Cooper, C., and Ellingboe, A. "Etching of SiO₂ in an electron cyclotron resonance argon plasma." *Appl. Phys. Lett.*, Vol. 56, No. 14, 1990, pp. 1311-1313.
 77. Moslehi M., and Davis, C. "Advanced epitaxial Si and Ge_xSi_{1-x} multiprocessing for semiconductor

- device technologies." *J. Mat. Res.*, Vol. 5, No. 6, 1990, pp. 1159-1162.
78. Miki, N., Kikuyama, H., Kawanabe, I., Miyashita, M., and Ohmi, T. "Gas-phase selective etching of native oxide." *IEEE Trans. Electron Devices*, Vol. 37, No. 1, 1990, pp. 107-115.
 79. Syverson, D., and Duranko, G. "Contamination aspects of anhydrous HF gas processing." *Solid State Technol.*, 1988, pp. 101-104.
 80. Deal, B., McNeilly, M., Kao, D., and de Larios, J. "Vapor phase wafer cleaning thin film, oxide etching, and thin film growth." *Proc. of the First International Symp. on Cleaning Technol. in Semiconductor Device Mfg. (The Electrochem. Soc.)*, Vol. 90-9, 1989, p. 121.
 81. Sedgwick, T., Berknblit, M., and Kuan, T. "Low-temperature selective epitaxial growth of silicon at atmospheric pressure." *Appl. Phys. Lett.*, Vol. 54, No. 26, 1989, pp. 2689-2691.
 82. Campbell, S., Leighton, J., Case, G., and Knutson, K. "Very thin silicon epitaxial layers grown using rapid thermal vapor phase epitaxy." *J. Vac. Sci. Technol.*, Vol. B7, No. 5, 1989, pp. 1080-1083.
 83. Regolini, J., and Bensahel, D. "Selective epitaxial silicon growth in the 650-1100 °C range in a reduced pressure chemical vapor deposition reactor using dichlorosilane." *Appl. Phys. Lett.*, Vol. 54, No. 7, 1989, pp. 658-659.
 84. Gibbons, J., King, C., Hoyt, J., Noble, D., Gronet, C., Scott, M., Rosner, S., Reid, G., Laderman, S., Nauka, K., Turner, J., and Kamins, T. "Si/Si_{1-x}Ge heterojunction bipolar transistors fabricated by limited reaction processing." *IEDM Tech. Dig.*, 1988, pp. 566-569.
 85. Yew, T., and Reif, R. "Silicon selective epitaxial growth at 800 °C using SiH₄/H₂ assisted by H₂/Ar plasma sputter." *Appl. Phys. Lett.*, Vol. 55, No. 10, 1989, pp. 1014-1016.
 86. Ohmi, T., Hashimoto, K., Morita, M., and Shibata, T. "In situ-doped epitaxial silicon film grown at 250 °C by an ultra-clean low-energy bias sputtering." *IEDM Tech. Dig.*, 1989, pp. 53-56.
 87. Ohmi, T. "Future trends and applications of ultra-clean technology." *IEDM Tech. Dig.*, 1989, pp. 49-52.
 88. Gao, Q., Hariu, T., and Ono, S. "Low-temperature surface cleaning of Si and successive plasma-assisted epitaxial growth of GaAs." *Jpn. J. Appl. Phys.*, Vol. 26, No. 10, 1987, pp. L1576-L1578.
 89. Sherman, A. "In situ removal of native oxide from silicon wafers." *J. Vac. Sci. Technol. B*, Vol. 8, No. 4, 1990, pp. 656-657.
 90. Miyake, K. "Removal of a thin SiO₂ layer by low-energy hydrogen ion bombardment at elevated temperatures." *Jpn. J. Appl. Phys.*, Vol. 28, No. 11, 1989, pp. 2376-2381.
 91. Yamada, H. "Low-temperature surface cleaning method using low-energy reactive ionized species." *J. Appl. Phys.*, Vol. 65, No. 2, 1989, pp. 775-781.
 92. Suehune, I., Kishimoto, A., Hamaoka, K., Honda, Y., Kan, Y., and Yamanishi, M. "Dependence of GaAs etch rate on the angle of incidence of a hydrogen plasma beam excited by electron cyclotron resonance." *Appl. Phys. Lett.*, Vol. 56, No. 24, 1990, pp. 2393-2395.
 93. Breaux, L., Anthony, B., Hsu, T., Banerjee, S., and Tasch, A. "Homoepitaxial films grown in Si(100) at 150 °C by remote plasma-enhanced chemical vapor deposition." *Appl. Phys. Lett.*, Vol. 55, No. 18, 1989, pp. 1885-1887.
 94. Burns, G. "Low-temperature native oxide removal from silicon using nitrogen trifluoride prior to low-temperature silicon epitaxy." *Appl. Phys. Lett.*, Vol. 53, No. 15, 1988, pp. 1423-1425.
 95. Saito, Y., Yamaoka, O., and Yoshida, A. "Plasmaless cleaning process of silicon surface using chlorine trifluoride." *Appl. Phys. Lett.*, Vol. 56, No. 12, 1990, pp. 1119-1121.
 96. Hirose, F., Suemitsu, M., and Miyamoto, N. "High quality silicon epitaxy at 500 °C using silane gas-source molecular beam technique." *Jpn. J. Appl. Phys.*, Vol. 28, No. 11, 1989, pp. L2003-L2006.
 97. Murota, J., Nakamura, N., Kato, M., and Miko-shiba, N. "Low-temperature silicon selective deposition and epitaxy on silicon using the thermal decomposition of silane under ultraclean environment." *Appl. Phys. Lett.*, Vol. 54, No. 11, 1989, pp. 1007-1009.
 98. Morar, J., Meyerson, B., Karisson, U., Himpfel, F., and McFeely, F. "Oxygen removal from Si via reaction with absorbed Ge." *Appl. Phys. Lett.*, Vol. 50, No. 8, 1987, pp. 463-465.
 99. Novak, R. "Anhydrous HF etching of native SiO₂: Applications to device fabrication." *Solid State Technol.*, March 1988, pp. 39-41.
 100. Deal, B., McNeilly, M., Kao, D., and deLarios, J. "Vapor phase wafer cleaning: Processing for the 1990s." *Solid State Technol.*, July 1990, pp. 73-77.
 101. Miki, N., Kikuyama, H., Maeno, M., Murota, J., and Ohmi, T. "Selective etching of native oxide by dry processing using ultra clean anhydrous hydrogen fluoride." *IEDM Tech. Dig.*, 1988, pp. 730-733.
 102. Nara, Y., Yamazaki, T., Sugii, T., Sugino, R., Ito, T., and Ishikawa, H. "Photochemical cleaning and epitaxy of Si." *SPIE Symp. on Advances in Semiconductors and Semiconductor Structures*, Panama City, 1988.
 103. Takahagi, T., Nagai, I., Ishitani, A., and Kuroda, H. "The formation of hydrogen passivated silicon single-crystal surfaces using ultraviolet cleaning and HF etching." *J. Appl. Phys.*, Vol. 64, No. 7, 1988, pp. 3516-3521.
 104. Vig, J. "UV/ozone cleaning of surfaces." *Proc. of the First International Symp. on Cleaning Technol.*

- in *Semiconductor Device Mfg. (The Electrochem. Soc.)*, Vol. 90-9, 1989, p. 105.
105. Itoh, T., Sugino, R., Watanabe, S., Nara, Y., and Sato, Y. "UV enhanced dry cleaning of silicon wafers." *Proc. of the First International Symp. on Cleaning Technol. in Semiconductor Device Mfg. (The Electrochem. Soc.)*, Vol. 90-9, 1989, p. 114.
 106. Logan, J. S., Costable, J., and Jones, F., "High rate radio frequency sputtering using in-phase plasma confinement." *J. Vac. Sci. Technol. A*, Vol. 8, 1990, p. 1835.
 107. Horwitz, C. M. "Silicon deposition in diode and hollow cathode systems." *J. Vac. Sci. Technol. B*, Vol. 7, 1989, p. 443.
 108. Horwitz, C. M., Boronkay, S., Gross, M., and Davies, K. "Hollow cathode etching and deposition." *J. Vac. Sci. Technol. A*, Vol. 6, 1988, p. 1837.
 109. Davies, K. E., and Horwitz, C. M. "Diode and hollow cathode etching in CF_3 ." *J. Vac. Sci. Technol. A*, Vol. 7, 1989, p. 2705.
 110. Shirakawa, T., Toyoda, H., and Sugai, H. "RF plasma production at ultralow pressures with surface magnetic confinement." *J. Vac. Sci. Technol. B*, Vol. 7, 1989, p. 443.
 111. Cook, J. M., Ibbotson, D. E., and Flamm, D. L. "Application of a low-pressure radio frequency discharge source to polysilicon etching." *J. Vac. Sci. Technol. B*, Vol. 8, 1990, p. 1.
 112. Kitagawa, A., Takeuchi, M., and Kasagi, M. "Enhanced growth of silicon dioxide films by parallel resonant RF plasmas." *Jpn. J. Appl. Phys.*, Vol. 7, 1990, p. L1178.
 113. Keller, J. H. "Novel RF induction plasma processing techniques." *Abstracts of 37th Amer. Vac. Soc. Symp.*, Toronto, 1990.
 114. Boswell, R. W., and Porteous, R. K. "Large volume high density RF inductively coupled plasma." *Appl. Phys. Lett.*, Vol. 50, 1987, p. 1130.
 115. Boswell, R. W., Perry, A. J., and Emami, M. "Multipole confined diffusion plasma produced by 13.56 MHz electrodeless source." *J. Vac. Sci. Technol. A*, Vol. 7, 1989, p. 3345.
 116. Thornton, J. A. "Magnetron sputtering: Basic physics and applications to cylindrical magnetrons." *J. Vac. Sci. Technol.*, Vol. 15, 1978, p. 171.
 117. Yeom, G. Y., and Kushner, M. J. "Magnetic field effects on cylindrical reactive ion etching of Si/SiO_2 in CF_4/H_2 plasmas." *J. Vac. Sci. Technol. A*, Vol. 7, 1989, p. 987.
 118. Bright, A. A., and Kaushik, S. "Magnetron ion etching with CF_4 based plasmas: Effects of magnetic field on plasma chemistry." *J. Vac. Sci. Technol. B*, Vol. 7, 1990, p. 542.
 119. Sheridan, T. E., Goeckner, M. J., and Goree, J. "Model of energetic electron transport in magnetron discharges." *J. Vac. Sci. Technol. A*, Vol. 8, 1990, p. 30.
 120. Wendt, A. E., and Lieberman, M. A. "Spatial structure of a planar magnetron discharge." *J. Vac. Sci. Technol. A*, Vol. 8, 1990, p. 902.
 121. Okano, H., Horiike, Y., and Yamazaki, T. "Etching characteristics of n^+ poly-Si and Al employing a magnetron plasma." *Jpn. J. Appl. Phys.*, Vol. 23, 1984, p. 482.
 122. Asmussen, J. "Electron cyclotron resonance microwave discharges for etching and thin film deposition." *J. Vac. Sci. Technol. A*, Vol. 7, 1989, p. 883.
 123. Shufflebotham, P. K., and Thompson, D. J. "Diagnostic techniques for plasma stability in electron cyclotron plasma processing." *J. Vac. Sci. Technol. A*, Vol. 8, 1990, p. 3713.
 124. Uetake, H., Matsuura, T., and Ohmi, T. "Anisotropic etching of n^+ polycrystalline silicon with high selectivity using a chlorine and nitrogen plasma in an ultraclean electron cyclotron resonance etcher." *Appl. Phys. Lett.*, Vol. 57, 1990, p. 596.
 125. Matsuo, S., and Adachi, Y. "Reactive ion etching using a broad beam ECR ion source." *Jpn. J. Appl. Phys.*, Vol. 21, 1982, p. L4.
 126. Burke, R. R., Pelletier, J., Pornot, C., and Vallier, L. "Distributed electron cyclotron resonance in silicon processing: epitaxy and etching." *J. Vac. Sci. Technol. A*, Vol. 8, 1990, p. 2931.
 127. Constantine, C., Johnson, D., Pearton, S. J., Chakrabarti, U. K., Emerson, A. B., Hobson, W. S., and Kinsella, A. P. "Plasma etching of III-V semiconductors in $CH_4/H_2/Ar$ electron cyclotron resonance discharges." *J. Vac. Sci. Technol. A*, Vol. 8, 1990, p. 596.
 128. Popov, O. A., and Waldron, H. "Electron cyclotron resonance plasma stream source for plasma enhanced chemical vapor deposition." *J. Vac. Sci. Technol. A*, Vol. 7, 1989, p. 914.
 129. Wei, J., Kawarada, H., and Suzuki, J. "Low temperature synthesis of diamond films using magneto-microwave plasma CVD." *Jpn. J. Appl. Phys.*, Vol. 29, 1990, p. L1483.
 130. Ono, T., Oda, M., Takahashi, C., and Matsuo, S. "Reactive ion stream etching utilizing electron cyclotron resonance plasma." *J. Vac. Sci. Technol. B*, Vol. 4, 1988, p. 696.
 131. Holber, W. M., and Forster, J. "Ion energetics in electron cyclotron resonance discharges." *J. Vac. Sci. Technol. A*, Vol. 8, 1990, p. 3720.
 132. Samukawa, S., Mori, S., and Sasaki, M. "Dependence of ECR plasma etching characteristics on sub magnetic field and substrate position." *Jpn. J. Appl. Phys.*, Vol. 29, 1990, p. 792.
 133. Hopwood, J., Reinhard, D. K., and Asmussen, J. "Charged particle densities and energy distributions in a multipolar electron cyclotron resonance plasma etching source." *J. Vac. Sci. Technol. A*, Vol. 8, 1990, p. 3103.
 134. Pichot, M., and Durandet, A. "Microwave multipolar plasmas excited by distributed electron cyclotron resonance: Concept and performance." *Rev. Sci. Instrum.*, Vol. 59, 1988, p. 1072.
 135. Tachi, S., Tsujimoto, K., and Okudaira, S. "Low-temperature reactive ion etching and microwave

- plasma etching of silicon." *Appl. Phys. Lett.*, Vol. 52, No. 8, 1988, p. 616.
136. Izumi, K., et al. *Trans. Inst. Elect. and Comm. Eng. Jpn. Sect. E (Japan)*, Vol. E62, No. 12, Dec. 1979, p. 921.
137. Burnham, M. E., and Wilson, S. R. *Proc. SPIE*, Vol. 530, 1985, p. 240.
138. Hemment, P. L. F. *Mater. Res. Soc. Symp. Proc.*, Vol. 53, 1986, p. 207.
139. Mao, B.-Y., et al. *IEEE Trans. Nucl. Sci.* Vol. NS-33, No. 6, 1986, p. 1702.
140. Davis, G. E., et al. *IEEE Trans. Nucl. Sci.*, Vol. NS-32, No. 6, 1985, p. 4432.
141. Colinge, J. P. *Silicon On Insulator Technology: Materials to VLSI*. Dordrecht: Kluwer Academic Publishers, 1991.
142. Bailey, W. E., Lu, H., Blake, T. G. W., Hite, L. R., Mei, P., Hurta, D., Houston, T. W., and Pollack, G. P. *1991 IEEE International SOI Conference Proceedings*, 1991, p. 134.
143. McKitterick, J. B., and Caviglia, A. L. *IEEE Trans. Elect. Dev.*, Vol. 36, No. 6, 1989, p. 1133.
144. Celler, G. K., et al. *1989 IEEE SOS/SOI Technology Conference Proceedings*, 1989, p. 139.
145. Mayer, D. *1989 IEEE SOS/SOI Technology Conference Proceedings*, 1989, p. 52. Hashimoto, K., et al. *IEDM Technical Digest*, 1985, p. 672.
146. Houston, T. W., Lu, H., Mei, P., Blake, T. G. W., Hite, L. R., Sundaresan, R., Matloubian, M., Bailey, W. E., Liu, J., Peterson, A., and Pollack, G. *IEEE SOS/SOI Technology Conference Proceedings*, 1989, Paper L1, p. 137.
147. Miki, H., Omura, Y., Ohmancuda, T., Kumon, M., Asada, K., Izumi, K., Sakai, T., and Sugano, T. *IEDM Tech. Dig.*, 1989, Paper 34.7.1, p. 906.
148. Yamazaki, K., et al. *Ext. Abstracts of the 8th Internat. Workshop on Future Electron Devices—Three Dimensional ICs and Nano-Meter Functional Devices, Kochi Prefecture, Japan*, 1990, p. 105.
149. Aibara, R., et al. *Ext. Abstracts of the 8th Internat. Workshop on Future Electron Devices—Three Dimensional ICs and Nano-Meter Functional Devices, Kochi Prefecture, Japan*, 1990, p. 113.

21

Optoelectronics

*Revised and Expanded by
Gregory E. Stillman*

The Optical Spectrum	21-2
Radiometry	21-3
Terms and Definitions	
Blackbody Radiation	
Interaction of Optical Waves With Matter	21-9
Reflectance	
Absorptance	
Transmittance	
Refraction	
Optical Sources	21-10
Tungsten Lamps	
Fluorescent Lamps	
Arc Lamps	
Light-Emitting Diodes	
High-Intensity Visible LEDs	
Superluminescent Diodes	
Semiconductor Lasers	
Optical Detectors	21-14
Terms and Figures of Merit	
Characterization of Detectors	
Ultimate Sensitivity of Detectors	
Thermal Detectors	
Quantum Detectors	
Spectral Response of Semiconductor Detectors	

Optoelectronics is the technological marriage of the fields of optics and electronics. It includes the generation and evaluation of electromagnetic radiation in the optical wavelength range and its conversion into electrical current or signals, the interaction of light with matter, radiometry, and the characteristics of sources and detectors.

THE OPTICAL SPECTRUM

The optical spectrum is generally defined to encompass electromagnetic radiation with wavelengths in the range from 10 nm to 10³ μm, or frequencies in the range from 300 GHz to 3000 THz (Fig. 1).

Other units are often used to describe the optical spectrum, as discussed below. For engineering calculations, three significant figures usually give sufficient accuracy. More precise values of the constants can be obtained from NBS/NIST data.

Wavelength λ:

$$1 \mu\text{m} = 10^{-3} \text{ mm} = 10^3 \text{ nm} = 10^4 \text{ \AA}$$

Frequency (ν or f):

$$1 \text{ cm}^{-1} \text{ (wave number)} = 30 \text{ GHz} = 3 \times 10^{10} \text{ Hz}$$

$$\nu \text{ (or } f) = \frac{c}{\lambda} = \frac{3 \times 10^{14}}{\lambda \text{ (}\mu\text{m)}} \text{ Hz}$$

$$\bar{\nu} \text{ (wave number)} = \frac{\nu}{c} = \frac{1}{\lambda} \text{ (cm}^{-1}\text{)}$$

Photon energy (E):

$$E = h\nu = \frac{hc}{\lambda} = \frac{1.24}{\lambda \text{ (}\mu\text{m)}} \text{ eV,}$$

where the velocity of light in free space is

$$c = 3 \times 10^{10} \text{ cm/s}$$

and Planck's constant *h* is

$$h = 6.63 \times 10^{-34} \text{ J}\cdot\text{s}$$

The optical spectrum is divided into three major categories, as follows.

Ultraviolet: Wavelengths shorter than those in the visible spectrum and longer than for X rays are collectively designated ultraviolet (UV). Ultraviolet is classified according to wavelength as extreme (100–2000 Å), far (2000–3000 Å), or near (3000–3700 Å). Ultraviolet is also sometimes designated as short-wave or long-wave.

Visible: Those wavelengths in the approximate range 3700–7500 Å can be perceived by the human eye and are therefore collectively designated as visible light. Visible light is classified according to the various colors its wavelengths elicit in the mind of a standard observer. The major color categories are violet (3700–4550 Å), blue (4560–4920 Å), green (4930–5770 Å), yellow (5780–5970 Å), orange (5980–6220 Å), and red (6230–7500 Å).

Infrared: Those wavelengths longer than those in the visible spectrum and shorter than microwaves are collectively designated infrared (IR). Infrared is classified according to its wavelength as near (0.75–3 μm), middle (3–6 μm), far (6–15 μm), and extreme or submillimeter (15 μm–1 mm), although different authors often use slightly different wavelength ranges for these classifications.

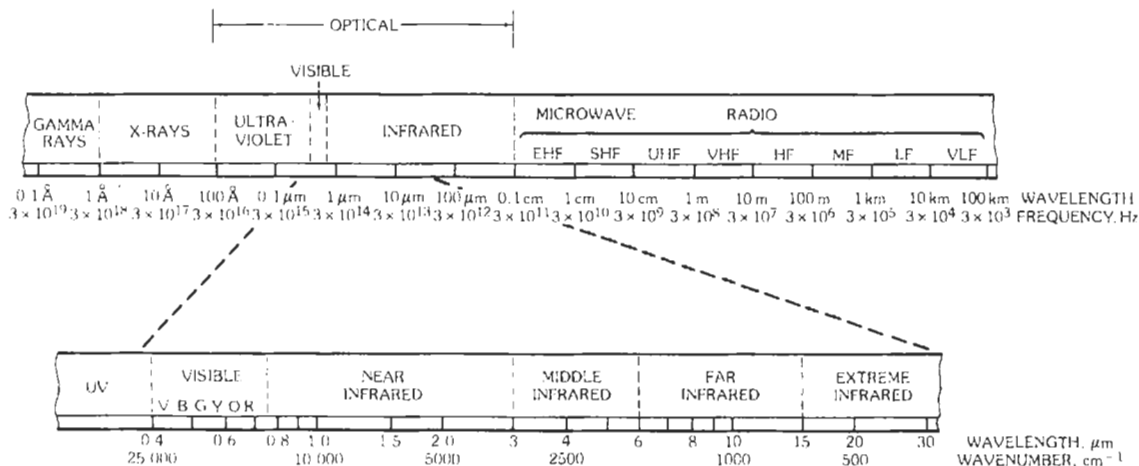


Fig. 1. The electromagnetic spectrum.

RADIOMETRY

Terms and Definitions

Radiometry is the science of measurement of optical radiation at any wavelength, based simply on physical measurements. Radiant energy cannot be measured quantitatively directly, but must always be converted into some other form such as thermal, electrical, or chemical. Radiometry applies over the entire electromagnetic spectrum, not just the optical region. The terms used to describe radiant power in radiometry are summarized in Table 1. All the terms used in radiometry are defined in terms of energy. When certain quantities are considered as a function of wavelength or frequency, the adjective "spectral" is used to modify the term, and the symbol for the quantity is followed by the appropriate spectral symbol, λ , ν , or $\bar{\nu}$, in parentheses. When the spectral concentration of a quantity is considered, it is also described by the adjective spectral, but in this case the symbol is subscripted with the appropriate spectral quantity; for example the spectral radiant exitance can be given as W_λ in units of $W \cdot m^{-2} \cdot \mu m^{-1}$.

Blackbody Radiation

Any surface of a body with a temperature greater than absolute zero ($T = 0$ K) is a source of radiation. The term "blackbody" applies to a thermal radiator that absorbs completely all incident electromagnetic radiation regardless of the wavelength, the direction of incidence, or the polarization. Such a radiator also has the maximum emission possible for any wavelength and in any direction for a thermal radiator in thermal equilibrium at a given temperature. Kirchhoff's law states that for any body (all materials) in an isothermal enclosure at temperature T , the ratio of the radiant exitance (emittance), W , to the absorptance, α , is equal to the radiant exitance of a blackbody, W_{bb} , at the same temperature:

$$W(T)/\alpha = W_{bb}(T)$$

Planck's law gives the blackbody spectral power for unpolarized radiation emitted at temperature T between the wavelengths λ and $\lambda + d\lambda$ as

$$W_\lambda(\lambda, T)d\lambda = 2\pi c^2 h \lambda^{-5} (e^{hc/\lambda kT} - 1)^{-1} d\lambda$$

$$= C_1 d\lambda / \lambda^5 (e^{C_2/\lambda T} - 1)$$

where,

$k = 1.38 \times 10^{-23}$ J·K⁻¹ is Boltzmann's constant,

$$C_1 = 2\pi c^2 h = 3.74 \times 10^4 \text{ W} \cdot \text{cm}^{-2} \cdot \mu\text{m}^4,$$

$$C_2 = hc/k = 1.44 \times 10^4 \mu\text{m} \cdot \text{K}.$$

Fig. 2 shows the variation of $W_\lambda(\lambda, T)$ with λ for various values of T . The integral of $W_\lambda(\lambda, T)d\lambda$ over all

wavelengths gives the Stefan-Boltzmann law for total blackbody radiant exitance $W_{bb}(T)$ in W/cm^2 ,

$$W_{bb}(T) = \int_0^\infty W_\lambda(\lambda, T)d\lambda = (2\pi^5 k^4 / 15c^2 h^3) T^4$$

$$= (\pi^4 C_1 / 15 C_2^4) T^4 = \sigma T^4$$

where $\sigma = 5.67 \times 10^{-12}$ W·cm⁻²·K⁻⁴ is the Stefan-Boltzmann constant. The blackbody spectrum has a pronounced maximum at a particular wavelength λ_m for a given temperature T , and the relationship

$$\lambda_m T = C_2 / 4.9651 = 2897.8 \mu\text{m} \cdot \text{K}$$

is known as Wien's displacement law.

Planck's law can also be specified in terms of frequency or wave number instead of wavelength, and these relations result in slightly different forms of Wien's displacement law. For calculations it is often convenient to use the dimensionless quantity

$$x = hc/\lambda kT = h\nu/kT = C_2/\lambda T$$

for which Planck's law becomes

$$W_\lambda(x, T)dx = (15\sigma T^4 / \pi^4) [x^3 dx / (e^x - 1)]$$

The energy flux per unit area for wavelengths lying between x_1 and x_2 can easily be calculated using the series expansions below:

$$W_{0-\lambda_0}(T) = \int_0^{\lambda_0} W_\lambda(\lambda, T)d\lambda$$

$$= (15\sigma T^4 / \pi^4) \int_{x_0}^\infty x^3 dx / (e^x - 1)$$

$$= (15\sigma T^4 / \pi^4) \sum_{m=1}^\infty (e^{-m x_0} / m^4)$$

$$[(m x_0)^3 + 3(m x_0)^2 + 6(m x_0) + 6]$$

where

$$x_0 = hc/\lambda_0 kT = C_2/\lambda_0 T$$

and the energy flux per unit area for wavelengths between λ_2 and λ_1 is given by

$$W_{\Delta\lambda} = (15\sigma T^4 / \pi^4) \sum_{n=1}^\infty (e^{-n x} / n^4) [(n x)^3 + 3(n x)^2 + 6n x + 6] \Big|_{x_2}^{x_1}$$

for $x_i \geq 2$.

TABLE 1. COMMON RADIOMETRIC TERMINOLOGY

Symbol	Term	Description	Units	Equation
U	Radiant energy	Energy transferred by electromagnetic waves	J	
u	Radiant energy density	Radiant energy per unit volume	$J \cdot \text{cm}^{-3}$	$u = dU/dV$
P	Radiant flux	Rate of transfer of radiant energy	W	$P = dU/dt$
W	Radiant emittance	Radiant flux emitted per unit area of a source	$W \cdot \text{cm}^{-2}$	$W = dP/dA$
Q	Radiant photon emittance	Number of photons emitted per second per unit area	$\text{Photons} \cdot \text{s}^{-1} \cdot \text{cm}^{-2}$	
J	Radiant intensity	Radiant flux per unit solid angle	$W \cdot \text{sr}^{-1}$	$J = dP/dQ$
N	Radiance	Radiant flux per unit solid angle per unit area	$W \cdot \text{cm}^{-2} \cdot \text{sr}^{-1}$	$N = W/\pi$
H	Irradiance	Radiant flux incident per unit area	$W \cdot \text{cm}^{-2}$	$H = dP/dA$
P_λ	Spectral radiant flux	Radiant flux per unit wavelength interval at a particular wavelength	$W \cdot \mu\text{m}^{-1}$	$P_\lambda = dP/d\lambda$
W_λ	Spectral radiant emittance	Radiant emittance per unit wavelength interval at a particular wavelength	$W \cdot \text{cm}^{-2} \cdot \mu\text{m}^{-1}$	$W_\lambda = dW/d\lambda$
Q_λ	Spectral radiant photon emittance	Radiant photon emittance per unit wavelength interval at a particular wavelength	$\text{Photons} \cdot \text{s}^{-1} \cdot \text{cm}^{-2} \cdot \mu\text{m}^{-1}$	$Q_\lambda = dQ/d\lambda$
J_λ	Spectral radiant intensity	Radiant intensity per unit wavelength interval at a particular wavelength	$W \cdot \text{sr}^{-1} \cdot \mu\text{m}^{-1}$	$J_\lambda = dJ/d\lambda$
N_λ	Spectral radiance	Radiance per unit wavelength interval at a particular wavelength	$W \cdot \text{cm}^{-2} \cdot \text{sr}^{-1} \cdot \mu\text{m}^{-1}$	$N_\lambda = dN/d\lambda$
H_λ	Spectral irradiance	Irradiance per unit wavelength interval at a particular wavelength	$W \cdot \text{cm}^{-2} \cdot \mu\text{m}^{-1}$	$H_\lambda = dH/d\lambda$
ϵ	(Radiant) emissivity	Ratio of radiant emittance of a source to that of a blackbody at the same temperature	(Numeric)	
α	(Radiant) absorptance	Ratio of absorbed radiant flux to incident radiant flux	(Numeric)	
θ	(Radiant) reflectance	Ratio of reflected radiant flux to incident radiant flux	(Numeric)	
τ	(Radiant) transmittance	Ratio of transmitted radiant flux to incident radiant flux	(Numeric)	

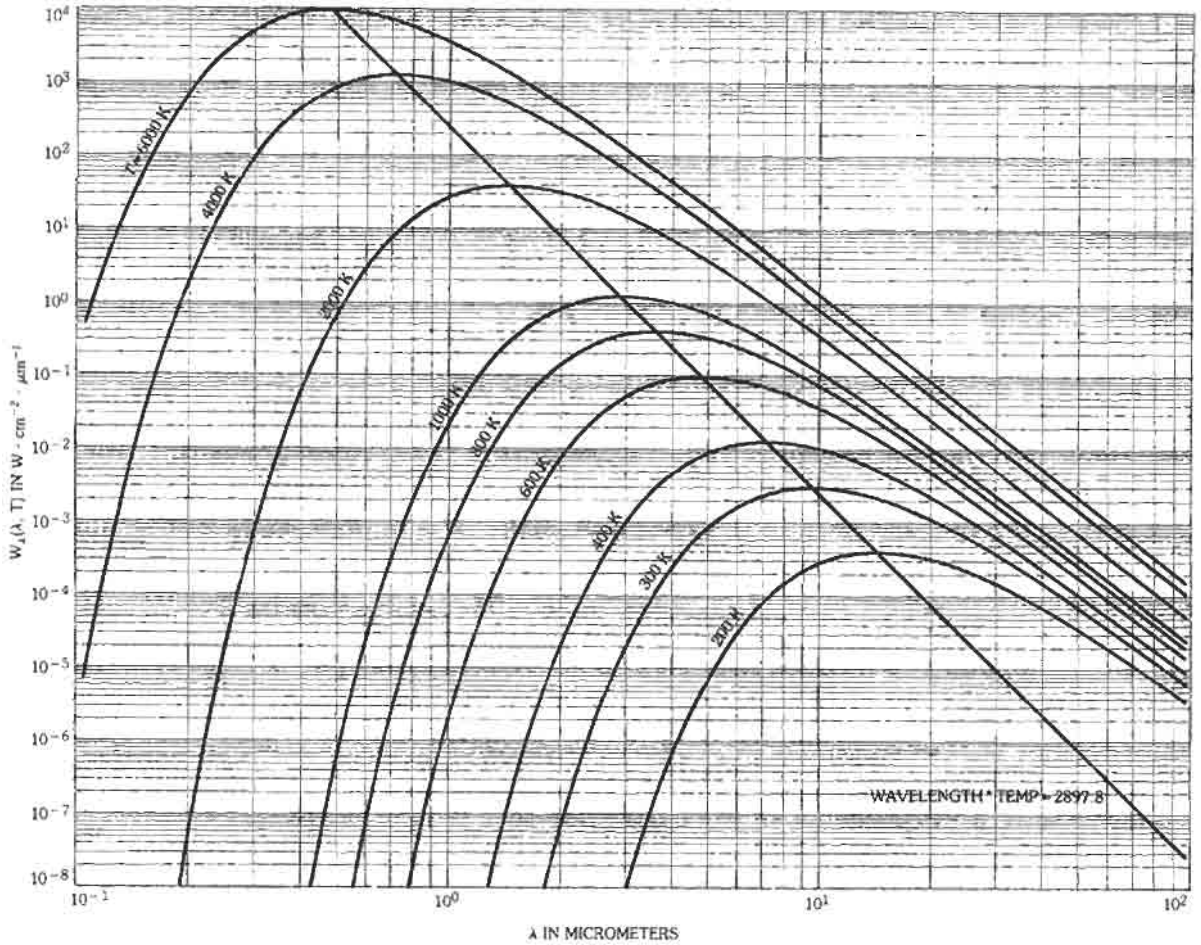


Fig. 2. Blackbody spectral radiant emittance, $W_\lambda(\lambda, T)$, for different blackbody temperatures.

A similar expression can be derived for values of $x_1 < 2$ using

$$W_{0-\lambda_0}(T) = \sigma T^4 \{1 - (15/\pi^4)x_0^3 [1/3 - x_0/8 + x_0^2/60 - x_0^4/5040 + x_0^6/272160 - x_0^8/13305600 + \dots]\}$$

which converges rapidly for $x_0 < 2$.

Similar relations for the number of quanta or photons emitted by a body at temperature T can be obtained by dividing $W_\lambda(\lambda, T)$ by the energy of the quanta at wavelength λ , hc/λ , as

$$Q_\lambda(\lambda, T)d\lambda = (C_3/\lambda^4)[d\lambda/(e^{c_2/\lambda T} - 1)]$$

where $C_3 = C_1/hc = 2\pi c$. The variation of $Q_\lambda(\lambda, T)$ with λ for several temperatures is shown in Fig. 3. The relation corresponding to Wien's displacement law

relating the wavelength of maximum quanta emission λ'_m to the temperature is

$$\lambda'_m T = 3669.73 \mu\text{m}\cdot\text{K}$$

for λ'_m in micrometers, and the Stefan-Boltzmann law for the total number of photons emitted per unit area in blackbody radiation is

$$Q_{bb}(T) = \sigma T^4/2.75kT = \sigma' T^3$$

where $\sigma' = 1.49 \times 10^{11} \text{ s}^{-1}\cdot\text{cm}^{-2}\cdot\text{K}^{-3}$. These relations are summarized in Table 2.

Since the blackbody radiant emission and radiant photon emission are both functions of λT , it is possible to plot a single normalized curve of blackbody radiant emission as a function of λT that makes possible a simple determination of the fraction of the peak emission at any value of λT , as well as the fraction of the total radiant emittance below λT (i.e., for values of λ_0

TABLE 2. BASIC LAWS AND CONSTANTS OF BLACKBODY RADIATION

Stefan-Boltzmann Law	Wien's Displacement Law		Planck's Equation
$W_{bb}(T) = \sigma T^4$	$\lambda_m T = b$	$W_{\lambda_m} = aT^5$	$W_\lambda(\lambda, T) = C_1 \lambda^{-5} (e^{C_2/\lambda T} - 1)^{-1}$
$Q_{bb}(T) = \sigma' T^3$	$\lambda_m' T = b'$	$Q_{\lambda_m} = a' T^4$	$Q_\lambda(\lambda, T) = C_3 \lambda^{-4} (e^{C_2/\lambda T} - 1)^{-1}$
$W_{bb}(T)$ = Total blackbody radiant emittance, $W \cdot cm^{-2}$, for temperature T			$W_\lambda(\lambda, T)$ = Blackbody spectral radiant emittance, $W \cdot cm^{-2} \cdot \mu m^{-1}$
$Q_{bb}(T)$ = Total blackbody photon emittance, photons $cm^{-2} \cdot s^{-1}$			W_{λ_m} = Maximum spectral radiant emittance, $W \cdot cm^{-2} \cdot \mu m^{-1}$
σ = Stefan-Boltzmann constant $= (5.67) \times 10^{-12} W \cdot cm^{-2} \cdot K^{-4}$			$Q_\lambda(\lambda, T)$ = Blackbody spectral photon emittance, photons $\cdot s^{-1} \cdot cm^{-2} \cdot \mu m^{-1}$
σ' = Photon Stefan-Boltzmann constant $= 1.49 \times 10^{11} s^{-1} \cdot cm^{-2} \cdot K^{-3}$			Q_{λ_m} = Maximum spectral radiant photon emittance, photon $\cdot s^{-1} \cdot cm^{-2} \cdot \mu m^{-1}$
λ_m = Wavelength of maximum radiant emittance			$c = 3 \times 10^{10} cm \cdot s^{-1}$
λ_m' = Wavelength of maximum radiant photon emittance			$C_1 = 2\pi c^2 h = 3.74 \times 10^4 W \cdot cm^{-2} \cdot \mu m^4$
λ = Radiant wavelength			$C_2 = ch/k = 1.44 \times 10^4 \cdot \mu m \cdot K$
$a = 1.2862 \times 10^{-15} W \cdot cm^{-2} \cdot \mu m^{-1} \cdot K^{-5}$			
$a' = 2.10098 \times 10^7 s^{-1} \cdot cm^{-2} \cdot \mu m^{-1} \cdot K^{-4}$			$C_3 = C_1/hc = 2\pi c$
$b = (2897.8 \pm 0.4) \mu m \cdot K$			
$b' = (3669.7 \pm 0.4) \mu m \cdot K$			
$h = (6.63) \times 10^{-34} W \cdot s$ Planck's constant			$k = 1.38 \times 10^{-23} J \cdot K^{-1}$ Boltzmann's constant

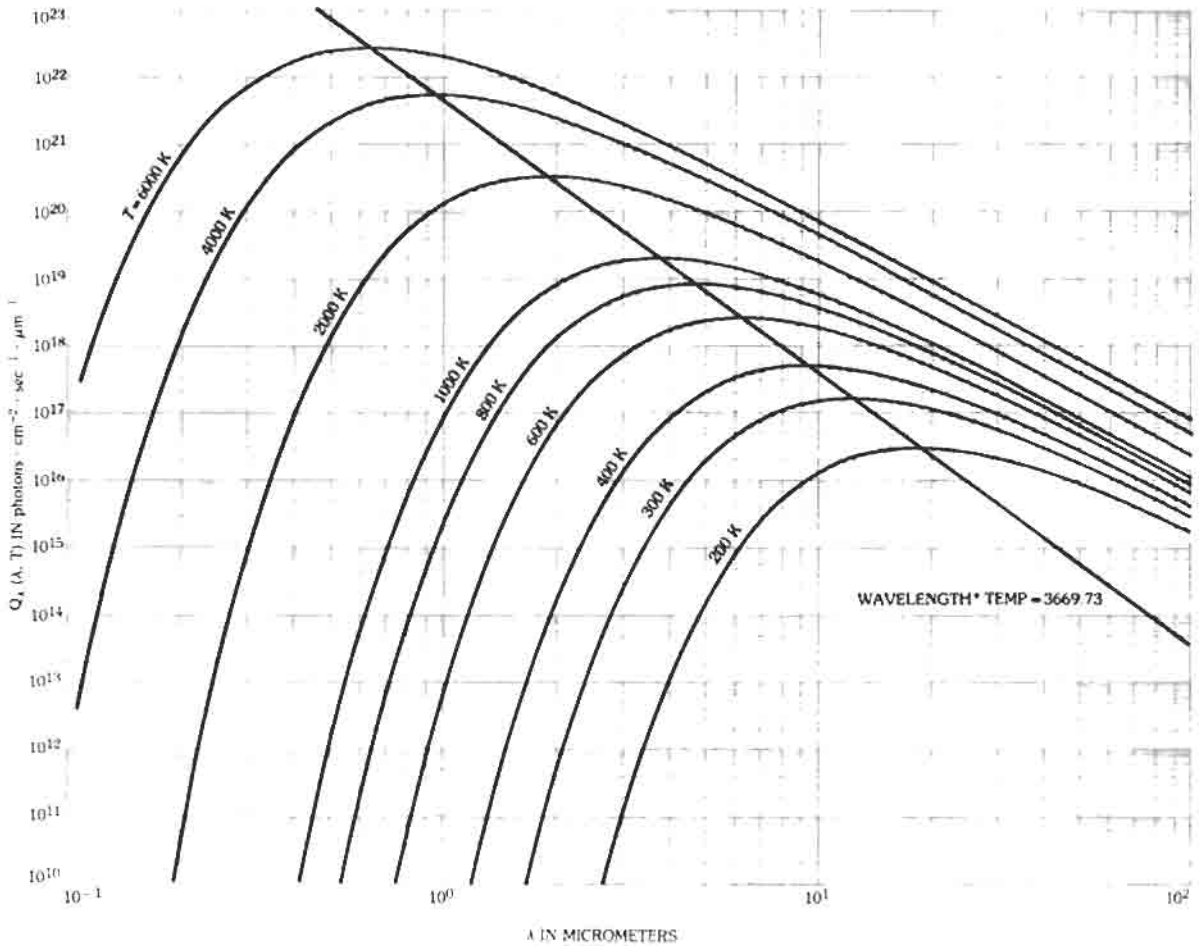


Fig. 3. Blackbody spectral radiant photon emittance, $Q_{\lambda}(\lambda, T)$, versus λ for different blackbody temperatures.

shorter than λ in λT). This type of plot is particularly useful in determining the effect of background radiation on the sensitivity of infrared detectors. A plot of this normalized photon emittance is given in Fig. 4. To use this curve, the product of the wavelength, λ , and the blackbody or background temperature in kelvins, T , is determined, and the photon emittance, in $\text{cm}^{-2}\cdot\text{sec}^{-1}\cdot\mu\text{m}^{-1}$ at the peak, given by

$$Q\lambda m = a'T^4$$

and the total photon emittance, in $\text{photons}\cdot\text{cm}^{-2}\cdot\text{sec}^{-1}$, for the blackbody or background at temperature T is determined from

$$Q(T) = b'T^5$$

as given in Table 2. Then, from values read from Fig. 4, the actual spectral emittance at wavelength λ for λT and

the total photon emittance ($\text{photons}\cdot\text{cm}^{-2}\cdot\text{sec}^{-1}$) for wavelengths shorter than λ in λT can be directly calculated. The latter gives the total photon flux capable of exciting free carriers in photon detectors with a cutoff wavelength of $\lambda_c = \lambda$.

A more direct indication of the total photon emittance from the threshold wavelength to zero wavelength, i.e., the total photon emittance with energies greater than that determined by the wavelength λ_c , is given in Fig. 5 for blackbody temperature ranging from 250 K to 1500 K.

The results given above apply only to a blackbody. The emissivity of actual bodies can be accounted for by multiplying $W_{\lambda}d\lambda$, for example, by $\epsilon(\lambda, T)$, the variation of the emissivity of a particular material with a particular surface preparation or condition with wavelength and temperature. If $\epsilon(\lambda, T) = \epsilon(T) = \text{constant} < 1$, the body is referred to as a "graybody." The relationship between the emissivity and spectral radiant emittance is shown in Fig. 6.

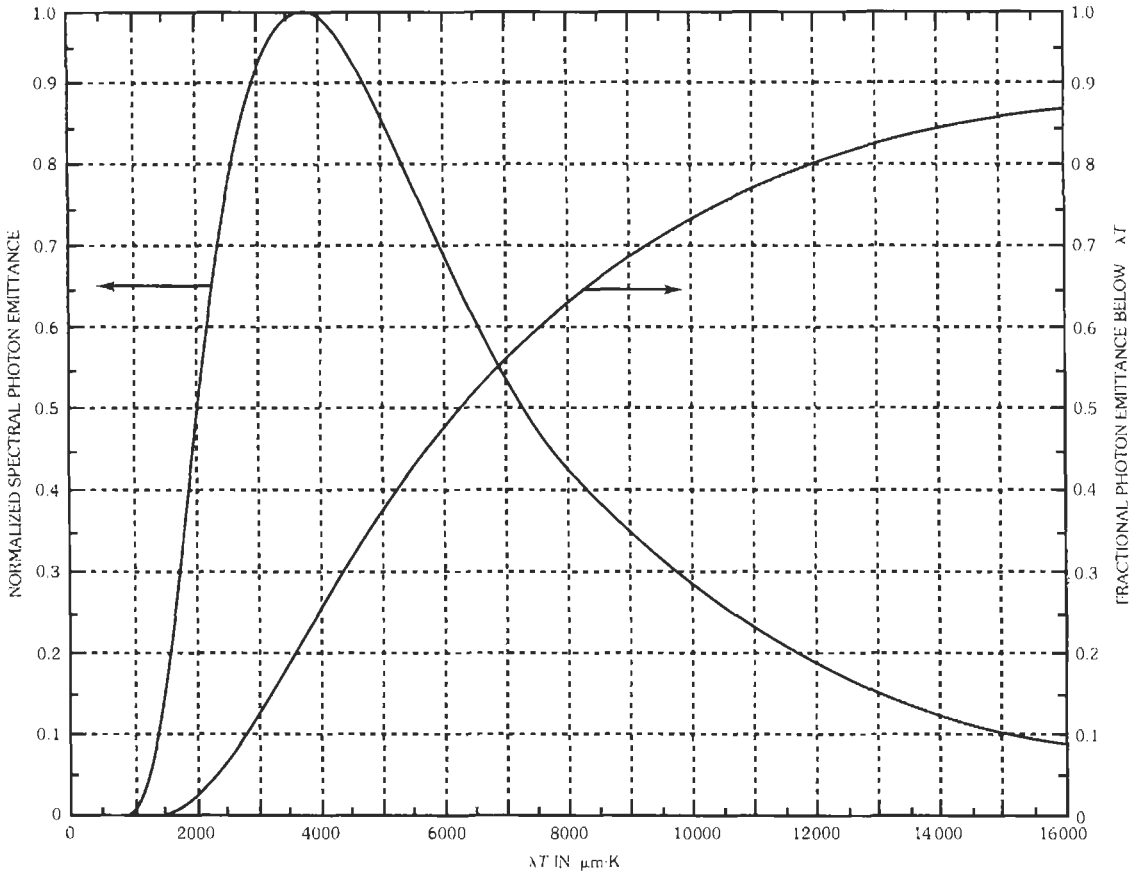


Fig. 4. Normalized photon emittance.

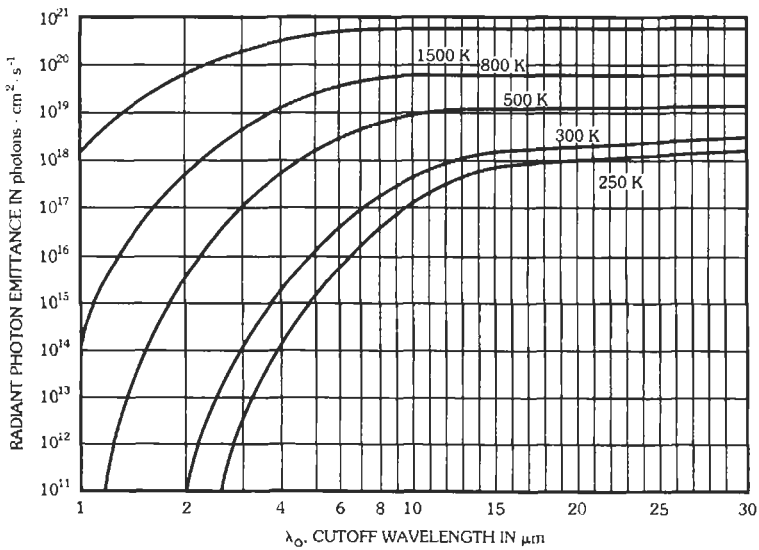


Fig. 5. Total radiant photon emittance for wavelengths shorter than λ_0 for different blackbody temperatures. (From Santa Barbara Research Center Wall Chart, copyright 1991.)

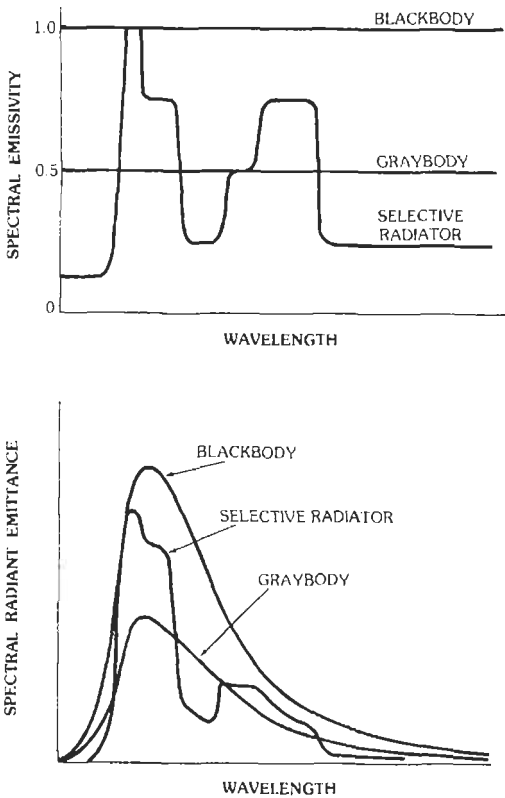


Fig. 6. Spectral emissivity and spectral radiant emittance of three types of radiators. (From Hudson, R. D. Jr., Infrared System Engineering, New York: John Wiley & Sons, Inc., 1969.)

INTERACTION OF OPTICAL WAVES WITH MATTER

An optical wave may interact with matter by being reflected, refracted, absorbed, or transmitted. The interaction normally involves two or more of these effects.

Reflectance

Some of the optical radiation impinging upon a surface is usually reflected from the surface. Reflectance varies according to the properties of the surface and the wavelength, and in real circumstances may range from more than 98% (smoked MgO at visible wavelengths) to less than 1% (lampblack at visible wavelengths). Reflection from a surface may be either diffuse, specular, or both. A diffuse reflector has a surface that is rough when compared to the wavelength of the impinging radiation. Lambert's law specifies a perfectly diffuse surface as one having a constant radiance independent of the viewing angle according to

$$N = W/\pi$$

where,

N is radiance ($W/m^2/sr$),
 W is the radiant excitation (W/m^2).

In practice, the reflectance of real reflectors varies with the cosine of the viewing angle. A specular reflector has a surface that is smooth when compared to the wavelength of the impinging radiation. A perfect specular reflector will reflect an oncoming beam without altering the divergence of the beam. A narrow beam of optical radiation impinging upon a specular reflector obeys two rules:

1. The angle of reflection is equal to the angle of incidence.
2. The incident ray and the reflected ray lie in the same plane as a normal line extending perpendicularly from the surface.

Fig. 7 illustrates both diffuse and specular reflection.

Absorptance

Some of the optical radiation impinging upon any substance is absorbed by the substance. Absorptance varies according to the properties of the substance and the wavelength, and in real circumstances may range from a low of less than 2 dB/km for certain ultrapure fused silica glasses to more than 98% for lampblack.

Transmittance

Some of the optical radiation impinging upon a substance is transmitted into the substance. The penetration depth may be slight, in which case the transmittance is zero. Certain ultrapure silica glasses may have a transmittance of greater than 75% at certain wavelengths over a distance of 1 km. The reflectance, absorptance, and transmittance of a substance must satisfy

$$\rho + \alpha + \tau = 1$$

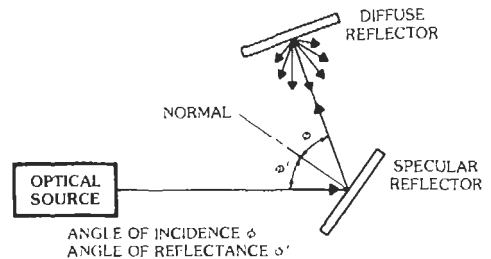


Fig. 7. Specular and diffuse (Lambertian) reflectance. (From Mims, F. M. III, Optoelectronics, Indianapolis: Howard W. Sams & Co., Inc., 1975, p. 19.)

where,

ρ is reflectance,
 α is absorptance,
 τ is transmittance.

Refraction

A ray of optical radiation passing from one medium to another is bent at the interface of the two mediums if the angle of incidence is unequal to 90° . The index of refraction for a substance is the sine of the angle of incidence divided by the sine of the angle of refraction. Refractive index varies with wavelength and ranges from 1.002914 (air at 656 nm) to 2.7 (crystalline titanium oxide).

OPTICAL SOURCES

Important optoelectronic sources are tungsten lamps, fluorescent lamps, glow-discharge lamps, electroluminescent diodes, and lasers.

Tungsten Lamps*

Fig. 8 shows the spectral output of a tungsten lamp at a temperature of 2800 K. A typical tungsten lamp is an efficient optical source, but only about 5% of its radiant flux falls within the visible wavelengths. The addition of a halogen gas to a tungsten lamp can increase the life and efficiency of the lamp. At high operating temperature, the halogen (iodine or bromine) vaporizes and combines chemically with the evaporated tungsten on the quartz envelope and in the gas. The resulting tungsten-halogen gas migrates back to the filament, where the very high temperatures decompose it so that tungsten is redeposited on the filament and the halogen repeats the cycle. This permits uniform and constant output throughout the lamp life.

Fluorescent Lamps†

A typical fluorescent lamp is a sealed glass tube filled with argon gas and containing a small amount of mercury. When an electrical discharge is established in the tube, ultraviolet radiation is produced, which causes a phosphor coating on the inside wall of the tube to fluoresce with a bright white glow. Fluorescent lamps provide high efficiency.

* Wolfe, W. L., and Zessis, G. J., eds. *The Infrared Handbook*. Ann Arbor, Mich.: Environmental Research Institute of Michigan, 1978.

† Grum, F., and Becherer, Richard J. *Optical Radiation Measurements*, Vol. 1, *Radiometry*. New York: Academic Press, Inc., 1979, p. 142.

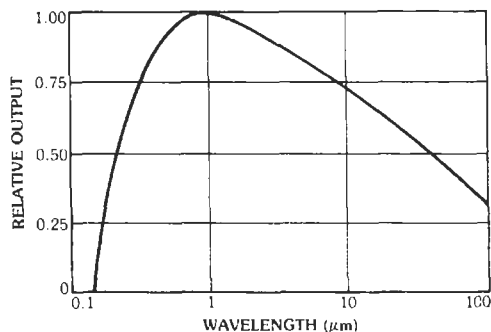


Fig. 8. Spectral output of a tungsten lamp at 2800 K.

Arc Lamps§

Arc lamps operate at considerably higher temperatures than other sources and, with the exception of certain lasers, are the most brilliant artificial sources. A representative arc lamp consists of a heavy-walled quartz envelope filled with gas (mercury, xenon, krypton, etc.) at a typical pressure of 20–40 atmospheres at operating temperature. The emission from a typical high-pressure arc lamp originates from the gaseous discharge, from the incandescence of the electrodes, and from the envelope.

Light-Emitting Diodes§

The operation of light-emitting diodes (LEDs) is dependent on the radiative recombination of electron-hole pairs at a forward-biased pn junction. The recombination of holes and electrons that occurs when current is injected into a pn junction results in the release of energy. This energy, which corresponds to the band gap, may be in the form of a photon (radiative recombination), a series of photons (lattice vibrations, non-radiative), or both. In some cases, the energy may be transferred to another electron. Light-emitting diodes are electroluminescent diodes in which radiative recombination is significant. Recombination radiation emitted by LEDs is peaked at or near the band-gap energy, and in practical devices ranges from $0.55 \mu\text{m}$ to $34 \mu\text{m}$. The spectral emission width of a representative device may be 25–30 nm, and this provides sufficient monochromaticity for the production of discrete wavelength bands (including the visible colors green, yellow, orange, and red).

Direct band-gap semiconductors make more efficient LEDs than indirect band-gap semiconductors (those in which electron-hole recombination is accompanied by emission of a photon and phonons). Silicon and germa-

§ Saleh, B. E. A., and Teich, M. C. *Fundamentals of Photonics*. New York: John Wiley and Sons (Wiley Interscience), 1991.

nium, both indirect band-gap materials, produce far too little recombination radiation for practical use, but other direct and indirect band-gap materials or direct-indirect band-gap alloys, particularly GaAs, GaP, and solid solution alloys of GaAs and GaP, GaAs and AlAs, and InAs and GaP, can be tailored to produce relatively efficient radiative recombination.

Gallium arsenide (GaAs) LEDs and similar direct band-gap LEDs in other materials exhibit internal quantum efficiencies (ratio of emitted photons per injected electrons) very nearly unity. Due to such factors as internal absorption, contact shadows, and refractive-index-induced surface reflectance, the external quantum efficiency of practical LEDs is usually much lower (less than about 0.1). Both quantum efficiency and recombination radiation wavelength are affected by temperature, and the external quantum efficiency of specially fabricated LEDs may exceed 0.4 at 20 K. Wavelength is directly related to junction temperature and typically varies about 0.25 nm per degree Celsius. Therefore, changes in either the ambient temperature or junction heating will alter the emitted wavelength.

A wide range of economical long-lived plastic- and metal-packaged LEDs is commercially available. It is common practice to install the diode chip inside a miniature directional reflector and encapsulate the entire assembly in an index-matching epoxy to enhance the radiation extraction efficiency. This fabrication technique improves the external efficiency by a factor of 2 or more. Visible emitters suitable for use as visual indicators and other display roles include: GaP and GaP:N (green—550 nm); GaAs_{0.25}P_{0.75}:N (yellow—610 nm); GaAs_{0.6}P_{0.4} (red—660 nm); Al_{0.3}Ga_{0.7}As (red—675 nm); and GaP:Zn, O (red—690 nm). Several of these materials, particularly GaAs_{0.6}P_{0.4}, are in widespread use as display devices for digital clocks, watches, calculators, and electronic instruments. The most common infrared-emitting materials have been GaAs (905 nm) and GaAs:Si (940 nm). Until recently, these materials offered much higher efficiencies than visible emitters. They are useful in such applications as optical communications and ranging, position sensing, object detection, and electro-optical isolation. Ordinary GaAs LEDs are characterized by a reasonable power output (about 1.5 mW @ 100 mA I_F) and fast turn-on time (about 1 ns), whereas special GaAs diodes doped with Si are more efficient (about 10 mW @ 100 mA I_F) but slower (t_{on} of about 300 ns). Since the radiant emission from electroluminescent diodes is generally linear with respect to applied current below the saturation region, they can be easily and directly pulse or analog modulated by simply controlling the forward bias. Electroluminescent diodes of GaAs, GaAs:Si, GaAs-Al_xGa_{1-x}As, and related materials emitting in the near infrared (approximately 750–950 nm) have found wide use in fiber-optical communication applications, but because of lower fiber losses in the 1.3- to 1.6- μ m wavelength range, InGaAsP and InGaAs LEDs

are being rapidly developed for use at these longer wavelengths.

Electroluminescent diodes normally emit radiation in a relatively broad pattern and therefore cannot be coupled to optical waveguide fibers as efficiently as the laser sources described below. Nevertheless, reasonable coupling efficiency can be obtained by utilizing epoxy lenses or special structures such as the Burrus diode structure shown in Fig. 9. This structure has a typical emission area of 2×10^{-5} cm² and a radiance of 100 W/sr/cm² at a forward bias of 150 mA. The output wavelength can be peaked anywhere between 750 and 905 nm by adjusting the aluminum concentration in the Al_xGa_{1-x}As alloy making up the diode. About 2 mW can be coupled into a suitable optical waveguide fiber when this structure is used. This and most other electroluminescent diodes exhibit operating lifetimes of more than 10⁴ hours at the relatively high current densities required to obtain significant optical output.

High-Intensity Visible LEDs

A typical package commonly used for LEDs is shown in Fig. 10. This device is a Hewlett Packard Optoelectronic Division high-power LED fabricated from AlInGaP and packaged with an epoxy lens. This material has very high luminous efficiency, and AlInGaP devices are capable of producing high light output over a wide range of drive currents. The package shown incorporates an advanced optical design that produces an extremely high peak intensity in a viewing angle as narrow as three degrees. The emission of these AlInGaP LEDs occurs in a narrow spectral range and can be

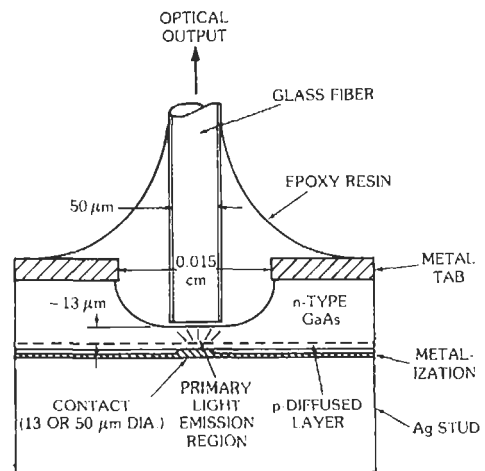


Fig. 9. GaAs electroluminescent diode specifically designed for direct coupling to an optical-fiber waveguide. (After Burrus, C. A., and Dawson, R. W., "Small-Area High-Current-Density GaAs Electroluminescent Diodes and a Method of Operation for Improved Degradation Characteristics," *Applied Physics Letters*, 1 August 1970, pp. 97–99.)

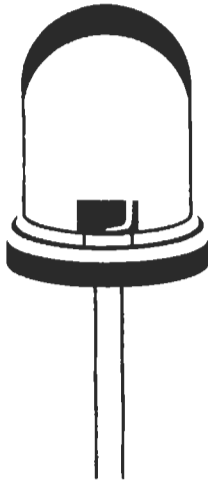


Fig. 10. Epoxy package for an AlInGaP diode with a high pack intensity in the visible (590 nm) wavelength region and a narrow viewing angle of three degrees.

tuned to wavelengths in the range 590 nm to 620 nm by adjusting the AlInGaP alloy composition during growth as shown in Fig. 11. These lamps are suitable for applications in optical communications, as emitters in optical couplers, as LEDs in high ambient lighting conditions, as replacements for low-power laser diodes, and even as alternatives to incandescent lamps.

Superluminescent Diodes

Specially fabricated GaAs-Al_xGa_{1-x}As double heterostructure and other electroluminescent diodes can be made to emit both spontaneous (quantum noise) and stimulated (amplified quantum noise) radiation. As in the case of a laser, stimulated emission implies a narrower spectral emission width and higher radiance than spontaneous emission. Fig. 12 shows the construction of a typical superluminescent diode (SLD). The

structure of the device in the figure is virtually identical to that of a stripe-geometry injection laser with the major exception being an incomplete upper electrode. The incomplete electrode eliminates current injection near one of the mirrors, and this effectively isolates the mirror from the optical wave propagating along the plane of the junction and suppresses laser action. A similar effect can be had by angling one mirror of an injection laser a few degrees away from the normal or by coating one of the end mirrors of an injection laser with an antireflective film. The SLD is less efficient than most other electroluminescent sources, but its narrow spectral width and high radiance make it well suited for optical-waveguide links. Pulsed radiation of 50 mW has been coupled into an optical fiber from an SLD source.

Semiconductor Lasers*

Several fundamental modifications of the basic p-n junction electroluminescent diode exist, and chief among these is the semiconductor injection laser. In its simplest form, the injection laser is a direct band-gap LED having an exceptionally flat and uniform junction (the active region) bounded on facing sides by two parallel mirrors perpendicular to the plane of the junction, which provide a Fabry-Perot resonant cavity. The mirrors are usually produced by cleaving the semiconductor chip along parallel planes to produce perfectly parallel and flat surfaces. The remaining two sides of the chip perpendicular to the junction plane are intentionally roughened during the sawing process, which separates bars of material into individual chips. This surface roughening suppresses off-axis lasing

* Saleh, B. E. A., and Teich, M. C., loc. cit.; Kressel, H., and Butler, J. K., *Semiconductor Lasers and Heterojunction LEDs* (New York: Academic Press, Inc., 1977); Casey, H. C., Jr. and Panish, M. B., *Heterostructure Lasers* (New York: Academic Press, Inc., 1978); Thompson, G. H. B., *Physics of Semiconductor Laser Devices* (Chichester: John Wiley & Sons, Inc., 1980).

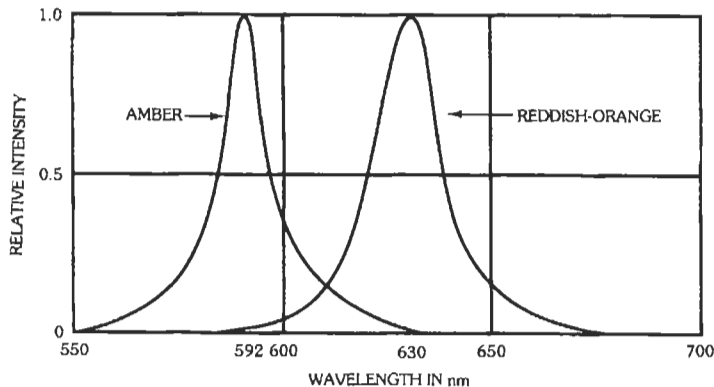


Fig. 11. Relative intensity versus wavelength for AlInGaP LEDs. (Courtesy Hewlett Packard Optoelectronic Devices.)

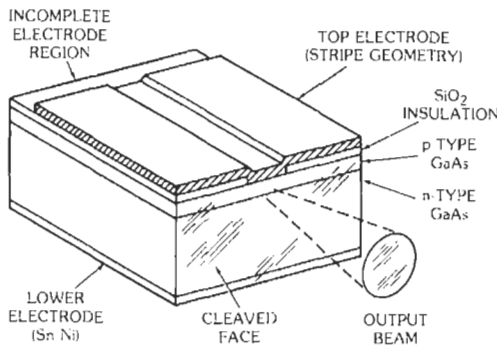


Fig. 12. Stripe-geometry semiconductor injection laser modified with an incomplete upper contact electrode to achieve superluminescent emission. (After Lee, T, Burrus, C. A. Jr., and Miller, B. I., "A Stripe-Geometry Double-Heterostructure Amplified-Spontaneous-Emission [Superluminescent] Diode," IEEE Journal of Quantum Electronics, August 1973, pp. 820-828.)

modes. The high index of refraction of the semiconductor provides sufficient reflectance at the semiconductor-air interface for the optical feedback necessary for laser action.

Below threshold, the injection laser behaves like a conventional LED, but as the current injection is increased, a threshold point occurs where the hole-electron population in the active region becomes inverted. Spontaneous recombination of holes and electrons then produces photons that stimulate in-phase recombination and photon emission by other holes and electrons, and lasing occurs as the optical gain in the active region overcomes absorption and other losses. The mirrors on either end of the active region provide the optical feedback necessary to sustain laser action, and a small fraction of the wave propagating between the mirrors emerges from each on each pass. One end facet on many commercial lasers is overcoated with a reflective gold film to cause all the radiation to emerge from only one end of the device and thus enhance collection efficiency.

The most common and best developed injection lasers utilize GaAs and AlGaAs (820-905 nm) and InGaAs and InP (1.3-1.55 μm) materials, though many other semiconductors have been used to produce wavelengths ranging from as short as 630 nm ($\text{Al}_x\text{Ga}_{1-x}\text{As}$) to as long as 34 μm (PbSnSe). The high current density required to achieve lasing in broad-area homojunction injection lasers (8000 A/cm^2 or more) precludes continuous operation above temperatures greater than about 77 K. Room-temperature operation of a homojunction laser requires short current pulses (10-100 amperes for typical diodes) no more than 200 ns wide and with a duty cycle of 0.1%. Peak pulse power outputs at 300 K of up to 100 watts from single diodes and several kilowatts from diode arrays can be obtained from commercial devices. Injection lasers whose GaAs active region is sandwiched between two

$\text{Al}_x\text{Ga}_{1-x}\text{As}$ layers to produce a double heterostructure which confines both the recombination and the optical wave to a very thin layer have been fabricated with threshold current densities of less than 1000 A/cm^2 . These devices are now commercially available and can be operated continuously at 300 K with appropriate heatsinking. Individual diodes may emit up to a few tens of milliwatts.

Injection-laser recombination radiation is characterized by a narrow spectral bandwidth, beam directionality, and moderate beam divergence. The beam directionality permits collection of up to 90% of the radiation by simple optics in most cases. The beam divergence, which may be 10-15° in the plane of the junction (half-power measurement) and 25-45° perpendicular to the plane of the junction, is a result of the very small emission region of the laser. Essentially a diffraction limited slit, the emission region and internal lasing modes both contribute to the various diffraction patterns and structures seen in the far-field pattern of the beam from most injection lasers.

Injection laser diodes at shorter wavelengths in the blue and green regions of the visible spectrum would be of considerable interest in many different applications. The wide band-gap II-VI compounds such as CdS, ZnSe, and ZnS have been studied extensively and have high luminescence efficiency when excited optically or with electron beams, but until 1991 injection lasers had not been obtained in wide band-gap II-VI semiconductors. In that year, however, heterostructure lasers utilizing the ZnSe-based single quantum well structure shown in Fig. 13 were reported.† The lattice matching of the different layers in heterostructure lasers is an important consideration for room-temperature semicon-

† Haase, M. A., Qui, J., DePuydt, J. M., and Cheng, H. "Blue-green laser diodes." *Appl. Phys. Lett.*, 59, (11), 1991, pp. 1272-1274.

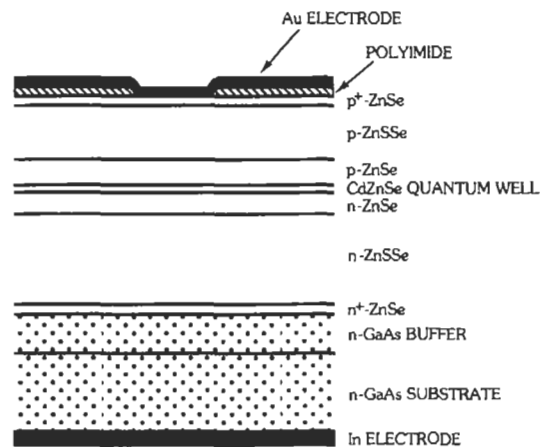


Fig. 13. Cross section of a blue-green laser diode.

ductor lasers, and the band-gap and lattice constant for important II-VI, III-V, and elemental semiconductors and semiconductor alloys are given in Fig. 14. The laser emission wavelength at 77 K was 490 nm, and this is the shortest wavelength ever generated by a semiconductor laser diode. More work will be required to optimize these devices and obtain room-temperature operation, but it is a very promising result.

OPTICAL DETECTORS

Detectors are designed to convert optical radiation into a current or voltage. Associated with both the "signal" radiation and the background radiation incident on the detector are fluctuations in the photon arrival rate, and these fluctuations induce fluctuations, or noise, referred to as signal noise and background noise, respectively, in the current or voltage produced by the detector. There are other sources of noise, such as thermal noise (Johnson or Nyquist noise), due to random thermal motion in any resistance in the circuit (including the detector itself), and generation-recombination noise in photoconductive detectors, due to random generation and recombination of free carriers, that can limit the sensitivity of an optical detector. There are several figures of merit used to evaluate the performance of detectors.

Terms and Figures of Merit

Noise equivalent power: (Abbreviated NEP.) The rms value of sinusoidally modulated radiant energy falling on the detector required to produce an rms-signal-to-rms-noise ratio of unity. The post-detection electrical bandwidth or noise bandwidth must be specified. Usually the NEP varies as the square root of the noise bandwidth, and the NEP is referenced to a 1-Hz bandwidth and expressed in units of $W/Hz^{1/2}$. The smaller the value of the NEP, the more sensitive is the detector. Noise equivalent power can be given for a particular wavelength (NEP_λ) or blackbody temperature NEP (T).

Detectivity: Because the NEP is smaller for higher sensitivity, it is common to define a figure of merit called the detectivity that is the reciprocal of the NEP reduced to a 1-Hz bandwidth, $D = 1/NEP$.

D^* : (Pronounced "dee star.") For some detector applications, especially those where the NEP and/or D is limited by background radiation, the NEP increases with the square root of the detector area, A , and to permit comparison of the performance of similar detectors with different areas, a normalized detectivity

$$D^* = A^{1/2} B^{1/2} / NEP \quad \text{cm Hz}^{1/2} \text{ W}^{-1}$$

is defined, where B is the electrical noise bandwidth.

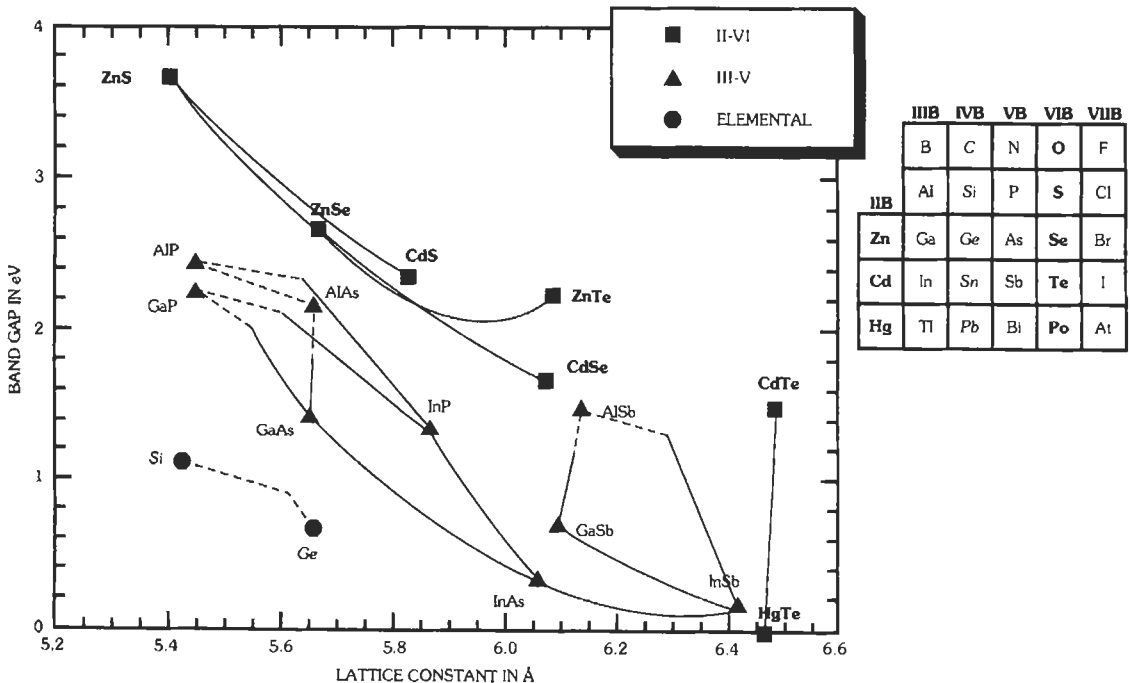


Fig. 14. Band gap and lattice constant of several II-VI, III-V, and elemental semiconductors and semiconductor alloys. (Courtesy Photonics Research Laboratory, 3M.)

Because of the way infrared detectors are characterized, it is common to specify the D^* as D^*_{bb} or D^*_λ to indicate measurement of the D^* with a blackbody or with a monochromatic source. Usually, D^* is given as $D^*_{bb}(T, f, \Delta f)$ for a blackbody temperature T , chopping frequency f , and bandwidth Δf , and $D^*_\lambda(\lambda, f, \Delta f)$ for measurement or evaluation at wavelength λ , chopping frequency f , and bandwidth Δf .

D^{**} : (Pronounced "dee-double star.") The detectivity of a background-limited detector depends not only on the temperature of the background but, for cooled detectors, also on the field of view. As the field of view is decreased, the background radiation is reduced, and the detectivity is increased, so long as the detectivity is limited by the background radiation. The theoretical improvement factor in the detectivity as a function of the cone angle θ is shown in Fig. 15. To permit the comparison of different detectors of the same kind under different background conditions, a figure of merit is defined as

$$D^{**} = (\Omega/2\pi)^{1/2} D^*$$

where for circular geometry the solid angle Ω subtended by the background for a detector with an angle of view of θ is $\Omega = 2\pi \sin^2(\theta/2)$. Thus,

$$D^{**} = D^* \sin(\theta/2) \text{ cm Hz}^{1/2} \text{ sr}^{1/2} \text{ W}^{-1}$$

For $\theta/2 = \pi/2$ or $\theta = \pi$, and $\Omega = 2\pi \text{ sr}$,

$$D^{**} = D^*$$

This variation assumes that the thermal radiation from

the cold shield which defines the field of view is negligible.

Responsivity: The responsivity is defined as the rms signal voltage or rms signal current per unit rms radiant power incident on the detector in units of volts/watt or amperes/watt. The voltage responsivity must be specified for a given load resistance or referenced to an open-circuit condition. The responsivity is also specified as R_{bb} or R_λ , denoting a blackbody responsivity or spectral responsivity at a particular wavelength λ as for D^* .

If R_λ is the responsivity in A/W for a detector at wavelength λ , then the signal current that results when an optical power $P_\lambda d\lambda$ at wavelength λ in interval $d\lambda$ is incident on the detector is

$$I_S = R_\lambda \cdot P_\lambda d\lambda \text{ amperes}$$

Quantum efficiency: The quantum efficiency η_λ of a photon or quantum detector is defined as

$$\eta_\lambda = \frac{\text{number of photoelectrons collected per second}}{\text{number of incident photons per second at wavelength } \lambda}$$

The number of photoelectrons collected per second is given by

$$I_S/q = R_\lambda P_\lambda d\lambda/q = q\eta_\lambda P_\lambda d\lambda/h\nu = \eta_\lambda P_\lambda d\lambda/(hc/\lambda)$$

where q is the charge of an electron and $h\nu$ and hc/λ are the photon energy at wavelength λ . Therefore, the quantum efficiency is given by

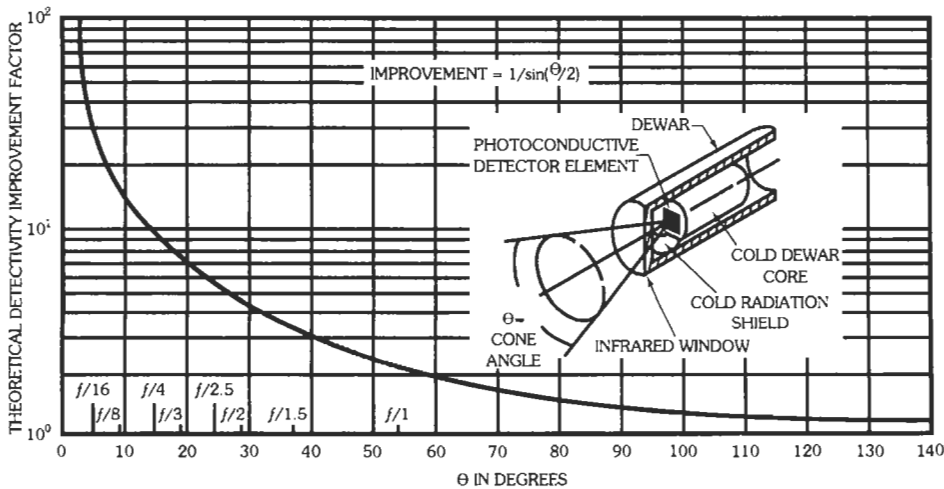


Fig. 15. Dependence of background-limited detectivity improvement on cooled-shield-limited field of view determined by cone angle θ or f number.

$$\eta_{\lambda} = R_{\lambda}hc/\lambda q$$

and $hc/q = 1.24 \times 10^{-6} \text{ W}\cdot\text{m}\cdot\text{A}^{-1}$, so

$$\eta_{\lambda} = 1.24 R_{\lambda}/\lambda \quad \text{for } \lambda \text{ in } \mu\text{m}$$

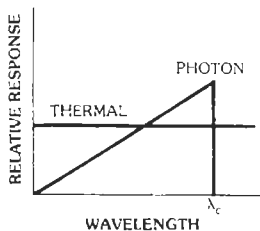
This relation permits a direct calculation of the quantum efficiency at wavelength λ from the measurement of the responsivity in A/W at the same wavelength. The term quantum efficiency usually refers to the interaction of the incident radiation with the detector, not including any gain mechanism in the device or in the electronic equipment used in the measurement of the signal current or voltage. Occasionally, however, external quantum efficiencies which include internal gain as in avalanche photodiodes are quoted, and in this case η can be significantly greater than unity.

Spectral response: The spectral responsivity variation for ideal photon detectors and thermal detectors is shown in Fig. 16. For photon detectors, the responsivity $R(\lambda)$, defined in terms of incident power, increases linearly with increasing wavelength up to the long-wavelength cutoff or the minimum energy for carrier excitation and then decreases abruptly to zero, whereas the responsivity for an ideal thermal detector is independent of wavelength. For photon detectors, a related figure of merit is the quantum efficiency η , the ratio of the number of electrons collected to the incident number of photons.

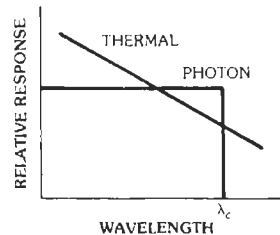
Response time: Another important measure of the performance of optical detectors is the response time. This is often specified in terms of a simple time constant, a rise time, or a cutoff modulation frequency.

Characterization of Detectors

The characterization of detectors is important because the results of such measurements make it possible to select the best possible detector for a particular application. The determination of NEP, D , D^* , or D^{**} described above each requires measurements of the responsivity, or the rms signal for a given incident power, and the noise of the detector being characterized.



(A) Constant flux per unit wavelength interval.



(B) Constant photons per unit wavelength interval.

Fig. 16. Idealized spectral response curves for photon and thermal detectors. (From Grum, F., and Becherer, R. J. *Optical Radiation Measurements, Vol. 1, Radiometry*. New York: Academic Press, Inc., 1979.)

Responsivity—The power responsivity of a photon detector can be determined with the use of a blackbody source by the following procedure:

1. The relative spectral power responsivity is measured at all wavelengths to which the detector responds by using a monochromator source and thermal detector that is assumed to have a flat or constant spectral power responsivity.
2. The absolute spectral power responsivity is then determined by either: (A) measuring the absolute power responsivity at one or a few wavelengths by using a laser or calibrated source of known power output, or (B) measuring the blackbody responsivity of the detector and, from this result and the relative spectral responsivity, determining the responsivity at the peak of the relative spectral response curve by means of the technique described below.

Fig. 17 shows a schematic diagram of an experimental arrangement for the measurement of the relative spectral power responsivity. The monochromator provides a nearly monochromatic source of radiation that is modulated by the chopper at some suitable frequency (in the range from about 10 Hz to about 1000 Hz) depending on the type of detector being characterized. This chopped radiation is divided by a 50% beam splitter that is flat over the wavelength range of interest so that equal intensities of radiation fall on both the detector being characterized and a thermocouple or other reference detector that is assumed to be spectrally flat. Alternatively, the radiation can be switched back and forth from the detector being characterized and the reference detector at a low frequency (about 10 Hz) compatible with the thermal detector by means of mirrored chopper blades. The signal from the thermal detector is then detected synchronously with the mirrored chopper. The output signals from the detector being measured and the thermocouple detector are then ratioed by the ratiometer to give a reading proportional to the power responsivity of the test detector. By recording this signal as the wavelength of the output radiation of the monochromator is scanned over the desired range, the relative spectral response of the test detector, $r(\lambda)$, can be obtained. Automated computer

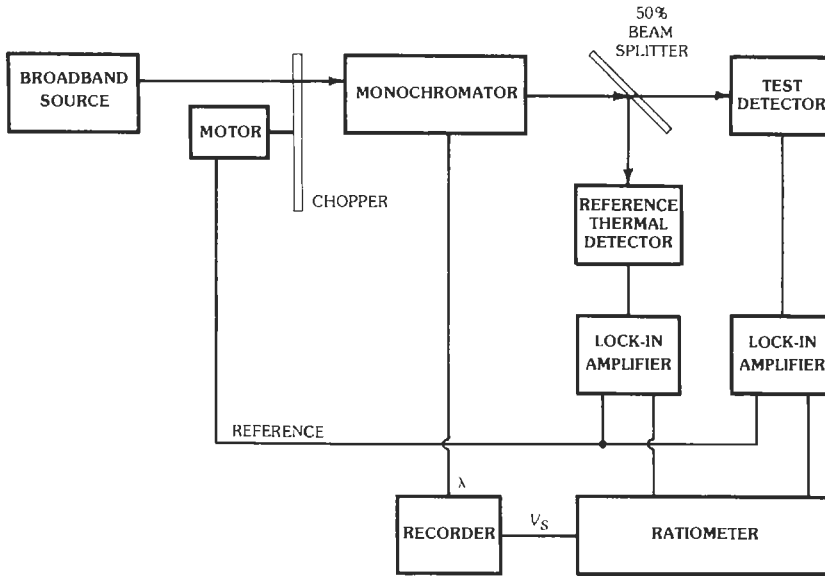


Fig. 17. Measurement of relative spectral power responsivity.

acquisition of data and calculated ratios can now often be used more easily than the analog ratiometer.

The absolute spectral responsivity of near-infrared detectors is usually determined by measuring the absolute responsivity at one or two wavelengths by using lasers with emission in the spectral range of interest and a calibrated output power. From these values of responsivity at specified wavelength, the relative spectral power responsivity curves can be made absolute.

For detectors in the wavelength range where calibrated lasers or reference detectors are not available, the spectral power responsivity is usually determined by using blackbody radiation. Fig. 18 shows an experimental arrangement for this process. The blackbody responsivity is

$$R_{bb}(f, T) = \frac{[(0.45 A_1 A_2) R_p / \pi d^2] \int_0^\infty r(\lambda) W(\lambda, T) d\lambda}{[(0.45 A_1 A_2) / \pi d^2] \int_0^\infty W(\lambda, T) d\lambda} = \frac{R_p \int_0^\infty r(\lambda) W(\lambda, T) d\lambda}{\sigma T^4}$$

$$R_{bb}(f, T) = \frac{R_p \int_0^\infty r(\lambda) P_{rms}(\lambda, T) d\lambda}{P_{rms}(T)}$$

where,

- R_p is the peak responsivity of the detector,
- $r(\lambda)$ is the relative spectral power responsivity of the detector ($0 \leq r(\lambda) \leq 1$),
- $P_{rms}(\lambda, T)$ is the rms value of the incident radiation at modulation frequency f and wavelength λ ,
- $P_{rms}(T)$ is the total rms blackbody radiation falling on the detector.

For a square-wave-modulated blackbody source of area A_1 at a distance d from the detector of area A_2 , as in the diagram of Fig. 18,

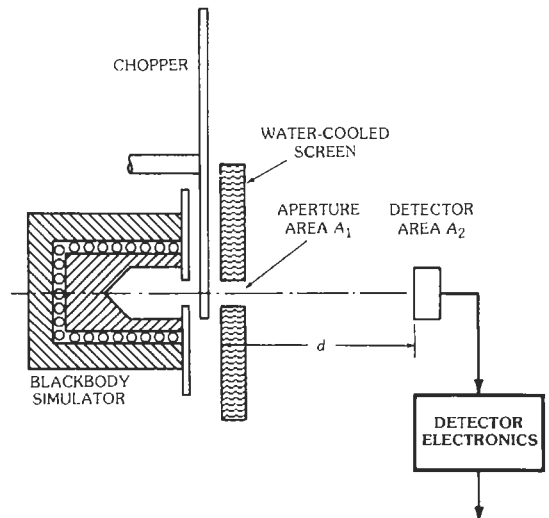


Fig. 18. Determination of spectral power responsivity with blackbody radiation.

But, the experimentally determined value of $R_{bb}(f, T)$ is

$$R_{bb}(f, T) = V_{rms}(f, T)/(0.45 A_1 A_2 / \pi d^2) \sigma T^4$$

where $V_{rms}(f, T)$ is the experimentally measured value of the rms voltage, for example, generated by the detector. The peak value of responsivity R_p is given in terms of the experimental blackbody responsivity $R_{bb}(f, T)$ by

$$R_p = \sigma T^4 R_{bb}(f, T) / \int_0^\infty r(\lambda) W(\lambda, T) d\lambda$$

$$= R_{bb}(f, T) / \gamma(T)$$

where $\gamma(T)$ is called the effectiveness factor. For an ideal thermal detector, $r(\lambda) = 1$, $\gamma = 1$, and $R_p = R_{bb}(f, T)$. The units of R_p are amperes/watt or volts/watt.

Noise Measurements—In order to determine the NEP and related figures of merit for the detectors described above, it is necessary to measure the rms noise voltage or current characteristics of the detector being evaluated. The rms noise can be measured with rms-responding meters, lock-in amplifiers, spectrum analyzers with special noise-measuring capability, or tuned voltmeters with average-responding meters. When the latter are used, it should be noted that the meter reading for the rms value of a sinusoid must be multiplied by 1.128 in order to obtain the rms value of Gaussian noise. The NEP can then be calculated from the responsivity and the rms noise as $NEP_\lambda = V_n / R_\lambda$, where V_n is the rms noise voltage, for example. The D , D^* , etc., can be calculated similarly.

Ultimate Sensitivity of Detectors*

There are many noise sources that can limit the detection capability of any detection system. These noise sources can generally be categorized as noise in the detector, noise in the amplifier, and fluctuations in the radiating background and the signal itself. The noise in the detector comes from fluctuations in the concentration or motion of the current carriers. The amplifier noise consists of thermal noise as well as excess noise due to the same type of fluctuations as in the detector. The background noise or photon noise is due to the quasi-random absorption of photons from the surroundings of the detector; for thermal detectors it also comes from the quasi-random emission of photons from the detector. It is this background noise that limits the ultimate performance of infrared detectors.

Johnson, Nyquist, or Thermal Noise—This is the random fluctuation in the current or voltage at the electrical contacts of a resistor due to the random thermal motion of electrons in the resistor. The rms values of these fluctuations are given by

$$\overline{v_j^2}^{1/2} = (4kTB R_L)^{1/2}$$

$$= 1.287 \times 10^{-10} R_L^{1/2} (T/300)^{1/2} B \text{ volts}$$

and

$$\overline{i_j^2}^{1/2} = (4kTB/R_L)^{1/2}$$

$$= 1.287 \times 10^{-10} R_L^{-1/2} (T/300)^{1/2} B \text{ amperes}$$

where,

R_L is the resistance in ohms,
 T is the absolute temperature,
 B is the electrical bandwidth.

This noise source is independent of current for a given temperature.

For photovoltaic detectors, the diode terminal current for small voltages is linear with voltage and can be written as

$$I_t = -I_p + (1/R_o) V$$

where,

I_t is the diode terminal current,
 I_p is the photogenerated current,
 V is the open-circuit photovoltage,
 $R_o = (\partial V / \partial I) |_{V=0}$.

When the diode is operated as an open-circuited photovoltaic detector, $I_t = 0$. Then,

$$V = I_p R_o = q \eta (P_\lambda / E_\lambda) R_o$$

where P_λ is the incident optical power at wavelength λ , E_λ is the photon energy at wavelength λ , and the rest of the parameters have meanings defined previously. For this case, the voltage responsivity $R_{V,\lambda}$ at the wavelength λ is given by

$$R_{V,\lambda} = V/P_\lambda = \eta q R_o / E_\lambda$$

The noise voltage for the unbiased open-circuit photovoltaic detector is just the Johnson noise of the incremental diode resistance,

$$\overline{v^2} = 4kTR_o \Delta f$$

The NEP is given by the power P_λ required to give a signal/noise ratio of unity, and if the photovoltaic detector noise is only due to the Johnson noise of R_o ,

* Vincent, John David. *Fundamentals of Infrared Detector Operation & Testing*. New York: John Wiley and Sons, 1990.

$$NEP = 2E_{\lambda} (kT)^{1/2} / \eta q R_o^{1/2}$$

and the detectivity D^*_{λ} is†

$$D^*_{\lambda} = A^{1/2} / NEP = \eta q (R_o A)^{1/2} / 2E_{\lambda} (kT)^{1/2}$$

The variation of D^*_{λ} with R_o for a Johnson-noise limited open-circuit photon detector, i.e., under sufficiently low background conditions, is given in Fig. 19.

Shot Noise—Shot noise results whenever the fluctuation arises because of a series of independent events occurring at random, such as the emission of electrons by a temperature-limited thermionic cathode, or a photocathode, or the crossing of a junction by electrical carriers (electrons or holes) in a pn junction or transistor. Shot noise also results when transitions occur between two energy levels as for generation and recombination of carriers in a semiconductor or when photons are emitted by a laser. The root-mean-square current fluctuations for the process can be expressed in terms of the average dc current as

$$\begin{aligned} \overline{i_s^2}^{1/2} &= (2qI_{dc}B)^{1/2} \\ &= 5.657 \times 10^{-10} (I_{dc})^{1/2} B^{1/2} \text{ amperes} \end{aligned}$$

where,

† Melngailis, I., and Harmon, T. C. "Single-Crystal Lead-Tin Chalcogenides," in *Semiconductors and Semimetals*, Vol. 5. Academic Press, 1970, pp. 111-163.

I_{dc} is the average dc current in amperes,
 B is the electrical bandwidth in hertz.

Generation-Recombination Noise—Although the noise that results from fluctuations in the instantaneous values of the free carrier concentration in semiconductors can be described as shot noise for some cases, for photoconductive detectors it is more meaningful to describe this noise source in terms of the generation, recombination, and trapping processes that cause the fluctuations in the free carrier concentrations. The G-R noise for a homogeneous n-type extrinsic semiconductor is

$$\overline{i_{GR}^2}^{1/2} = \{(4I_{dc}^2/N_0)[\tau_0/(1 + \omega^2\tau_0^2)]B\}^{1/2}$$

amperes

where,

I_{dc} is the average dc current,
 N_0 is the total number of free carriers in the semiconductor,
 τ_0 is the carrier lifetime.

1/f Noise—There are several types of noise that increase as the frequency decreases. Because of the different, and not necessarily well understood, sources of this noise, there is no simple expression that describes 1/f noise accurately. However, empirical expressions like

$$\overline{i_{1/f}^2}^{1/2} = (KI_{dc}^{\alpha} B/f^{\beta})^{1/2}$$

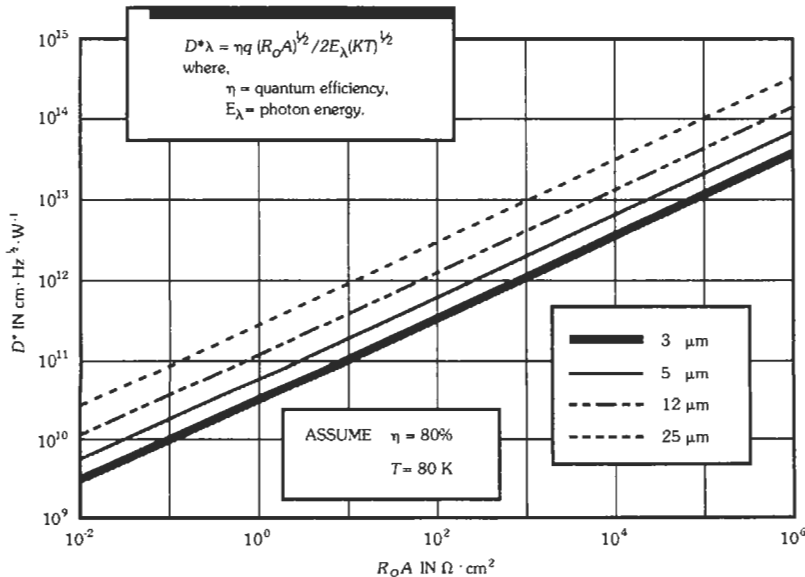


Fig. 19. D^*_{λ} for Johnson-noise-limited photon detectors where R_o is the differential resistance at zero bias.

where K is a proportionality factor, $\alpha \approx 2$, and $\beta \approx 1$ are frequently used. In this expression, I_{dc} is the average current, and B is the bandwidth.

Photon Noise—Also called radiation noise, photon noise results from fluctuations in the rate of absorption (and for thermal detectors, emission) of photons by the detector. For photons from a laser or nonblackbody source, this noise can be calculated accurately by the shot noise of the average current induced by the photon flux. For blackbody radiation, the noise is greater than that estimated by the shot-noise formula because blackbody photons, described by Bose-Einstein statistics, do not arrive randomly in time but, instead, are clustered together. Thus, the mean square fluctuation in the arrival rate of photons, which for independent events would simply be the average photon arrival rate, \bar{n} , is increased by the factor $e^{h\nu/kT}/(e^{h\nu/kT} - 1)$. For practical applications, this distinction is generally not important.

Temperature Noise—Temperature noise is the fluctuation in temperature of a thermal detector that is in contact with its surroundings by both conduction and radiation. If the detector is only in contact with its surroundings through radiation, the temperature noise reduces to blackbody photon noise due to both absorption and emission of photons.

Amplifier Noise—There is always the thermal Johnson noise of some equivalent load resistance associated with a detector circuit, but in addition, with real amplifiers there is some additional added noise called amplifier noise. This noise, as well as the Johnson noise of the load resistor, is usually referred to the input of the amplifier. The rms noise voltage referred to the amplifier input is described by the equation

$$\overline{v_a^2}^{1/2} = [4k(T_L + T_A)R_L B]^{1/2}$$

where,

- T_L is the noise temperature of the equivalent load resistance, R_L ,
- T_A is the amplifier noise temperature,
- B is the electrical bandwidth.

The noise figure, F , is another useful measure of amplifier performance and is the ratio of the signal-to-noise ratio at the amplifier input to the signal-to-noise ratio at the output for a room-temperature load resistor ($T_L = 290\text{ K}$). In terms of T_A and T_L ,

$$F = 1 + T_A/T_L$$

The noise figure is usually expressed as a logarithmic function,

$$NF = 10 \log_{10} F \text{ decibels}$$

In this notation,

$$\overline{v_a^2} = 4kFT_L R_L B$$

It should be noted that F for a given amplifier usually depends on both R_L and the frequency, f .

Ultimate Detector Sensitivity—The ultimate sensitivity of a detector is limited by the shot noise of the signal itself. For photon detectors this is given by

$$NEP_\lambda = 2hcB/\lambda\eta(\lambda)$$

where $\eta(\lambda)$ is the detector quantum efficiency at wavelength λ . In practice, this limit is seldom if ever achieved, because of the various noise sources described above. However, because of the large amount of background radiation present in much of the infrared spectral region, it is common for the sensitivity of good infrared detectors to be limited by the photon noise due to background radiation. The mean square noise power spectral density of the photon flux from blackbody radiation at temperature T incident on a thermal detector of area A is given by

$$S_p(f) = \int_0^\infty \frac{2A2\pi h^2 \nu^4}{c^2} \frac{e^{h\nu/kT} d\nu}{(e^{h\nu/kT} - 1)^2} = 8AkT\sigma T^4$$

and for a thermal detector, in equilibrium with and at the same temperature as the background (and with unity emissivity), this fluctuation must be doubled because of the fluctuations due to the photons emitted by the detector; so

$$S_p(T) = 16 AkT\sigma T^4$$

Thus, the mean square noise power in a bandwidth B is

$$\overline{P_b^2} = 16AkT\sigma T^4 B$$

and the D^* for such a thermal detector is

$$D^* = A^{1/2} B^{1/2} / \overline{P_b^2}^{1/2} = 1/4(kT\sigma T^4)^{1/2} \\ = 5.5171 \times 10^{-11} (T/300)^{5/2} \text{ cm Hz}^{-1/2} \text{ W}^{-1}$$

The D^* for an ideal photon detector with a long-wavelength cutoff of $\lambda_0 = c/\nu_0$ for a signal at ν_s within the spectral range of the detector is given by

$$D_\lambda^* = \frac{c\eta(\nu_s)}{2h\nu_s \pi^{1/2} \left[\int_{\nu_0}^\infty \frac{\eta(\nu)\nu^2 \exp(h\nu/kT) d\nu}{[\exp(h\nu/kT) - 1]^2} \right]^{1/2}}$$

and this can be readily evaluated for the peak detectivity of an ideal detector for ν_s and ν_0 , assuming $\eta(n)$ is independent of wavelength, by using a computer or a programmable calculator and the series solution

$$D_{\lambda_0}^* = \frac{c[\eta(\nu_0)]^{1/2}}{2\pi^{1/2} h^{1/2} \nu_0^2 (kT)^{1/2}} U \quad \text{where } U \text{ is}$$

$$\left\{ \sum_{n=1}^{\infty} \exp(-nh\nu_0/kT) [1 + 2kT/mh\nu_0 + 2(kT/mh\nu_0)^2] \right\}^{-1/2}$$

The background-limited $D^*_{\lambda_0}$ at the peak wavelength, λ_0 , is shown in Fig. 20 for several values of the background temperature. This result applies to photoemission and photovoltaic detectors where only the background fluctuations contribute to the noise. For $T = 300 \text{ K}$, the minimum $D^*_{\lambda_0}$ occurs at about $14 \mu\text{m}$. The rapid increase in $D^*_{\lambda_0}$ for detectors with short cutoff wavelengths occurs because these detectors respond to less of the photon noise of this background. The small increase in $D^*_{\lambda_0}$ for longer wavelength occurs because the increase in the photon noise received from the background is more than compensated by the increase in the number of photons per watt of signal. For photoconductive photon detectors, there is the additional fluctuation due to fluctuations in the recombination of carriers in addition to the excitation of free carriers, so for photoconductive detectors the photon noise is larger by $\sqrt{2}$ and the background-limited D^* is smaller by $\sqrt{2}$.

The various types of thermal and photon or quantum detectors are described briefly below.

Thermal Detectors

Thermal detectors respond to optical-radiation-induced temperature variations and are therefore well-suited for broadband detection throughout the optical spectrum. Thermal detectors include the bolometer,

thermocouple, thermopile, thermopneumatic cell, and pyroelectric detector.

Bolometer—The bolometer changes its resistance in response to thermal energy resultant from impinging radiant energy. The most common bolometric detector is the thermistor.

Thermocouple—A thermocouple is a junction of two dissimilar metals that, upon absorbing thermal energy, produces an emf.

Thermopile—The thermopile is an array of thermocouples. Miniature thermopiles made with thick film deposition techniques are commonly used in infrared detection applications.

Thermopneumatic Cell—The thermopneumatic cell senses the presence of thermal energy by means of a sealed cell that expands or contracts in response to variations in applied radiant energy. The magnitude of the expansion can be detected by means of interferometric techniques.

Pyroelectric Detector—The pyroelectric detector is a temperature-sensitive current source. In a typical detector, a thin wafer of a ferroelectric crystal such as triglycine sulfate or lithium tantalate forms a capacitor whose capacitance is altered by thermal energy. Because of the ferroelectric effect, this detector will not detect constant unmodulated incident radiation.

Quantum Detectors

Quantum detectors respond to variations in the number of incident photons. Quantum detectors have a

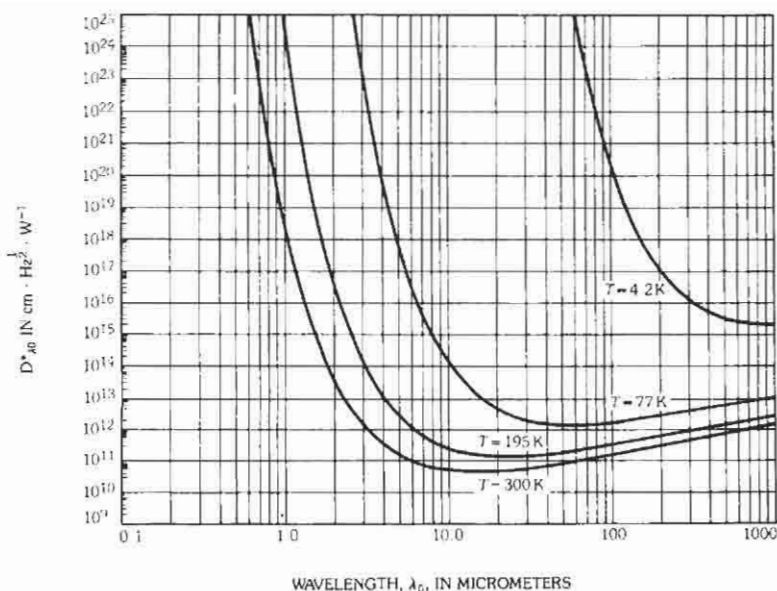


Fig. 20. Background-limited $D^*_{\lambda_0}$ at peak wavelength, λ_0 , for ideal photon detection with threshold wavelength λ_0 and different background temperatures T .

more limited spectral sensitivity range than thermal detectors but are generally characterized by relatively fast response time and high sensitivity. Quantum detectors include photoemissive devices, photovoltaic cells, photoconductive cells, and photoelectromagnetic cells.

Photoemissive Devices—If photons with sufficient energy impinge on a photocathode, electrons are emitted. Such electrons are known as photoelectrons. For an input flux of fixed relative spectral distribution, the number of photoelectrons is proportional to the intensity of the input flux, whereas the energy of the photoelectrons is independent of this intensity. The maximum energy of emitted electrons expressed in volts, V , depends on the wavelength, λ , and the temperature. At absolute zero, according to Einstein's law

$$e(V + \phi) = hc/\lambda$$

where,

- e = electron charge = 1.6×10^{-19} coulomb,
- ϕ = work function in volts,
- h = Planck's constant = 6.6×10^{-34} joule-second,
- c = velocity of light in meters/second,
- λ = wavelength in meters.

If a threshold wavelength, λ_0 , is defined by

$$e\phi = hc/\lambda_0$$

then V is seen to be zero (except for thermal velocities) at the wavelength λ_0 ; for $\lambda > \lambda_0$, there is no photoelectric emission at absolute zero. At temperatures above absolute zero, there is always a finite probability of some photoemission at all wavelengths due to the thermalization of the electron distribution.

Photocathode Response to Monochromatic Radiation—The output current dI_λ in amperes, generated by a photocathode subjected to a monochromatic input flux dW_λ in watts, is given by

$$dI_\lambda = s_\lambda dW_\lambda$$

where s_λ is the monochromatic radiant responsivity of the photocathode in amperes/watt defined by this equation. Similarly, the number of electrons/second, dn_λ , generated by an input flux of dN_λ photons/second is given by

$$dn_\lambda = \eta_\lambda dN_\lambda$$

where η_λ is the monochromatic quantum efficiency of the photocathode in electrons/photon defined by this equation.

The monochromatic radiant responsivity, s_λ , in A/W, is related to the monochromatic quantum efficiency, η_λ , at wavelength λ , in μm , by

$$s_\lambda = e\lambda\eta_\lambda/hc = 0.81\lambda/\eta_\lambda$$

Typical values of the monochromatic radiant responsivity, s_λ , and corresponding monochromatic quantum efficiency, η_λ , as a function of wavelength, λ , are shown in Fig. 21 for some commonly used photocathodes, designated by their JEDEC registered "S numbers." Table 3 gives typical peak responsivities for the various surfaces, and Table 4 indicates the general composition and other properties of the common surfaces.

Vacuum Photodiodes—The combination of a photocathode and an anode electrode for collecting the emitted photocurrent in an evacuated envelope is called a vacuum photodiode. A positive anode potential sufficient to assure collection of all emitted photoelectrons (that is, to "saturate" the diode phototube) is normally required, the tube then acting as a constant-current generator (Fig. 22). The power-supply potential, V_B , must assure sufficient anode potential in the presence of a voltage drop in the load resistor, R_L .

Under these conditions the total anode output current, I_a , neglecting all noise fluctuations, is given by the equation

$$I_a = I_s + I_b + I_d + I_L$$

where,

- I_s = emitted photocathode signal current,
- I_b = emitted photocathode photocurrent due to stray background flux,
- I_d = photocathode thermionic dark current,
- I_L = residual dark current (leakage, etc.).

The instantaneous value of signal current I_s will follow the instantaneous signal flux input magnitude from direct current up to an upper frequency limit (commonly 0.2–2 GHz) set by the transit-time spread of the electrons crossing the gap between cathode and anode, and including induced displacement currents during transit.

For steady-state or slowly varying input flux rates, the total noise current output, i_n , from the diode and load resistor is given by

$$\overline{i_n^2} = 2qB(I_s + I_b + I_d) + (4kTB/R_L) + i_L^2$$

where,

- q = the electron charge = 1.6×10^{-19} coulomb,
- B = noise-current measurement bandwidth,
- k = Boltzmann's constant = 1.38×10^{-23} joule per K,
- T = absolute temperature of the load resistor (K),
- R_L = load resistance (ohms),
- i_L = residual dark noise current (from leakage, stray pickup, etc.).

To increase the absolute level of the noise voltage generated by noise current i_n so that tube noise predominates over the noise voltage of the subsequent amplifiers, and to suppress load-resistor noise relative to tube

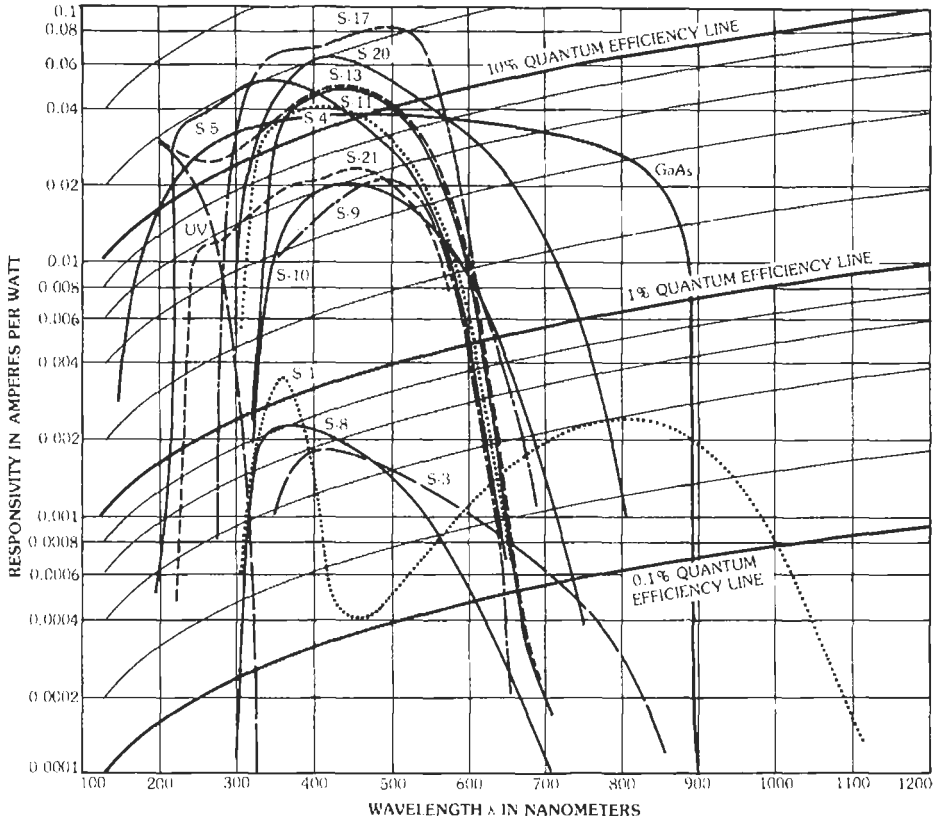


Fig. 21. Typical absolute spectral response characteristics of various photocathodes.

noise, large values of the load resistor (of the order of 10^7 – 10^9 ohms) are commonly used when response time is not a limitation.

TABLE 3. TYPICAL PEAK PHOTOCATHODE RESPONSIVITIES

S Number	Radiant Responsivity		Quantum Efficiency	
	$S_{\lambda_{max}}$ ($A \cdot W^{-1}$)	λ_{max} (nanometers)	$\eta_{\lambda_{max}}$ (electron photon ⁻¹)	λ_{max} (nanometers)
S1	0.0025	800*	0.004	770*
S3	0.0019	420	0.0058	400
S4	0.042	420	0.13	380
S5	0.052	330	0.21	320
S8	0.0024	360	0.0082	350
S9	0.023	490	0.056	480
S10	0.021	440	0.061	420
S11	0.048	450	0.14	400
S13	0.048	450	0.20	<200
S17	0.083	500	0.26	350
S20	0.066	420	0.20	400
S21	0.024	460	0.07	370

* Neglecting short wavelength peak

For a plane-parallel vacuum photodiode, the space-charge-limited output current, $I_{a(max)}$, in amperes for a given applied cathode-to-anode potential difference, V , in volts is given by

$$I_{a(max)} = 2.33 \times 10^{-6} AV^{3/2}/d^2$$

where,

A = uniformly emitting emission area (meters²),
 d = anode-cathode spacing (meters).

In practice, linear output currents up to approximately half of this maximum limit can be obtained.

For a plane-parallel vacuum phototube, the output anode current, I_A , as a function of time, t , for an ultrashort exciting light pulse is given by

$$I_A = (2QRC/T^2)[(t/RC) + \exp(-t/RC) - 1] \quad \text{for } 0 < t < T$$

$$I_A = (2QRC/T^2)[(T/RC) + \exp(-T/RC) - 1] \exp[-(t - T)/RC] \quad \text{for } t \geq T$$

and

TABLE 4. CHARACTERISTICS OF STANDARD PHOTOSURFACES

S Number ¹	Principal Photocathode Components ²	Entrance Window Material	Photocathode Supporting Substrate ³	Typical Luminous Responsivity ⁴ ($\mu\text{A}/\text{lumen}$)	Typical Photocathode Dark Current ⁵ at 25 °C (A/cm^2)
S1	Ag-O-Cs	Visible-light-transmitting glass ⁶	Entrance window or opaque material ⁷	25	$10^{-11} - 10^{-13}$
S3	Ag-O-Rb	Visible-light-transmitting glass ⁶	Opaque material ⁷	6.5	10^{-12}
S4	Cs-Sb	Visible-light-transmitting glass ⁶	Opaque material ⁷	40	10^{-14}
S5	Cs-Sb	Ultraviolet-transmitting glass	Opaque material ⁷	40	10^{-14}
S8	Cs-Bi	Visible-light-transmitting glass ⁶	Opaque material ⁷	3	$10^{-14} - 10^{-15}$
S9	Cs-Sb	Visible-light-transmitting glass ⁶	Entrance window	30	10^{-14}
S10	Ag-Bi-O-Cs	Visible-light-transmitting glass ⁶	Entrance window	40	$10^{-13} - 10^{-14}$
S11	Cs-Sb	Visible-light-transmitting glass ⁶	Entrance window	60	$10^{-14} - 10^{-15}$
S13	Cs-Sb	Fused silica	Entrance window	60	$10^{-14} - 10^{-15}$
S17	Cs-Sb	Visible-light-transmitting glass ⁶	Opaque reflecting material ⁷	125	$10^{-14} - 10^{-15}$
S19	Cs-Sb	Fused silica	Opaque material ⁷	40	10^{-14}
S20	Sb-K-Na-Cs	Visible-light-transmitting glass ⁶	Entrance window	150	$10^{-15} - 10^{-16}$
S21	Cs-Sb	Ultraviolet-transmitting glass	Entrance window	30	10^{-14}
UV ⁸	Cs-Te	Sapphire	Opaque material ⁷	0	—

Notes:

1. The S number is the designation of the spectral response characteristic of the device and includes the transmission of the device window material.
2. Principal components of the photocathode are listed without regard to order of processing or relative proportions.
3. When the supporting substrate is the entrance window, an intermediate semitransparent electrically conductive layer may be used.
4. Corresponding to the specific absolute response curves shown in Fig. 21 using a 2854 K color-temperature tungsten-lamp test source.
5. Specific dark current excludes direct-current leakage.
6. Lime glass and Kovar sealing borosilicate glass are commonly used for visible-light-transmitting glass.
7. The opaque material used as the supporting substrate for photocathodes in which the input radiation is incident on the same side as the emitted photoelectrons is usually metallic in nature.
8. An S number designation has not yet been assigned to this experimental "solar blind" photoemissive surface.

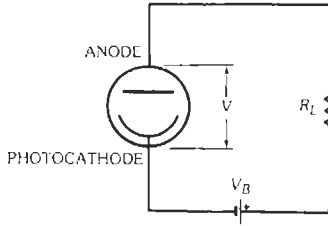


Fig. 22. Photodiode circuit.

$$T = d(2m/qV)^{1/2}$$

$$= (3.37 \times 10^{-6})d/(V)^{1/2} \text{ seconds}$$

where,

- T = transit time of the charge from cathode to anode (seconds),
- V = cathode-to-anode potential (volts),
- C = total capacitance including external circuit capacitance,
- R = load resistance,
- Q = total charge,
- m = electron mass.

Photomultipliers—The combination of a photocathode and a secondary-emission electron multiplier is called a photomultiplier. Emitted photoelectrons from the photocathode are directed under the influence of a suitable electrode, often called the “focus electrode,” to the surface of a secondary-emitting electrode, the “first dynode.” Subsequently emitted secondary electrons, increased in number by the effective secondary-emission ratio, σ_1 , are then directed to the secondary emitting surface of a subsequent dynode by an appropriate electric field for further multiplication. Continuing this process for n successive dynodes and collecting the multiplier charge at an output electrode called the “anode” or “collector” leads to a charge or current amplification, G , given by

$$G = \sigma^n, \text{ for } \sigma_1 = \sigma_2 = \sigma_3 \dots = \sigma$$

Gains as high as 10^5 – 10^9 are commonly achieved in 10-stage to 16-stage multipliers.

Disregarding all noise fluctuations of the output current, and assuming operation within the usual linear-response region, a photomultiplier acts as a constant-current source generating an output current, I_o , given by

$$I_o = I_s + I_b + I_d = GeI_{ks} + GeI_{kb} + GeI_{kd} + I_{ad}$$

where,

- I_s = anode signal current due to an incident signal flux to be detected,
- I_b = anode current due to any background flux simultaneously present on the photocathode,

- I_d = anode dark current,
- G = current gain of the electron multiplier,
- ϵ = collection efficiency (ratio of current entering the electron multiplier to emitted photocathode current),
- I_{ks} = photocathode signal current due to an incident signal flux to be detected,
- I_{kb} = photocathode current generated by any background flux simultaneously present on the photocathode,
- I_{kd} = photocathode dark current,
- I_{ad} = component of anode dark current I_d not originating from the photocathode.

The output signal current, I_s , follows the instantaneous value of the input signal flux from direct current up to an upper frequency limit (typically 20–200 megahertz) established by the response time or “transit-time spread” (typically 1–10 nanoseconds).

Noise fluctuations of the output current in photomultipliers can be divided into two classes: dark noise, occurring in the absence of input flux, and noise-in-signal, including “quantum” noise resulting from the inherent quantum nature of the input flux as well as uncontrolled fluctuations of that flux. The presence of an appreciable, in fact often predominant, noise-in-signal current component in photomultipliers depending on the instantaneous signal current magnitude requires caution in applying noise concepts to photomultipliers and may lead to erroneous conclusions regarding photomultiplier behavior, particularly for a modulated flux input.

For a steady-state unmodulated flux input, the total noise current, i_n , in the load resistance, R , is given by

$$i_n^2 = 2eGK\Delta f(I_s + I_b + GeI_{kd}) + i_r^2 + (4kT\Delta f/R)$$

where,

- e = electron charge,
- K = photomultiplier noise factor,
- Δf = noise bandwidth of the noise-current measuring circuits (Hz),
- i_r = residual photomultiplier anode dark noise current, excluding dark-current emission from the photocathode (amperes),
- $(4kT\Delta f/R)^{1/2}$ = Johnson-Nyquist noise current in load resistance R (amperes),
- k = Boltzmann’s constant, 1.38×10^{-23} joule/K,
- T = absolute temperature of the load resistance (kelvins).

For photomultipliers with a constant gain per stage, σ , in the first few stages of the electron multiplier, the noise factor, K , may be estimated from

$$K = \sigma/(\sigma - 1)$$

Photomultipliers as Scintillation and Single-Electron Counters—In combination with suitable

scintillating material, typically thallium-activated NaI crystals, photomultipliers are extensively used to detect the single flashes of light generated by the scintillating material on bombardment by a single triggering particle, typically X rays or a gamma ray from a nuclear disintegration process. If the scintillating material generates an average of N photons per disintegration incident on the effective photocathode of peak quantum efficiency, Y_{\max} , the resultant average charge pulse, Q_A , appearing in the anode circuit (disregarding all photons or electrons producing no output charge) will be given by the equation

$$Q_A = NY_{\max} \alpha G e$$

where α is a spectral matching factor describing the relative match between the scintillator spectral output and cathode sensitivity.

Because of the random statistical fluctuations of cathode quantum efficiency Y_{\max} and electron-multiplier gain G , as well as in the number of effective photons, N , generated by the scintillator, the anode charge, Q_A , will vary in magnitude from pulse to pulse, introducing ambiguity in the determination of the average magnitude of N , which in turn is used to determine the energy of the triggering input particle, for example the gamma ray. The ratio of the spread of the amplitude of individually observed values of charge Q_A at half maximum to the most probable value $Q_{A(\max)}$ is called the "energy resolution" of the photomultiplier-plus-scintillator combination and is commonly 7-10% minimum (see Fig. 23).

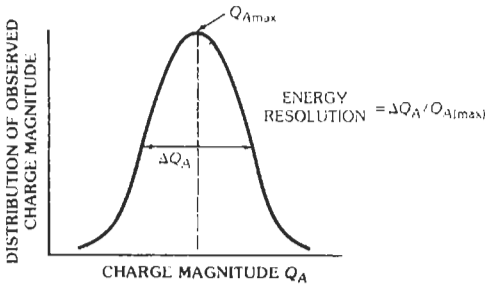


Fig. 23. Energy resolution of photomultiplier.

If the input flux has no time-coherent groups of photons, as it does in scintillation detection, photoelectrons are emitted singly at random emission times from the photocathode and also generate an average output charge, Q_A , given by

$$Q_A = Gq$$

where all photons or electrons generating no output pulse are disregarded in measuring G and computing

Q_A . Assuming sufficiently large gain G and sufficiently low generation of dark pulses of similar charge amplitude, the individual anode pulses of charge amplitude Gq can be detected and counted individually, the photomultiplier then acting as a single-electron counter.

Photoconductivity—Photoconductivity is the increase in electrical conductivity of a material that takes place when the material is illuminated with infrared, visible, or ultraviolet light.

The absorption of light is a quantum process in which electrons are excited to higher energy levels. Ordinarily, the excited electrons are more mobile than unexcited

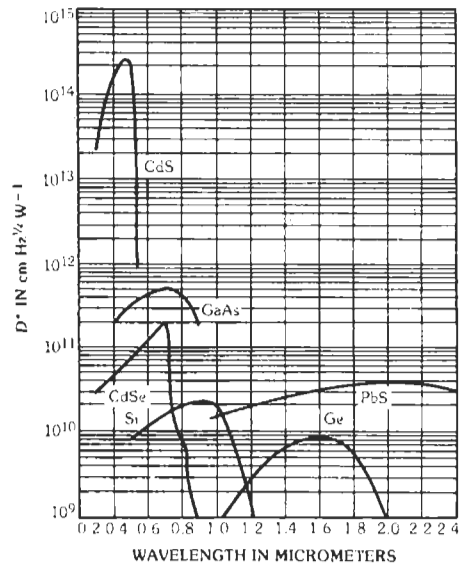


Fig. 24. Detectivity for some typical photoconductors at room temperature.

electrons. Photoconductivity is commonly analyzed in terms of the number and mobility of the excited electrons in an electron conduction band and of holes in a lower-energy valence band. To maintain a steady current, both types of current carriers must be generated in the volume of the material, or else charge carriers must enter the photoconductor at one of the electrodes. Many high-resistivity photoconductors make "ohmic" contacts with their electrodes. These serve as practically unlimited reservoirs of mobile electrons, free to enter the photoconductor volume. Even in these photoconductors, the steady dark current is usually limited to a low value by a build-up of a space-charge-potential barrier in the photoconductor.

At the same time that mobile photoelectrons are excited (thermally or optically) in the interelectrode volume, positive charges must also be generated; these

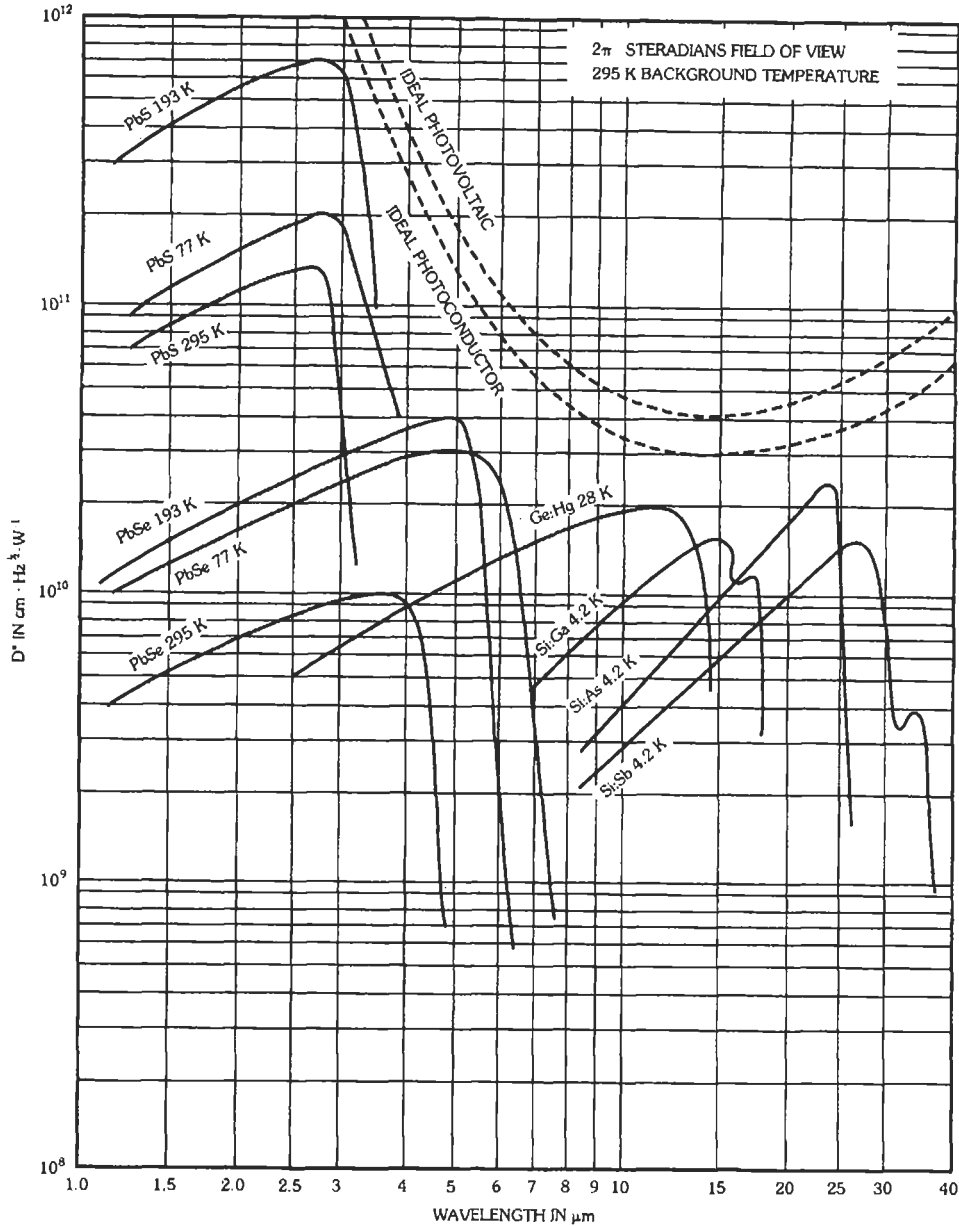


Fig. 25. Spectral detectivities of several thin-film and doped germanium and silicon extrinsic photoconductive detectors. (From Santa Barbara Research Center.)

compensate the charge of the photoelectrons in such a way that a "photocurrent" can be superimposed on the small space-charge-limited "dark" current originating at the electrode. If the positive charges are immobile, then long after the photoelectron has passed through the photoconductor into the anode, the immobile positive charges may remain to support a photocurrent of electrons, drawing on the reservoir of electrons at the

cathode. This will continue until the immobile "holes" or impurity centers are neutralized by recombination with some of the mobile electrons. Since the recombination lifetime may be much longer than the electron transit time between electrodes, the number of "photoelectrons" transported across the photoconductor may be much larger than the rate of generation of photoelectrons in the photoconductor volume. This ratio of

photocurrent to generation rate is called the photoconductive gain.

The photoconductive gain of a pure material can often be greatly increased by addition of localized traps lying near the conducting band. Since these are in thermal equilibrium with the conducting band, they serve as an additional reservoir of the charge carriers. This can increase both the response time and the sensitivity by a large factor.

Practically all materials are photoconductors in the sense that light of the correct wavelengths will generate current carriers. However, in many materials the photoconductivity is not detectable by ordinary measurements, either because of very short carrier lifetimes or because of a large dark current. The useful photoconductors, characterized by comparatively long lifetimes and low dark currents, have most of their charge carriers immobile (in the dark). Light of the proper energy can excite these carriers through the forbidden energy regions into the conduction bands. The long-wavelength limit of photoconductivity at low temperatures is given approximately by

$$\lambda_{max} = hc/E_g$$

where,

- E_g is the forbidden band gap,
- h is Planck's constant,
- c is the velocity of light.

For wavelengths longer than 5 micrometers, this equation gives a band gap smaller than 1/4 eV. Photoconductors with such small energy gaps are usually cooled to reduce the dark conductivity due to thermal excitation of carriers across the gap.

Values of D^* , a commonly used figure of merit for photoconductors, for several typical photoconductors at room temperature are shown in Fig. 24. The photoconductors with long-wavelength cutoffs will be considerably more sensitive if they are cooled below room temperature.

Photodiodes†—When photons irradiate a semiconductor having a band gap less than the energy of the photons, hole-electron pairs are produced. By forming a pn junction in a semiconductor, those hole-electron pairs can be collected. The resulting photodiode may have a quantum efficiency (ratio of collected photoelectrons to impinging photons) of from 0.2 to 0.7 or higher with antireflection coating.

Photodiodes have two primary operating modes, the photovoltaic mode and the reverse-biased mode that is often referred to incorrectly as the photoconductive mode. In the photovoltaic mode, the unbiased junction is illuminated to stimulate the production of hole-

electron pairs. Charge separation then occurs in the field of the junction, and a current and/or voltage results. Since photovoltaic cells produce a short-circuit photocurrent that is linear with respect to the radiation incident upon the pn junction, they are well suited for low-level light detection and measurement. This is especially true for narrow-band-gap materials and devices where the leakage current of reverse-biased diodes is often high enough to contribute undesirable noise. In the figures where the D^* limit for ideal photovoltaic detectors is shown, this refers to short-circuit photovoltaic detectors or reverse-biased photodiode detectors. If the photovoltaic detectors are operated in the open-circuit mode, the limiting $D^* \lambda$ is the same as that shown for ideal photoconductive detectors, since in this case the recombination of carriers in the open-circuit photovoltaic detector contributes noise just as in photoconduction.

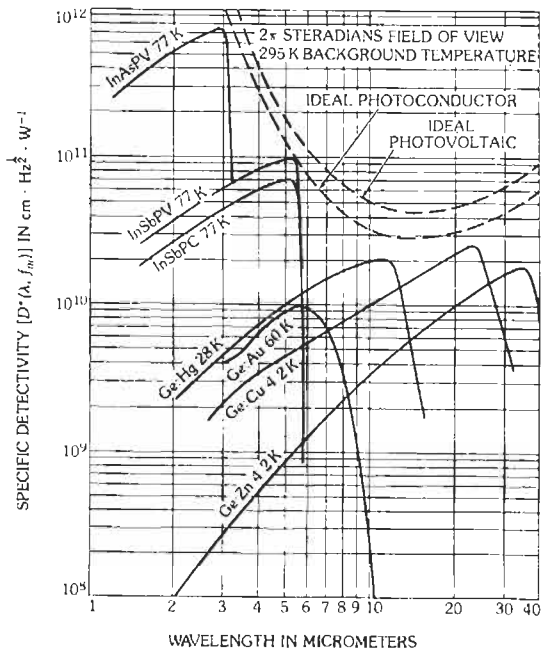


Fig. 26. Spectral detectivities of several semiconductor detectors. (Santa Barbara Research Center data.)

The major application for photovoltaic cells, however, is the conversion of solar radiation into electrical power. The operating specification of most interest in this application is the power conversion efficiency (ratio of the power produced by the cell to the power in the incident photon flux). Cell efficiency is affected by such factors as surface reflectance, absorption between the surface and the pn junction, and transmissivity of the junction region. Power efficiencies of commercial solar

† Stillman, G. E., and Wolfe, C. M., "Avalanche Photodiodes," in *Semiconductor and Semimetals*, Vol. 12 (New York: Academic Press, Inc., 1977; pp. 291-393).

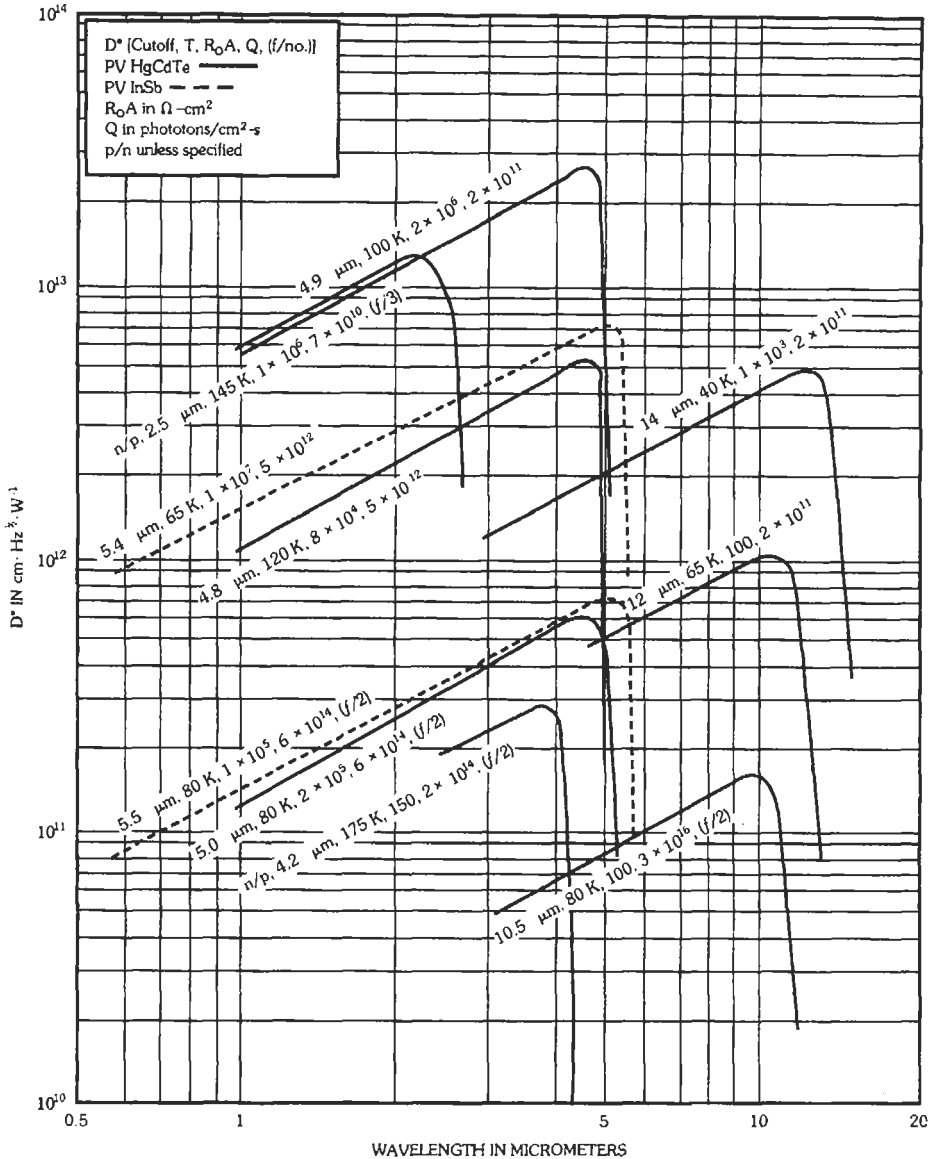


Fig. 27. Spectral response of detectors.

cells range from 5% to a maximum of about 15%. The best developed solar-cell material is silicon, but considerable work has been expended on gallium arsenide, indium phosphide, and cadmium sulfide. Fabrication advances now permit the production of long silicon ribbons, and this should eventually result in significant cost reductions for silicon solar-energy converters.

In the photoconductive mode, a photodiode (usually of special construction) is reverse biased and exposed to optical radiation. The resulting hole-electron pairs separate at the junction and establish a current that consists

of a small and constant reverse leakage (the dark current) and the photocurrent (the signal). The photocurrent is linear with respect to the incident radiation. The dark current establishes a minimum noise level below which operation is impractical. An important advantage of photoconductive operation over the photovoltaic mode is response time. A photodiode resembles a parallel-plate capacitor. Unbiased operation is slowed by high junction capacitance, whereas reverse-biased operation reduces junction capacitance and enhances response time.

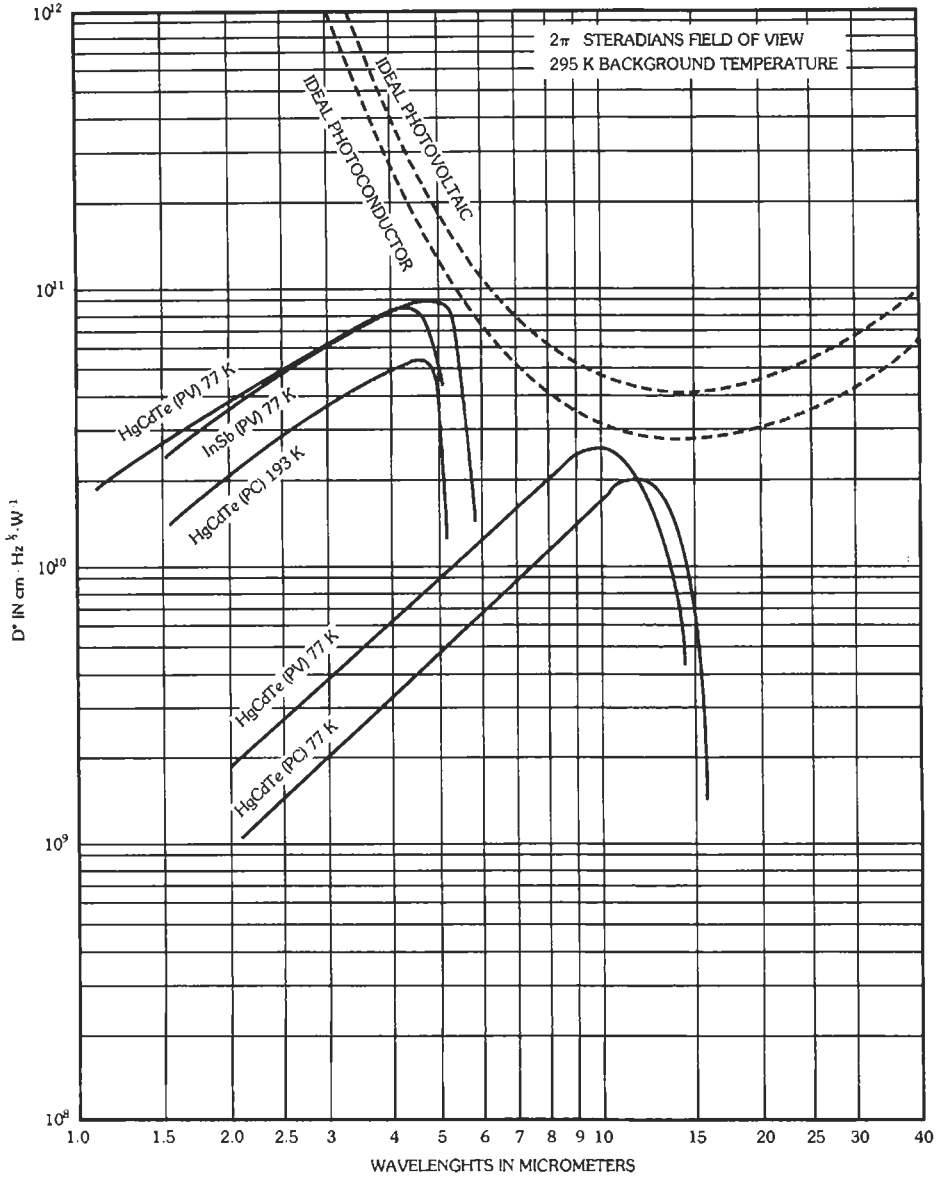


Fig. 28. Spectral response of detectors.

A photodiode of major importance is the avalanche photodiode. This device is a modified photodiode biased just below the avalanche breakdown region. Photons impinging upon the junction region create hole-electron pairs that initiate an avalanche multiplication of carriers across the junction. The resulting internal gain and response time make avalanche photodiodes suitable for operation in wide-bandwidth applications, where the s/n ratio is limited by the thermal noise of the load resistor and amplifier noise of wide-bandwidth amplifiers, formerly reserved for the photomultiplier tube.

Spectral Response of Semiconductor Detectors

A knowledge of detector spectral response and sensitivity or D^*_λ is essential in choosing a detector for a particular application. The spectral response of various detectors is summarized in Figs. 25 through 28. The operating temperatures can have a strong influence on the directivity of detectors, and the variation of D^* with temperature for several different detectors is shown in Fig. 29.

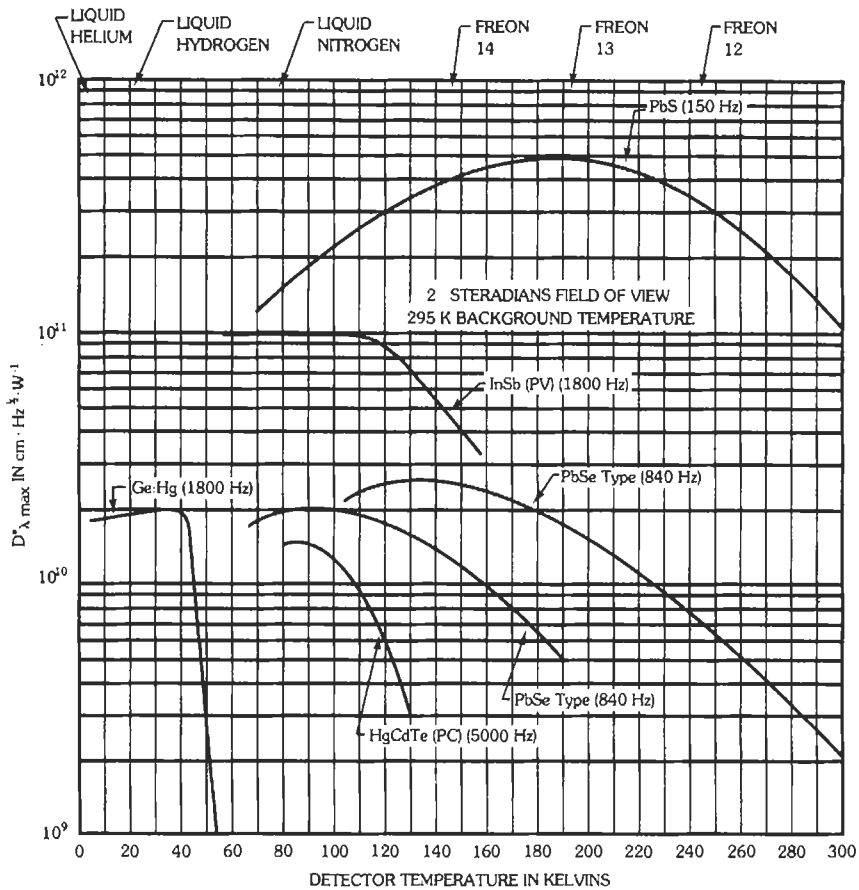


Fig. 29. Variation of $D^*_{\lambda_{max}}$ for several different photoconductive and photovoltaic semiconductor detectors with operating temperature for a 295 K 2π steradians background condition. The chopping frequency used for the measurements is shown in parentheses.

22

Optical Communications

Chester S. Gardner and George Papen

Optical Detectors and Noise	22-2
Modulation Techniques	22-3
Pulse-Position Modulation	
Binary Orthogonal Signaling	
Optical Transmitters	22-5
Incoherent Sources	
Coherent Sources	
Optical Receivers	22-8
Direct-Detection Receivers	
Heterodyne-Detection Receivers	
Background Radiation	22-12
Atmospheric Effects	22-13
Absorption	
Scattering	
Turbulence	
Fiber Waveguide Transmission	22-19
Terminology	
Losses	
Dispersion	
Fiber Bandwidth	
Splicing	
Connectors	
Sample Calculation	

In recent years, advances in device technology have significantly enhanced the capabilities of optical communication systems. High-data-rate systems have been constructed for space applications, and fiber-optic systems are finding wide use in telephone and data communications. In this chapter, the design of both fiber and free-space optical communication systems is discussed. Systems aspects are emphasized in this chapter. More detailed information on device technology can be found in Chapter 21 on Optoelectronics.

OPTICAL DETECTORS AND NOISE

Optical detectors can be classified into two general categories. Photomultiplier tubes (PMT) and vacuum photodiodes are photoemissive devices that emit electrons into a vacuum or gas. Avalanche photodiodes (APD), pin photodetectors, and phototransistors are solid-state devices in which the excited charge is transported in the solid by holes and electrons. Photomultiplier tubes, APDs, and pin photodetectors are the most widely used in optical communication systems.

Photomultiplier tubes are typically used when high sensitivity and bandwidth are required. Because their internal gain can be as high as 10^7 , the output signal level is large compared to the thermal noise of the succeeding amplifiers. Receivers employing PMTs are, therefore, shot-noise limited. The disadvantages of PMTs are that they are relatively bulky and fragile compared to solid-state detectors and they require high-voltage power supplies. Solid-state detectors are small, reliable, and inexpensive, and they exhibit high quantum efficiency through the visible to near-infrared portions of the spectrum. Avalanche photodiodes are characterized by high gain and bandwidth.

Optical detectors are square-law devices because they respond to intensity rather than amplitude. The detection process involves the interaction of incident photons with the detector material. In the case of vacuum photodiodes and PMTs, the incident photons generate photoelectrons that are ejected from the photocathode and flow to the anode. Electron-hole pairs are created in the lattice structure of solid-state detectors by the incident photons. The quantum efficiency (η) of a detector is the fraction of incident photons that produce a photoelectron or electron-hole pair. The quantum efficiencies of photocathodes in PMTs range from a few tenths of a percent to about 30 percent. Silicon photodetectors can have quantum efficiencies exceeding 80 percent. In either case, each detected photon causes an impulse of current to flow in the detector load resistor. The photocurrent is equal to the number of electrons emitted per second times the electron charge ($e = 1.6 \times 10^{-19}$ C). Photodetector efficiency is sometimes expressed in terms of responsivity, which is defined as the photocurrent divided by the incident optical power.

Responsivity (R) can be written in terms of the detector quantum efficiency, electron charge, and photon energy (hf)

$$R = \eta e/hf \tag{Eq. 1}$$

where,

h is Planck's constant (6.63×10^{-34} J · s),
 f is the optical frequency.

The detector output current can be modeled by using the following formula:

$$i(t) = \sum_{j=1}^{N(t)} G_j g(t - \tau_j) + i_{Th}(t) \tag{Eq. 2}$$

where,

- $N(t)$ = number of photons that have been detected during the time interval $(-\infty, t)$,
- τ_j = arrival time of the j th photon,
- G_j = detector gain for the j th photon,
- $g(t)$ = detector impulse response,
- $i_{Th}(t)$ = detector thermal noise current.

In the above expression, N is a random process whose statistics are related to the optical intensity. If the intensity is deterministic, N is Poisson distributed

$$\text{Prob}(N = k) = (\mu^k/k!)e^{-\mu} \tag{Eq. 3}$$

where,

$$\mu = (\eta/hf) \int_{-\infty}^t P(\tau) d\tau, \tag{Eq. 4}$$

$P(t)$ is the optical power.

The expected value and variance of N equal μ .

$$E(N) = \text{Var}(N) = (\eta/hf) \int_{-\infty}^t P(\tau) d\tau \tag{Eq. 5}$$

The internal amplification processes of PMTs and APDs are noisy so that G is also random. This effect introduces additional noise into the detection process. Amplification noise is most significant in APDs and should be considered in analyzing receiver performance. However, if the gain is assumed to be constant, the mean and variance of the detector output current are given by Campbell's theorem

$$E[i(t)] = G(\eta/hf)g(t) * P(t) \tag{Eq. 6}$$

$$\text{Var}[i(t)] = G^2(\eta/hf)g^2(t) * P(t) + \sigma_{Th}^2 \tag{Eq. 7}$$

The first term on the right-hand side of Eq. 7 is called shot noise. It is caused by the statistical nature of the

photon detection process. The term σ_{Th}^2 is the thermal noise variance, which is given by

$$\sigma_{Th}^2 = 4kT_L B/R_L \quad (\text{Eq. 8})$$

where,

- R_L = detector load resistance,
- k = Boltzmann's constant,
- T_L = equivalent temperature of the detector load resistor (kelvins),
- B = detector bandwidth.

Typically, the detector bandwidth will be larger than or equal to the bandwidth of $P(t)$. In this case, the convolutions in Eqs. 6 and 7 can be approximated by

$$g(t) * P(t) \cong P(t) \int_{-\infty}^{\infty} g(\tau) d\tau = eP(t) \quad (\text{Eq. 9})$$

$$g^2(t) * P(t) \cong P(t) \int_{-\infty}^{\infty} g^2(\tau) d\tau = 2e^2BP(t) \quad (\text{Eq. 10})$$

In deriving Eqs. 9 and 10, it was assumed that $g(t)$ is an ideal low-pass filter, with bandwidth B . Consequently, the mean and variance of the output current are given by

$$E[i(t)] = (\eta eG/hf) (P_s + P_b) \quad (\text{Eq. 11})$$

$$\text{Var}[i(t)] = (2\eta e^2 G^2 B/hf) (P_s + P_b) + 4kT_L B/R_L \quad (\text{Eq. 12})$$

where,

- P_s is the signal power,
- P_b is the optical background noise power.

For atmospheric systems operating during the daytime, P_b can be quite large. For fiber systems, P_b is negligible.

Eq. 12 is valid for PMTs and detectors with no internal gain ($G = 1$). For APDs, the random variation of the internal gain must be considered when calculating the variance of the output current. The details are complicated and need not be repeated here. The variance of the output current for an APD is given by*

$$\text{Var}[i(t)] = [2\eta e^2 G^2 F(G)B/hf] (P_s + P_b) + 4kT_L B/R_L \quad (\text{Eq. 13})$$

where,

$$F(G) = \rho G + (2 - 1/G)(1 - \rho), \quad (\text{Eq. 14})$$

- G is the mean gain,
- ρ is the ratio of hole collision ionization probability to electron collision ionization probability.

For silicon APDs, ρ is between 0.02 and 0.03, and G is typically less than 200. For noise-free avalanche gain, $F(G)$ would be equal to 1.

The current variance is proportional to the noise power at the detector output. The signal power is proportional to the square of the mean signal current. Therefore, the signal-to-noise power ratio (SNR) at the detector output is given by

$$\text{SNR} = \frac{[(\eta eG/hf)P_s]^2}{(2\eta e^2 G^2 B/hf) (P_s + P_b) + 4kT_L B/R_L} \quad (\text{Eq. 15})$$

Thermal-noise-limited detection occurs when the internal gain of the detector and the optical signal power are low. In this case, thermal noise dominates the detection process, and the SNR becomes

$$\text{SNR} = [(\eta eG/hf)P_s]^2 / (4kT_L B/R_L) \quad (\text{Eq. 16})$$

Thermal noise limited

If the internal gain of the detector is high so that detection is shot-noise limited, the SNR can be written as

$$\text{SNR} = (\eta/2hfB) [P_s^2 / (P_s + P_b)] \quad (\text{Eq. 17})$$

Shot noise limited

MODULATION TECHNIQUES

An optical carrier can be modulated with any of the conventional amplitude, phase, and frequency techniques employed at radio frequencies. However, intensity modulation and polarization modulation are the most widely used methods for optical systems. In digital communications, the polarization state of the optical carrier can be used to represent the value of a data bit. Bits can be represented as either right or left circular polarization or as any two orthogonal linear polarization states. Lithium niobate (LiNbO_3) and lithium tantalate (LiTaO_3) are two of the more commonly used electro-optic crystals for polarization modulators. Intensity modulation can be used for analog and digital signals. Analog signals are sometimes used to modulate RF subcarriers by conventional phase or frequency modulation. The optical intensity is then modulated by the RF subcarrier. Intensity modulation can be accomplished directly in laser diodes and LEDs by varying the diode forward current. External acousto-optic and electro-optic modulators can also be used to vary the optical intensity continuously.

Pulse-Position Modulation

Pulse-position modulation (PPM) is a very effective digital intensity-modulation technique. In M -ary PPM, each time slot is divided into M equal intervals. A single pulse of constant energy is transmitted during one of the M intervals. During a single time slot, one of M

* Reference 1.

different messages can be transmitted. If $M = 2^k$, this corresponds to a binary word of k bits. This modulation technique is a form of block coding. In general, the optimum detection strategy can be quite complicated. However, when the optical detector is shot-noise limited, the maximum-likelihood receiver reduces to a particularly simple form. For each time slot, the receiver counts the number of detected photons in each of the M intervals. The transmitted message is assumed to be the message corresponding to the interval with the largest count. The probability of word error (PWE) for this case is*

$$PWE = 1 - \{ \exp[-(N_s + MN_b)] / M - \sum_{k=1}^{\infty} [(N_s + N_b)^k / k!] \exp[-(N_s + N_b)] \times \left[\sum_{j=1}^{k-1} (N_b^j / j!) e^{-N_b} \right]^{M-1} (1/aM)[(1+a)^M - 1] \}$$

(Eq. 18)

where,

$$a = (N_b^k / k!) / \sum_{i=0}^{\infty} N_b^i / i!$$

(Eq. 19)

* Reference 3.

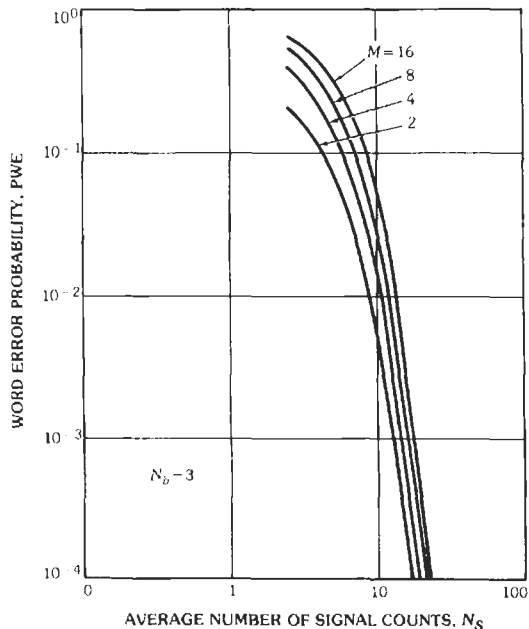


Fig. 1. Word error probability for block-coded pulse-position modulation. (From Gagliardi, R., and Karp, S. Optical Communications. New York: John Wiley & Sons, Inc., 1976, p. 264.)

$$N_s = (\eta/hf)J_s \tag{Eq. 20}$$

$$N_b = (\eta/hf)J_b \tag{Eq. 21}$$

In these expressions, N_s is the expected signal photon count; N_b is the expected background noise count per pulse interval; and J_s and J_b are, respectively, the signal and background noise energies per pulse interval. In Fig. 1, PWE is plotted versus N_s for several values of M . The plots indicate that PWE increases as M increases. It is misleading to use these plots to compare M -ary systems at different values of M . As M increases, more bits are transmitted per pulse. Also, for a fixed time-slot width, the data rate increases with M . Comparisons of M -ary systems should be made only on a bit-error-rate basis for a fixed data rate. In Fig. 2, the equivalent probability of bit error (PBE) is plotted versus N_s for several values of M . Notice that for N_s sufficiently large, the PBE continually decreases with increasing M . When N_s is above the crossover point, there is a significant advantage to increasing M . However, it should be noted that the peak power of the optical transmitter and the bandwidths of the modulator and receiver must also increase.

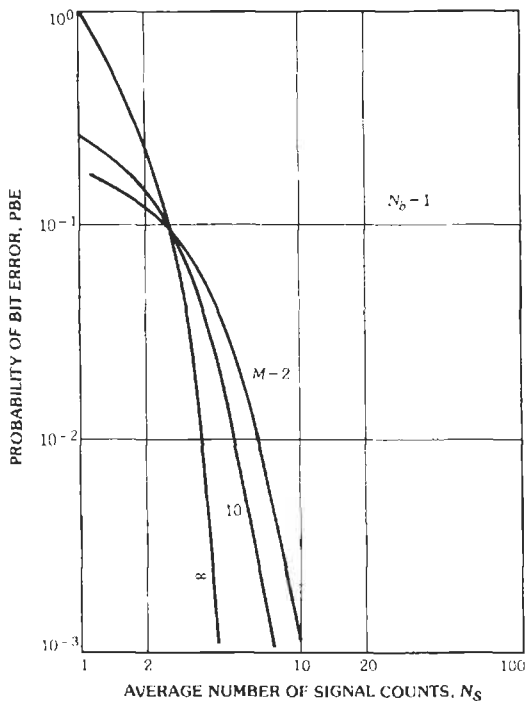


Fig. 2. Equivalent bit-error probability for block-coded pulse-position modulation. (From Gagliardi, R., and Karp, S. Optical Communications. New York: John Wiley & Sons, Inc., 1976, p. 271.)

Binary Orthogonal Signaling

When $M = 2$, we have binary modulation. In this case, the PBE is given by*

$$PBE = \frac{1}{2} \{1 + Q[\sqrt{2N_b}, \sqrt{2(N_s + N_b)}] - Q[\sqrt{2(N_s + N_b)}, \sqrt{2N_b}]\} \quad (Eq. 22)$$

where,

$$Q(a,b) = \int_b^\infty \exp[-(a^2 + x^2)/2] I_0(ax) x dx \quad (Eq. 23)$$

is Marcum's Q function. This expression for the PBE is also valid for other types of binary orthogonal signaling such as polarization modulation.† In Fig. 3, the PBE is plotted versus N_s for several values for N_b . For thermal-noise-limited detection, the PBE can be expressed in terms of the error function and the SNR:

$$PBE = \frac{1}{2} [1 - \text{erf}(\sqrt{\text{SNR}}/2)] \quad (Eq. 24)$$

* Reference 4.
† Reference 5.

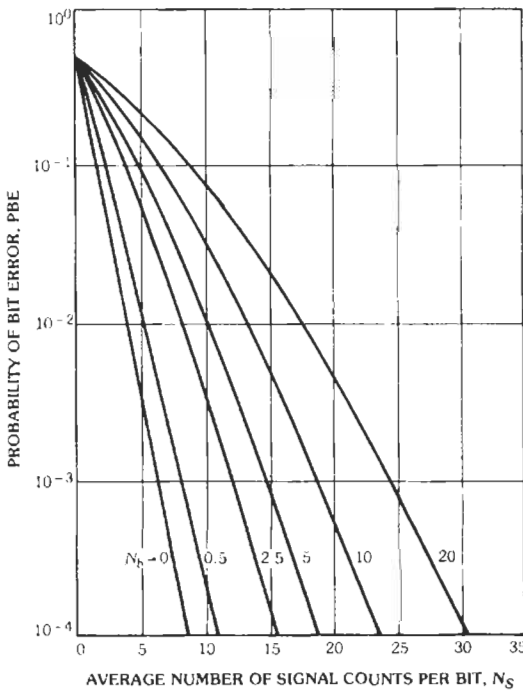


Fig. 3. Bit-error probability for shot-noise-limited binary orthogonal signaling. (From Pratt, W. K. Laser Communications Systems. New York: John Wiley & Sons, Inc., 1969, p. 209.)

where,

$$\text{SNR} = [G(eN_s/\tau_B)]^2 / (4kT_L/R_L\tau_B) \quad (Eq. 25)$$

τ_B is the pulse interval width.

Eq. 25 was obtained from Eq. 16 by replacing the detector bandwidth by $1/\tau_B$ and P_s by N_s/τ_B . The PBE for the thermal-noise-limited detection is plotted in Fig. 4. In general, the signal levels required to realize a given PBE are much higher for thermal-noise-limited detection than for shot-noise-limited detection. This is expected because, by definition, the thermal noise power is much larger than the shot noise power in a thermal-noise-limited receiver.

OPTICAL TRANSMITTERS

When an incoherent source such as an LED is used in an optical communication system, the transmitting optics are designed to image the source onto the receiving telescope. For systems employing lasers, a beam expander is often used for the transmitting tele-

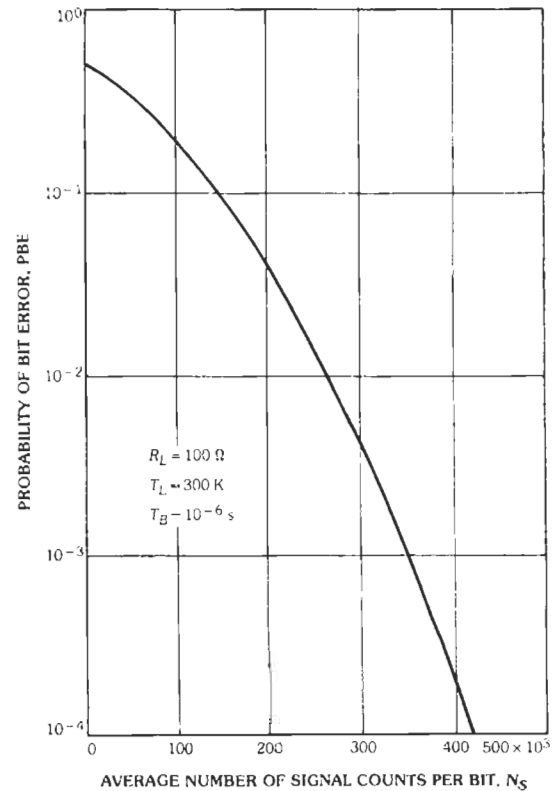


Fig. 4. Bit-error probability for thermal-noise-limited binary orthogonal signaling. (T_B = signal pulse width). (From Pratt, W. K. Laser Communications Systems. New York: John Wiley & Sons, Inc., 1969, p. 210.)

scope to decrease beam divergence. The transmitter and receiver can share the same telescope to reduce cost and weight. A variety of telescope configurations have been used in optical communication systems. The most common are the Cassegrain reflector and the simple refracting telescope. Because of the high cost of low-loss lens materials, reflecting optics are utilized almost exclusively in systems operating at the middle and far IR wavelengths. Both reflecting and refracting telescopes are used for near IR and visible systems.

Incoherent Sources

A simple imaging telescope for an LED source is diagrammed in Fig. 5. The objective lens is designed to project an image of the LED onto the receiver. The image intensity depends on the image magnification and the power radiated by the LED. The image magnification (*m*) is the ratio of the image distance (*z*) to the object distance (*Z_o*)

$$m = z/Z_o \tag{Eq. 26}$$

The image and object distances are related through the lens-maker's equation

$$1/f_L = 1/z + 1/Z_o \tag{Eq. 27}$$

where *f_L* is the lens focal length. By solving Eq. 27 for the object distance, we obtain

$$Z_o = f_L / (1 - f_L/z) \tag{Eq. 28}$$

Usually, the image distance is large compared to the lens focal length so that *Z_o* = *f_L* and *m* = *z/f_L*. The image diameter (*d_i*) is equal to the LED diameter (*d_{LED}*) multiplied by the magnification

$$d_i = (z/f_L)d_{LED} \tag{Eq. 29}$$

The total optical power in the image is equal to the power collected by the transmitting telescope objective multiplied by the atmospheric transmittance. To compute the power collected by the objective, we need to know the radiation pattern of the LED. If the LED has no lens, the radiation is Lambertian and is proportional

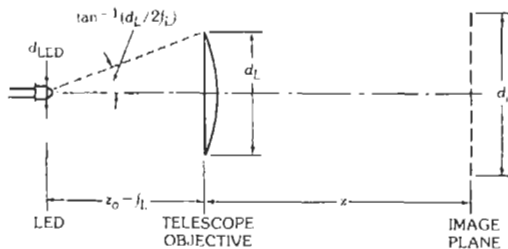


Fig. 5. Transmitting telescope for an LED source.

to the cosine of the angle between the observation point and the optical axis. Many LEDs are packaged with small lenses to collimate the radiation and reduce the beamwidth. In most cases, the radiation pattern can be fairly accurately modeled by the formula

$$P_L(\theta, \phi) = [(\xi + 1)/2\pi] P_L \cos^\xi \theta \tag{Eq. 30}$$

where,

$$\xi = -\ln 2 / \ln \cos(\theta_B/2) \tag{Eq. 31}$$

- θ* and *φ* are the spherical coordinate angles,
- P_L* is the total power emitted by the LED into the hemisphere,
- θ_B* is the full-width at half-maximum (FWHM) beamwidth.

For a Lambertian source, *θ_B* = 120° and *ξ* = 1. For LEDs, *θ_B* can be as small as 10° and *ξ* can approach 200.

The power collected by the objective lens is calculated by integrating the LED radiation pattern over the solid angle of the lens.

$$P_T = \int_0^{2\pi} \int_0^{\tan^{-1}(d_L/2f_L)} \{(\xi + 1)/2\pi\} P_L \cos^\xi \theta \cdot \sin \theta \, d\theta \, d\phi$$

$$= P_L [1 - (1 + d_L^2/4f_L^2)^{-(\xi+1)/2}] \tag{Eq. 32}$$

The quantity *P_T* is the total power transmitted by the telescope, and *d_L* is the diameter of the telescope objective lens. The maximum value of *P_T* is *P_L*, the total power emitted by the LED. The bracketed term is a factor describing the efficiency of the LED-telescope combination. The efficiency increases with decreasing beamwidth (increasing *ξ*) and decreasing *f/l* number (*f_L/d_L*) of the objective lens. This is illustrated in Fig. 6, where telescope efficiency is plotted versus the LED

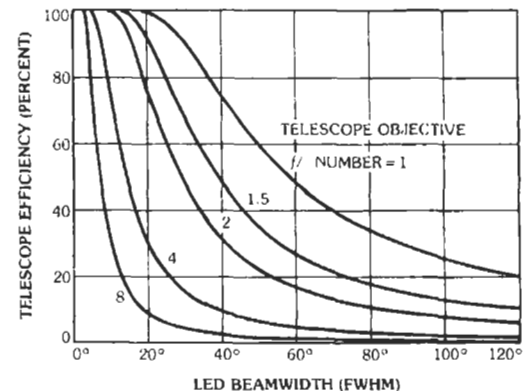


Fig. 6. Variation of the transmitting-telescope efficiency with LED beam width.

beamwidth for several values of the objective-lens f/l number. For a Lambertian source, Eq. 32 reduces to

$$P_T = P_L d_L^2 / (d_L^2 + 4f_L^2) \quad (\text{Eq. 33})$$

If we assume that the image is uniformly bright, then the optical power density or signal intensity (I_s) in the receiver plane is equal to the total power in the image divided by the image area. For a circular LED source, the signal intensity is given by

$$I_s = (4f_L^2 / \pi d_{LED}^2) (T_a P_T / z^2) \quad (\text{Eq. 34})$$

where T_a is the atmospheric transmittance. As expected, the signal intensity decreases inversely with the square of the distance between the transmitter and receiver. Although the telescope efficiency decreases with increasing f_L , the f_L^2 dependence in Eq. 34 dominates so that I_s increases with increasing f_L . However, this does not always mean that it is better to have a large focal length. The magnification and image size decrease with increasing f_L . Smaller images require more precise pointing and tracking. The maximum permitted pointing error (θ_E) is approximately equal to the arctangent of the image radius divided by the image distance

$$\theta_E = \tan^{-1}(d_i/2z) = \tan^{-1}(d_{LED}/f_L) \quad (\text{Eq. 35})$$

This is the condition for the edge of the image to lie on the optical axis of the receiving telescope. The transmitting telescope must be pointed with an accuracy less than the angle given by Eq. 35. Small changes in the pointing angle can be obtained by moving the LED perpendicular to the optical axis of the transmitting telescope.

In direct analogy with conventions established in microwave-communication systems analysis, it is sometimes convenient to express the signal power density in terms of the transmitter or source power and an effective antenna gain. In this case, Eq. 34 is written as

$$I_s = (G_T/4\pi)(T_a P_L/z^2) \quad (\text{Eq. 36})$$

where,

$$G_T = (16f_L^2/d_{LED}^2)[1 - (1 + d_L^2/4f_L^2)^{-(\xi+1)/2}] \quad (\text{Eq. 37})$$

Coherent Sources

For spatially coherent laser sources, the transmitting telescope is usually a beam expander that is designed to reduce the beam divergence. The spatial mode structure of laser resonators has been discussed extensively in the literature.* The fundamental mode for a radially symmetric cavity has a Gaussian cross section

$$I_L(x,y,z) = [2P_L/\pi\omega^2(z)] \exp[-2(x^2 + y^2)/\omega^2(z)] \quad (\text{Eq. 38})$$

where,

I_L is the beam intensity,
 P_L is the total power in the beam,
 $\omega(z)$ is the beamwidth at the e^{-2} intensity point.

In Eq. 38, the optical axis of the laser is assumed to be the z axis. The laser is located at the origin, and the beam is propagating in the positive z direction. The factor $\omega(z)$ is a function of the beamwidth and phase-front curvature at the laser. If the receiver is in the far field, then $\omega(z)$ can be written as

$$\omega(z) = z \tan \theta_L \quad (\text{Eq. 39})$$

where θ_L is the beam divergence angle. Usually, θ_L is a few milliradians or less. The smallest beam is obtained whenever the laser radiation is focused onto the receiver plane. If the receiver is in the far field, this is equivalent to collimating the beam. For a collimated beam, $\omega(z)$ is given by

$$\omega(z) = \omega_0(1 + \lambda^2 z^2 / \pi^2 \omega_0^4)^{1/2} \quad (\text{Eq. 40})$$

or

$$\omega(z) = \omega_0[1 + (\lambda z / \pi \omega_0^2)^2]^{1/2}$$

where,

λ is the optical wavelength,
 ω_0 is the beamwidth at the output of the beam expander.

The divergence angle for a collimated beam is

$$\theta_L = \lambda / \pi \omega_0 \quad (\text{Eq. 41})$$

In the far field, the on-axis signal intensity is calculated by substituting Eq. 40 into Eq. 38 and evaluating the result at $x = y = 0$:

$$I_s = (2\pi\omega_0^2/\lambda^2) (T_a P_L/z^2) \quad (\text{Eq. 42})$$

where T_a is the atmospheric transmittance. This expression is similar to Eq. 34 for incoherent sources. The beam divergence decreases and the on-axis intensity increases with increasing beamwidth, ω_0 .

It is sometimes convenient to express the on-axis intensity in terms of an equivalent antenna gain:

$$I_s = (G_T/4\pi) (T_a P_L/z^2) \quad (\text{Eq. 43})$$

where,

$$G_T = 8\pi^2 \omega_0^2 / \lambda^2 \quad (\text{Eq. 44})$$

Because a Gaussian beam is infinite in extent, Eq. 44 is the maximum gain that would be obtained with an infinitely large beam expander. However, in practice, the beam will be truncated by the finite extent of the transmitting aperture. In addition, if a Cassegrain

* Reference 6.

reflecting telescope is used to expand the beam, the central obscuration caused by the secondary mirror will also reduce the effective gain. It is necessary to match the beamwidth to the telescope to minimize the gain loss.* This is illustrated in Fig. 7, where the relative far-field axial gain is plotted versus the ratio of the aperture radius to beamwidth for several values of obscuration ratio ($\gamma = b/a$). The radii of the primary and secondary mirrors are denoted, respectively, by a and b . The gain is plotted relative to the gain for an unobscured uniformly illuminated aperture, $4\pi^2 a^2/\lambda^2$. A uniformly illuminated aperture has the maximum possible on-axis gain. Thus, the data in Fig. 7 represent the efficiency of the transmitting antenna. Truncation losses dominate when a/ω_0 is small, whereas obscuration losses dominate when a/ω_0 is large. Maximum efficiency is obtained whenever a/ω_0 satisfies the equality

$$a/\omega_0 = 1.12 - 1.30\gamma^2 + 2.12\gamma^4 \quad (\text{Eq. 45})$$

The maximum efficiency is plotted versus obscuration ratio in Fig. 8. The loss in efficiency over an ideal uniformly illuminated aperture increases from approxi-

* Reference 7.

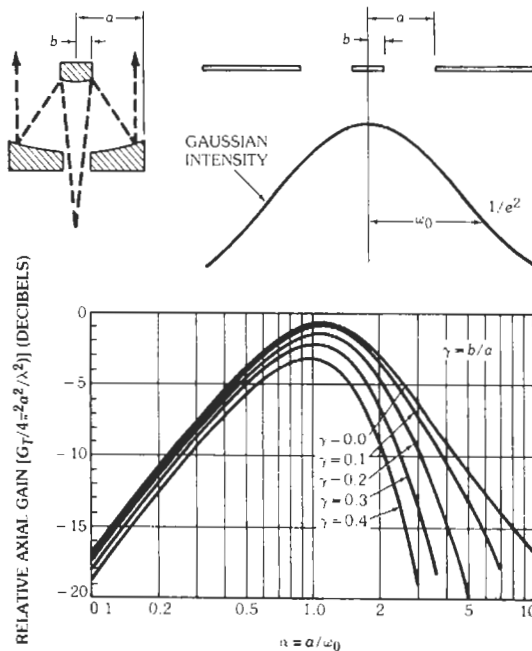


Fig. 7. Far-field axial gain of a centrally obscured transmitting telescope as a function of the ratio of the telescope radius to the laser beam width. (From Klein, B. J., Degnan, J. J. "Optical Antenna Gain. 1. Transmitting Antennas," Appl. Opt., 13, September 1974, p. 2137.)

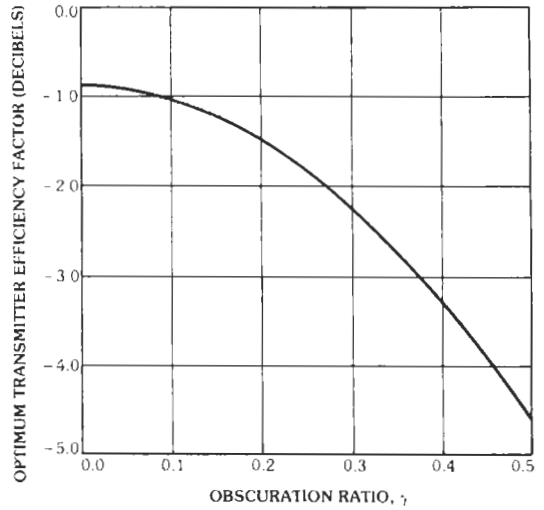


Fig. 8. Optimum telescope efficiency relative to an unobscured uniformly illuminated aperture ($4\pi^2 a^2/\lambda^2$) as a function of obscuration ratio. (From Klein, B. J., and Degnan, J. J. "Optical Antenna Gain. 1. Transmitting Antennas," Appl. Opt., 13, September 1974, p. 2139.)

mately 1 dB to almost 5 dB as the obscuration ratio increases from 0 to 0.5. As a consequence, these effects must be taken into account, particularly for systems designed with a low gain margin. In some cases, the spider support structure for the secondary mirror can also have a measurable effect on the axial gain.* However, the effects are usually significant only when the obscuration ratio is small ($\gamma < 0.2$).

The pointing and tracking requirements become more severe as the beamwidth increases and beam divergence decreases. Small pointing errors can result in substantial loss of signal. This is illustrated in Fig. 9, where the loss due to transmitter pointing error is plotted versus pointing error for several values of the obscuration ratio. These curves correspond to the optimum antenna configuration where Eq. 45 is satisfied. The loss is approximately 10 dB when the pointing error is $\lambda/2a$ radians.

Because of the geometry of the emitting stripe in a semiconductor laser diode, the output beam has an elliptical cross section that is highly asymmetric. Therefore, the transmitting telescope must be designed to expand the minor axis of the beam more than the major axis to produce a symmetric cross section. This can be accomplished by using an anamorphic prism pair or cylindrical telescope.

OPTICAL RECEIVERS

Receiving telescopes are designed to focus the optical signal onto the photodetector and to reject as much of

* Reference 8.

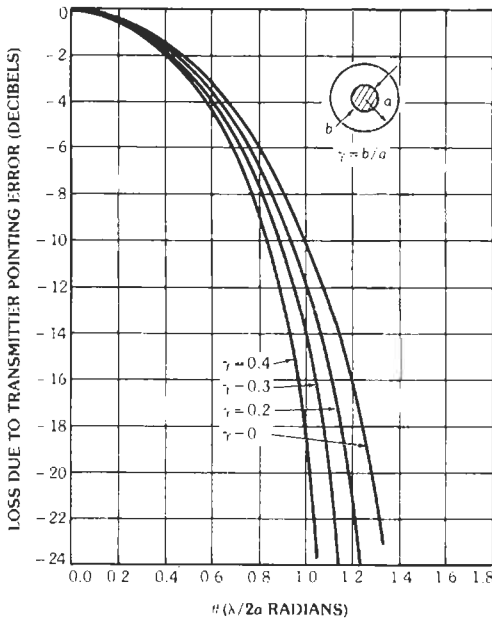


Fig. 9. The decibel loss due to transmitter pointing error as a function of the angular pointing error for optimum antenna configuration (maximum far-field gain) for several values of the obscuration ratio. (From Klein, B. J. and Degnan, J. J. "Optical Antenna Gain. I. Transmitting Antennas." Appl. Opt., 13, September 1974, p. 2139.)

the background radiation as is practical. Interference filters are employed to eliminate background radiation that is not the same wavelength as the optical signal. Field stops are used to reject radiation that is not emitted from the region surrounding the transmitting telescope. Direct-detection receivers respond to the signal intensity and are the most widely used in communication systems. In heterodyne-detection receivers, the optical signal is combined with a local-oscillator beam, and then both signals are focused onto the same detector. Heterodyne-detection receivers respond to signal amplitude and are used primarily in the far infrared to overcome limitations imposed by poor detector sensitivity.

Direct-Detection Receivers

A simple direct-detection receiver is diagrammed in Fig. 10. Because of cost, the objective is usually a reflector such as a Cassegrain system whenever the aperture is greater than about 10 cm. The field-stop iris limits the field of view (FOV) of the telescope to a small region surrounding the transmitter. The FOV should be as small as possible to minimize background noise. However, because the pointing and tracking constraints are more severe with a smaller FOV, there is a trade-off. Even for point-to-point communication systems where the locations of the transmitter and receiver are fixed,

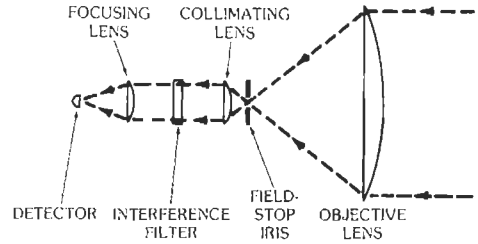


Fig. 10. Simple direct-detection receiver.

atmospheric turbulence and mechanical vibrations can cause the apparent position of the transmitter to change randomly. In these cases, the receiver FOV will have to be adjusted to prevent inadvertent loss of signal.

The planar (θ_R) and solid-angle (Ω_R) fields of view can be calculated from simple geometric considerations. If f_R and d_I denote, respectively, the focal length of the receiver objective and the field-stop iris diameter, then we have

$$\theta_R = \tan^{-1}(d_I/2f_R) \tag{Eq. 46}$$

$$\Omega_R = 4\pi \sin^2(\theta_R/2) \tag{Eq. 47}$$

These expressions for the FOV are useful for calculating the amount of background noise received by the telescope.

Interference filters are multilayered thin-film devices. They are constructed with dielectric and metallic layers and can have bandwidths as small as 1 Å. Usually, interference filters are designed to work with collimated light. If the optical signal is not adequately collimated, the peak transmission of the filter is reduced, and the center wavelength is shifted to a shorter wavelength. The center wavelength can also be shifted to shorter wavelengths by tilting the filter from normal to oblique incidence. For small angles, the shift is without distortion of the passband or reduction of the peak transmittance. Consequently, a receiver can be tuned over a limited wavelength range by tilting the filter. In terms of the angle of incidence, ϕ , the center wavelength at small angles of incidence ($\phi < 5^\circ$) is given by the equation

$$\lambda = \lambda_{\max}[1 - (n_o/n_e)^2 \sin^2 \phi]^{1/2} \tag{Eq. 48}$$

where,

- n_o is the refractive index of air,
- n_e is the effective refractive index of the filter spacer.

Typical values for n_e range from about 1.4 to 2.1.

The center wavelength is also temperature dependent because the layer thicknesses and refractive indices change with temperature. The center wavelength shifts to longer wavelengths as the temperature increases. The thermal coefficient of wavelength shift is a function of

wavelength and filter construction. For visible and near IR wavelengths, the coefficient varies from approximately 0.1 to 0.3 Å/°C. To maintain proper tuning, it is usually necessary to mount the filter in a temperature-controlled oven for bandwidths below 10 Å FWHM. The peak transmittance of interference filters is a decreasing function of the bandwidth. Typical values range from over 70% for bandwidths on the order of 100 Å to less than 10% for bandwidths on the order of 1 Å.

The received signal power is calculated by multiplying the signal power density times the receiver area (A_R) and efficiency (η_R):

$$P_S = \eta_R A_R (G_T / 4\pi) (T_a P_L / z^2) \quad (\text{Eq. 49})$$

where,

- P_L is LED or laser power,
- G_T is the transmitting antenna gain,
- T_a is the atmospheric transmittance.

Efficiency η_R includes the effects of reflective and absorptive losses in the receiver optics including the interference filter. Low-loss optical materials and anti-reflection coatings can reduce losses to a few percent or less for each lens or mirror. Because losses are wavelength dependent, the optics should be coated for the specific wavelength of interest to obtain optimum performance.

Equation 49 can be expressed in terms of an equivalent receiving antenna gain (G_R):

$$P_S = (\lambda / 4\pi z)^2 T_a \eta_T \eta_R G_T G_R P_L \quad (\text{Eq. 50})$$

where,

$$G_R = 4\pi A_R / \lambda^2 \quad (\text{Eq. 51})$$

A transmitting-antenna efficiency factor (η_T) has been included in Eq. 50 to account for losses in the transmitter optics. The transmitting antenna gain is given by Eq. 37 for LED sources and by Eq. 44 for laser sources with a Gaussian cross section. The factor $(\lambda / 4\pi z)^2$ is called the free-space loss and represents the effect of path length on the signal strength. The SNR for a direct-detection receiver is calculated by substituting Eq. 50 into Eq. 15. The background noise power can be computed from the data presented in the section on background radiation.

Heterodyne-Detection Receivers

A simple heterodyne-detection receiver telescope is diagrammed in Fig. 11. The local oscillator (LO) beam is combined with the optical signal, and both are focused onto the same detector. Because the objective and collimating lenses demagnify the signal beam, the receiving telescope illustrated in Fig. 11 is optically equivalent to combining the signal and LO in front of

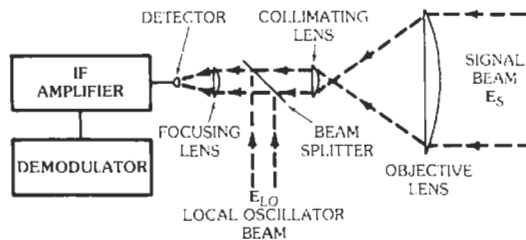


Fig. 11. Simple heterodyne-detection receiver.

the objective lens. Optical mixing of the signal and LO beams produces an intermediate-frequency (IF) signal in the detector output that is proportional to the signal amplitude. The IF is the difference between the signal and LO frequencies.

The electric field vectors for the signal (E_S) and LO (E_{LO}) can be written in the form

$$E_S = A_S \cos(2\pi f_s t + \phi_s) \quad (\text{Eq. 52})$$

$$E_{LO} = A_{LO} \cos(2\pi f_{LO} t + \phi_{LO}) \quad (\text{Eq. 53})$$

where,

- A_S and A_{LO} are the vector amplitudes,
- f_s and f_{LO} are the optical frequencies,
- ϕ_s and ϕ_{LO} are the phases of the signal and LO, respectively.

The amplitudes and phases of both beams can depend on time and position. The detector responds to the total incident power, which is calculated by integrating over the receiver aperture the magnitude squared of the total field:

$$P_D = \int_{\text{Receiver Aperture}} |\overline{E_S + E_{LO}}|^2 dx dy = P_S + P_{LO} + S_{IF} \quad (\text{Eq. 54})$$

The bar in Eq. 54 represents a time average over an interval that is short compared to the inverse of the detector bandwidth but is long compared to the optical periods of the signal and LO. Terms P_S and P_{LO} are, respectively, the total signal and LO powers focused onto the detector by the telescope. Term S_{IF} is the IF signal, which is given by

$$S_{IF} = \int_{\text{Receiver Aperture}} A_S \cdot A_{LO} \cos(2\pi f_{IF} t + \phi_s - \phi_{LO}) dx dy \quad (\text{Eq. 55})$$

where,

$$f_{IF} = f_s - f_{LO} \quad (\text{Eq. 56})$$

The receiver sensitivity and S_{IF} are maximum when the polarization, amplitude distribution, and phase of the LO beam are matched to the signal beam.* If the receiver is in the far field and the effects of atmospheric turbulence are neglected, the signal beam can be approximated as a uniform plane wave. Then for uniform LO and signal beams, A_S and A_{LO} are constant throughout the receiver aperture and can be moved outside the integral sign in Eq. 55. To obtain maximum sensitivity, the signal and LO beams must be properly aligned. Small pointing errors can significantly reduce the IF signal strength. To illustrate, the phase fronts of the signal beam are assumed to be tilted with respect to the LO. The LO phase is constant, and the signal phase is

$$\phi_S = \phi_0 + (2\pi/\lambda)y \sin \theta_R \quad (\text{Eq. 57})$$

For simplicity, the coordinate system was chosen so that the tilt (θ_R) is with respect to the y -axis. Angle θ_R may be regarded as the angular pointing error of the receiver.

After substituting Eq. 57 into Eq. 55 and carrying out the integration for an unobscured circular aperture, we obtain

$$S_{IF} = \frac{2J_1[(2\pi a/\lambda) \sin \theta_R]}{(2\pi a/\lambda) \sin \theta_R} A_R A_S \cdot A_{LO} \cos(2\pi f_{IF} t + \phi_0 - \phi_{LO}) \quad (\text{Eq. 58})$$

where,

- A_R is the receiver aperture area (πa^2),
- a is the aperture radius,
- J_1 is a Bessel function of the first RESET TO ADD ONE LINEkind.

The factor involving the Bessel function is a gain-reduction factor related to pointing error. When θ_R is zero (i.e., perfect LO and signal beam alignment), this factor is one. The IF signal is proportional to the LO amplitude. By increasing the LO power, the IF signal can be made much larger than the receiver thermal noise and the signal and background shot noises. Typically, the limiting noise for heterodyne detection is shot noise contributed by the LO. There is an IF signal contribution due to background noise. However, for typical IF filter bandwidths (B_{IF} less than approximately 1 GHz), this background-noise contribution is small compared to the IF laser signal and can be neglected.

The detector output is processed by an IF amplifier and filter and final demodulator. The expected current at the IF-amplifier output due to the optical signal is calculated by substituting Eq. 58 into Eq. 11 and is

$$E[i_{IF}(t)] = (\eta e G/hf) S_{IF} \quad (\text{Eq. 59})$$

If the polarization of the LO is matched to the signal, the IF signal power is given by

$$P_{IF} = 2(\eta e G/hf)^2 \left(\frac{2J_1(2\pi a/\lambda) \sin \theta_R}{(2\pi a/\lambda) \sin \theta_R} \right)^2 P_{LO} P_S \quad (\text{Eq. 60})$$

where,

$$P_S = A_R A_S^2/2 \quad (\text{Eq. 61})$$

$$P_{LO} = A_R A_{LO}^2/2 \quad (\text{Eq. 62})$$

The noise power is calculated from Eq. 12 with B replaced by the IF bandwidth (B_{IF}) and P_{LO} added to the signal and background shot-noise power. Therefore, when the pointing error is zero, the IF SNR is given by

$$\text{SNR}_{IF} = \frac{2(\eta e G/hf_s)^2 P_{LO} P_S}{(2\eta e^2 G^2 B_{IF}/hf_s)(P_S + P_{LO} + P_b) + 4k T_L B_{IF}/R_L} \quad (\text{Eq. 63})$$

If the LO power is large, the thermal noise power and the signal plus background shot-noise power will be negligible compared to the LO shot-noise power. In this case, the SNR reduces to

$$\text{SNR}_{IF} = \eta P_S/hf_s B_{IF} \quad (\text{Eq. 64})$$

The signal power can be calculated from Eq. 50, which was derived for a direct-detection receiver. The background noise power can be computed from the data presented in the section on background radiation.

The pointing error can have a significant effect on the IF SNR. This is illustrated in Fig. 12, where the IF SNR loss in decibels is plotted versus normalized pointing error for an unobscured circular aperture and uniform LO beam. To prevent significant loss, the pointing error must be very small. The condition for negligible loss is

$$\theta_R \ll \lambda/2\pi a \quad (\text{Eq. 65})$$

Equation 65 is the condition for pointing within the diffraction-limited FOV of the telescope. Signals that lie outside the diffraction-limited FOV will be rejected by the telescope. Consequently, the pointing and tracking requirements for heterodyne-detection receivers are quite severe.

In practice, atmospheric turbulence will cause the amplitude and phase of the signal beam to fluctuate randomly in space and time. The effect reduces the IF signal power, particularly when the receiving aperture is large. The IF signal strength is maximum when the LO and signal beams add coherently across the aperture.

* Reference 3.

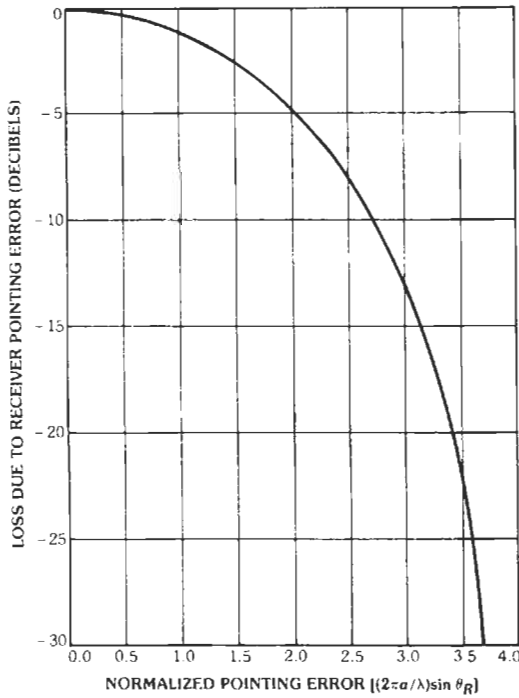


Fig. 12. Intermediate-frequency signal power loss due to receiver pointing error.

Turbulence reduces the spatial coherence of the signal beam and causes the signal phase to vary randomly across the aperture. As the diameter of the aperture increases, the additional signal may add incoherently or out of phase with the LO beam, resulting in either no increase or a reduction in the SNR. This phenomenon is discussed in more detail in the section on atmospheric effects.

BACKGROUND RADIATION

During the day, the dominant background-noise source is the sun, whereas at night, radiation from man-made sources, the moon, the stars, and the planets becomes important. Background radiation may enter the receiving telescope directly, when the source falls within the receiver FOV, or indirectly, when the radiation is reflected or scattered. The spectral distribution of the background radiation depends on many factors, including the spectral distribution of the source, absorption characteristics of the atmosphere, and reflection characteristics of the earth and objects near the transmitter. Uniform background sources are conveniently described by their spectral radiance, $N(\lambda)$, which is defined as the power radiated at wavelength λ per unit bandwidth into a unit solid angle per unit

source area.* The total background noise power collected by the receiving telescope is

$$P_b = \begin{cases} N(\lambda) \Delta \lambda A_R \Omega_{\text{FOV}} & \text{if } \Omega_{\text{FOV}} \in \Omega_S \\ N(\lambda) \Delta \lambda A_R \Omega_S & \text{if } \Omega_S \in \Omega_{\text{FOV}} \end{cases} \quad (\text{Eq. 66})$$

where,

$\Delta \lambda$ is the optical bandwidth of the receiver,
 A_R is the receiver area,
 Ω_{FOV} is the receiver solid-angle FOV,
 Ω_S is the solid angle subtended by the source when viewed from the receiver.

When $\Omega_{\text{FOV}} \in \Omega_S$, the source completely fills the receiver FOV, and only a fraction of the total power within the receiver optical bandwidth collected by the aperture is focused onto the detector. This is the case for sky background. When $\Omega_S \in \Omega_{\text{FOV}}$, the complete image of the source is focused onto the detector. In this case, Ω_S must be known to evaluate Eq. 66. This is not particularly convenient for small sources such as stars and planets. Consequently, the radiation from point sources is usually described in terms of the spectral irradiance, $H(\lambda)$, which is the power per unit bandwidth per unit receiver area.* The spectral irradiance is related to spectral radiance by

$$H(\lambda) = \Omega_S N(\lambda) \quad (\text{Eq. 67})$$

Expressed in terms of the spectral irradiance, the received background noise power is

$$P_b = H(\lambda) \Delta \lambda A_R \quad (\text{Eq. 68})$$

Often the spectral radiance of a background source can be approximated by a blackbody curve

$$N(\lambda) = (2hc^2/\lambda^5) [\exp(hc/\lambda kT_b) - 1]^{-1} \quad (\text{Eq. 69})$$

where,

h is Planck's constant (6.63×10^{-34} J · s),
 k is Boltzmann's constant (1.38×10^{-23} J/K),
 T_b is the source temperature,
 c is the velocity of light.

Most earth objects have temperatures between 200 K and 300 K so that their spectral radiance peaks within the wavelength region 9 to 15 μm . The moon has an equivalent blackbody temperature of 373 K. Radiation from the sun and stars peaks at the near IR and visible wavelengths. The equivalent blackbody temperature of the sun is 5900 K.

* References 2 and 4.

Absorption in the atmosphere of the earth can significantly alter the spectral radiance of a background source. This is illustrated in Fig. 13, where the spectral irradiance of the sun outside the atmosphere of the earth and at sea level is plotted. Also noted in the figure are the atmospheric constituents that are responsible for some of the major absorption bands. Fig. 14 illustrates the measured spectral radiance of the clear daytime sky. (The zenith angle of the sun is 45°, and the visibility is excellent.) The radiance for sunlit clouds is approximately one order of magnitude larger. On a clear day, the color temperature of the sky is approximately 20 000 K to 25 000 K.

The spectral irradiances of planets and stars also approximate blackbody radiation curves. Fig. 15 shows the calculated spectral irradiances outside the atmosphere of the earth from the planets. It is likely that many stars of widely differing temperatures will be within the receiver FOV. The probable spectral irradiance from a one-square-degree star field near the galactic plane is plotted in Fig. 16.

If the background radiation is due to the reflection of the sun and sky, knowledge of the reflectance of the

terrain surrounding the transmitter is required. In Fig. 17, the typical reflectances of water, snow, soil, and vegetation are plotted versus wavelength.

ATMOSPHERIC EFFECTS

Absorption, scattering, and turbulence affect any optical signal that propagates through the atmosphere of the earth. Absorption by minor constituents such as water vapor, carbon dioxide, and ozone can significantly attenuate the signal beam. Rayleigh scattering by air molecules and Mic scattering by larger particles such as aerosols, dust, and clouds also attenuate the signal. Turbulence is caused by small temperature fluctuations in the atmosphere that give rise to random variations in the refractive index. The optical signal is scattered by these refractive inhomogeneities. Portions of the signal that are scattered by different inhomogeneities interfere at the receiver aperture. The amplitude and phase of the total signal vary randomly in space and time. This loss of signal coherence results in reduced gain and fading in both direct-detection and heterodyne-detection receivers.

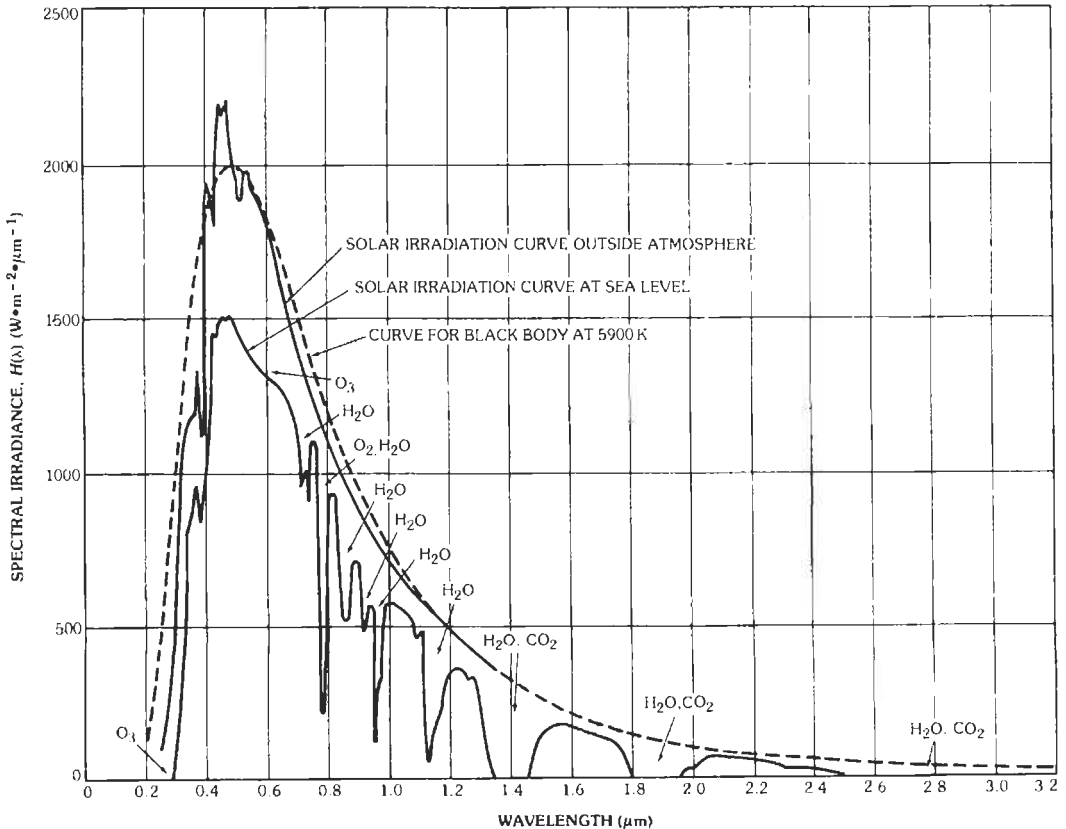


Fig. 13. Solar spectral irradiance with sun at zenith. (From Gait, P. R. "Solar Spectral Irradiance," in Handbook of Geophysical and Space Environment, Sec. 16.1, S. L. Valley, ed. Cambridge, Mass., 1965)

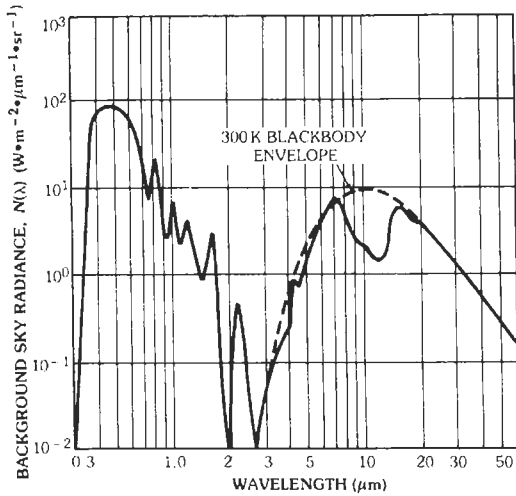


Fig. 14. Diffuse component of the typical background radiance from the sky viewed at sea level. (From Pratt, W. K. Laser Communications Systems. New York: John Wiley & Sons, Inc., 1968, p. 121.)

The transmittance of the atmosphere (T_a) over a path length z for radiation of wavelength λ is given by Bouguer's law (references 9 and 10)

$$T_a = \exp \left[- \int_0^z \alpha(\lambda, r) dr \right] \quad (\text{Eq. 70})$$

where $\alpha(\lambda, r)$ is the attenuation or extinction coefficient and the integral is taken over the propagation path. The attenuation coefficient has units of inverse meters. The integral of $\alpha(\lambda, r)$ in Eq. 70 is called the optical thickness of the propagation path. Optical thickness is dimensionless. Because the atmosphere attenuates through absorption and scattering by both gases and particles, the attenuation coefficient can be written as the sum of four terms

$$\alpha = \alpha_{g,a} + \alpha_{g,s} + \alpha_{p,a} + \alpha_{p,s} \quad (\text{Eq. 71})$$

where the subscripts $g, p, a,$ and s denote, respectively, gases, particles, absorption, and scattering. Absorption by gases and scattering by particles are the most significant at optical wavelengths.

Absorption

The main absorbing gases in the atmosphere for ultraviolet, visible, and infrared wavelengths in order of importance are $\text{H}_2\text{O}, \text{CO}_2, \text{O}_3, \text{N}_2\text{O}, \text{CO}, \text{O}_2, \text{CH}_4,$ and N_2 . The variation of the molecular absorption coefficient is very complicated due to the presence of

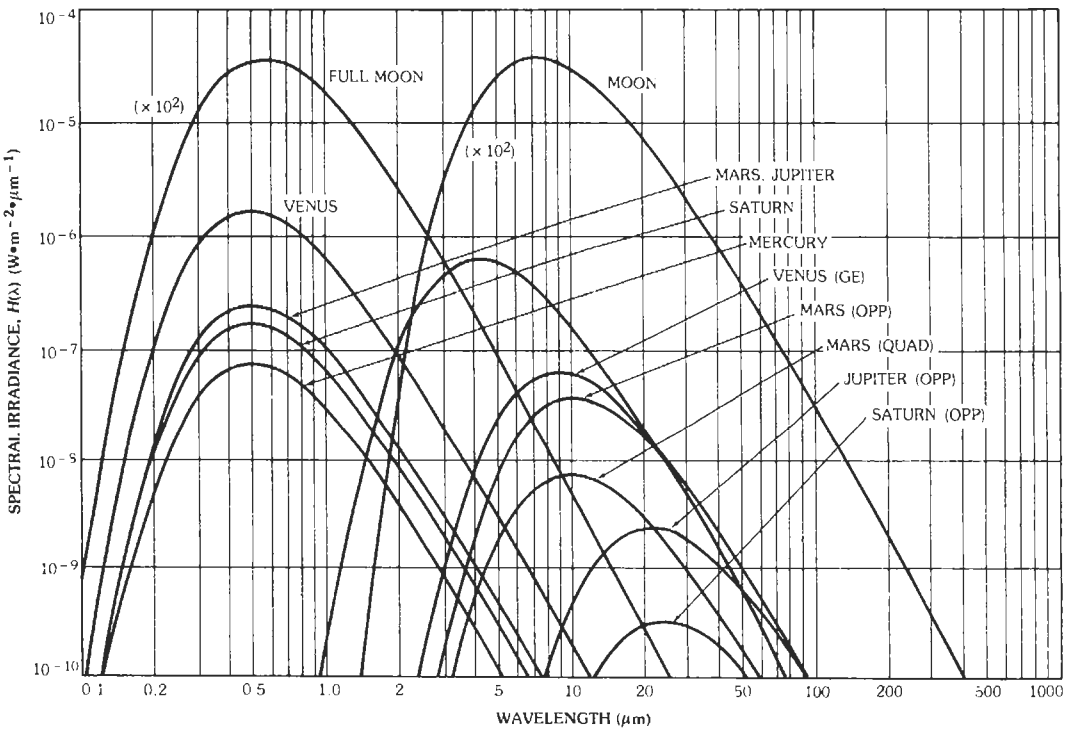


Fig. 15. Calculated planetary and lunar spectral irradiance outside the terrestrial atmosphere. (From Pratt, W. K. Laser Communications Systems. New York: John Wiley & Sons, Inc., 1968, p. 123.)

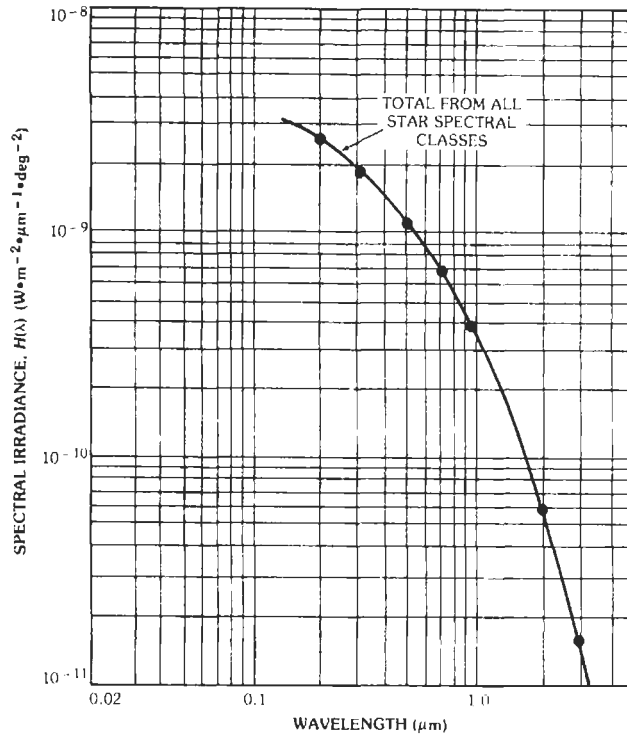


Fig. 16. Probable spectral irradiance from a one-square-degree star field near the galactic plane. (From RCA Electro-Optics Handbook. RCA Commercial Engineering, 1974, p. 69.)

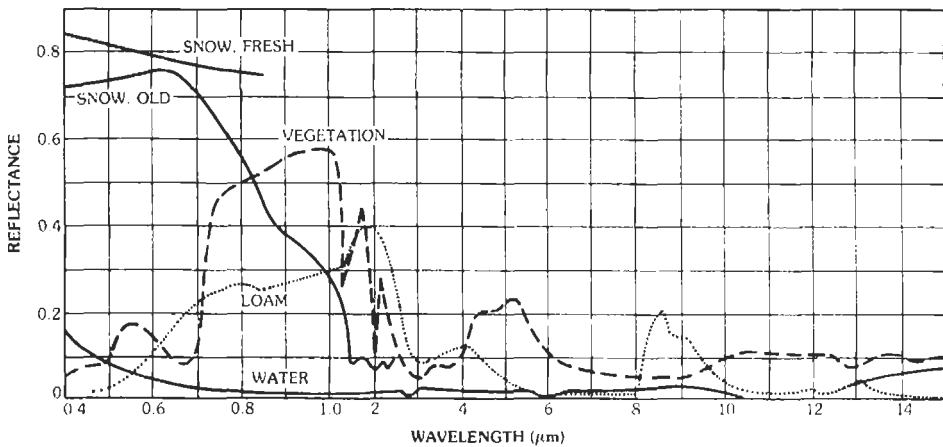


Fig. 17. Typical reflectance of water surface, snow, dry soil, and vegetation. (From Handbook of Optics. New York: McGraw-Hill Book Co., 1978, pp. 14-54.)

various absorption band complexes. Strong absorption lines occur most frequently at ultraviolet ($\lambda < 0.3 \mu\text{m}$) and infrared ($\lambda > 0.9 \mu\text{m}$) wavelengths, where they can severely limit the effective range of an optical communication system. Computer programs have been developed to calculate the detailed absorption spectrum for various

atmospheric models.* Attenuation coefficients have also been tabulated for some of the more common laser wavelengths.† Fig. 18 shows the spectral transmittance

* Reference 11.
 † Reference 12.

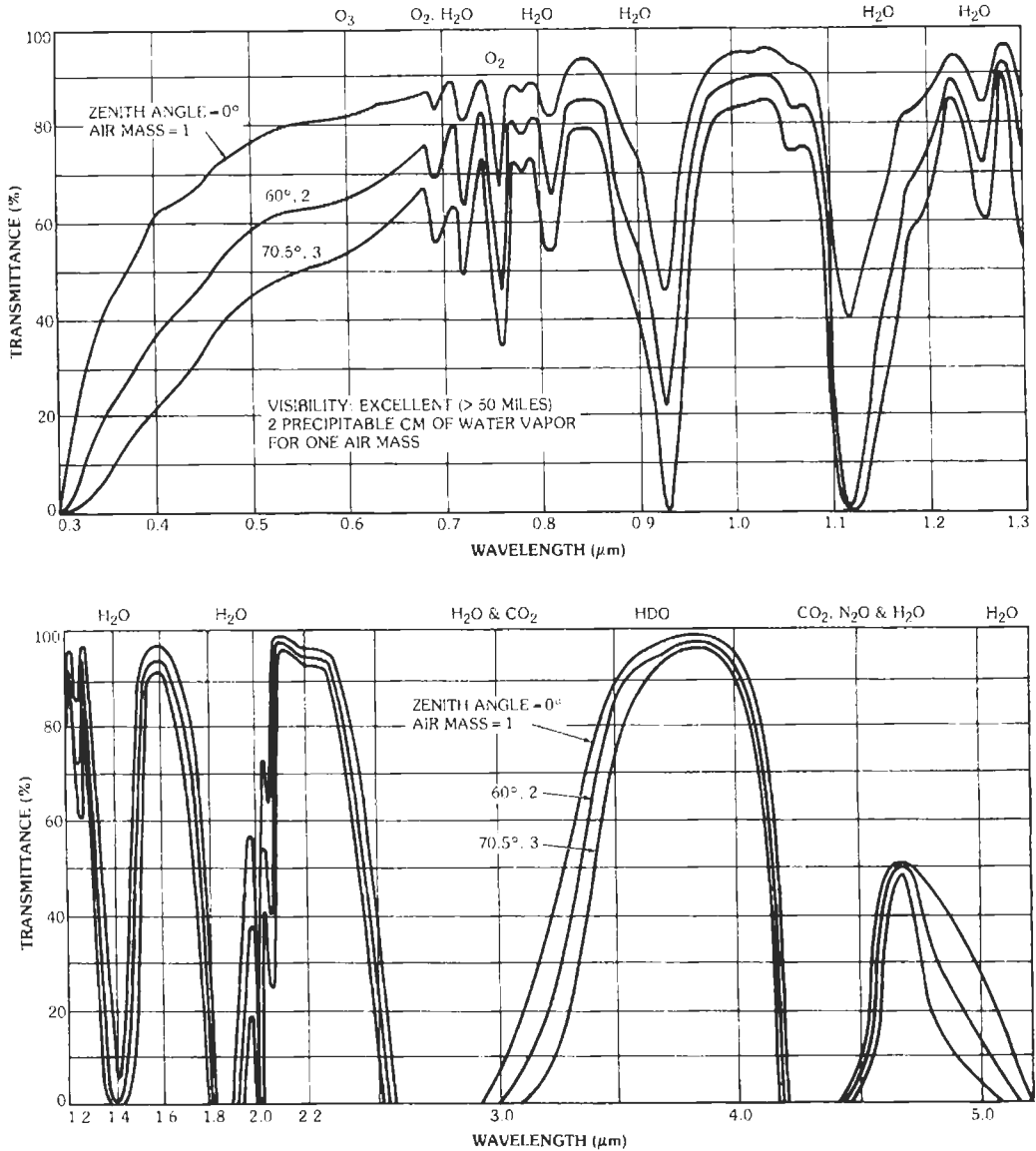


Fig. 18. Spectral transmittance from the atmosphere of the earth for varying optical air masses. (From RCA Electro-Optics Handbook. RCA Commercial Engineering, 1974, p. 83.)

through the entire atmosphere from sea level to space for several zenith angles. Fig. 19 shows the transmittance over a 1000-foot horizontal path at sea level between 0.5 μm and 25 μm. These figures illustrate the fact that certain wavelengths are strongly attenuated even for relatively short propagation path lengths.

Scattering

The attenuation coefficient due to scattering by gas molecules is dominated by the elastic component.

Because the size of the gas molecules is small compared to the optical wavelength, the attenuation coefficient can be expressed in terms of the Rayleigh cross section σ_R

$$\alpha_{g,s} = N_g \sigma_R \quad (\text{Eq. 72})$$

where N_g is the number density of the gas molecules. For optical wavelengths, σ_R at altitudes below 100 km is given by*

* Reference 10.

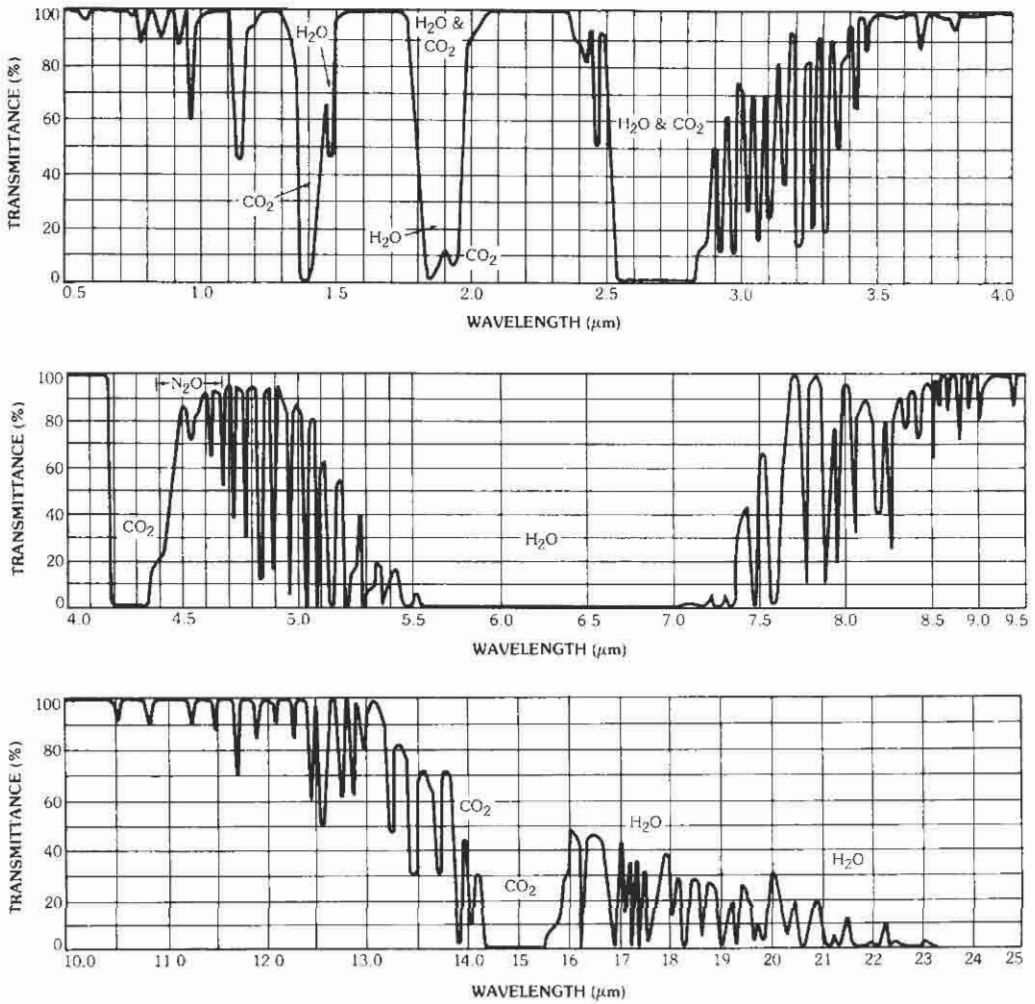


Fig. 19. Transmittance of 1000-ft horizontal air path at sea level containing 5.7 mm precipitable water at 79 °F. (From RCA Electro-Optics Handbook. RCA Commercial Engineering, 1974, p. 84.)

$$\sigma_R = 4.59[\lambda (\mu\text{m})/0.55]^{-4} \times 10^{-27} \text{ cm}^2 \quad (\text{Eq. 73})$$

Extensive tabulations of the Rayleigh attenuation coefficient based on several model atmospheres are included in reference 12.

Particulate and aerosol scattering depend on many factors, including the size distribution, shape, and composition (refractive index) of the scatterers. Fig. 20 shows the wavelength variation of the attenuation coefficient ($\alpha_{p,s}$) for various atmospheric conditions. At visible wavelengths, the attenuation coefficient is approximately*

$$\alpha(\lambda) \approx (3.912/V_M)(0.55 \mu\text{m}/\lambda)^{0.585} V_M^{1/3} \quad (\text{Eq. 74})$$

* Reference 10.

where V_M (km) is the meteorological visibility. Because precipitation particles are large compared to the optical wavelength, the corresponding attenuation coefficient does not depend on wavelength. Also, the dependence on the microstructure of the particles is negligible compared with the dependence on precipitation intensity. The attenuation coefficient for rains is approximately*

$$\alpha \approx 0.21 r^{0.74} \text{ km}^{-1} \quad (\text{Eq. 75})$$

where r is the rainfall rate (mm/hr).

Turbulence

Scattering by atmospheric turbulence reduces the spatial and temporal coherence of the signal beam. The

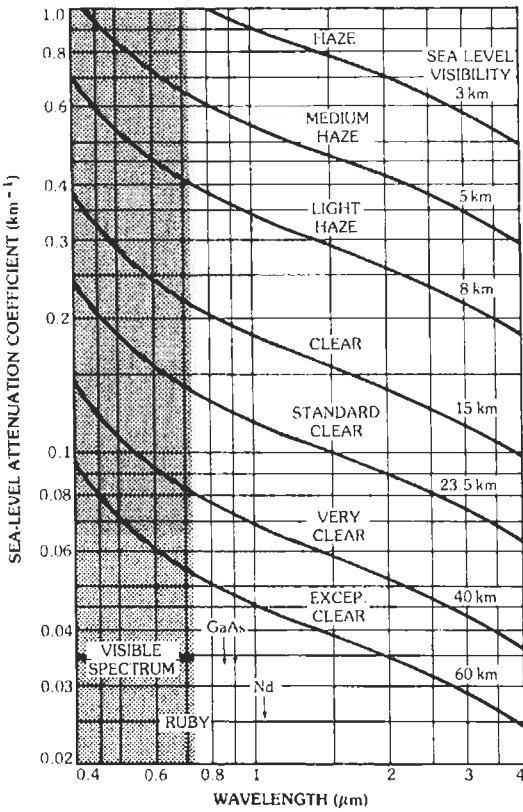


Fig. 20. Approximate variation of attenuation coefficient with wavelength at sea level for various atmospheric conditions. Absorption by water vapor and carbon dioxide is neglected. (From RCA Electro-Optics Handbook, RCA Commercial Engineering, 1974, p. 89.)

effect causes fading or scintillation in the signal power and a reduction in the SNR in heterodyne-detection receivers. In weak turbulence, the probability density of the intensity fluctuations is very nearly log-normal*

$$p(I) = [1/(2\pi)^{1/2} I \sigma] \exp \{-[\ln(I/\langle I \rangle) + (1/2) \sigma^2]^2 / 2\sigma^2\} \quad (\text{Eq. 76})$$

where,

$$\sigma^2 = \ln(1 + \sigma_I^2 / \langle I \rangle^2) \quad (\text{Eq. 77})$$

σ_I^2 is the intensity variance,
 $\langle I \rangle$ is the mean intensity.

In very strong turbulence, the intensity probability density approaches the exponential distribution

$$p(I) = (I/\langle I \rangle) \exp(-I/\langle I \rangle) \quad (\text{Eq. 78})$$

* Reference 13.

The intensity variance is a function of wavelength, path length, and turbulence strength. For weak turbulence, the intensity variance is given by

$$\sigma_I^2 / \langle I \rangle^2 = \sigma_1^2 \quad (\text{Eq. 79})$$

where,

$$\sigma_1^2 = 1.23(2\pi/\lambda)^{7/6} C_n^2 z^{11/16} \quad (\text{Eq. 80})$$

C_n^2 is the refractive index structure parameter,
 z is the path length.

The value of C_n^2 varies with altitude, meteorological conditions, and time of day. Typically, C_n^2 is maximum near noon, reaching values on the order of 10^{-14} to $10^{-13} \text{ m}^{-2/3}$. At night, the value ranges between 10^{-16} and $10^{-15} \text{ m}^{-2/3}$. Eq. 79 is valid for values of σ_1^2 up to approximately 0.3. In stronger turbulence, the normalized intensity variance is less than σ_1^2 . For very strong turbulence, $\sigma_I^2 / \langle I \rangle^2$ approaches a value near one.

Eqs. 76 through 79 describe the statistics of the signal-intensity fluctuations at a point. They can also be used to characterize the signal-power fluctuations for a direct-detection receiver if the aperture diameter is smaller than the intensity correlation length, ρ_c . The intensity correlation length is given approximately by

$$\rho_c \approx \begin{cases} (\lambda z)^{1/2} & \sigma_1^2 \approx 0.3 \\ 0.36(\lambda z)^{1/2} / (\sigma_1^2)^{3/5} & \sigma_1^2 \gg 1 \end{cases} \quad (\text{Eq. 81})$$

Because of aperture averaging, power fluctuations are reduced in direct-detection receivers if the aperture area is large compared to the intensity coherence area ($\pi\rho_c^2$). This is illustrated in Fig. 21, where the ratio of the signal-power variance for a receiver of diameter $2a$

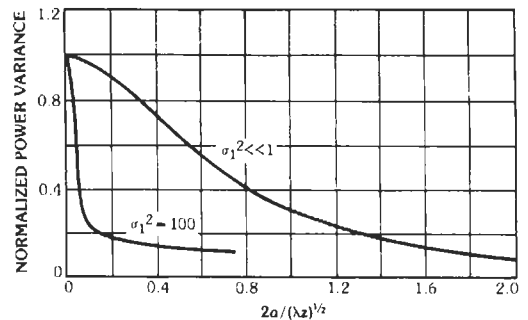


Fig. 21. Aperture averaging in direct-detection receivers for strong and weak turbulence. The received power variance is normalized to the variance for a point aperture; a = aperture radius. (From Fante, R. L. "Electromagnetic Beam Propagation in Turbulent Media," Proc. IEEE, 63, December 1975, p. 1680.)

to the variance for a point aperture is plotted versus aperture diameter. In both strong and weak turbulence, there is a significant reduction in the power fluctuations whenever the aperture diameter exceeds ρ_c .

Temporal fluctuations of signal intensity are caused by the movement and breakup of refractive inhomogeneities within the propagation path. The intensity temporal power spectrum has a cutoff frequency (f_c) that is related to the wind velocity normal to the propagation path (V_{\perp}) and the intensity correlation length:

$$f_c \approx V_{\perp} / \rho_c \tag{Eq. 82}$$

For weak turbulence, f_c is typically 10 to 100 Hz, whereas for strong turbulence, the cutoff frequency is approximately an order of magnitude higher. The inverse of f_c is the intensity coherence time and can be used to estimate the fade time.

Heterodyne-detection receivers are most affected by the signal phase fluctuations. Areas of the signal beam that are out of phase with the LO will add destructively at the detector and reduce the IF signal amplitude. This results in a reduction of the IF SNR. For plane-wave sources that are propagating horizontally, the phase correlation length (r_0) is given by the equation*

$$r_0 = [2.36 (2\pi/\lambda)^2 C_n^2 z]^{-3/5} \tag{Eq. 83}$$

The loss of phase coherence reduces the effective area of the receiver aperture. This is illustrated in Fig. 22, where the effective aperture area is plotted versus aperture diameter ($2a$). The effective area never exceeds the phase coherence area (πr_0^2). Because r_0 can be as small as 10 to 20 cm, this effect can severely limit receiver performance. Spatial diversity can be used to overcome some of these problems.†

* Reference 14.

† Reference 15.

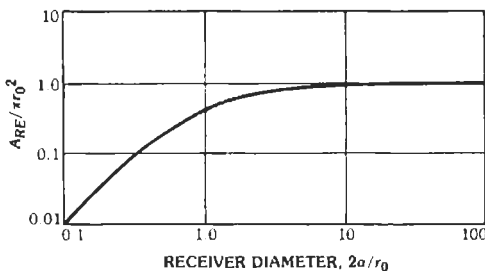


Fig. 22. Effective heterodyne receiver area A_{RE} ; a is the aperture radius, and r_0 is the phase correlation length of received signal. (From Gagliardi, R., and Karp, S. Optical Communications. New York: John Wiley & Sons, Inc., 1976, p. 199.)

FIBER WAVEGUIDE TRANSMISSION

In designing fiber-optic communication systems, there are two principle characteristics of optical fibers, attenuation and dispersion, that determine how much information can be sent and how far it can be sent. All current fiber-optic systems are either dispersion or attenuation limited. These characteristics, coupled with the transmitter and receiver characteristics, allow one to specify the type of fiber, transmitter, and receiver required for a particular application. Detailed design of fiber-optic systems is covered in specific fiber-optic handbooks. A representative sampling is given in the references.

Terminology

In order to specify the type of components required for a particular application, it is useful to be familiar with some of the important terminology and theory that is used in fiber optics.

There are two types of optical fiber predominately in use: single-mode step-index fiber and graded-index fiber. Single-mode step-index fiber is primarily used in high-bandwidth long-haul telecommunication systems. Graded-index fiber is primarily used in short-haul low-bandwidth systems such as local area networks (LANs).

The two types of fiber differ in how they guide the light. Single-mode step-index fiber has a small core region of a high-index glass (n_1) surrounded by a cladding region of lower index (n_2) (Fig. 23A). The normalized difference in the indices $\Delta = (n_1 - n_2)/n_1$ is usually small (0.001–0.02). The light is guided because it experiences total internal reflection at the core-cladding interface. (See “Snell’s Law” below.)

Graded-index fiber differs from step-index fiber in that it has a continuously varying index-of-refraction profile from the center of the core to the cladding region (Fig. 23B). The functional form of the index profile is typically quadratic and is chosen to compensate for modal dispersion. In Fig. 23B, all the rays, regardless of the launch angle, take the same time to travel the length of the fiber.

Fig. 23C shows a multimode step-index fiber. Rays with different launch angles travel at different effective velocities, producing modal dispersion.

Snell’s Law—For large core diameters, the light-guiding properties of optical fiber can be understood by the use of Snell’s law and ray optics. Snell’s law (refer to Fig. 24A) can be derived from the boundary conditions at the dielectric interface and is given by

$$n_1 \sin \theta_1 = n_2 \sin \theta_2 \tag{Eq. 84}$$

* The choice of n_1 or n_2 in the denominator varies.

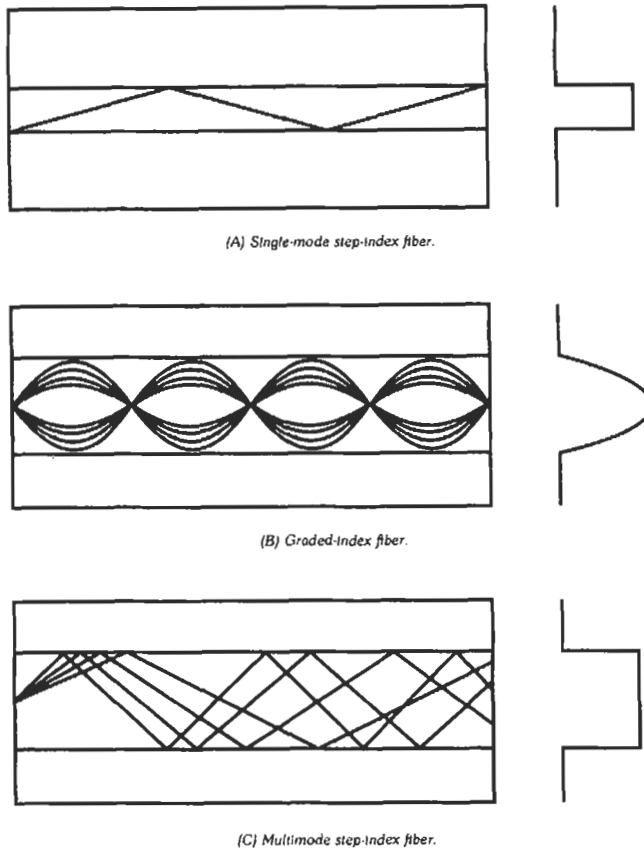


Fig. 23. Types of optical fiber. In each case, the index-of-refraction profile is shown on the right.

where the angles are measured from the normal to the interface. For light traveling from a higher-index material (n_1) to a lower-index material (n_2), the ray is refracted away from the normal. The critical angle θ_c is the incident angle where the refracted ray travels along the interface. When this occurs, $\sin \theta_2 = \sin(\pi/2) = 1$ (Fig. 24B), and Snell's Law yields

$$\theta_c = \sin^{-1}(n_2/n_1) \quad (\text{Eq. 85})$$

At angles greater than the critical angle, the ray is totally internally reflected (Fig. 24C) and is guided along the core.

Numerical Aperture—The numerical aperture of a fiber is a measure of the light-collecting ability of the fiber. For a large-core step-index fiber, it is the maximum angle θ_m that a ray can make with the axis of the fiber and still experience total internal reflection (Fig. 25). From the geometry in Fig. 25,

$$\sin \theta_m = \sqrt{n_1^2 - n_2^2} \approx n_1 \sqrt{2\Delta} \quad (\text{Eq. 86a})$$

$$\text{NA} = \sqrt{n_1^2 - n_2^2} \approx n_1 \sqrt{2\Delta} \quad (\text{Eq. 86b})$$

where $\Delta = (n_1 - n_2)/n_2$. The value of the numerical aperture changes slightly from the calculated values in Eqs. 86 for single-mode step-index fiber because a ray analysis is no longer strictly valid and diffractive effects have to be considered.

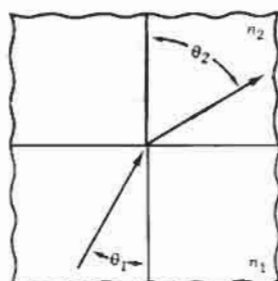
The numerical aperture of a graded-index fiber is defined slightly differently and is a function of the distance from the center of the core

$$\text{NA}(r) = \sqrt{n^2(r) - n_2^2} \quad (\text{Eq. 87})$$

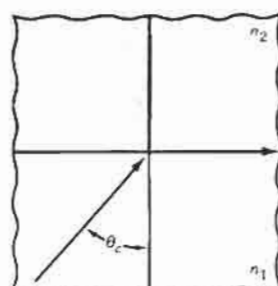
where $n(r)$ is the index profile (typically quadratic). The numerical aperture is now a function of r and is maximum for a thin beam along the axis. Values of the numerical aperture vary from ≈ 0.13 for single-mode step-index fiber to ≈ 0.3 for large-core graded-index fiber.

Group Index and Group Delay—The index of refraction n of a plane wave $\exp[-j(\omega t - \beta z)]$ in a medium is given by

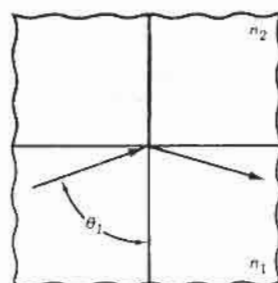
$$n = \beta/k \quad (\text{Eq. 88})$$



(A) Angle of Incidence is less than critical angle



(B) Angle of incidence equals critical angle



(C) Total internal reflection; incident angle exceeds critical angle

Fig. 24. Snell's law of reflection.

where,

- ω is the temporal frequency,
- β is the propagation constant,
- $k = \omega/c$,
- c is the speed of light in free space.

The group index is the speed at which a pulse travels in a medium. For a pulse with a temporal width τ much less than the inverse of the optical frequency ($\tau \ll 1/\omega$), the group index N is given by

$$N = \partial\beta/\partial k = \partial(nk)/k = n + k(dn/dk) = n + \omega(dn/d\omega) = n [-\lambda(dn/d\lambda)] \quad (\text{Eq. 89})$$

The total time for a pulse to travel a distance L is given by

$$t = L/v_g = NL/c \quad (\text{Eq. 90})$$

and the differential delay δt of the pulse is given by

$$\delta t = L/v_g = (L/c) \delta N \quad (\text{Eq. 91})$$

These formulas are useful for dispersion calculations.

V-Number—Most commercial single-mode step-index fiber has a small index contrast Δ . Solutions to the wave equation for these "weakly guiding" geometries yield modes that are commonly referred to as LP modes. An important parameter that characterizes these solutions is the V number defined as

$$V = k \cdot a \cdot \text{NA} \quad (\text{Eq. 92})$$

where,

- $k = \omega/c$,
- a is the core radius,
- NA is the numerical aperture ($\text{NA} = \sqrt{n_1^2 - n_2^2} = n_1 \sqrt{2\Delta}$).

Note that V depends on λ through k . By using weakly guiding theory,* it can be shown that for $V < 2.40483$ (first zero of $J_0(x)$), only one mode will propagate. The dependence of β on V for the LP_{01} mode is shown in Fig. 26.

Losses

The dominant loss mechanisms in silica fiber are absorption (ultraviolet, infrared, and impurity) and

* Reference 16.

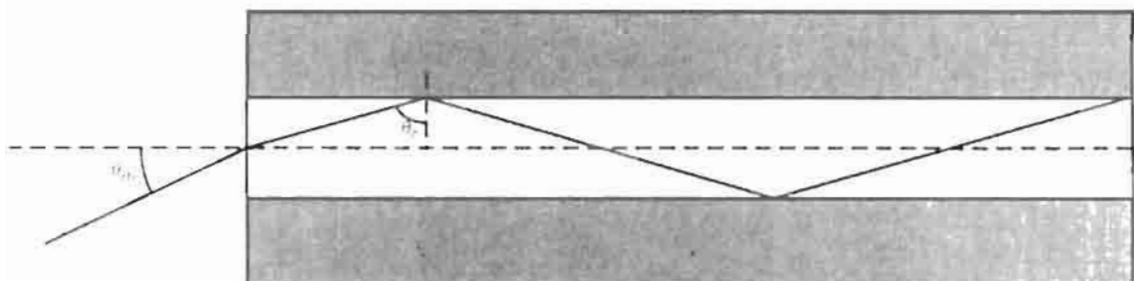


Fig. 25. Rays being reflected or refracted at the core-clad boundary.

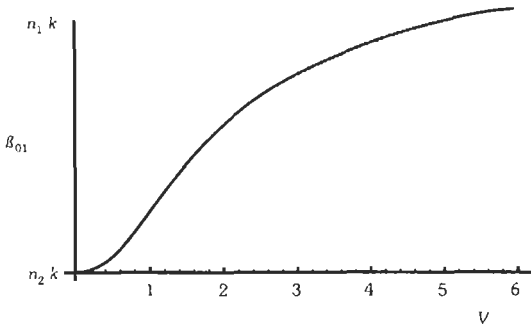


Fig. 26. The dependence of the propagation constant β on the V parameter for the LP_{01} mode in a single-mode fiber.

Rayleigh scattering. Ultraviolet absorption and Rayleigh scattering limit short-wavelength operation, and infrared absorption is dominant for wavelengths above about $1.6 \mu\text{m}$. Impurity absorption from residual OH effects adds losses between 1.3 and $1.5 \mu\text{m}$. For commercial fiber, it is Rayleigh scattering, which is proportional to $1/\lambda^4$, and infrared absorption that set the limits on losses within the wavelength regions of interest. Modern optical fiber operates at very close to these theoretical limits. The attenuation minimum for silica fiber occurs at about $1.55 \mu\text{m}$ and is $< 0.2 \text{ dB/km}$. The attenuation minimum at $1.55 \mu\text{m}$ is the reason why this wavelength is used in long-haul systems. Typical absorption curves for various fibers are shown in Fig. 27.

Dispersion

Dispersion is the wavelength or modal dependence of propagation constant β . In fiber optic systems, one is generally interested in pulse propagation through the fiber. In general, the pulse can be decomposed into modes, and for each mode the pulse can be further decomposed into wavelength components. Each of the modes and each wavelength component within each mode will travel at a slightly different velocity. The result is that the pulse broadens as it travels down the fiber, and thus the amount of information that can be transmitted is reduced.

Modal Dispersion—The dependence of the propagation constant on the mode number is called modal dispersion. It can be eliminated by using single-mode step-index fiber. However, in graded-index fiber that supports many modes, an incident pulse will excite several modes. Each mode will travel at a different group velocity. It can be shown* that a quadratic index profile minimizes the modal dispersion and that the modal dispersion is given by

$$\delta t = (Ln_1/c)(\Delta^2/4) \tag{Eq. 93}$$

where,

* Reference 17.

δt is the temporal spread of the pulse,
 L is the total distance traveled,
 $\Delta = (n_1 - n_2)/n_1$,
 n_1 is the peak index at the center of the core.

Material Dispersion—Material dispersion is caused by the wavelength dependence of the index of refraction. It affects pulse broadening because each wavelength component of an initial pulse travels at a different group velocity.

If the initial pulse has a spectral width $\delta\lambda$, then if we use differentials $\delta\lambda/\delta N = d\lambda/dN$. The differential time δt for a pulse to travel a length L is $\delta t = \delta NL/c$. Using these relations, we can write the material response time as

$$\delta t = |D_\lambda| \delta\lambda L \tag{Eq. 94}$$

where D_λ is the material dispersion coefficient given by

$$D_\lambda = (-\lambda/c)(d^2n/d\lambda^2) \tag{Eq. 95}$$

The units of the material dispersion coefficient are usually $\text{ps/nm} \cdot \text{km}$. For silica fiber, the term $d^2n/d\lambda^2$ crosses zero near $\lambda \approx 1.3 \mu\text{m}$. This is the zero material dispersion wavelength and is one of the reasons that $1.3\text{-}\mu\text{m}$ light is used in fiber optic systems.

As can be seen by examining Eq. 95, the effect of material dispersion can also be minimized by using a transmitter with a very narrow wavelength spectrum so that all of the wavelength components lie within a narrow band. This is one of the reasons why narrow-spectral-width lasers (typically distributed feedback lasers) are used in long-haul systems.

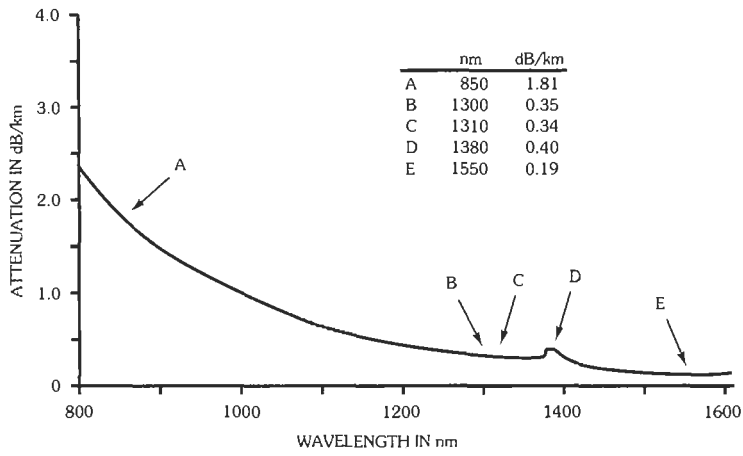
Waveguide Dispersion—In the absence of material and modal dispersion, there is still a dependence of the propagation constant β on the wavelength. This is called waveguide dispersion and can be understood by considering the distribution of the field within the core and cladding regions (Fig. 28). For a single mode, if the wavelength is such that the mode is well guided, then a large fraction of the energy is carried in the core region. As the wavelength changes and approaches cutoff for that mode, more and more energy is carried in the cladding region, and a different effective group index results.

Waveguide dispersion is determined by using the expression for the group index, Eq. 89

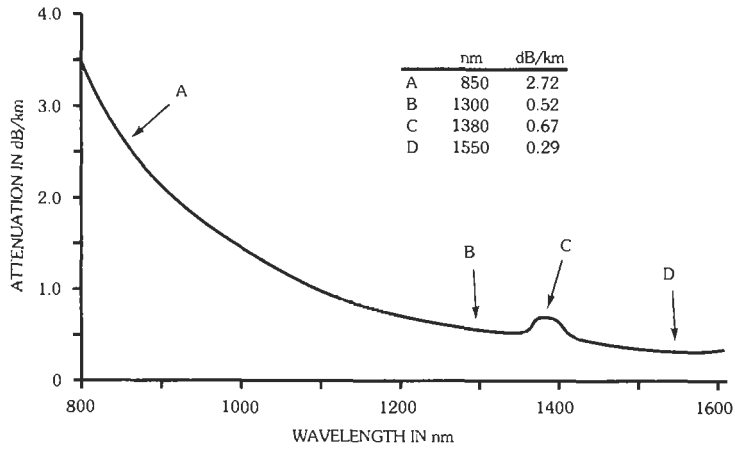
$$N = d\beta/dk = (d\beta/dV)(dV/dk) = (d\beta/dV) a \cdot \text{NA} \tag{Eq. 96}$$

where $V = k \cdot a \cdot \text{NA}$ has been used. Now by using $\delta t = (dN/d\lambda) L/c \delta\lambda$ and changing variables from λ to V , the waveguide dispersion can be written as

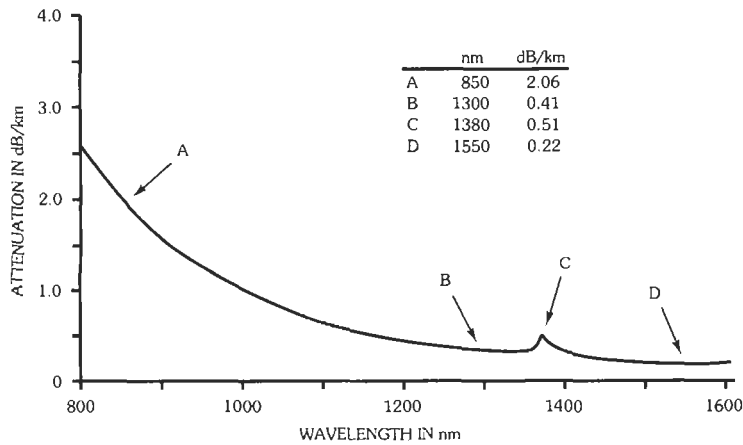
$$\delta t = |D_w| \delta\lambda L \tag{Eq. 97}$$



(A) Single-mode step-index fiber.



(B) 50- μ m-core graded-index fiber.



(C) Dispersion-shifted single-mode fiber.

Fig. 27. Typical attenuation versus wavelength for fibers. (Courtesy Corning, Inc.)

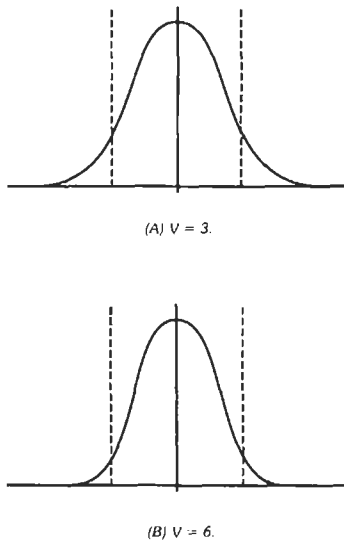


Fig. 28. Field distribution for the LP_{01} mode. The different distributions of the field within the core region yield different effective indices and produce waveguide dispersion.

where

$$D_w = -(1/2\pi c) V^2 (d^2\beta/dV^2) \quad (\text{Eq. 98})$$

The dependence of D_w on λ is through the dependence of V on λ . The units of D_w are the same as those for material dispersion.

Dispersion-Flattened and Dispersion-Shifted Fibers—The combined effect of material and waveguide dispersion is commonly called chromatic dispersion because both of these effects depend on the wavelength of the light.* The total chromatic dispersion is a function of both of D_λ and D_w . The material dispersion coefficient is generally fixed because most fiber is made of silica glass. However, the waveguide dispersion is a function of the geometry of the fiber. By changing the guiding properties of the fiber, the combined effect of the material and waveguide dispersion can be shifted from the material dispersion minimum of $1.3 \mu\text{m}$ to the attenuation minimum of $1.55 \mu\text{m}$ (see Fig. 29A). One can also produce two chromatic dispersion minimal results in a fiber than has a flat dispersion response between the two chromatic dispersion zeros (Figs. 29B and 29C).

Fiber Bandwidth

The bandwidth of a fiber is limited by dispersion; depending on the type of fiber, this may be modal or

chromatic dispersion. For all the dispersive mechanisms, the amount of pulse spreading is a function of the distance L traveled. Thus, fiber bandwidth is generally specified as the product of the bandwidth times the distance (i.e., $\text{MHz} \cdot \text{km}$ or $\text{GHz} \cdot \text{km}$).

Single-Mode Step-Index Fibers—The bandwidth of a single-mode step-index fiber is limited by material and waveguide dispersion. Modal dispersion is not present. However, both material and waveguide dispersion depend on the spectral width of the source. The data sheets of most manufacturers of single-mode step-index fiber will fit the total chromatic dispersion $\text{ps/nm} \cdot \text{km}$ to a polynomial within the wavelength range from about 1.3 to about $1.55 \mu\text{m}$. They also list a range of values for the zero chromatic dispersion wavelength. From these values and the spectral width of the source, an overall bandwidth distance product ($\text{GHz} \cdot \text{km}$) can be determined.

Graded-Index Fibers—The core diameter for graded-index fiber is much larger than for single-mode step-index fiber (50 and 62.5 , 85 and $100 \mu\text{m}$ are standard core sizes) and allows for easier alignment and splicing. However, even with the quadratic index profile, modal dispersion is usually the dominant dispersive mechanism because of the difficulty in manufacturing an optimal index profile. In general, modal dispersion is only weakly dependent on the spectral width of the laser. Thus, in contrast to single-mode step-index fiber, the bandwidth distance product of graded-index fiber is specified independently of the spectral width of the source. The actual bandwidth of graded-index fiber varies from run to run because of the sensitivity of the bandwidth on the index profile. Typical available bandwidths and attenuation are given in Table 1. A frequency response curve of a $50\text{-}\mu\text{m}$ -core graded-index fiber at $1.3 \mu\text{m}$ is shown in Fig. 30.

Splicing

Whenever two sections of fiber are to be permanently joined, a splice is required. There are several types of splicing technology, with fusion splicing and mechanical splicing being two of the more common technologies. In fusion splicing, the ends of the fiber are permanently fused. In mechanical splicing, the fibers are aligned by using a mechanical sleeve, and then epoxy is applied to produce an index match between the two ends. A sample histogram of splice losses for a typical commercially available mechanical splice for a single-mode fiber is shown in Fig. 31. Losses for fusion splicing and other types (such as the AT&T rotary splice) are in the range $0.02\text{--}0.07 \text{ dB}$.

Connectors

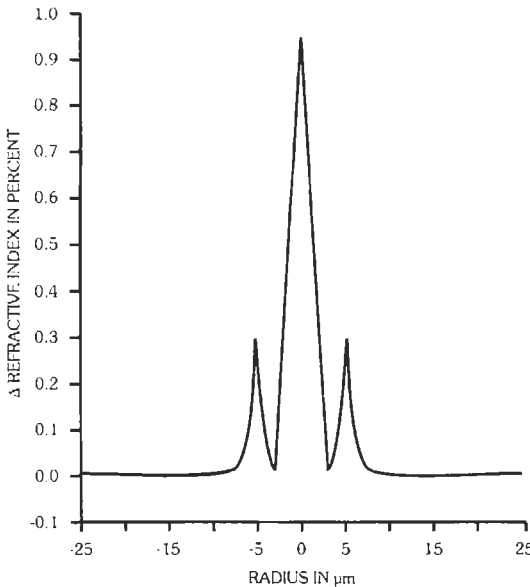
Connectors are used for temporary connections and whenever a system may need to be reconfigured. This is

* Other references restrict chromatic dispersion to refer only to material dispersion.

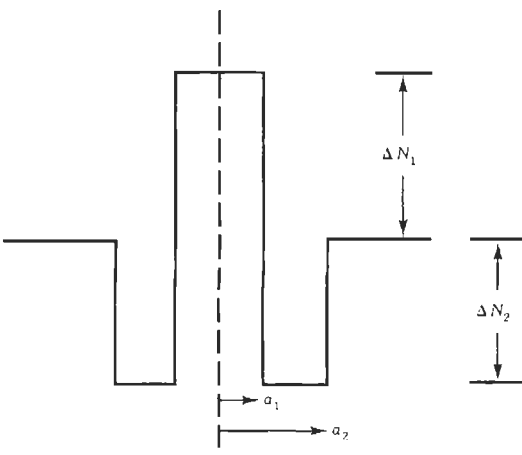
typically the case with fiber-optic LANs. There are many different kinds of connectors for both multimode and single-mode fiber. A typical loss histogram of two of the more common multimode fiber connectors (ST and SMA) is shown in Fig. 32.

Sample Calculation

If the frequency response of the fiber and the type of pulse for the data are known, the maximum data rate can be determined for graded-index fiber. This maximum data rate is approximately independent of the spectral width of the source.



(A) Index profile for commercial dispersion-shifted fiber. (Courtesy Corning, Inc.)



(B) Index profile for "W"-profile dispersion-flattened fiber.

For example, if the data is assumed to consist of Gaussian pulses of the form

$$f(t) = [1/(\sigma\sqrt{2\pi})] \exp(-t^2/2\sigma^2) \quad (\text{Eq. 99})$$

then integrating Eq. 99 from -2σ to 2σ yields 95% of the total energy. If we define 4σ as the pulse interval, the pulse rate R_m is

$$R_m = 1/4\sigma \quad (\text{Eq. 100})$$

The relationship between the optical 3-dB point (Fig. 30) and σ is given by*

$$\text{BW (3 dB opt.)} = 0.19/\sigma \quad (\text{Eq. 101})$$

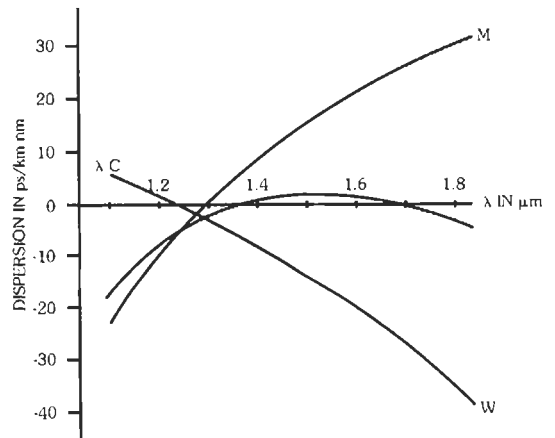
and is shown in Fig. 33. Note that while Eqs. 100 and 101 are strictly true only for a Gaussian pulse, Eq. 100 is approximately true if the rms values are used for the pulse width. (See sections 2.4.2 and 2.4.3 of reference 18.)

From Eqs. 100 and 101, the maximum pulse rate for the fiber shown in Fig. 30 is given by

$$R_m = 925/(0.19 \times 4) \approx 1.2 \text{ Gb/s} \cdot \text{km} \quad (\text{Eq. 102})$$

at $1.3 \mu\text{m}$. Thus for a 10-km link, the maximum data rate is roughly 120 Mb/s.

* Reference 18.



(C) Dispersion characteristics for "W"-profile fiber. W is the waveguide dispersion. M is the material dispersion. (From Okamoto, K., Edahiro, T., Kawana, A., and Miya, T. "Dispersion minimisation in single mode fibers over a wide spectral range." Electronics Letters 15, 1979, pp. 729-731.)

Fig. 29. Waveguide profiles of dispersion-shifted and dispersion-flattened fiber.

TABLE I. ATTENUATION AND BANDWIDTH CELLS FOR GRADED-INDEX FIBERS

Core/Cladding (μm)	Attenuation (dB/km)			Bandwidth (MHz · km)		
	850 nm only	1300 nm only	850/1300 nm	850 nm only	1300 nm only	850/1300 nm
50/125	≤ 2.4 to ≤ 3.0	≤ 0.6 to ≤ 1.2	$\leq 2.4/0.5$ to $\leq 2.7/1.2$	≥ 400 to ≥ 600	≥ 400 to ≥ 1500	$\geq 400/400$ to $\geq 600/1000$
62.5/125	≤ 3.0 or ≤ 3.2	≤ 0.7 or ≤ 0.9	$\leq 3.0/0.7$ or $\leq 3.2/0.9$	≥ 160 to ≥ 200	≥ 200 to ≥ 1000	$\geq 160/200$ to $\geq 200/600$
100/140	≤ 3.5 to ≤ 4.5	≤ 1.5 to ≤ 2.0	$\leq 3.5/1.5$ to $\leq 4.5/2.0$	≥ 100 to ≥ 300	≥ 100 to ≥ 500	$\geq 100/100$ to $\geq 200/200$

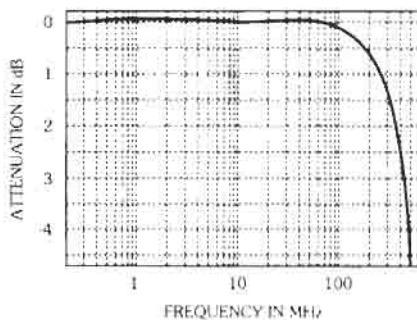


Fig. 30. Representative bandwidth measurement of a 2.2-km spool of 50- μm -core graded-index fiber at 1.3 μm . The optical 3-dB point is approximately 420 MHz, yielding a bandwidth of approximately 925 MHz · km.

REFERENCES

References Cited in Text

1. Personick, S. D. "Photodetectors for Fiber Systems." *Fundamentals of Optical Fiber Communications*. M. K. Barnoski, ed. New York: Academic Press, Inc., 1981.
2. *RCA Electro-Optics Handbook*. Tech. Series EOH-11, RCA Commercial Engineering, 1974.
3. Gagliardi, R. M., and Karp, S. *Optical Communications*. New York: John Wiley & Sons, Inc., 1976.
4. Pratt, W. K. *Laser Communications Systems*. New York: John Wiley & Sons, Inc., 1969.
5. Peterson, G. D., and Gardner, C. S. "Cross-Correlation Interference Effects in Multiaccess Op-

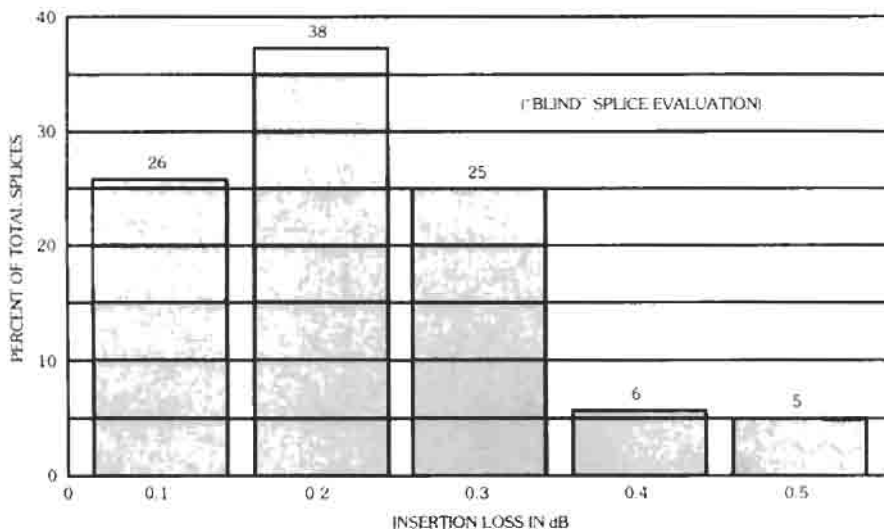


Fig. 31. Loss histogram of commercial mechanical single-mode-fiber splice at 1.3 μm . (Courtesy PSI Telecommunications, Inc.)

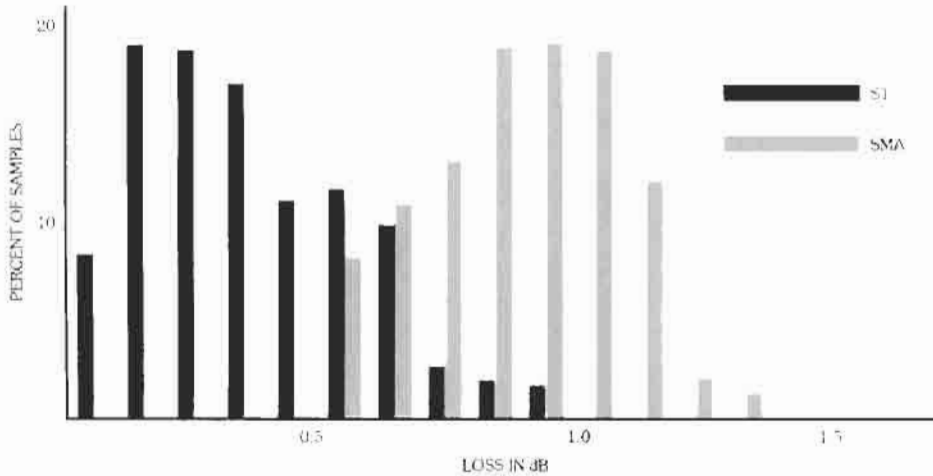


Fig. 32. Loss histogram for ST and SMA connectors. (From reference 19.)

tical Communications." *IEEE Trans. Aerospace Elect. Sys.*, AES-17, March 1981, pp. 199-207.

6. Kogelnik, H., and Li, T. "Laser Beams and Resonators." *Proc. IEEE*, 54, October 1966, pp. 1312-1329.
7. Klein, B. J., and Degnan, J. J. "Optical Antenna Gain. 1: Transmitting Antennas." *Appl. Opt.*, 13, September 1974, pp. 2134-2141.
8. Klein, B. J., and Degnan, J. J. "Optical Antenna Gain. 3: The Effect of Secondary Support Struts on Transmitter Gain." *Appl. Opt.*, 15, April 1976, pp. 977-979.
9. McCartney, E. J. *Optics of the Atmosphere*. New York: John Wiley & Sons, Inc., 1976.

10. Hinkley, E. D., ed. *Laser Monitoring of the Atmosphere*. New York: Springer-Verlag, 1976.
11. McClatchey, R. A., et al. "AFCRL Atmospheric Absorption Line Parameter Compilation." Air Force Camb. Res. Lab. Environ. Res. Pap. 434, AFCRL-TR-73-0096 (1973).
12. Driscoll, W. G., ed. *Handbook of Optics*. New York: McGraw-Hill Book Co., 1978.
13. Fante, R. L. "Electromagnetic Beam Propagation in Turbulent Media." *Proc. IEEE*, 63, December 1975, pp. 1669-1692.
14. Fried, D. L., and Mevers, G. E. "Evaluation of r_0 For Propagation Down Through the Atmosphere." *Appl. Opt.*, 13, November 1974, pp. 2620-2622.
15. Churnside, J. H., and McIntyre, C. M. "Heterodyne Receivers for Atmospheric Optical Communications." *Appl. Opt.*, 19, 15 February 1980, pp. 582-590.
16. Gloge, D. "Weakly guiding fibers." *Applied Optics*, Vol. 10, No. 10, 1971, pp. 2252-2258.
17. Saleh, B. E. A., and Teich, M. C. *Fundamentals of Photonics*. New York: John Wiley & Sons, Inc., 1991.
18. Gowar, J. *Optical Communication Systems*. Englewood Cliffs, NJ: Prentice Hall, 1984.
19. Buckman, A. B. *Guided-Wave Photonics*. Orlando, FL: Saunders College Publishing, 1992.

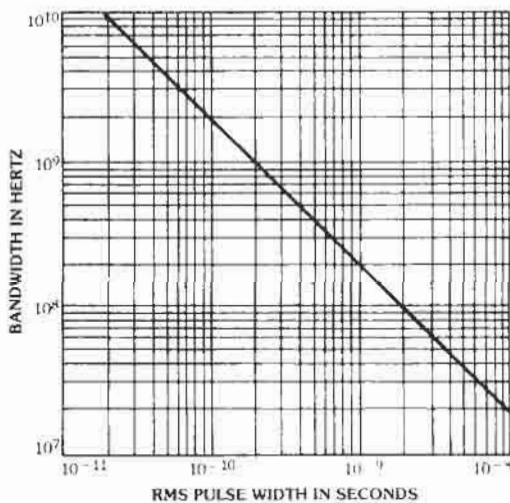


Fig. 33. Relationship between rms pulse width for a Gaussian pulse and the optical 3-dB point of a fiber.

Additional References on Fiber Optics

Ainslie, B. J., and Day, C. R. "A review of single-mode fibers with modified dispersion characteristics." *IEEE J. Lightwave Tech.*, 4:967-979, 1986.

Allard, F. C. *Fiber Optics Handbook*. New York: McGraw-Hill, 1990.

Corning Glass. *Engineering Handbook*. New York: Corning Inc., 1990.

Hecht, J. *Understanding Fiber Optics*. Indianapolis: Howard W. Sams & Co., 1987.

Hentschel, C. *Fiber Optics Handbook*. Hewlett Packard Company, Germany, 1988.

Hoss, R. J. *Fiber Optic Communications*. Englewood Cliffs, NJ: Prentice Hall, 1990.

Murata, H. *Handbook of Optical Fibers and Cables*. New York: Marcel Dekker, 1988.

Yeh, C. *Handbook of Fiber Optics: Theory and Applications*. Orlando, FL: Academic Press, 1990.

23

Analog Communications

*Revised by
Michael B. Pursley*

PART 1—CONTINUOUS MODULATION

Analytic Signal Representation of Modulated Waveforms 23-3

Amplitude Modulation 23-4

- Double Sideband (DSB)
- Conventional Amplitude Modulation (AM)
- Vestigial Sideband
- Single Sideband (SSB)

Demodulation, or Detection, of Amplitude Modulation 23-5

- DSB Detection
- AM Detection
- SSB Detection

Comparison of Amplitude-Modulation Systems 23-6

Exponential Modulation 23-6

- Frequency Spectrum of Single-Tone Angular Modulation
- Multitone Angle Modulation
- Spectral Distribution of an FM/FM Signal
- Bandwidth Considerations in Multitone FM
- Signal-to-Noise Improvement of FM Systems
- Signal-to-Noise Improvement Through De-Emphasis

PART 2—PULSE MODULATION

Sampling 23-11

- Sampling in the Time Domain
- Sampling in the Frequency Domain
- Sampling of a Bandpass Function ($B_0, B_0 + B$)

Pulse-Amplitude Modulation (PAM) *23-13*

 Spectra of Amplitude-Modulated Pulses
 Signal-to-Noise Ratio in PAM

Pulse-Time Modulation (PTM) *23-13*

 Spectra of Time-Modulated Pulses
 Signal-to-Noise Improvement Ratio in PTM

Pulse-Code Modulation (PCM) *23-16*

 Quantization Noise in a PCM System
 False-Pulse Noise in a PCM System

Delta Modulation (DM) *23-18*

 Signal-to-Noise Ratio in DM

Modulation is a process whereby certain characteristics of a wave (often called a carrier) are varied or selected in accordance with a message signal. Modulation can be divided into continuous modulation, in which the modulated wave is always present, and pulsed modulation, in which no signal is present between pulses. Digital data modulation is discussed in Chapter 24.

PART 1—CONTINUOUS MODULATION

In continuous modulation* the modulated carrier can be given by the expression $s(t) = A(t) \cos \theta(t)$, where $A(t)$ is the *instantaneous amplitude* and $\theta(t)$ is the *instantaneous phase*. For a sinusoidal carrier of angular frequency ω_c , this expression reduces to $s(t) = A(t) \cos[\omega_c t + \phi(t)]$, where $\phi(t)$ is the carrier phase. When the instantaneous amplitude $A(t)$ is varied linearly by the message function and the carrier phase is constant, the process is called *amplitude modulation*; when the carrier phase angle $\phi(t)$ is modulated by the message function, the process is called *angular or phase modulation*.

The concept of a rotating vector can be used to represent a sinusoidal vector modulated in both amplitude and phase as shown in Fig. 1, where $s(t)$ is represented as the projection of a rotating vector on a fixed reference axis.

$$s(t) = A(t) \cos[\omega_c t + \phi(t)]$$

$$= \text{Re}\{A(t) \exp[j(\omega_c t + \phi(t))]\}$$

In these expressions, $A(t)$ represents the envelope of the modulated carrier, and $\phi(t)$ is the modulated phase. The vector rotates with an instantaneous angular frequency $\omega_i(t)$ given by

$$\omega_i(t) = \omega_c + [d\phi(t)/dt]$$

In amplitude modulation, only the amplitude changes, and the general expression reduces to

$$s(t) = \text{Re}\{A(t) \exp(j\phi_0) \cdot \exp(j\omega_c t)\}$$

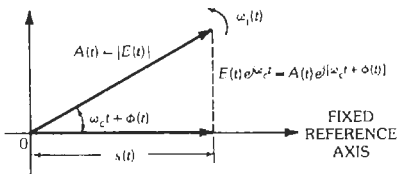


Fig. 1. Fixed-reference vector diagram. (From P. F. Panter, *Modulation, Noise, and Spectral Analysis*, Fig. 2-7, © 1965, McGraw-Hill Book Co.)

In phase modulation, only the phase changes so that

$$s(t) = \text{Re}\{A_c \exp[j\phi(t)] \cdot \exp(j\omega_c t)\}$$

where A_c is constant.

ANALYTIC SIGNAL REPRESENTATION OF MODULATED WAVEFORMS

A real signal

$$s(t) = A(t) \cos[\omega_c t + \phi(t)]$$

may be expressed either as

$$s(t) = \text{Re}\{A(t) \exp[j(\omega_c t + \phi(t))]\}$$

or as

$$s(t) = \text{Re}\{\psi(t)\}$$

where $\psi(t)$ is the analytic signal defined by

$$\psi(t) = s(t) + j\hat{s}(t)$$

The function $\hat{s}(t)$ is the Hilbert transform of $s(t)$, namely

$$\hat{s}(t) = \pi^{-1} \int_{-\infty}^{\infty} \frac{s(\tau)}{t - \tau} d\tau$$

Basically, the analytic signal $\psi(t)$ is a complex function of a real variable whose real and imaginary parts form a Hilbert pair. The analytic signal is simply a formalized version of the "rotating vector" discussed above. If $S(j\omega)$ is the Fourier transform of $s(t)$, then $\Psi(j\omega)$, the Fourier transform of $\psi(t)$, can be written in terms of $S(j\omega)$ as

$$\Psi(j\omega) = 2S(j\omega) \quad \omega > 0$$

$$= S(j\omega) \quad \omega = 0$$

$$= 0 \quad \omega < 0$$

Also, $\hat{S}(j\omega)$, the Fourier transform of $\hat{s}(t)$, is given by the equation

$$\hat{S}(j\omega) = -j(\text{sgn} \omega)S(j\omega)$$

where

$$\text{sgn} x = 1 \quad x > 0$$

$$= 0 \quad x = 0$$

$$= -1 \quad x < 0$$

and $\text{sgn} x$ is the signum function.

* P. F. Panter, *Modulation, Noise, and Spectral Analysis* (New York: McGraw-Hill Book Co., 1965; Ch. 5 and 6).

AMPLITUDE MODULATION

In amplitude modulation, the frequency components of the modulating signal are translated to occupy a different position in the spectrum. It is essentially a multiplication process in which the time functions that describe the modulating signal and carrier are multiplied together. The following amplitude-modulation systems are discussed.

- (A) Double-sideband suppressed carrier (DSB-SC), also called DSB
- (B) Conventional amplitude modulation (AM)
- (C) Vestigial sideband
- (D) Single sideband (SSB)

Double Sideband (DSB)

In DSB modulation, the message signal $g(t)$, whose Fourier transform is $G(j\omega)$, is considered to have zero dc component. The product

$$e(t) = A_c g(t) \cos \omega_c t$$

represents a double-sideband suppressed-carrier signal, and A_c = amplitude of unmodulated carrier. The radio-frequency envelope follows the waveform of the modulating signal $g(t)$ as shown in Fig. 2. The spectral components of the DSB signal $e(t)$ are given by its Fourier transform

$$E(j\omega) = \frac{1}{2}G[j(\omega - \omega_c)] + \frac{1}{2}G[j(\omega + \omega_c)]$$

as shown in Fig. 3. Note that the upper and lower sidebands are translated symmetrically $\pm\omega_c$ about the origin.

Conventional Amplitude Modulation (AM)

In amplitude modulation, a dc term is added to the modulating signal $g(t)$. The resulting waveform, shown

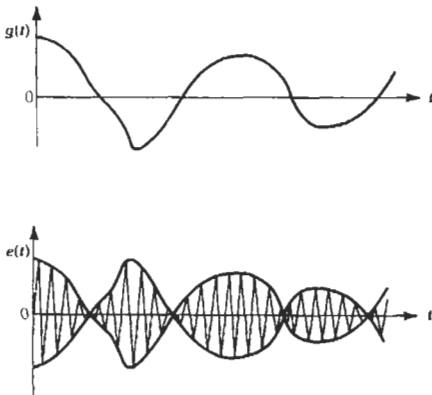


Fig. 2. Double-sideband waveforms. (From P. F. Panter, Modulation, Noise, and Spectral Analysis, Fig. 5-3, © 1965, McGraw-Hill Book Co.)

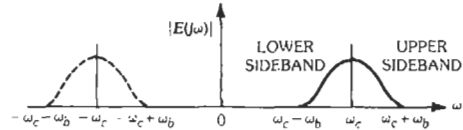
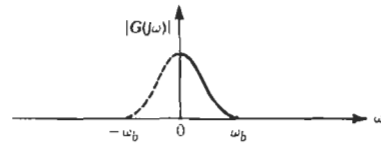


Fig. 3. Baseband signal and double-sideband spectra. (From P. F. Panter, Modulation, Noise, and Spectral Analysis, Fig. 5-2, © 1965, McGraw-Hill Book Co.)

in Fig. 4, is given by

$$\begin{aligned} e(t) &= [A_0 + as(t)]\cos \omega_c t \\ &= A_0 [1 + m_a s(t)]\cos \omega_c t \end{aligned}$$

where,

- a = maximum amplitude of modulating function, $g(t) = as(t)$, $|s(t)| \leq 1$,
- $m_a = a/A_0$ = modulation index or degree of modulation, $0 \leq m_a \leq 1$,
- A_0 = amplitude of unmodulated carrier, $|m_a s(t)| \leq 1$, to ensure an undistorted envelope.

Vestigial Sideband

Vestigial-sideband modulation is derived from a DSB signal by passing the output of the product modulator through a filter whose transfer function is $H_v(j\omega)$, as shown in Fig. 5. The transfer function $H_v(j\omega)$ of the filter treats the two sidebands of the DSB signal in such a manner as to attenuate one sideband differently from the other. The process of vestigial-sideband modulation by the use of the filter network $H_v(j\omega)$ may be replaced by the equivalent vestigial system shown in Fig. 6, where the transfer functions $H_i(j\omega)$ and $H_q(j\omega)$ are given by

$$H_i(j\omega) = \frac{1}{2}\{H_v[j(\omega - \omega_c)] + H_v[j(\omega + \omega_c)]\}$$

$$H_q(j\omega) = (1/2j)\{H_v[j(\omega - \omega_c)] - H_v[j(\omega + \omega_c)]\}$$

Single Sideband (SSB)

Single-sideband transmission may be produced in the same manner as vestigial sideband by using a high-pass filter $H_s(j\omega)$ which completely eliminates all signals on one side of the carrier frequency. The transfer function $H_s(j\omega)$ of the ideal high-pass filter is defined by

$$H_s(j\omega) = \left[\frac{1}{2} + \frac{1}{2} \operatorname{sgn}(\omega - \omega_c) \right] + \left[\frac{1}{2} - \frac{1}{2} \operatorname{sgn}(\omega - \omega_c) \right]$$

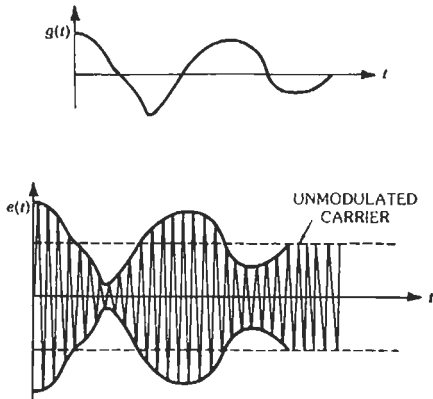


Fig. 4. Amplitude modulation. The modulating signal is at top and the modulated carrier at bottom. (From P. F. Panter, Modulation, Noise, and Spectral Analysis, Fig. 5-4, © 1965, McGraw-Hill Book Co.)

where $\text{sgn}\omega$ is the signum function. The output spectrum $E_s(j\omega)$ is given by

$$E_s(j\omega) = H_s(j\omega)E(j\omega) = \frac{1}{2}G[j(\omega - \omega_c)]\left[\frac{1}{2} + \frac{1}{2}\text{sgn}(\omega - \omega_c)\right] + \frac{1}{2}G[j(\omega + \omega_c)]\left[\frac{1}{2} - \frac{1}{2}\text{sgn}(\omega + \omega_c)\right]$$

The SSB signal can also be regarded as the resultant of quadrature modulation of a carrier by a pair of signals in phase quadrature. The modulated wave

$$e_s(t) = s(t) \cos\omega_c t - \sigma(t) \sin\omega_c t$$

represents an upper-sideband signal with no spectral components below the carrier angular frequency ω_c , where $s(t)$ is an arbitrary message function and $\sigma(t)$ its harmonic conjugate.

This equation can be written in the form

$$e_s(t) = [s^2(t) + \sigma^2(t)]^{1/2} \cos\{\omega_c t + \tan^{-1}[\sigma(t)/s(t)]\} = \alpha(t) \cos[\omega_c t + \phi(t)]$$

regarding the single-sideband signal as a hybrid amplitude-modulated and phase-modulated wave. The envelope $\alpha(t)$ and phase $\phi(t)$ are related by the analytic signal

$$\psi(t) = s(t) + j\sigma(t) = \alpha(t) \exp[j\phi(t)]$$

where $\sigma(t) = s(t)$, the Hilbert transform of $s(t)$. The amplitude and phase of the complex signal $\psi(t)$ are identical to the envelope and phase of the single-sideband wave. The Fourier transform of the analytic signal $\psi(t)$ is

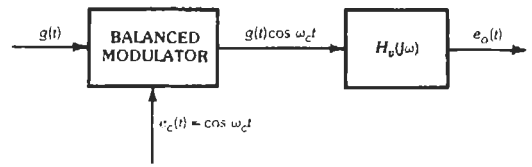


Fig. 5. Vestigial-sideband transmission system. (From P. F. Panter, Modulation, Noise, and Spectral Analysis, Fig. 5-7, © 1965, McGraw-Hill Book Co.)

$$\Psi(j\omega) = S(j\omega) + jS(j\omega) = S(j\omega) + S(j\omega) = 2S(j\omega), \quad \omega > 0 = S(j\omega) - S(j\omega) = 0, \quad \omega < 0$$

Thus, a study of single sideband can be made through the analytic signal without reference to the arbitrary carrier frequency ω_c .

DEMODULATION, OR DETECTION, OF AMPLITUDE MODULATION

The process of separating the modulating signal from a modulated carrier is called demodulation or detection. In DSB or SSB detection, the detector must be supplied with a carrier wave that is synchronized with the wave used at the transmitter. This method of detection is called coherent or synchronous detection. In conventional amplitude-modulation systems, coherent detection is not necessary, and the modulating signal may be recovered by the use of envelope detection; e.g., the modulated carrier is applied to a half-wave rectifier whose output is then filtered to provide the desired modulating signal.

DSB Detection

In DSB reception, the incoming signal $e_r(t)$ is multiplied by a locally generated signal that is phase-synchronized with the carrier component of the received

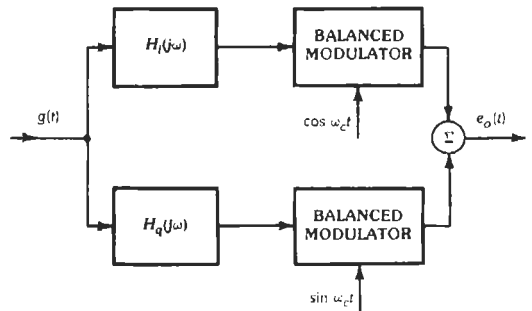


Fig. 6. Equivalent vestigial-sideband transmission system. (From P. F. Panter, Modulation, Noise, and Spectral Analysis, Fig. 5-8, © 1965, McGraw-Hill Book Co.)

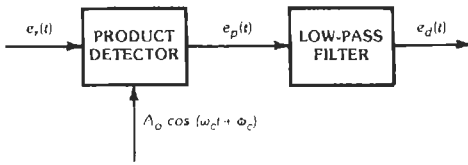


Fig. 7. Block diagram of double-sideband (DSB) receiver. (From P. F. Panter, Modulation, Noise, and Spectral Analysis, Fig. 6-1. © 1965, McGraw-Hill Book Co.)

signal $e_r(t)$, as shown in Fig. 7. The detected output after filtering is given by

$$e_d(t) = kg(t) \cos(\phi_c - \phi_0), \quad k = \text{constant}$$

where $(\phi_c - \phi_0)$ represents the phase difference between the transmitted carrier and the locally generated oscillator. When the local carrier is in phase with the incoming carrier, the detected signal is maximum. The output signal-to-noise ratio $(S/N)_o$ is related to the input signal-to-noise ratio $(S/N)_i$ by the expression

$$\frac{(S/N)_o}{(S/N)_i} = 2 \cos^2(\phi_c - \phi_0)$$

where the noise in each case is measured in a band occupied by the signal. This represents a maximum improvement of 3 decibels when the local oscillator is in phase with the incoming carrier.

AM Detection

Synchronous Detection:

$$\frac{(S/N)_o}{(S/N)_i} = \frac{2m_a^2 \langle g \rangle^2(t) \cos^2(\phi_c - \phi_0)}{1 + m_a^2 \langle g \rangle^2(t)}, \quad |g(t)| \leq 1$$

where $\langle g \rangle^2(t)$ equals the mean-square value of the message function, which is maximum for $m_a = 1$ and $\phi_c = \phi_0$.

Envelope Detection: In case of a carrier much stronger than the noise (high input carrier-to-noise ratio), we have

$$\frac{(S/N)_o}{(S/N)_i} = \frac{2m_a^2 \langle g \rangle^2(t)}{1 + m_a^2 \langle g \rangle^2(t)}$$

which is identical to the case of synchronous detection with $\phi_c = \phi_0$.

In case of poor input carrier-to-noise ratio, the message function $g(t)$ may be lost in the noise, which results in a threshold effect. This effect exists only in envelope detection and does not exist if synchronous or coherent detection is used.

SSB Detection

$$(S/N)_o / (S/N)_i = \cos^2(\phi_c - \phi_0)$$

where the signal component of the output is measured by the correlation of the detected output with the transmitted signal.

COMPARISON OF AMPLITUDE-MODULATION SYSTEMS

For equal power in the sidebands, the output signal-to-noise power ratios are identical.

For the same average total transmitted power, the following relations hold.

$$(S/N)_o(\text{DSB}) / (S/N)_o(\text{AM}) = 1 + r^{-1}$$

where r equals the ratio of the mean-square power of the message function to its peak power, and

$$(S/N)_o(\text{DSB}) / (S/N)_o(\text{SSB}) = 1$$

For equal peak power

$$(S/N)_o(\text{DSB}) / (S/N)_o(\text{AM}) = 4$$

for any waveform of the modulating signal.

To compare the merits of SSB versus DSB and AM on the basis of the signal-to-noise ratio, the waveform of the modulating signal must be specified. This is illustrated in Fig. 8 for a modulating signal $\sin^v x$, $0 \leq v \leq 1$.

EXPONENTIAL MODULATION

In exponential or angular modulation, the carrier analytic signal $A_c \exp[j(\omega_c t + \phi_c)]$ is multiplied by the transformed message function $\exp[j\psi(t)]$ to produce an angle-modulated carrier analytic signal.

$$\begin{aligned} e(t) &= \text{Re}\{A_c \exp[j(\omega_c t + \phi_c)] \cdot \exp[j\psi(t)]\} \\ &= \text{Re}\{A_c \exp[j\psi(t)]\} \end{aligned} \quad (1)$$

where,

A_c = amplitude of unmodulated carrier,

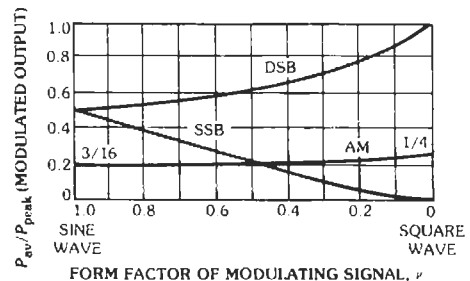


Fig. 8. Average-to-peak power relations as a function of modulating signal. (After W. K. Squires and E. Bedrosian, "The Computation of Single-Sideband Peak Power," Proceedings of the IRE, Vol. 48, p. 124, Fig. 2; January 1960.)

ω_c = angular frequency of unmodulated carrier,
 ϕ_c = carrier phase angle,
 $\phi(t) = [\omega_c t + \phi_c + \psi(t)]$,
 = instantaneous phase angle modulated by the message function, $g(t)$.

Expanding Eq. 1 in powers of $\psi(t)$, we have

$$e(t) = \text{Re}\{A_c \exp[j(\omega_c t + \phi_c)] \times [1 + j\psi(t) - (1/2!)\psi^2(t) - j(1/3!)\psi^3(t) + \dots]\}$$

When $|\psi(t)|_{\max} \gg 1$, we have nonlinear modulation, since the carrier is multiplied by higher powers of $\psi(t)$. In case $|\psi(t)|_{\max} \ll 1$, the exponential modulation is approximately linear and is given by

$$e(t) \cong \text{Re}\{A_c [1 + j\psi(t)] \exp[j(\omega_c t + \phi_c)]\}$$

Note that for amplitude modulation we have

$$e_{AM}(t) = \text{Re}\{A_c [1 + m_a g(t)] \exp[j(\omega_c t + \phi_c)]\}$$

Expressing Eq. 1 in the real form, we obtain

$$e(t) = A_c [\cos \omega_c t + \phi_c + \psi(t)]$$

where for phase modulation

$$\psi(t) = m_p g(t), \quad m_p = \text{constant}$$

and for frequency modulation

$$\psi(t) = m_f \int_0^t g(\tau) d\tau, \quad m_f = \text{constant}$$

The instantaneous frequency $\omega_i(t)$ is defined by

$$\omega_i(t) = \left[\frac{d\phi(t)}{dt} \right] = \left[\omega_c + \frac{d\psi(t)}{dt} \right]$$

In *phase modulation*, the instantaneous phase of the modulated signal varies proportionally with the modulating signal $g(t)$

$$e_{PM}(t) = A_c \cos[\omega_c t + m_p g(t)]$$

where ϕ_c has arbitrarily been set to zero.

For single-tone sinusoidal modulation, $g(t) = \cos \omega_m t$, we have

$$e_{PM}(t) = A_c \cos(\omega_c t + m_p \cos \omega_m t)$$

where $m_p = \Delta \theta$, and the peak phase deviation is independent of ω_m .

The instantaneous frequency

$$\begin{aligned} \omega_i(t) &= d\phi(t)/dt \\ &= \omega_c - m_p \omega_m \sin \omega_m t \end{aligned}$$

and the peak frequency deviation $\Delta \omega = m_p \omega_m$ is proportional to the modulating frequency ω_m .

In *frequency modulation*, the instantaneous frequency of the modulated signal is proportional to $g(t)$

$$\omega_i(t) = \omega_c + m_f g(t)$$

or

$$e_{FM}(t) = A_c \cos \left[\omega_c t + m_f \int_0^t g(\tau) d\tau \right]$$

For single-tone sinusoidal modulation

$$\omega_i(t) = \omega_c + \Delta \omega \cos \omega_m t$$

$$e_{FM}(t) = A_c \cos[\omega_c t + (m_f/\omega_m) \sin \omega_m t]$$

The peak frequency deviation $\Delta \omega \equiv m_f$ is independent of ω_m , whereas the peak phase deviation $\Delta \theta = \Delta \omega/\omega_m$ is inversely proportional to ω_m ; $\Delta \theta$ (in radians) is the modulation index, often denoted by β . For broad-band application, $\Delta \omega \ll \omega_c$ and $\beta \gg 1$.

Frequency Spectrum of Single-Tone Angular Modulation

Small Phase Deviation (Narrow-Band PM):

$$e(t) = A_c \cos(\omega_c t + \beta \sin \omega_m t), \quad \beta \ll 1$$

$$\begin{aligned} e(t) &\cong A_c (\cos \omega_c t - \beta \sin \omega_m t) \\ &= \underbrace{A_c \cos \omega_c t}_{\text{carrier}} - \underbrace{\frac{1}{2}(A_c \beta) \cos(\omega_c - \omega_m)t}_{\text{lower sideband}} \\ &\quad + \underbrace{\frac{1}{2}(A_c \beta) \cos(\omega_c + \omega_m)t}_{\text{upper sideband}} \end{aligned}$$

The corresponding equation for AM is

$$\begin{aligned} e_{AM}(t) &= A_c \cos \omega_c t + \frac{1}{2}(A_c m_a) \cos(\omega_c - \omega_m)t \\ &\quad + \frac{1}{2}(A_c m_a) \cos(\omega_c + \omega_m)t \end{aligned}$$

Large Phase Deviation (Wideband PM):

$$\begin{aligned} e(t) &= A_c \cos(\omega_c t + \beta \sin \omega_m t), \quad \beta \gg 1 \\ &= A_c [\cos \omega_c t \cos(\beta \sin \omega_m t) - \sin \omega_c t \sin(\beta \sin \omega_m t)] \\ &= A_c [\cos \omega_c t \sum_{n=-\infty}^{\infty} J_n(\beta) \cos n \omega_m t - \sin \omega_c t \sum_{n=-\infty}^{\infty} J_n(\beta) \sin n \omega_m t] \end{aligned}$$

The waveform for wideband modulation is given by the equation

$$e(t) = A_c \{ J_0(\beta) \cos \omega_c t - J_1(\beta) [\cos(\omega_c - \omega_m)t - \cos(\omega_c + \omega_m)t] + J_2(\beta) [\cos(\omega_c - 2\omega_m)t + \cos(\omega_c + 2\omega_m)t] - J_3(\beta) [\cos(\omega_c - 3\omega_m)t - \cos(\omega_c + 3\omega_m)t] + \dots \}$$

$$= A_c \sum_{n=-\infty}^{\infty} J_n(\beta) \cos(\omega_c + n\omega_m)t$$

as shown in Fig. 9.

In practical application, the required bandwidth is finite, for—beyond a certain frequency range from the carrier, depending on the magnitude of β —the sideband amplitudes, which are proportional to $J_n(\beta)$, are negligibly small (see Fig. 10). Note that at $\beta = 2.404$, $J_0(\beta) = 0$ and the carrier amplitude is zero.

The average power in an angle-modulated wave is constant

$$P = \frac{1}{2}(A_c^2) \sum_{-\infty}^{\infty} J_n^2(\beta)$$

$$= \frac{1}{2}(A_c^2)$$

Multitone Angle Modulation

For two-tone angle modulation with ω_1 and ω_2 , the instantaneous frequency is given by

$$\omega_i(t) = \omega_c + \Delta \omega_{c1} t \cos \omega_1 t + \Delta \omega_{c2} \cos \omega_2 t$$

where $\Delta \omega_{c1}$ and $\Delta \omega_{c2}$ denote the corresponding frequency deviations of ω_c , and the FM signal is

$$e(t) = A_c \cos(\omega_c t + \beta_1 \sin \omega_1 t + \beta_2 \sin \omega_2 t)$$

where $\beta_1 = \Delta \omega_{c1} / \omega_1$ and $\beta_2 = \Delta \omega_{c2} / \omega_2$.

The spectral components are as follows.

(A) Carrier:

$$J_0(\beta_1) J_0(\beta_2) A_c \cos \omega_c t$$

(B) Sidebands due to ω_1 :

$$J_n(\beta_1) J_0(\beta_2) A_c \cos(\omega_c \pm n\omega_1)t, \quad n = 1, 2, 3$$

(C) Sidebands due to ω_2 :

$$J_m(\beta_2) J_0(\beta_1) A_c \cos(\omega_c \pm m\omega_2)t, \quad m = 1, 2, 3$$

(D) Beat frequencies at $\omega_c \pm n\omega_1 \pm m\omega_2$:

$$J_n(\beta_1) J_m(\beta_2) A_c \cos(\omega_c \pm n\omega_1 \pm m\omega_2)t$$

Spectral Distribution of an FM/FM Signal

Let ω_c = carrier angular frequency, ω_s = subcarrier, and ω_m = modulating angular frequency. The instantaneous frequency of the carrier wave is

$$\omega_i(t) = \omega_c + \Delta \omega \cos[\omega_s t + \phi_s + \beta_s \sin(\omega_m t + \phi_m)]$$

where,

$\Delta \omega$ = peak frequency deviation of carrier,
 $\beta_s = \Delta \omega_s / \omega_m$ = peak phase deviation of subcarrier.

The spectral distribution is given by

$$e(t) = A_c \sum_{p=-\infty}^{\infty} \sum_{q=-\infty}^{\infty} J_p(\beta) J_q(p\beta_s)$$

$$\times \cos[(\omega_c + p\omega_s + q\omega_m)t + \phi_c + p\phi_s + q\phi_m]$$

where $\beta = \Delta \omega / \omega_s$ = peak phase deviation of carrier.

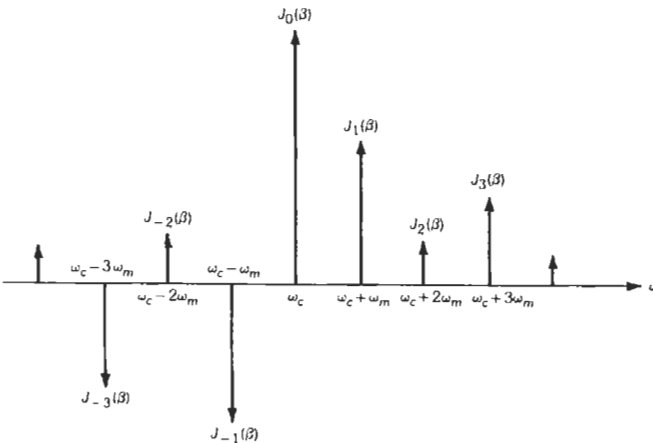


Fig. 9. Composition of FM wave into sidebands. (From P. F. Panter, Modulation, Noise, and Spectral Analysis, Fig. 7-6, © 1965, McGraw-Hill Book Co.)

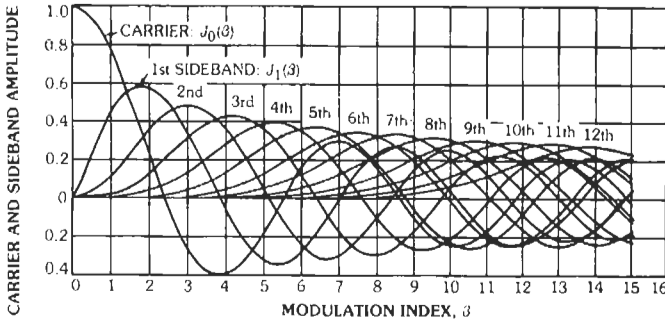


Fig. 10. Plot of Bessel functions of first kind as a function of argument β . (From P. F. Panter, Modulation, Noise, and Spectral Analysis, Fig. 7-8, © 1965, McGraw-Hill Book Co.)

Bandwidth Considerations in Multitone FM

An estimate of the IF bandwidth required for transmission of FM carrier by a complex modulating signal is given by

$$\beta_{IF} = 2(\Delta F + 2f_m) = 2\Delta F(1 + 2/\beta)$$

where,

ΔF = peak frequency deviation for the system,
 f_m = highest baseband frequency (see Fig. 11).

Signal-to-Noise Improvement in FM Systems

The performance of a conventional FM receiver in the presence of random fluctuation noise is commonly judged on the basis of the variation of the output signal-to-noise $(S/N)_o$ power ratio as a function of the carrier-to-noise power ratio $(C/N)_i$ measured at the input to the limiter. This relationship is shown graphically in Fig. 12. The threshold of full improvement occurs when $(C/N)_i$ is about 12 decibels. For all values of the carrier greater than the threshold, the output $(S/N)_o$ is proportional to the input $(C/N)_i$. The signal-

to-noise improvement ratio for a single-channel FM system is given by

$$(S/N)_o/(C/N)_i = (\Delta\Phi)^2, \text{ using a phase detector}$$

where $\Delta\Phi$ = peak phase deviation, and

$$(S/N)_o/(C/N)_i = 3(\Delta F/f_m)^2 = 3\beta^2, \text{ using a frequency discriminator}$$

where,

ΔF = peak frequency deviation,
 f_m = highest modulating frequency.

The signal-to-noise improvement ratio for a particular channel of a multiplex system is given by

$$(S/N)_o/(C/N)_i = (B/B_c)(\Delta F_m/f_n)^2$$

where,

$2B$ = IF bandwidth,
 B_c = channel bandwidth,

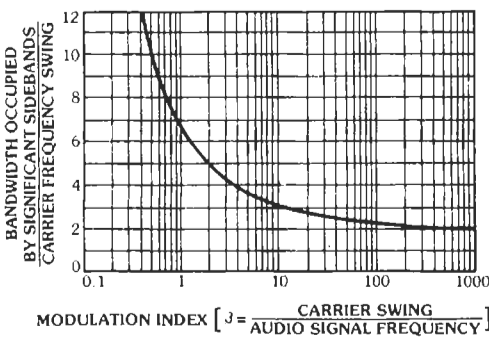


Fig. 11. Significant bandwidth (normalized) vs modulation index β . (From C. E. Tibbs and G. G. Johnstone, Frequency Modulation Engineering, John Wiley & Sons, Inc., New York. Courtesy of Chapman & Hall, Ltd., London.)

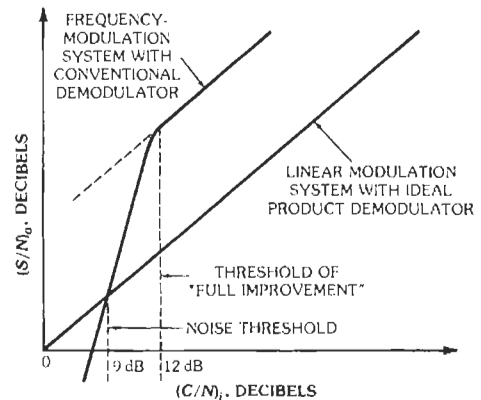


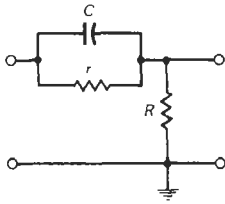
Fig. 12. Noise performance of conventional FM receiver. (From P. F. Panter, Modulation, Noise, and Spectral Analysis, Fig. 14-2, © 1965, McGraw-Hill Book Co.)

ΔF_m = peak channel frequency deviation,
 f_n = midband channel frequency in the n th channel.

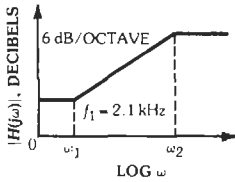
In the nonlinear region, when the noise is larger than the carrier there exists a signal-suppression effect, the average amplitude of the discriminator output is reduced, and for $(C/N)_i \ll 1$ we have $(S/N)_o \propto (C/N)_i^2$.

Signal-to-Noise Improvement Through De-Emphasis

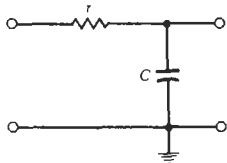
The $(S/N)_o$ ratio of the high frequency end of the baseband can be increased by passing the modulating signal (at the transmitting end) through a pre-emphasis network (Figs. 13A and B), which emphasizes the higher signal frequencies, and then passing the output of the discriminator through a de-emphasis network (Figs. 13C and D) to restore the original signal-power distribution. Typical pre-emphasis and de-emphasis cir-



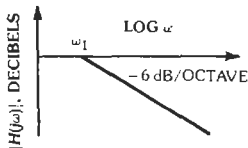
(A) Pre-emphasis network ($\tau \gg R, rC = 75 \mu s$).



(B) Asymptotic response ($\omega_1 = 1/rC, \omega_2 = 1/RC$).



(C) De-emphasis network ($rC = 75 \mu s$).



(D) Asymptotic response ($f_1 = 2.1 \text{ kHz}$).

Fig. 13. Pre-emphasis and de-emphasis networks. (From P. F. Panter, Modulation, Noise, and Spectral Analysis, Figs. 14-6 and 14-7, © 1965, McGraw-Hill Book Co.)

cuit responses for general time constants τ are shown in Fig. 14.

The improvement factor ρ_{FM} is given by

$$\rho_{FM} = \frac{(2\pi f_m \tau)^3}{3(2\pi f_m \tau - \tan^{-1} 2\pi f_m \tau)}$$

where f_m denotes the highest baseband frequency.

For narrow-band FM

$$\rho_{FM} \rightarrow 1$$

For wideband FM, f_m is large, and

$$\rho_{FM} \rightarrow (2\pi f_m \tau)^2/3$$

The mean $(S/N)_o$ ratio for FM with pre-emphasis is given by

$$(S/N)_o = \rho_{FM} \cdot 3\beta^2(C/N)_i$$

PART 2—PULSE MODULATION

In pulse-modulation systems, the unmodulated carrier is usually a series of regularly recurrent pulses; information is conveyed by modulating some parameter of the transmitted pulses such as the amplitude, duration, time of occurrence, or shape of pulse. This type of modulation is based on the "sampling principle," which states that a continuous message waveform that has a spectrum of finite width could be recovered from a set of discrete instantaneous samples whose rate is higher than twice the highest signal frequency. This discrete set of periodic samples of the message function is used to modulate some parameter of the carrier pulses. In *pulse-amplitude modulation* (PAM), the series of periodically recurring pulses is modulated in amplitude by the corresponding instantaneous samples of the message function. In *pulse-time modulation* (PTM), the instantaneous samples of the message function are used to vary the time of occurrence of

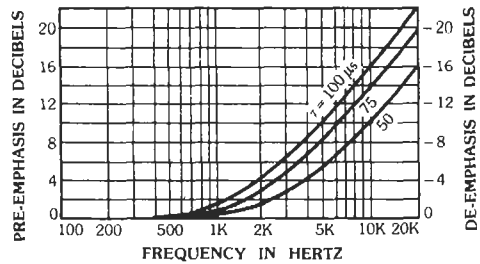


Fig. 14. Pre-emphasis and de-emphasis circuit response, for time constants of $\tau = 50, 75,$ and $100 \mu s$. (From C. E. Tibbs and G. G. Johnstone, Frequency Modulation Engineering, John Wiley & Sons, Inc., New York. Courtesy of Chapman & Hall, Ltd., London.)

some parameter of the pulsed carrier. Pulse-duration, pulse-position, and pulse-frequency modulation are particular forms of pulse-time modulation. In *pulse-duration modulation* (PDM), the time of occurrence of either the leading or trailing edge of each pulse (or both) is varied from its unmodulated position by the samples of the modulating wave. This is also called *pulse-length* or *pulse-width modulation* (PWM). In *pulse-position* (or *phase*) modulation (PPM), the samples of the modulating wave are used to vary the position in time of a pulse, relative to its unmodulated time of occurrence. Pulse-position modulation is essentially the same as PDM, except that the variable edge is now replaced by a short pulse. In *pulse-frequency modulation* (PFM), the samples of the message function are used to modulate the frequency of the series of carrier pulses.

The pulse-modulation systems enumerated so far are examples of uncoded pulse systems. In *pulse-code modulation* (PCM), the modulating signal waveform is sampled at regular intervals as in conventional pulse modulation. However, in PCM, the samples are first quantized into discrete steps; i.e., within a specified range of expected sample values, only certain discrete levels are allowed, and these are transmitted over the system by means of a code pattern of a series of pulses.

Another example of a code-modulation system is *delta modulation*. As in PCM, the range of signal amplitudes is quantized, and binary pulses are produced at the sending end at regular intervals. However, in delta-modulation systems, instead of the absolute quantized signal amplitude being transmitted at each sampling, the transmitted pulses carry the information corresponding to the derivative of the amplitude of the modulating signal.

SAMPLING

Sampling in the Time Domain

If a signal $f(t)$ is sampled at regular intervals of time and at a rate higher than twice the highest significant signal frequency, then the samples contain all the information of the original signal. The function $f(t)$ may be reconstructed from these samples by the use of a low-pass filter. The reconstruction equation is

$$f(t) = \alpha \sum_{n=-\infty}^{\infty} f(n\alpha/2B) \frac{\sin 2\pi B (t - n\alpha/2B)}{2\pi B (t - n\alpha/2B)},$$

$0 < \alpha \leq 1$

where $f(t)$ is band-limited to B hertz, and the samples are taken at sampling intervals $\alpha/2B$ seconds apart.

Sampling in the Frequency Domain

A time-limited signal $f(t)$ that is zero outside the range $t_1 < t < t_2$ is completely determined by the values

of the spectrum function $F(j\omega)$ at the angular-frequency sampling points given by

$$\omega_n = n[2/(t_2 - t_1)]$$

The function $f(t)$ expressed in terms of its sampling values in the frequency domain is given by the reconstruction equation

$$f(t) = \sum_{n=-\infty}^{\infty} (t_2 - t_1)^{-1} F \left(j \frac{2\pi n}{t_2 - t_1} \right) \times \exp [j2\pi n t / (t_2 - t_1)]$$

Sampling of a Bandpass Function ($B_0, B_0 + B$)

The reconstruction equation for $f(t)$ in terms of its sampled values is

$$f(t) = 2BT \sum_{n=-\infty}^{\infty} f(nT) \frac{\sin \pi B(t - nT)}{\pi B(t - nT)} \times \cos 2\pi B_c(t - nT)$$

where $B_c = B_0 + (B/2)$, the center frequency of the bandpass signal, and the permissible values of T are given by

$$m/2B_0 \leq T \leq [(m + 1)/2(B_0 + B)], \quad m = 0, 1, 2, \dots$$

provided $B_0 \neq 0$.

The minimum sampling frequency for a band-limited signal of width B is illustrated in Fig. 15.

PULSE-AMPLITUDE MODULATION (PAM)

In PAM, the samples of the message function are used to amplitude modulate the successive carrier pulses. When the modulated pulses follow the amplitude variation of the sampled time function during the sampling interval, the process is called *natural sampling*

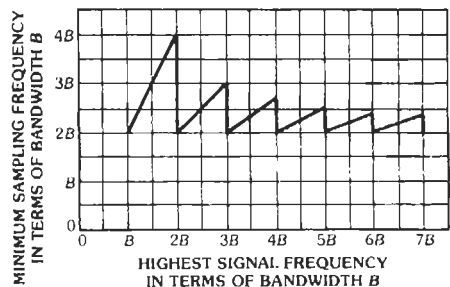


Fig. 15. Minimum sampling frequency for band of width B . (From P. F. Panter, Modulation, Noise, and Spectral Analysis, Fig. 17-13, © 1965, McGraw-Hill Book Co.)

or *top sampling*. In contrast with natural sampling is instantaneous or square-topped sampling, where the amplitude of the pulses is determined by the instantaneous value of the sampled time function corresponding to a single instant (i.e., center or edge) of the sampling interval. Pulse-amplitude modulation can be instrumented by two distinct methods. The first produces a variation of the amplitude of a pulse sequence about a fixed nonzero value or pedestal and constitutes double-sideband amplitude modulation (Figs. 16A and B). In the second method, the pedestal is zero, and the output signal consists of double-polarity modulated pulses and constitutes double-sideband suppressed-carrier modulation (Figs. 16C and D).

Spectra of Amplitude-Modulated Pulses

Double-Polarity AM Pulses, Natural (or Top) Sampling—In the process of natural sampling (or exact scanning), the modulated pulses follow the sampled time function during the sampling interval. The unit sampling function (Fig. 16E) consists of a train of unmodulated periodic pulses of unit amplitude given by

$$p_T(t) = (\tau/T) \sum_{n=-\infty}^{\infty} \frac{\sin(n\pi\tau/T)}{n\pi\tau/T} \exp(jn\omega_0 t)$$

where,

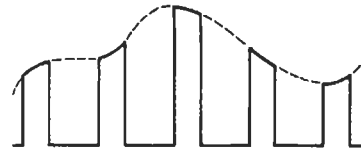
- $\omega_0 = 2\pi f_0 = 2\pi/T$ is the fundamental angular frequency of the pulse train,
- τ is the duration of the pulse,
- τ/T is the duty cycle.

Double-polarity AM pulses are obtained by multiplying the message signal $f(t)$ by the unit sampling function $p_T(t)$. In case of sinusoidal modulation, $f(t) = A \cos(\omega_m t + \phi)$, and the waveform of the AM pulses is given by

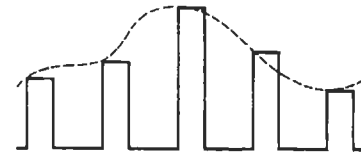
$$\begin{aligned} f_{s1}(t) &= f(t)p_T(t) \\ &= (\tau/T)A \cos(\omega_m t + \phi) \\ &\quad + (\tau/T)A \sum_{n=1}^{\infty} \frac{\sin(n\pi\tau/T)}{n\pi\tau/T} \\ &\quad \times \cos [(n\omega_0 \pm \omega_m)t \pm \phi] \end{aligned}$$

In the general case, the message function $f(t)$ is band-limited, and its spectrum is $F(j\omega)$. The output spectrum is

$$\begin{aligned} F_{s1}(j\omega) &= (\tau/T)F(j\omega) + (\tau/T) \sum_{n=1}^{\infty} \frac{\sin(n\pi\tau/T)}{n\pi\tau/T} \\ &\quad \times \{F[j(\omega - n\omega_0)] + F[j(\omega + n\omega_0)]\} \\ &= (\tau/T) \sum_{n=-\infty}^{\infty} \frac{\sin(n\omega_0\tau/2)}{n\omega_0\tau/2} F[j(\omega - n\omega_0)] \end{aligned}$$



(A) Single-polarity pulses.



(B) Single-polarity flat-top pulses.



(C) Double-polarity pulses.



(D) Double-polarity flat-top pulses.



(E) Unit sampling function.

Fig. 16. Various shapes of amplitude-modulated pulses. (From H. S. Black, Modulation Theory. Courtesy D. Van Nostrand Co., Inc., Princeton, N.J.)

The spectrum of the double-polarity AM pulses consists of the original modulation spectrum and an infinite number of upper and lower sidebands around ω_0 and its harmonics.

Double-Polarity AM Pulses, Instantaneous (or Square-Top) Sampling—In case of sinusoidal modulation, the output waveform is given by

$$\begin{aligned} f_{s2}(t) &= (\tau/T)A \sum_{n=-\infty}^{\infty} \frac{\sin[\pi(\tau/T)(n\omega_0 + \omega_m)/\omega_0]}{[\pi(\tau/T)(n\omega_0 + \omega_m)/\omega_0]} \\ &\quad \times \cos[(n\omega_0 + \omega_m)(t - \frac{1}{2}\tau) + \phi] \end{aligned}$$

In the general case, the output spectrum is

$$F_{s2}(j\omega) = (\tau/T) \frac{\sin(\omega\tau/2)}{\omega\tau/2} \sum_{n=-\infty}^{\infty} F[j(\omega - n\omega_0)]$$

Single-Polarity AM Pulses, Natural Sampling—
For sinusoidal modulation

$$f_{s3}(t) = [1 + m_a \cos(\omega_m t + \phi)] \times (\tau/T) \sum_{n=-\infty}^{\infty} \frac{\sin(n\pi\tau/T)}{n\pi\tau/T} \exp(jn\omega_0 t) = p_T(t) + f_{s1}(t)$$

where m_a is the modulation index. In the general case

$$F_{s3}(j\omega) = P(j\omega) + F_{s1}(j\omega)$$

where $P(j\omega)$ is the Fourier transform of $p_T(t)$.

Single-Polarity AM Pulses, Instantaneous Scanning—For sinusoidal modulation

$$f_{s4}(t) = p_T(t) + f_{s2}(t)$$

In the general case

$$F_{s4}(j\omega) = P(j\omega) + F_{s2}(j\omega)$$

Signal-to-Noise Ratio in PAM

The signal to noise ratio for PAM is

$$(S/N)_i = \frac{1}{2} m_a^2 P/N_0 B$$

where,

- P = average power of unmodulated radio-frequency pulse train,
- N_0 = noise-power density in watts/hertz,
- B = channel (RF) bandwidth.

Also

$$(S/N)_o = \frac{1}{2} m_a^2 (\tau/T) P/N_0 f_m$$

where f_m = top frequency of message function. If the receiver is blocked between pulses to eliminate the noise in the interpulse period, the $(S/N)_o$ at the output of a low-pass filter is

$$(S/N)_o = \frac{1}{2} m_a^2 P/N_0 f_m$$

which is identical to the result obtained for conventional CW carrier amplitude modulation.

In practice, PAM provides a poorer signal-to-noise ratio than conventional AM, because the receiver is unblocked for rather longer than the pulse-duration time owing to the sloping sides of the pulse.

PULSE-TIME MODULATION (PTM)

The improvement in signal-to-noise ratio obtained by the use of time-modulated pulses of constant amplitude instead of amplitude-modulated pulses led to the devel-

opment of systems using pulse-duration and pulse-position modulation. The sampling associated with pulse modulation may be either natural or uniform (periodic). Natural sampling may be defined as a process of sampling in which the time of sampling coincides with the time of appearance of the time-modulated pulse as shown in Fig. 17A. In the process of natural sampling, the pulse duration τ_n corresponds to the value of the modulating signal $M(t_n)$ at that instant, and consequently the sampling intervals t_n are not equal but depend on the modulation level. Uniform sampling may be defined as a process of sampling in which the variation in the parameter of the pulse is proportional to the modulating signal at uniformly spaced sampling times. This is illustrated in Fig. 17B, where the width of the pulses is proportional to the modulating values $M(t_n)$ which are sampled at equal intervals $t_n = nT_r$ and are independent of the modulation process.

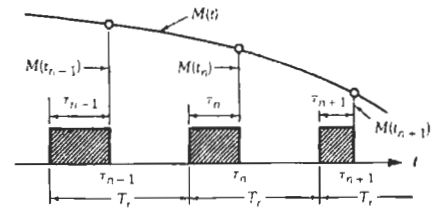
Spectra of Time-Modulated Pulses

The spectra of PTM pulses can be derived with reference to Fig. 18, where the two cosine waves, A and B , of angular frequency ω_r are displaced relative to each other by an amount τ , the width of the unmodulated pulse. The positive and negative steps that give rise to the pulse train are assumed to occur at the peaks of waveforms A and B , respectively. In the absence of modulation, the time of occurrence of the positive and negative steps is given by

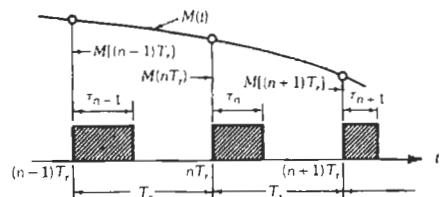
$$\omega_r(t + \pi/2) = 2n\pi$$

and

$$\omega_r(t - \pi/2) = 2n\pi$$



(A) Natural sampling.



(B) Uniform sampling.

Fig. 17. Pulse-duration modulation. (From P. F. Panter, Modulation, Noise, and Spectral Analysis, Fig. 18-14, © 1965, McGraw-Hill Book Co.)

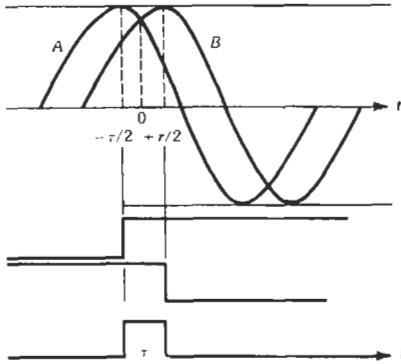


Fig. 18. Modulation process (modified). (From P. F. Panter, *Modulation, Noise, and Spectral Analysis*, Fig. 17-14, © 1965, McGraw-Hill Book Co.)

With natural modulation, the time of occurrence of the positive and negative steps is given by

$$\omega_r(t + \pi/2) + \beta \sin(\omega_m t + \phi) = 2n\pi$$

and

$$\omega_r(t - \pi/2) + \beta \sin(\omega_m t + \phi) = 2n\pi$$

Similarly, with uniform modulation, the time of occurrence or the position of the leading and trailing edges of the pulses is determined by

$$\omega_r(t + \pi/2) + \beta \sin(\omega_m t + \phi) = 2n\pi$$

$$\omega_r(t - \pi/2) + \beta \sin(\omega_m t - \tau + \phi) = 2n\pi$$

where,

ω_m is the modulating frequency,
 β is the modulation index.

Pulses whose moments of occurrence satisfy these equations are said to be time modulated. In pulse-frequency modulation, $\beta = \Delta \omega / \omega_m$, and in pulse-phase (or pulse-position) modulation, β is constant independent of the modulating frequency.

Pulse-Frequency Modulation, Natural Sampling—A useful expression for an infinite train of unmodulated pulses is in the form

$$p_T(t) = (A/2\pi j) \sum_{k=-\infty}^{\infty} k^{-1} \{ \exp[jk\omega_r(\pi/2)] - \exp[-jk\omega_r(\pi/2)] \} \exp(jk\omega_r t)$$

where,

A is the amplitude of the pulses,
 ω_r is the pulse repetition frequency.

Frequency modulation can be taken into account by

substituting for $\omega_r \pi/2$ in the expressions for the leading and trailing edges in the last equation, the expressions

$$\frac{1}{2}(\omega_r \tau) + \beta \sin(\omega_m t + \phi)$$

and

$$\frac{1}{2}(\omega_r \tau) - \beta \sin(\omega_m t + \phi)$$

The frequency-modulated pulse train is then

$$p_m(t) = (A/2\pi j)$$

$$\times \sum_{k=-\infty}^{\infty} k^{-1} \{ \exp[jk\omega_r(\pi/2) + k\beta \sin(\omega_m t + \phi)] - \exp[-jk\omega_r(\pi/2) - k\beta \sin(\omega_m t + \phi)] \} \times \exp(jk\omega_r t)$$

$$= \frac{A\omega_r \tau}{2\pi} + \frac{A\omega_r \tau}{\pi} \sum_{k=1}^{\infty} \frac{\sin[k\omega_r(\pi/2)]}{k\omega_r(\pi/2)}$$

$$\times \{ J_0(k\beta) \cos k\omega_r t + \sum_{n=1}^{\infty} J_n(k\beta) \times \{ \cos[(k\omega_r + n\omega_m)t + n\phi] + (-1)^n \cos[(k\omega_r - n\omega_m)t - n\phi] \} \}$$

This expression may be compared with that for the spectrum of a frequency-modulated continuous wave given by

$$e_{FM}(t) = AJ_0(\beta) \cos \omega_r t$$

$$+ A \sum_{n=1}^{\infty} J_n(\beta) \{ \cos[(\omega_r + n\omega_m)t + n\phi] + (-1)^n \cos[(\omega_r - n\omega_m)t - n\phi] \}$$

The conclusions reached are as follows.

- (A) With pulse-frequency modulation using natural sampling, the direct-current component of the pulse spectrum has no sideband of the modulating frequency.
- (B) The k th harmonic of the pulse-repetition frequency is frequency modulated, the modulation index being $k\beta$.

Pulse-Frequency Modulation, Uniform Sampling—In this type of modulation, the displacement of waveform B of Fig. 18 from its unmodulated position at any instant of time t will depend on the value of the modulating voltage at $(t - \tau)$. The expression for the modulated pulse train becomes

$$p_m(t) = (A/2\pi j)$$

$$\begin{aligned} & \times \sum_{k=1}^{\infty} k^{-1} \{ \exp[jk\omega_r(\pi/2) + k\beta \sin(\omega_m t + \phi)] \} \\ & - \exp\{-j[k\omega_r(\pi/2) - k\beta \sin(\omega_m t - \tau + \phi)]\} \\ & \times \exp(jk\omega_r t) \\ & = \frac{A\omega_r\tau}{2\pi} + A \left(\frac{\Delta\omega}{2\pi} \right) \tau \frac{\sin[\omega_m(\pi/2)]}{\omega_m(\pi/2)} \\ & \quad \times \cos[\omega_m t + \phi - (\omega_m \pi/2)] \\ & + \frac{A\omega_r\tau}{\pi} \sum_{k=1}^{\infty} \left(J_0(k\beta) \frac{\sin[k\omega_r(\pi/2)]}{k\omega_r(\pi/2)} \right. \\ & \times \cos k\omega_r t + \sum_{k=1}^{\infty} J_n(k\beta) \left\{ \frac{\sin(k\omega_r - n\omega_m)(\pi/2)}{k\omega_r(\pi/2)} \right. \\ & \quad \times \cos[(k\omega_r + n\omega_m)t + n\phi - n\omega_m(\pi/2)] \\ & \quad + (-1)^n \frac{\sin(k\omega_r - n\omega_m)(\pi/2)}{k\omega_r(\pi/2)} \\ & \quad \left. \left. \times \cos[(k\omega_r - n\omega_m)t - n\phi + n\omega_m(\pi/2)] \right\} \right) \end{aligned}$$

The conclusions reached are as follows.

- (A) The direct-current component of the pulse spectrum has a sideband of the modulating frequency of amplitude

$$(A\Delta\omega\pi/2\pi)\{\sin[\omega_m(\pi/2)]/\omega_m(\pi/2)\}$$

Modulation can therefore be recovered by means of a low-pass filter.

- (B) The upper and lower sidebands of the *k*th harmonic of the pulse-repetition frequency are not equal in amplitude, whereas in the case of natural sampling they are equal.

Pulse-Position (or Pulse-Phase) Modulation—

The waveform of pulse-phase modulation can be directly derived from that for pulse-frequency modulation by substituting $\omega_r\tau_d$ for β , where $\omega_r\tau_d$ represents the peak phase deviation of waveforms *A* and *B*, which is constant independent of the modulation frequency ω_m . The resulting waveform is the following:

Natural sampling:

$$\begin{aligned} P_m(t) &= \frac{A\omega_r\tau}{2\pi} + \frac{A\omega_r\tau}{\pi} \sum_{k=1}^{\infty} \frac{\sin[k\omega_r(\pi/2)]}{k\omega_r(\pi/2)} \\ & \times \left(J_0(k\omega_r\tau_d) \cos k\omega_r t + \sum_{n=1}^{\infty} J_n(k\omega_r\tau_d) \right. \\ & \quad \times \{ \cos[(k\omega_r + n\omega_m)t + n\phi] \\ & \quad \left. + (-1)^n \cos[(k\omega_r - n\omega_m)t - n\phi] \} \right) \end{aligned}$$

Note that each pulse-repetition-frequency harmonic is phase-modulated, with peak deviation equal to $k\omega_r\tau_d$. Also, there is no sideband accompanying the direct-current component of the pulse spectrum, and hence modulation cannot be recovered by means of a low-pass filter.

Uniform sampling:

$$\begin{aligned} p_m(t) &= \frac{A\omega_r\tau}{2\pi} + \frac{A\omega_r\omega_m\tau_d\tau}{2\pi} \frac{\sin[\omega_m(\pi/2)]}{\omega_m(\pi/2)} \\ & \times \cos[\omega_m t + \phi - (\omega_m \pi/2)] + \frac{A\omega_r\tau}{\pi} \sum_{k=1}^{\infty} \left(J_0(k\omega_r\tau_d) \right. \\ & \quad \times \frac{\sin[(k\omega_r(\pi/2)]}{k\omega_r(\pi/2)} \cos k\omega_r t + \sum_{n=1}^{\infty} J_n(k\omega_r\tau_d) \\ & \quad \times \left\{ \frac{\sin(k\omega_r + n\omega_m)(\pi/2)}{k\omega_r(\pi/2)} \cos[(k\omega_r + n\omega_m)t \right. \\ & \quad \left. + n\phi - n\omega_m(\pi/2)] + (-1)^n \frac{\sin(k\omega_r - n\omega_m)(\pi/2)}{k\omega_r(\pi/2)} \right. \\ & \quad \left. \left. \times \cos[(k\omega_r - n\omega_m)t - n\phi + n\omega_m(\pi/2)] \right\} \right) \end{aligned}$$

This is an equation very similar to that for pulse-frequency modulation.

Signal-to-Noise Improvement Ratio in PTM

In PDM, the noise manifests itself as jitter in the leading and trailing edges of the recovered pulses, and the slopes of the pulse edges influence noise reduction. A PPM system is affected by noise in the same manner as a PDM system. For trapezoidal pulses (Figs. 19 and 20), the *S/N* power ratio at the demodulator output is

$$(S/N)_o = \frac{1}{2} (t_0/\tau_r)^2 (A_c/\sigma)^2$$

The peak pulse power to mean noise power ratio is

$$(C/N)_i = (A_c/\sigma)^2$$

Hence

$$(S/N)_o = \frac{1}{2} (t_0/\tau_r)^2 (C/N)_i$$

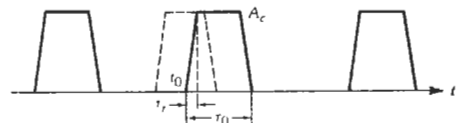


Fig. 19. Pulse-position modulation of trapezoidal pulses. (From P. F. Panter, Modulation, Noise, and Spectral Analysis, Fig. 18-26, © 1965, McGraw-Hill Book Co.)

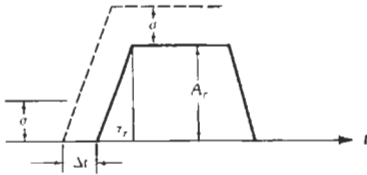


Fig. 20. Variation in pulse position due to the presence of noise or interference. (From P. F. Panter, *Modulation, Noise, and Spectral Analysis*, Fig. 18-27, © 1965, McGraw-Hill Book Co.)

The $(S/N)_0$ can be improved by decreasing the pulse rise τ_r , or correspondingly by widening the transmission bandwidth. For $B \cong 1/\tau_r$,

$$(S/N)_0 = \frac{1}{2} t_0^2 B^2 (C/N)_i$$

For $B \cong 1/2\tau_r$,

$$(S/N)_0 = 2 t_0^2 B^2 (C/N)_i$$

As in the case of FM, the $(S/N)_0$ ratio cannot be improved indefinitely by widening the bandwidth, because the noise power introduced at the receiver increases with bandwidth and eventually becomes comparable to the signal and "takes over" the system. A threshold level thus also exists just as in the FM case. This threshold level is usually taken as $A_c/\sigma = 2$, or $(C/N)_i = 4$ (6 dB).

PULSE-CODE MODULATION (PCM)

In PCM, several pulses are used as a code group to describe the quantized amplitude of a single sample. For example, a code group of n on-off pulses (binary code) can represent 2^n discrete amplitudes or levels, including zero level. In general, in an s -ary PCM system, the number of quantized amplitude levels the code group can express (including zero level) is given by

$$M = s^n$$

If a stands for 0 or 1, the binary notation with n digits, a_1, a_2, \dots, a_n , represents the number

$$a_1 2^0 + a_2 2^1 + a_3 2^2 + \dots + a_n 2^{n-1}$$

In the ternary number system, a stands for the pulse amplitude 0, 1, 2, and the code group of n digits represents the number

$$a_1 3^0 + a_2 3^1 + \dots + a_n 3^{n-1}$$

Table 1 shows how the 64 numbers from 0 through 63 are represented in binary, quaternary, and octonary notation.

TABLE 1. ENCODING INTO BINARY, QUATERNARY, AND OCTONARY NUMBERS*

Decimal No.	Binary No.	Quaternary No.	Octonary No.
0	000000	000	00
1	000001	001	01
2	000010	002	02
3	000011	003	03
4	000100	010	04
5	000101	011	05
6	000110	012	06
7	000111	013	07
8	001000	020	10
9	001001	021	11
10	001010	022	12
11	001011	023	13
12	001100	030	14
⋮	⋮	⋮	⋮
62	111110	332	76
63	111111	333	77

*From P. F. Panter, *Modulation, Noise, and Spectral Analysis*, Table 20-1, © 1965, McGraw-Hill Book Co., New York.

Quantization Noise in a PCM System

Representing the message signal by certain discrete allowed levels or steps is called quantizing. It inherently introduces an initial error in the amplitude of the samples, giving rise to quantization noise.

Uniform Spacing of Levels—In this case, the quantizing interval, or step, Δv is constant, and the quantizing noise power is given by

$$N_q = (\Delta v)^2 / 12$$

assuming that the quantization noise is uniformly distributed between $\pm \Delta v/2$. Assuming that the amplitudes of the samples are uniformly distributed, the signal power recovered from the quantized samples is

$$S_q = [(M^2 - 1)/12](\Delta v)^2$$

where M is the number of discrete levels assigned to the message signal. The ratio of the signal power to the quantizing noise power is

$$S_q/N_q = M^2 - 1 \cong M^2, \quad M \gg 1$$

Nonuniform Spacing of Levels—Quantization noise can be reduced by the use of nonuniform spacing of levels, to provide smaller steps for weaker signals and coarser quantization near the peak of large signals. Quantization noise can be minimized by an optimum level distribution that is a function of the probability density of the signal. The optimum level spacing Δv_k is given by

$$[p(v_k)]^{1/3} \Delta v_k = k/M, \quad k = \text{constant}$$

With optimum level spacing, the total minimum error power is

$$(N_q)_{\min} = (2/3M^2) \left\{ \int_0^V [p(v)]^{1/3} dv \right\}^3$$

where $p(v)$ is the probability density of the message signal, and the nonuniform levels are symmetrically disposed about zero level in the amplitude range $(-V, V)$.

In practice, nonuniform quantization is realized by compression, followed by uniform quantization as in Fig. 21. The logarithmic compression curve shown in Fig. 22 renders the distortion largely independent of the signal and is relatively easy to obtain in practice.

$$u = k \log[1 + (\mu v/V)]$$

where,

- v = input voltage,
- u = output voltage,
- μ = compression parameter,
- k = undetermined constant.

If the maximum values of the input and the compressed signals are adjusted to be equal, this gives

$$u = \frac{V \log(1 + \mu v/V)}{\log(1 + \mu)}, \quad 0 \leq v \leq V$$

and

$$u = \frac{-V \log(1 - \mu v/V)}{\log(1 + \mu)}, \quad -V \leq v \leq 0$$

The quantizing noise power with logarithmic compression is

$$N_q = (\alpha^2/12)(V^2 + \mu^2 S)$$

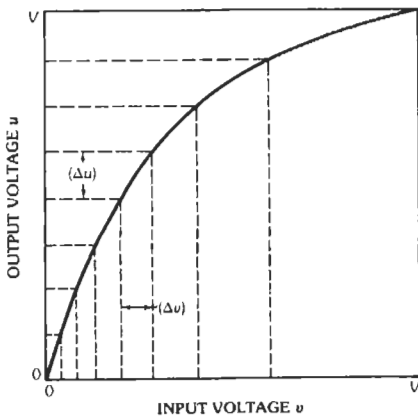


Fig. 21. Compression characteristic of "compressor." (From P. F. Panter, Modulation, Noise, and Spectral Analysis, Fig. 20-10, © 1965, McGraw-Hill Book Co.)

where

$$\alpha = \frac{2 \log(1 + \mu)}{\mu M}$$

and

$$S = \int_{-V}^V v^2 p(v) dv = \text{average signal power}$$

False-Pulse Noise in a PCM System

In addition to quantization noise, a PCM system is characterized by *false-pulse noise*, which originates primarily at the receiving end of the system and is caused by noise spikes breaking through the threshold. This type of noise decreases rapidly as the signal power is increased above threshold. The effect of the false pulses introduced in the code group is to introduce an error in the decoded samples. The mean-square error introduced in the decoded signal is defined as the false-pulse noise. The output signal-to-noise ratio at the decoder is

$$(S/N)_o = (1/4pq) - 1$$

where,

- p is the probability of sending out state one and receiving state zero and vice versa,
- q is the probability that no transmission fault occurs ($p + q = 1$).

The output signal-to-noise ratio drops from infinity with a noiseless channel ($p = 0, q = 1$) to zero in the case of an infinitely large channel noise ($p = q = 1/2$).

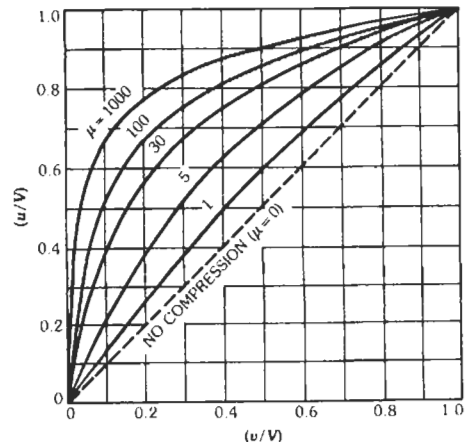


Fig. 22. Logarithmic compression characteristics. (From P. F. Panter, Modulation, Noise, and Spectral Analysis, Fig. 20-11, © 1965, McGraw-Hill Book Co.)

The output signal-to-noise ratio for K links in tandem is given by

$$(S/N)_{o,(K)} = \frac{(q - p)^{2K}}{1 - (q - p)^{2K}}$$

These expressions for the $(S/N)_o$ in a PCM system are given in terms of the probability of false pulses in the code group due to channel noise. The following expressions relate the output signal-to-noise ratio to the input carrier-to-noise ratio for one link.

$$(S/N)_o = (\pi/8)^{1/2} (V_0/2\sigma) \exp[\frac{1}{2}(V_0/2\sigma)^2]$$

For unipolar or on-off binary system

$$(S/N)_o = (\pi/16)^{1/2} [(C/N)_i]^{1/2} \exp[\frac{1}{4}(C/N)_i]$$

For bipolar binary system

$$(S/N)_o = (\pi/8)^{1/2} [(C/N)_i]^{1/2} \exp[\frac{1}{2}(C/N)_i]$$

For K links in tandem

$$(S/N)_{o,(K)} = (\text{erfx})^{2K} / [1 - (\text{erfx})^{2K}]$$

where $x = V_0/2\sigma$. For high $(c/n)_i$, $x \gg 1$, and

$$\text{erfx} \cong 1 - (2/\pi)^{1/2} \frac{\exp(-x^2/2)}{x}$$

The rapid improvement in $(S/N)_o$ for small increases in

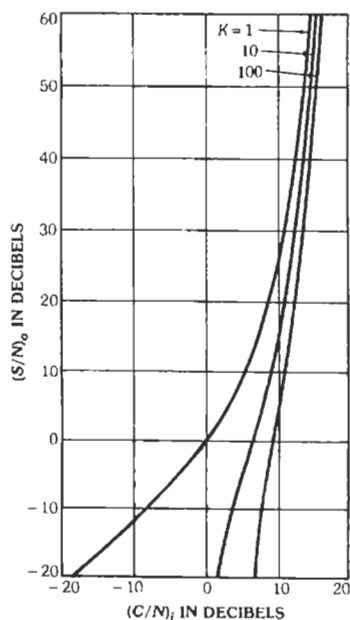


Fig. 23. Output signal-to-noise ratio for PCM. (From H. F. Mayer, "Principles of Pulse Code Modulation," Advan. Electron., Vol. 3, © 1951, Academic Press, Inc., New York.)

$(C/N)_i$ is illustrated in Fig. 23 for various links in tandem.

DELTA MODULATION (DM)

In a DM system, instead of the absolute signal amplitude being transmitted at each sampling, only the changes in signal amplitude from sampling instant to sampling instant are transmitted. As shown in Fig. 24, the transmitted pulse train $e_2(t)$ of positive and negative pulses at the output of the encoder can be assumed to be generated at a constant clock rate. The transmitted pulses from the pulse generator are positive if the change in signal amplitude is positive; otherwise the transmitted pulses are negative. In the decoder, the delta-modulated pulse train $e_2(t)$ is integrated into the voltage $e_1(t)$, which consists of the original message function plus noise components due to sampling. These are eliminated by a low-pass filter so that the reconstructed signal of the final output is a close replica of the original modulating signal $e_0(t)$.

Signal-to-Noise Ratio in DM

The difference between the original and reconstructed signals gives rise to a "quantizing noise" that can be decreased by increasing the "sampling frequency," which in DM is made equal to the pulse frequency. The quantized noise power when single integration is used is given by

$$N_0 = \frac{2}{3}(f_m/f_s)(\Delta v)^2$$

where,

- f_m = highest modulating frequency,
- Δv = height of unit step in volts.

A DM system has no fixed maximum signal amplitude limitation but overloads when the slope of the

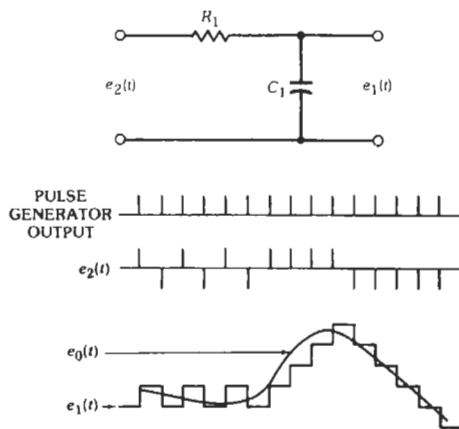


Fig. 24. Delta-modulation waveforms for single integration. (From P. F. Panter, Modulation, Noise, and Spectral Analysis, Fig. 22-2, © 1965, McGraw-Hill Book Co.)

signal is too large. The largest slope the system can reproduce is one that changes by one level or step every pulse interval, so that the maximum signal power depends on the type of signal. The signal power in the calculation of signal-to-noise ratio is taken as the power of the sinusoidal tone that is just below the overload point. The maximum amplitude of such a sinusoidal signal of frequency f that can be transmitted with single integration without overloading is

$$A = f_s (\Delta v) / 2\pi f$$

The average signal power is

$$S_0 = f_s^2 (\Delta v)^2 / 4\pi^2 f^2$$

so that the signal-to-noise ratio for single integration is

$$(S/N)_0 = \frac{3}{2} r^3 (f_m / \pi f)^2$$

where $r = f_s / 2f_m =$ bandwidth expansion factor. The signal-to-noise ratio for double integration is

$$(S/N)_0 = \frac{3}{2} r^5 (f_m / \pi f)^4$$

Thus, the improvement in signal-to-noise ratio varies with f_s^3 for the system with single integration, whereas it varies with f_s^5 for double integration.

24 Digital Communications

Michael B. Pursley

Baseband Signal Sets 24-2

Signal Sets for RF Channels 24-4

- Amplitude-Shift Keying (ASK)
- Binary Phase-Shift Keying (BPSK)
- Quadriphase Shift Keying (QPSK)
- Offset Quadriphase Shift Keying (OQPSK)
- Minimum-Shift Keying (MSK)
- Binary Frequency-Shift Keying (FSK)
- Continuous-Phase Frequency-Shift Keying (CPFSK)
- M*-ary Signaling and Multiple Frequency Shift Keying (MFSK)
- Spectra of PSK and MSK Signals

Optimum Receivers for Digital Communications 24-13

- Receivers for Binary Baseband Data Transmission
- Coherent Receivers for Binary and Quaternary RF Signals
- Noncoherent Receivers

Error Probabilities for Digital Communication Systems 24-18

- Coherent Systems With Additive White Gaussian Noise Channels
- Noncoherent Systems With Additive White Gaussian Noise Channels
- Noncoherent Systems With Nonselective Fading Channels

Spread-Spectrum Communications 24-24

- Direct-Sequence Spread-Spectrum Communications
- Frequency-Hop Spread-Spectrum Communications

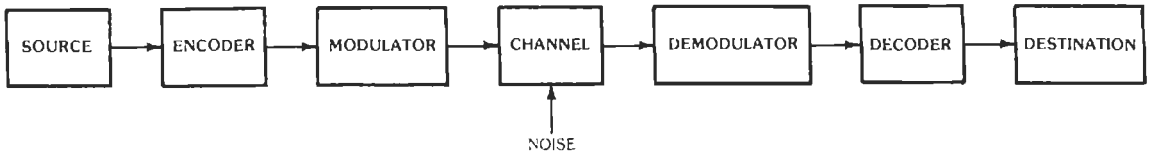


Fig. 1. General model for a communication system.

A general model for a digital communication system is illustrated in the block diagram of Fig. 1. The source, encoder, and modulator are part of the transmitter, and the demodulator, decoder, and destination are in the receiver. In this chapter, we are primarily concerned with the *digital data channel*, which consists of the modulator, channel, and demodulator; the encoder and decoder are discussed in Chapter 25. To focus attention on the digital data channel, consider the model of Fig. 2. In this model, the source and the encoder are combined into a single element called the message source, and the decoder is absorbed into the destination.

The information to be conveyed to the destination consists of a sequence of elements called *messages*. The *message set* is a set of M elements that are indexed by the integers $0, 1, \dots, M - 1$. The modulator produces signals that are used to represent these messages on the channel. During each transmission interval, the message source produces one of the M messages, and the corresponding waveform is transmitted over the channel to the receiver.

Rather than attempting to reproduce the transmitted waveform (as in an analog communication system), the goal of the demodulator in a digital communication system is to determine which of the M messages was sent. It accomplishes this by processing the received signal, which is a distorted, noisy version of the transmitted signal. The demodulator first produces a real number or vector called the *decision statistic*. This statistic is the input to a decision device (e.g., a threshold device), which bases its decision on a predetermined strategy referred to as a *decision rule*. Decision rules are selected according to various criteria such as average probability of error.

The set of allowable decisions can be different from the message set, such as when we permit the receiver to

erase symbols. For example, the binary erasure channel has a message set $\{0, 1\}$, but the set of allowable decisions is $\{0, 1, e\}$, where “ e ” denotes the erasure. This is in contrast to the binary “hard-decision” channel, which has the set $\{0, 1\}$ as both the message set and the set of allowable decisions.

The number of waveforms available to represent the messages is typically reflected in the names of the elements associated with the communications system (including the system itself). For $M = 2$, the signal sets, modulation techniques, and communications systems are referred to as *binary*; for $M = 3$, they are called *ternary*; and for $M = 4$, they are called *quaternary*. In this chapter, the term *M-ary* refers to nonbinary signaling (i.e., it is implicit that $M > 2$).

All of the digital modulation schemes described here use time-limited waveforms. A waveform $v(t)$ is *time limited* if $v(t) = 0$ for values of t outside some finite interval, and we say it is of *duration* T if this interval is $[0, T]$. The energy is perhaps the most important parameter of the transmitted signal. In general, the *energy* E_v in the waveform $v(t)$ is

$$E_v = \int_{-\infty}^{\infty} v^2(t) dt$$

If $v(t)$ has duration T ,

$$E_v = \int_0^T v^2(t) dt$$

BASEBAND SIGNAL SETS

A baseband signal set is a collection of M baseband signals $\{s_k : k = 0, 1, \dots, M - 1\}$, each of which is of duration T . The k th signal $s_k(t)$ represents the message

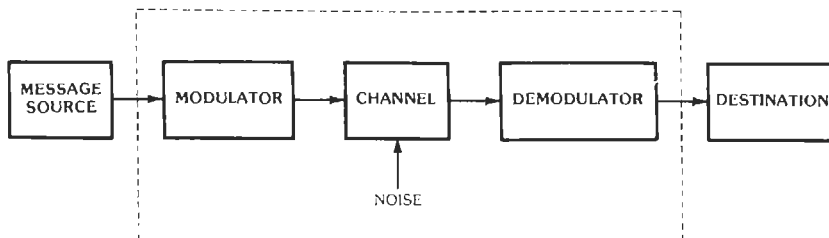


Fig. 2. Model for a digital data channel.

k , and a sequence of messages is sent to the receiver by transmitting the corresponding sequence of signals.

The energy in the k th signal is denoted by E_k . If all of the signals in the set have the same energy, the common value is denoted by E . Given two signals $s_k(t)$ and $s_j(t)$, their *inner product* (s_k, s_j) is defined by

$$(s_k, s_j) = \int_{-\infty}^{\infty} s_k(t)s_j(t)dt$$

Two types of signal sets are of particular interest: orthogonal signal sets and binary antipodal signal sets. In general, a signal set is said to be an *orthogonal* set if $(s_k, s_j) = 0$ for all $k \neq j$. A binary signal set is *antipodal* if $s_0(t) = -s_1(t)$ for all t in the interval $[0, T]$. Antipodal signals have equal energy E , and their inner product is $(s_0, s_1) = -E$.

Many of the baseband signal sets of interest can be defined in terms of a single waveform $v(t)$. To describe such sets, it is convenient to define a set of data variables $\{d_k : k = 0, 1, \dots, M - 1\}$ which are real numbers that represent the messages. Given a set of data variables and a pulse waveform $v(t)$ of duration T , a baseband signal set can then be defined as follows.

The k th signal in the set is given by

$$s_k(t) = Ad_kv(t)$$

where the constant A represents the signal amplitude due to amplification and attenuation (e.g., antenna gains and propagation losses). In many cases, A is the peak signal voltage at the receiver. Since $v(t)$ has duration T , each signal in the set $\{s_k : k = 0, 1, \dots, M - 1\}$ also has duration T . The energy in the k th signal is

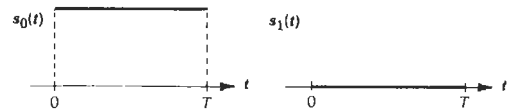
$$E_k = A^2d_k^2E_v$$

Each antipodal signal set can be defined in this way with $M = 2$, $d_1 = -d_0$, and $E_1 = E_0$.

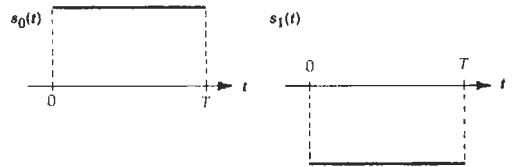
Four examples of binary baseband signal sets are illustrated in Fig. 3. Each of these signal sets can be obtained from a single waveform $v(t)$. For the signal sets shown in Figs. 3A and B, the waveform is a rectangular pulse of duration T . This waveform arises so frequently that it is convenient to give it a special notation. We denote the *rectangular pulse* of duration T by $p_T(t)$; it is defined by

$$p_T(t) = \begin{cases} 1 & \text{for } 0 \leq t < T \\ 0 & \text{otherwise} \end{cases}$$

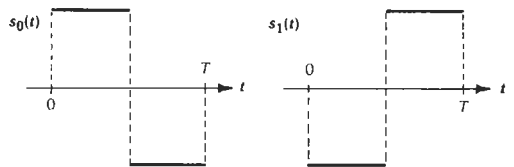
Thus for the signal sets shown in Figs. 3A and B, $v(t) = p_T(t)$. The basic waveform for the baseband signals shown in Fig. 3C is the split-phase or Manchester-coded pulse. For the signal set shown in Fig. 3D, the basic waveform is the *sine pulse*, which is one-half of a full period of a sine wave. That is, $v(t) = \sin(\pi t/T)$ for $0 \leq t < T$.



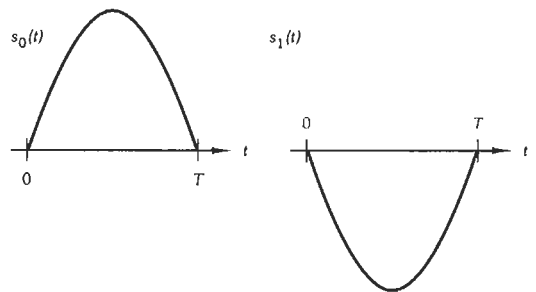
(A) Basic waveform: rectangular pulse.



(B) Basic waveform: rectangular pulse.



(C) Basic waveform: split-phase pulse.



(D) Basic waveform: sine pulse.

Fig. 3. Four examples of baseband signal sets.

The data symbols for the signal set of Fig. 3A are $d_0 = 1$ and $d_1 = 0$, so this is an orthogonal signal set. For the other three signal sets, $d_0 = +1$ and $d_1 = -1$; thus each of these is an antipodal signal set.

There are binary and M -ary baseband signal sets that cannot be described in terms of a single waveform of duration T . If M is a power of 2, for example, a set of M orthogonal signals can be obtained by letting the signals be sequences of pulses (each pulse is of duration T/M) with amplitudes determined by the rows of an M by M Hadamard matrix.* Other M -ary orthogonal signal sets of this type can be designed from maximal-length linear feedback shift-register sequences. A biorthogonal signal set* consists of the signals from an orthogonal set

* Reference 3.

together with all of their negatives. Two signals from a biorthogonal set are either orthogonal or antipodal.

Of course, the transmitter does not send a single message; it sends a sequence of messages. In order to send a sequence of messages, a corresponding sequence of signals must be transmitted. The message k is sent in the n th time interval by transmitting s_k during this interval. The baseband modulation process is described mathematically as follows. If m_n denotes the message produced by the source during the n th interval $[nT, (n + 1)T]$, the sequence of messages to be sent to the receiver is

$$(m_n) = \dots, m_{-1}, m_0, m_1, m_2, \dots$$

If the source output in the n th interval is k , $m_n = k$ and the transmitted signal $s(t)$ is equal to $s_k(t - nT)$ for $nT \leq t < (n + 1)T$. The signal $s(t)$ is actually a sequence of replicas of the basic waveforms, which suggests the representation

$$s(t) = \sum_{n=-\infty}^{\infty} s_{m_n}(t - nT)$$

Because the signals available for use in other intervals are simply replicas of those for $[0, T]$, most of the key features of a particular baseband modulation technique can be described and analyzed by considering only the signal set $\{s_k: k = 0, 1, \dots, M - 1\}$, which is defined on the interval $[0, T]$.

For signal sets that are derived from a single waveform $v(t)$ on $[0, T]$, the above description can be simplified by taking advantage of the fact that for such signal sets

$$s_k(t) = Ad_k v(t)$$

The transmitted signal is defined as follows. First, the message sequence (m_n) is represented by the corresponding sequence (b_n) of data variables; the latter sequence is defined by $b_n = d_k$ if $m_n = k$. Next, the transmitted signal can be written as

$$s(t) = \sum_{n=-\infty}^{\infty} Ab_n v(t - nT)$$

It follows that the signal set can be described and analyzed by considering only the waveform $v(t)$ and the set of data variables $\{d_k: k = 0, 1, \dots, M - 1\}$.

As an example, consider the waveform $v(t) = \sin(\pi t/T)$, which is shown in Fig. 3D. For binary communications, the message set is $\{0, 1\}$. In order to use antipodal signals, we can let $d_0 = +1$ and $d_1 = -1$, so that a binary 0 is transmitted as a positive sine pulse and a binary 1 is transmitted as a negative sine pulse. Message sequences for which $m_0 = 0, m_1 = 1, m_2 = 1$, and $m_3 = 0$ correspond to data sequences with $b_0 = +1, b_1 = -1, b_2 = -1$, and $b_3 = +1$. The resulting signal $s(t)$ is shown in Fig. 4.

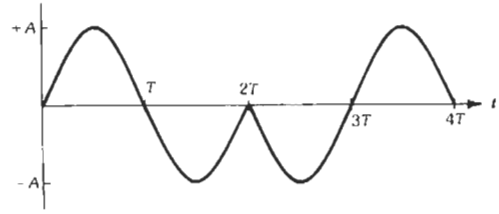


Fig. 4. Transmitted signal $s(t)$ for the message sequence 0110.

SIGNAL SETS FOR RF CHANNELS

The most general RF signals that are considered in this section are of the form

$$s(t) = A a(t) \cos[2\pi f_c t + \theta(t) + \phi]$$

where $a(t)$ and $\theta(t)$ are baseband signals of the type described in the preceding section. The signal $a(t)$ represents *amplitude modulation*, and $\theta(t)$ represents *phase modulation*. The variable A is a real number that represents the signal amplitude at the receiver. Physically, the parameter A may be the peak signal voltage, such as when $|a(t)| = 1$. The variable ϕ is the phase of the RF carrier at the time $t = 0$ in the absence of any phase modulation (i.e., $\theta(t) = 0$). This phase angle may or may not be known to the receiver.

A variety of RF signal sets can be obtained by employing the baseband waveforms of Fig. 3, or other elementary waveforms, as amplitude and phase modulation in the formulation given above. The most common RF signal sets that can be generated in this manner are described below.

The signal energy values given in the descriptions that follow are based on the assumption that either $2f_c T$ is an integer and both $|a(t)|$ and $\theta(t)$ are constant on $[0, T]$ (in which case the value is exact) or $f_c T \gg 1$ (in which case the value is a good approximation). These conditions imply that the double-frequency component of the square of the signal is negligible. It is also assumed that the signals $a(t)$ and $\theta(t)$ use baseband waveforms of duration T . For such signals the *energy per data pulse* is defined as follows. If the k th message is being sent in the interval $[0, T]$, the energy is

$$E_k = \int_0^T s^2(t) dt$$

For several of the modulation schemes, $\theta(t)$ is equal to a constant θ_k for $0 \leq t < T$ if the k th message is sent during the interval $[0, T]$. In this case, the energy is

$$E_k = (A^2/2) \int_0^T a^2(t) dt$$

Because of the assumption concerning $f_c T$, the energy does not depend on the value of θ_k . It does depend on

the energy in the waveform $a(t)$, so E_k may depend on the message k that is being sent.

Amplitude-Shift Keying (ASK)

The signal structure for amplitude-shift keying (ASK) is

$$s(t) = A a(t) \cos(2\pi f_c t + \phi)$$

where $a(t)$ is a sequence of baseband pulses whose amplitudes are modulated to represent the messages. This type of signaling is also commonly referred to as digital AM or multi-amplitude signaling. The baseband signal $a(t)$ can be written as

$$a(t) = \sum_{n=-\infty}^{\infty} b_n v(t - nT)$$

where $b_n = d_k$ if the k th message is being sent during the n th interval. The set $\{d_k; k = 0, 1, \dots, M-1\}$ of data variables is the set of amplitudes for the sequence of pulses. The simplest case is binary ASK modulation with $d_0 = 1$ and $d_1 = 0$; this is known as *on-off keying* (OOK).

Another important special case results if $M = 2$, $v(t)$ is the rectangular pulse of duration T , $d_0 = +1$, and $d_1 = -1$. This is just the amplitude-modulation representation for BPSK (see next subsection).

In general, the energy per data pulse for ASK depends on the message being sent. For the k th message, the energy is

$$E_k = (A d_k)^2 E_v / 2$$

where E_v is the energy in the waveform $v(t)$, as defined at the beginning of this chapter.

A modification of ASK is *quadrature ASK* (QASK), which is also known as quadrature AM (QAM). The QASK signal is of the form

$$s(t) = A \{a_1(t) \cos[2\pi f_c t + \phi] + a_2(t) \sin[2\pi f_c t + \phi]\}$$

The baseband signals $a_1(t)$ and $a_2(t)$ are sequences of pulses of duration T_s with amplitudes from the set $\mathcal{D} = \{d_k; k = 0, 1, \dots, M-1\}$. The QASK signal $s(t)$ can be written as

$$s(t) = s_I(t) + s_Q(t)$$

where

$$s_I(t) = A a_1(t) \cos[2\pi f_c t + \phi]$$

is the *in-phase* component of $s(t)$ and

$$s_Q(t) = A a_2(t) \sin[2\pi f_c t + \phi]$$

is the *quadrature* component. Each of these two components of the QASK signal is an ASK signal with pulse duration T_s .

The simplest, commonly used QASK signal set is 4-QASK in which the set \mathcal{D} contains only the two values $+1$ and -1 . The in-phase and quadrature signals for 4-QASK are each binary antipodal ASK signals. If $a_1(t)$ and $a_2(t)$ are sequences of binary pulses derived from a single binary source (e.g., odd-numbered bits go to the in-phase channel and even-numbered bits go to the quadrature channel), and if the source produces binary digits at the rate of one bit every T seconds, then $T_s = 2T$ for 4-QASK. The in-phase and quadrature binary ASK signals each have a data rate of $1/T_s$ bits per second, so the total data rate for the 4-QASK signal is $2/T_s = 1/T$ bits per second. A single binary ASK signal would have to use pulses of duration $T = T_s/2$ in order to provide the same data rate, but such a binary ASK signal would therefore require twice the bandwidth of the 4-QASK signal.

If each of these two signals $a_1(t)$ and $a_2(t)$ is allowed to take on more than two values, higher-order QASK signal sets are obtained. For example, letting $\mathcal{D} = \{-3, -1, +1, +3\}$ gives the 16-QASK signal set. The 16-QASK signal set is illustrated in Fig. 5, where points in signal space are determined by the amplitudes for the pulses in the in-phase and quadrature components of the signal. This signal set provides a data rate of $4/T_s$ bits per second, and so a binary ASK signal set requires four times the bandwidth of 16-QASK in order to provide the same data rate. Greater bandwidth savings relative to binary ASK can be achieved by letting \mathcal{D} be larger, and this is now a standard approach for such applications as high-speed modems* for telephone lines. The application of trellis-coded modulation† has further improved the performance of such modems.

Binary Phase-Shift Keying (BPSK)

One of the most commonly used binary signal sets is obtained by shifting the phase of the RF carrier by $+\pi/2$ radians or $-\pi/2$ radians, depending on whether the data bit is a 0 or a 1. Binary phase-shift keying can be viewed as binary phase modulation or binary amplitude modulation; the only requirement is that, during each signaling interval, $s(t)$ is one of two sinusoidal signals that differ in phase by π radians.

Suppose that the binary digit k is to be sent during the n th time interval. Viewed as phase modulation, the corresponding BPSK signal is

$$A \cos[2\pi f_c t + d_k(\pi/2) + \phi]$$

for $nT \leq t < (n+1)T$ where $|d_0 - d_1| = 2$. The two most common choices for d_k are $d_0 = 0$ and $d_1 = 2$ or

* Reference 19.

† References 18 and 23.

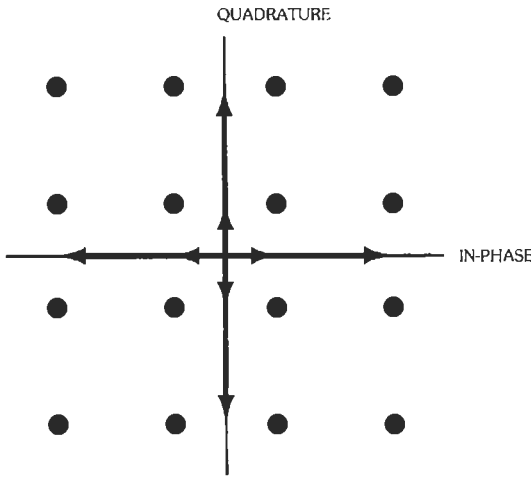


Fig. 5. A 16-QASK signal constellation.

$d_0 = +1$ and $d_1 = -1$. Based on the latter choice for d_k , the BPSK signal can be defined by

$$s(t) = A a(t) \cos[2\pi f_c t + \theta(t) + \phi]$$

where $a(t) = 1$ for all t and

$$\theta(t) = (\pi/2) \sum_{n=-\infty}^{\infty} b_n p_T(t - nT)$$

The data symbol b_k is d_k if k is to be sent in the n th time interval. This formulation corresponds to phase modulation by a sequence of rectangular pulses of duration T with amplitudes $(\pi/2)b_n$. The resulting signal $s(t)$ is illustrated in Fig. 6A for the data sequence $b_0 = -1$, $b_1 = +1$, and $b_2 = +1$. This illustration is for $f_c T = 2$ and $\phi = 0$.

Viewed as binary amplitude modulation, the BPSK signal can be expressed in a form that reveals its connection with binary ASK. Let $\theta(t) = 0$ for all t and use the amplitude to convey the information by defining

$$a(t) = \sum_{n=-\infty}^{\infty} b_n p_T(t - nT)$$

Thus, 0 is sent by transmitting

$$s(t) = A \cos[2\pi f_c t + \phi]$$

and 1 is sent by transmitting

$$s(t) = -A \cos[2\pi f_c t + \phi]$$

This formulation of binary PSK corresponds to amplitude modulation by a sequence of rectangular pulses of duration T with amplitudes $+1$ and -1 . However, it is equivalent to phase modulation of the form

$$s(t) = A \cos[2\pi f_c t + d_k(\pi/2) + \phi]$$

with $d_0 = 0$ and $d_1 = 2$.

Although the modulator implementations suggested by these two formulations are quite different, the resulting signal sets are mathematically equivalent. In particular, it is clear from either of these formulations that the BPSK signals representing the binary digits 0 and 1 form an antipodal signal set. Notice also that for each of these representations the energy per data pulse is $E = A^2 T/2$.

Quadrature Phase Shift Keying (QPSK)

Binary PSK can be generalized by allowing the phase waveform to take on more than two values. The QPSK signal is of the same basic form as BPSK, except that quaternary modulation is employed instead of binary modulation. Thus the transmitted signal for the interval $[nT, (n + 1)T]$ can be written as

$$s(t) = A \cos[2\pi f_c t + \theta_k + \phi]$$

where k denotes the message to be sent in the n th interval ($k = 0, 1, 2, \text{ or } 3$). The four different values for θ_k must be separated by $\pi/2$ radians. For example, we can let $\{\theta_k: k = 0, 1, 2, 3\}$ be the set $\{0, \pi/2, \pi, 3\pi/2\}$ or the set $\{\pi/4, 3\pi/4, 5\pi/4, 7\pi/4\}$. As with BPSK, the values themselves are unimportant, but the differences are constrained.

The QPSK signal can be expressed as

$$s(t) = A \cos[2\pi f_c t + \theta(t) + \phi]$$

The phase modulation $\theta(t)$ is a sequence of rectangular pulses of duration T and amplitudes $(\pi/2)b_n$ with $b_n = 0, 1, 2, \text{ or } 3$. This signal is illustrated in Fig. 6B for the data sequence $b_0 = 0$, $b_1 = 2$, and $b_2 = 3$. As in Fig. 6A, $f_c T = 2$ and $\phi = 0$.

The QPSK signal can also be viewed as quadrature amplitude modulation by writing

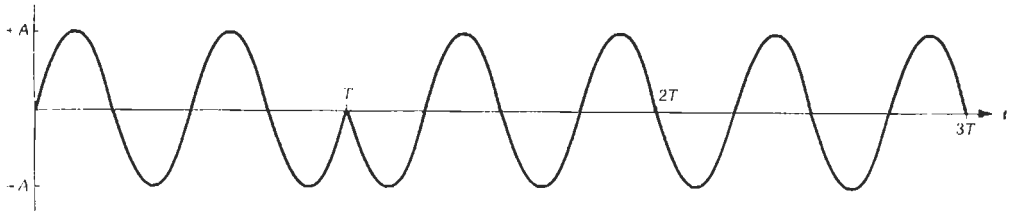
$$s(t) = A \cos[2\pi f_c t + \theta(t) + \phi]$$

as

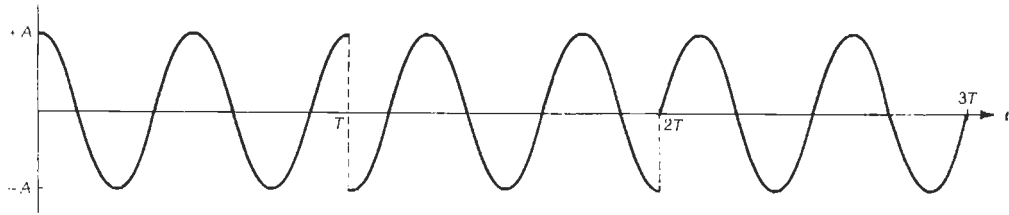
$$s(t) = A \{ \cos[\theta(t)] \cos[2\pi f_c t + \phi] - \sin[\theta(t)] \sin[2\pi f_c t + \phi] \}$$

If $\theta(t)$ takes on values in the set $\{\pi/4, 3\pi/4, 5\pi/4, 7\pi/4\}$, then $\cos[\theta(t)]$ is either $+1/\sqrt{2}$ or $-1/\sqrt{2}$, and $\sin[\theta(t)]$ takes on the same two values. Thus $s(t)$ can be written as

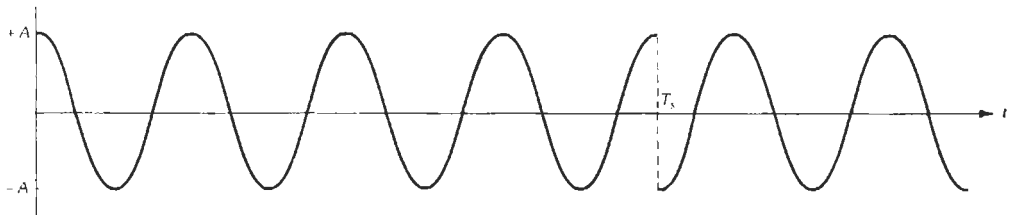
$$s(t) = (A/\sqrt{2}) \{ a_1(t) \cos[2\pi f_c t + \phi] + a_2(t) \sin[2\pi f_c t + \phi] \}$$



(A) Binary PSK, data rate $1/T$ bits per second.



(B) Example of QPSK, data rate $2/T$ bits per second.



(C) Example of QPSK, data rate $1/T$ bits per second.

Fig. 6. Various forms of phase-shift keying.

where $a_1(t)$ and $a_2(t)$ are sequences of positive and negative unit-amplitude rectangular pulses of duration T . This formulation of QPSK clearly exhibits its relationship to QASK. If we define the *in-phase* and *quadrature* signals by

$$s_I(t) = (A/\sqrt{2}) a_1(t) \cos[2\pi f_c t + \phi]$$

and

$$s_Q(t) = (A/\sqrt{2}) a_2(t) \sin[2\pi f_c t + \phi]$$

respectively, the QPSK signal can be written as

$$s(t) = s_I(t) + s_Q(t)$$

Each of the signals $a_1(t)$ and $a_2(t)$ has the same number of pulses per unit of time as the signal $a(t)$ that appears in the amplitude-modulation formulation of BPSK. Consequently, QPSK signaling has twice the data rate of BPSK. This can also be observed from the original quaternary-phase-modulation formulation by noticing that each pulse of duration T takes on one of four values, and this gives $\log_2(4) = 2$ bits of data per pulse.

The quadrature-amplitude-modulation formulation of QPSK displays the signal as the sum of two quadrature binary ASK signals. A QPSK signal can also be viewed as two quadrature BPSK signals. For example, we can write the QPSK signal as

$$s(t) = (A/\sqrt{2}) \{ \cos[2\pi f_c t + \theta_1(t) + \phi] + \sin[2\pi f_c t + \theta_2(t) + \phi] \}$$

where $\theta_1(t)$ and $\theta_2(t)$ take on the values 0 and π .

The power in the QPSK signal is $A^2/2$, just as for BPSK, and so the energy per pulse is $E = A^2T/2$. Since each pulse represents two bits of data, the energy per bit for QPSK is $E_b = A^2T/4$, which is one-half of the energy per bit for BPSK for the same pulse duration T . Because the pulse shape and duration are the same for BPSK and QPSK, the two signals have the same bandwidth.

For most purposes, it is better to compare BPSK and QPSK signals that have *the same energy per data bit and the same data rate* (rather than the same pulse duration as above). To do this, simply replace T by T_s in the QPSK expressions, and let $T_s = 2T$. That is, the signal $s(t)$ is given by

$$s(t) = A \cos[2\pi f_c t + \theta(t) + \phi]$$

as before, but now the phase modulation $\theta(t)$ is a sequence of rectangular pulses of duration $T_s = 2T$, rather than T . The amplitudes of these pulses are still given by $(\pi/2)b_n$ with $b_n = 0, 1, 2$, or 3 ; however, the sequence (b_n) is now transmitted at the rate of one quaternary symbol every T_s seconds. This signal is illustrated in Fig. 6C for $b_0 = 0$ and $b_1 = 2$, with $f_c T = 2$ and $\phi = 0$.

Because a QPSK pulse of duration T_s provides two bits of data, the data rate of a QPSK signal that uses pulses of duration T_s is

$$R_{\text{QPSK}} = 2/T_s = 1/T$$

bits per second. If the digital data consists of a single sequence of binary digits, half of the data bits can be modulated onto the in-phase component and half onto the quadrature component.

Since $1/T$ is also the data rate for a BPSK signal with pulse duration T , these two forms of PSK modulation can now be compared on the basis of the same data rate. In this case, the energy per pulse for QPSK is given by

$$E = A^2 T_s / 2$$

which is equivalent to an energy per bit of

$$E_b = A^2 T / 2$$

the same as BPSK. The pulse shapes for QPSK and BPSK are the same, but the pulse duration for QPSK is $T_s = 2T$, twice that for BPSK. This means the QPSK signal has one-half the bandwidth of the BPSK signal.

In summary, if QPSK and BPSK are compared for the same data rate, BPSK requires twice the bandwidth of QPSK. If the two signals have the same data rate and equal power, they also have the same energy per data bit.

Offset Quadrature Phase Shift Keying (OQPSK)

The OQPSK signal can be defined by replacing $\theta_1(t)$ with $\theta_1(t + T)$ in the quadrature-phase-modulation formulation of QPSK. That is, OQPSK can be represented in terms of two quadrature BPSK signals as

$$s(t) = (A/\sqrt{2}) \{ \cos[2\pi f_c t + \theta_1(t + T) + \phi] + \sin[2\pi f_c t + \theta_2(t) + \phi] \}$$

where $\theta_1(t)$ and $\theta_2(t)$ are sequences of pulses of duration $T_s = 2T$ as in the preceding subsection. Since the amount of the offset is one-half the pulse duration, the phase transitions of the in-phase and quadrature components of the OQPSK signal are separated in time by T seconds.

The OQPSK signal can be represented as quadrature amplitude modulation with the baseband signal of the

in-phase component offset by T seconds relative to the baseband signal of the quadrature component. Thus, the amplitude-modulation representation of OQPSK is

$$s(t) = (A/\sqrt{2}) \{ a_1(t + T) \cos[2\pi f_c t + \phi] + a_2(t) \sin[2\pi f_c t + \phi] \}$$

where $a_1(t)$ and $a_2(t)$ are as in the preceding subsection. Notice that the baseband signal for the in-phase component can change polarity at times nT for odd integers n only, but the transitions of the baseband signal for the quadrature component are at times nT for even values of n only. This is an important feature of OQPSK for band-limited channels with nonlinearities. A band-limited version of QPSK or BPSK has an envelope that may go to zero when the in-phase and quadrature signals switch polarity at the same time; the offset in OQPSK leads to a more nearly constant envelope for the band-limited signal.

If the OQPSK signal is expressed in the form

$$s(t) = A \cos[2\pi f_c t + \theta(t) + \phi]$$

the phase $\theta(t)$ can change every T seconds, whereas the phase shifts in QPSK occur only at $2T$ -second intervals. However, for OQPSK the magnitudes of the phase shifts are limited to 0 or $\pi/2$ only. Phase shifts for QPSK can have magnitudes $0, \pi/2, \text{ or } \pi$.

Minimum-Shift Keying (MSK)

The MSK signal can be considered as a special case of offset quadrature amplitude modulation in which the baseband waveform is the sine pulse $v(t) = \sin(\pi t/2T)p_{2T}(t)$. Thus the MSK signal can be written as

$$s(t) = A \{ a_1(t + T) \cos[2\pi f_c t + \phi] + a_2(t) \sin[2\pi f_c t + \phi] \}$$

where $a_1(t)$ and $a_2(t)$ are sequences of sine pulses of duration $2T$. That is, for each $i(i = 1 \text{ or } i = 2)$

$$a_i(t) = \sum_{n=-\infty}^{\infty} b_{i,n} v(t - 2nT)$$

where $(b_i) = \dots, b_{i,0}, b_{i,1}, b_{i,2}, \dots$ is a binary data sequence. Minimum-shift keying can also be viewed as continuous-phase frequency-shift keying (discussed in a later subsection).

Binary Frequency-Shift Keying (FSK)

Binary FSK signals are of the form

$$s(t) = A \cos[2\pi f_c t + \theta(t)]$$

where the modulation $\theta(t)$ is defined as follows. If 0 is the bit to be transmitted in the n th interval,

$$\theta(t) = 2\pi f_d t + \phi_0$$

and if 1 is sent,

$$\theta(t) = -2\pi f_d t + \phi_1$$

Thus, the signal

$$s(t) = A \cos[2\pi(f_c + f_d)t + \phi_0]$$

represents a 0 for the n th interval, and the signal

$$s(t) = A \cos[2\pi(f_c - f_d)t + \phi_1]$$

represents a 1. The quantity f_d is called the *frequency deviation*, and the parameter $h = 2f_d T$ is called the *deviation ratio* or *modulation index* for the FSK signal set. The FSK signals may be generated by switching between two oscillators or by applying a binary baseband signal at the input of a voltage-controlled oscillator (VCO).

In general, the phase angles ϕ_k are arbitrary, and ϕ_0 need not be related to ϕ_1 in any way. If the FSK signals are obtained by switching between two oscillators, one at frequency $f_c + f_d$ and one at frequency $f_c - f_d$, the phase angles ϕ_0 and ϕ_1 represent the phases of these oscillators at time $t = 0$. Alternatively, the FSK signals may be generated by applying a baseband signal to a VCO (or another frequency-modulation circuit), in which case the phase angles ϕ_0 and ϕ_1 may be related.

The signals in FSK systems are often referred to as ‘‘tones,’’ and the two signals are distinguished by calling one of them ‘‘mark’’ and the other ‘‘space.’’ The convention followed here is to refer to the signal at frequency $f_c + f_d$ as the space and the signal at frequency $f_c - f_d$ as the mark (the binary digit 0 is transmitted as a space, and the binary digit 1 is transmitted as a mark).

If $2f_c T$ and $2f_d T$ are integers, the mark and space signals form an *orthogonal* signal set for *all* values of ϕ_0 and ϕ_1 . It is often the case that $f_c T \gg 1$. Under this condition, the signals are approximately orthogonal if $2fdT$ is an integer or if $f_d T \gg 1$. If $2f_c T$ and $2f_d T$ are integers or if $(f_c - f_d)T \gg 1$, the energy per data bit is $A^2 T/2$.

For certain applications, it is necessary to generalize the above formulation of FSK by introducing a phase angle $\phi(n)$ that depends on the interval in which the signal is transmitted. To send a 0 in the n th interval, the transmitted signal for $nT \leq t < (n + 1)T$ is

$$s(t) = A \cos[2\pi(f_c + f_d)t + \phi(n) + \phi_0]$$

and to send a 1 in this same interval, the signal is

$$s(t) = A \cos[2\pi(f_c - f_d)t + \phi(n) + \phi_1]$$

This generalization is required in order to characterize

FSK signals in frequency-hopped spread-spectrum systems, for example, where the phase $\phi(n)$ represents a phase shift introduced by noncoherent frequency hopping.* This generalization is also required for the description of certain continuous-phase FSK modulation schemes.

Continuous-Phase Frequency-Shift Keying (CPFSK)

Continuous-phase FSK signals are of the form

$$s(t) = A \cos[2\pi f_c t + \theta(t)]$$

just as for other types of FSK signals. However, for CPFSK the phase modulation $\theta(t)$ is a continuous function of t . In order to describe CPFSK modulation, it is sufficient to consider only the case in which $\phi_0 = \phi_1 = 0$. This is because the phase angles ϕ_0 and ϕ_1 can always be absorbed in the phase modulation $\theta(t)$ for CPFSK signals.

The signal $\theta(t)$ is defined as follows. Suppose b_n is the data symbol to be sent in the n th interval, and b_n is either +1 or -1. The phase modulation is

$$\theta(t) = 2\pi b_n f_d t + \phi(n)$$

for $nt \leq t < (n + 1)T$. The phase angles $\phi(n)$ are such that the phase is continuous from one interval to the next. In order to make the phase continuous at time nT , the phase angles $\phi(n)$ and $\phi(n - 1)$ must satisfy

$$2\pi b_n f_d nT + \phi(n) = 2\pi b_{n-1} f_d nT + \phi(n - 1)$$

which is equivalent to the condition

$$\phi(n) = 2\pi(b_{n-1} - b_n)f_d nT + \phi(n - 1)$$

In other words, $\phi(n) = \phi(n - 1)$ if the two successive data symbols b_{n-1} and b_n are the same, but $\phi(n)$ and $\phi(n - 1)$ differ by $4\pi f_d nT$ radians if these two data symbols are different.

The simplest CPFSK signal is obtained for the case in which the modulation index $h = 2f_d T$ is an integer. For this case, $4\pi f_d nT$ is an integer multiple of 2π , and $\phi(n) = \phi(n - 1)$ modulo 2π , regardless of whether the successive data symbols are the same or different. But phase angles that differ by integer multiples of 2π are the same for our purposes, so we can let $\phi(n) = \phi$, where ϕ is a constant phase (independent of n). No phase changes are required in order to make $\theta(t)$ continuous provided the modulation index is an integer.

Another modulation index of considerable interest for CPFSK is 1/2. For $h = 1/2$, the frequency deviation is $f_d = 1/4T$, and the condition for continuous phase is

$$\phi(n) = \pi n[(b_{n-1} - b_n)/2] + \phi(n - 1)$$

* Reference 29.

If the successive data bits b_{n-1} and b_n are the same, the phase angles $\phi(n-1)$ and $\phi(n)$ are also the same. If the successive data bits are different, $b_{n-1} - b_n$ is either +2 or -2, and so the phase angles must differ by an integer multiple of π radians. It follows that if $\phi(0) = 0$ and $2f_c T$ is an integer, or if $f_c T \gg 1$, the mark and space signals are orthogonal. It is important to realize that orthogonality is obtained for $h = 1/2$ only if the phases in successive intervals are controlled.

From the relationship

$$\theta(t) = 2\pi b_n f_d t + \phi(n)$$

which holds for $nT \leq t < (n+1)T$, it follows that

$$\theta(nT) = 2\pi b_n f_d nT + \phi(n)$$

Making the substitution for $\phi(n)$ from the last equation into the expression for $\theta(t)$, we find that

$$\theta(t) = 2\pi b_n f_d (t - nT) + \theta(nT)$$

for $nT \leq t < (n+1)T$. This fact and the continuity condition imply

$$\theta((n+1)T) = 2\pi b_n f_d T + \theta(nT)$$

If we denote $\theta(kT)$ by θ_k and let $h = 1/2$ in this expression, it can be written as

$$\theta_{n+1} = (\pi/2)b_n + \theta_n$$

Similarly, the expression for $\theta(t)$ becomes

$$\theta(t) = (\pi/2T)b_n(t - nT) + \theta_n$$

which is valid for $nT \leq t < (n+1)T$.

The CPFSK signal is

$$s(t) = A \{ \cos[\theta(t)] \cos[2\pi f_c t] - \sin[\theta(t)] \sin[2\pi f_c t] \}$$

The above relationships for $\theta(t)$, θ_n , and θ_{n+1} imply that if n is an even integer and $h = 1/2$,

$$-\sin \theta(t) = -\sin[\theta_{n+1}] \sin[\pi(t - nT)/2T]$$

for $nT \leq t < (n+2)T$. Similarly, if n is an odd integer and $h = 1/2$,

$$\cos \theta(t) = -\cos[\theta_{n+1}] \sin[\pi(t - nT)/2T]$$

for $nT < t < (n+2)T$. Define the data sequence (B_n) by

$$B_n = -\sin \theta_{n+1}$$

if n is an even integer and by

$$B_n = -\cos \theta_{n+1}$$

if n is an odd integer. The data sequence (B_n) is related to the original data sequence (b_n) by

$$b_n = -B_n \cos \theta_n$$

if n is an even integer and by

$$b_n = B_n \sin \theta_n$$

if n is an odd integer.

From the above development, we conclude that for $h = 1/2$ the CPFSK signal $s(t)$ is of the form

$$s(t) = A \{ a_1(t+T) \cos[2\pi f_c t] + a_2(t) \sin[2\pi f_c t] \}$$

The baseband signals $a_1(t)$ and $a_2(t)$ are given by

$$a_1(t) = B_{2k-1} \sin[\pi(t - 2kT)/2T]$$

and

$$a_2(t) = B_{2k} \sin[\pi(t - 2kT)/2T]$$

for $2kT \leq t < 2(k+1)T$. A comparison of these last three equations with analogous equations in the subsection on minimum-shift keying shows that CPFSK with modulation index $h = 1/2$ is the same as MSK.

M-ary Signaling and Multiple Frequency Shift Keying (MFSK)

The general M -ary RF signal set has signals of the form

$$s_k(t) = A a_k(t) \cos[2\pi f_c t + \theta_k(t) + \phi_k]$$

for $k = 0, 1, \dots, M-1$. The signals $a_k(t)$ and $\theta_k(t)$ are baseband signals with bandwidths much smaller than the carrier frequency f_c . The messages are indexed by the integers $0, 1, \dots, M-1$, and in order to send the k th message during the n th interval, the signal $s(t) = p_T(t - nT) s_k(t)$ is transmitted. (The rectangular pulse of duration T is defined by $p_T(u) = 1$ for $0 \leq u < T$ and $p_T(u) = 0$ otherwise.)

The RF signal set is characterized by the inner products

$$(s_k, s_j) = \int_0^T s_k(t) s_j(t) dt$$

Since the amplitude modulation $a_k(t)$ and phase modulation $\theta_k(t)$ are narrow-band signals, the inner products are (at least approximately) given by

$$(s_k, s_j) = A^2 \int_0^T a_k(t) a_j(t) \cos[\theta_{k,j}(t) + \phi_{k,j}] dt$$

where

$$\theta_{k,j}(t) = \theta_k(t) - \theta_j(t)$$

and

$$\phi_{k,j} = \phi_k - \phi_j$$

If $(s_k, s_j) = 0$ for all $j \neq k$, the signals form an *orthogonal* signal set. One type of orthogonal signal set is obtained by letting $\theta_k(t) = 0$ for all t and using a set of orthogonal baseband signals for the amplitude modulation $a_k(t)$.

Multiple frequency shift keying (MFSK), which is also known as M -ary frequency shift keying, is a direct extension of binary FSK to modulation with more than two frequency tones. It is also a special case of M -ary signaling in which $a_k(t) = 1$ and

$$\theta_k(t) = (f_k - f_c)t$$

The resulting signals are of the form

$$s_k(t) = A \cos[2\pi f_k t + \phi_k]$$

In order to send the k th message in the n th interval, the transmitted signal is

$$s(t) = A p_T(t - nT) \cos[2\pi f_k t + \phi_k]$$

If the frequencies f_k are such that $f_k T = n_k$ for some integers n_0, n_1, \dots, n_{M-1} , the M different signals are orthogonal on each interval $[nT, (n + 1)T]$. The signals are at least approximately orthogonal if $|f_k - f_j|T \gg 1$ for each choice of $j \neq k$. The phases $\phi_0, \phi_1, \dots, \phi_{M-1}$ are arbitrary and, in general, not related to each other in any way.

Spectra of PSK and MSK Signals

By introducing appropriate random time delays and phase angles, we can model the signal $s(t)$ as a wide-sense stationary random process. The Fourier transform of the autocorrelation function for this random process is called the *power spectral density*, or simply the *spectrum*. The spectrum of a signal gives a measure of the distribution of its power as a function of frequency. For example, the bandwidth of a signal is usually defined in terms of its spectrum (e.g., see reference 15). The four digital modulation techniques considered are three forms of phase-shift keying (PSK) and MSK. The three PSK schemes are BPSK, QPSK, and OQPSK. The data rate in each case is assumed to be $R = 1/T$ bits per second.

The spectrum of the BPSK signal is

$$S_1(f) = (A^2 T/4)\{G_1[(f - f_c)T] + G_1[(f + f_c)T]\}$$

where the function G_1 is defined by

$$G_1(x) = [\text{sinc}(x)]^2 = [\sin(\pi x)/\pi x]^2$$

The spectrum of a QPSK signal with data rate $R = 1/T$ bits per second is given by

$$S_2(f) = (A^2 T/4)\{G_2[(f - f_c)T] + G_2[(f + f_c)T]\}$$

where

$$G_2(x) = 2 G_1(2x)$$

The spectrum of the OQPSK signal is the same as the spectrum of the QPSK signal. The spectrum of the MSK signal is

$$S_3(f) = (A^2 T/4)\{G_3[(f - f_c)T] + G_3[(f + f_c)T]\}$$

where the function G_3 is defined by

$$G_3(x) = \{[(4/\pi) \cos(2\pi x)]/(1 - 16x^2)\}^2$$

The power in each of the PSK and MSK signals is $P = A^2/2$. Since the power is equal to the integral of the power spectral density,

$$\int_{-\infty}^{\infty} S_i(f) df = (A^2/2) \int_{-\infty}^{\infty} G_i(x) dx = P$$

for each i . The comparisons below of the spectra of PSK and MSK signals are made on the basis of equal power in the signals. It is also common to compare the signals on the basis of equal values of $S_i(f_c)$, the spectral density at the carrier frequency (e.g., reference 15), but in this case the signals have different power levels.

Notice that the spectra $S_1(f)$, $S_2(f)$, and $S_3(f)$ are given in the form

$$S_i(f) = (A^2 T/4)\{G_i[(f - f_c)T] + G_i[(f + f_c)T]\}$$

so that comparisons of these spectra can be made by comparing the corresponding functions $G_1(x)$, $G_2(x)$, and $G_3(x)$. In fact, if $f_c T \gg 1$ the two components of the spectrum do not overlap significantly; that is, for each integer i ,

$$G_i[(f - f_c)T] \gg G_i[(f + f_c)T]$$

for $f > 0$, and the reverse inequality holds for $f < 0$. It follows that $S_i(f)$ is approximately

$$S_i(f) = G_i[(f - f_c)T]$$

for $f > 0$ and

$$S_i(f) = G_i[(f + f_c)T]$$

for $f < 0$. Consequently, for $f_c T \gg 1$, the shape of spectrum $S_i(f)$ for $f > f_c$ is just the shape of $G_i(x)$ for $x > 0$. Also notice that $G_i(x) = G_i(-x)$.

The spectra of BPSK, QPSK, OQPSK, and MSK are compared in Figs. 7 and 8. In Fig. 7, $G_i(x)$ is shown as a

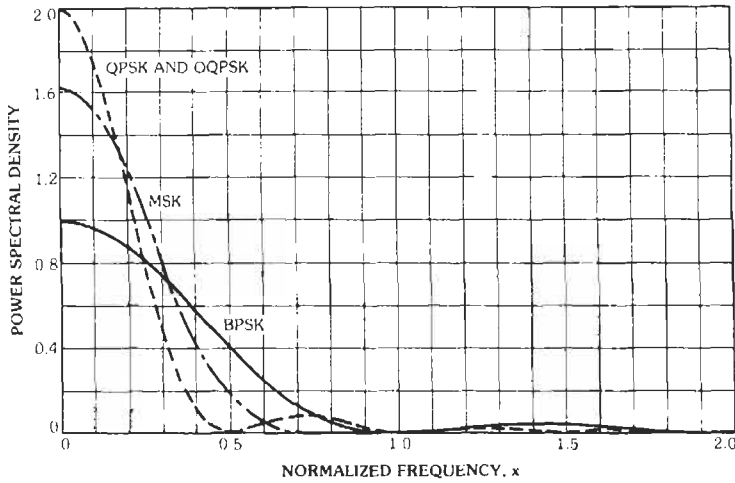


Fig. 7. Power spectral densities for PSK and MSK.

function of x for all three values of i . The parameter x can be thought of as *normalized frequency*, since it is frequency (in hertz) divided by the data rate R (in bits per second). Note that x is therefore a dimensionless parameter.

Specifications on signal spectra (e.g., bandwidths) are often given in terms of the power spectral densities in decibels (dB). Fig. 8 shows $G_i(x)$ in decibels as a function of x . Since $G_i(x)$ is a power density, the conversion is $10 \log_{10}[G_i(x)]$. The *3-dB bandwidth* of $S_i(f)$ is the width of a frequency band outside of which the power spectral density is at least 3 dB below $S_i(f_c)$, its value at the center of the band. Typically, bandwidths are normalized by dividing by R , which is

equivalent to specifying the bandwidth in terms of the normalized frequency. Since $10 \log_{10}(0.5)$ is approximately -3 dB, the 3-dB bandwidth is also termed the *half-power bandwidth*. In terms of the function G_i , the 3-dB, or half-power, bandwidth is the largest value of y for which $G_i(y/2) = 0.5 G_i(0)$.

Analogous definitions are given for other power-density levels (e.g., the 10-, 35-, or 50-dB bandwidths) in reference 15. The 3-dB bandwidths for the PSK and MSK signals are 0.88 for BPSK, 0.44 for QPSK and OQPSK, and 0.59 for MSK; and the 35-dB bandwidths are 35.12, 17.56, and 3.24, respectively.* The *null-to-*

* Reference 15.

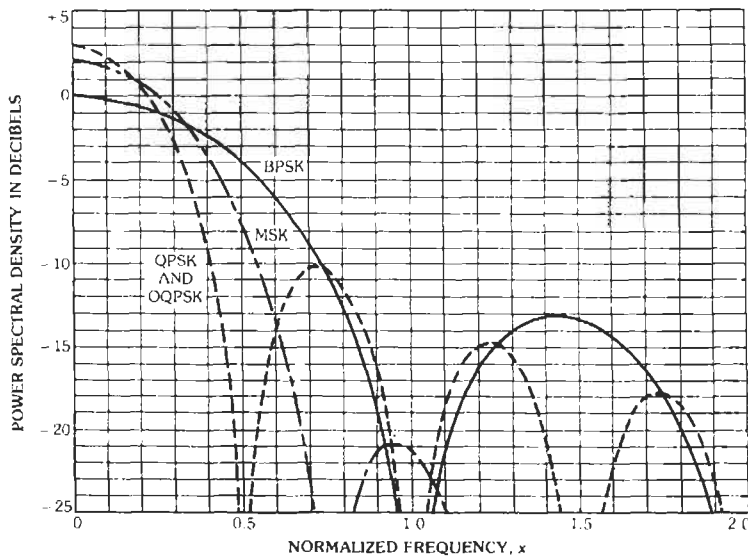


Fig. 8. Power spectral densities in decibels.

null bandwidth is the smallest value of γ for which $G_i(\gamma/2) = 0$. As can be seen in Figs. 7 and 8, the null-to-null bandwidths are 2.0 for BPSK, 1.0 for QPSK and OQPSK, and 1.5 for MSK. Notice that although MSK has lower side lobes and smaller 35-dB bandwidth than QPSK, QPSK has a narrower main lobe (i.e., smaller null-to-null bandwidth) and a smaller 3-dB bandwidth. These and related issues are discussed in greater detail in references 15 through 22.

OPTIMUM RECEIVERS FOR DIGITAL COMMUNICATIONS

Receivers for Binary Baseband Data Transmission

The general model for a binary baseband data transmission system with an additive Gaussian noise channel is shown in Fig. 9. The model of Fig. 9 is employed for coherent RF communications also, because it is the low-pass equivalent of a binary RF digital communications system with an additive Gaussian noise channel and a coherent receiver. The block diagram for the RF system is shown in Fig. 10. The input,

$$s(t) = \sqrt{2}s_k(t) \cos[2\pi f_c t + \phi]$$

for the RF channel corresponds to the input $s_k(t)$ for the baseband channel.

The principal elements of the channel are a binary signal set $\{s_0, s_1\}$, an additive Gaussian noise process $X(t)$, a linear filter with impulse response $h(t)$, a sampler, and a threshold device. The binary digit k is sent during the n th time interval by transmitting the signal $s_k(t - nT)$ for $nT < t < (n + 1)T$. For convenience, we consider the first time interval (i.e., $n = 0$) in all that follows. For this time interval, the channel output $Y(t)$ is the sum of the signal $s_k(t)$ plus the noise $X(t)$.

The channel output is the input to the linear time-invariant filter. The filter output $Z(t)$ is the sum of a filtered version of the signal, which is denoted by $\hat{s}_k(t)$, and a filtered version of the noise. The signal component of the filter output is given by

$$\hat{s}_k(t) = \int_{-\infty}^{\infty} s_k(u) h(t - u) du$$

the *convolution* of the functions s_k and h , which we denote by $\hat{s}_k = s_k * h$.

The output of the filter is sampled at time T_0 to give the decision statistic $Z(T_0)$. Our notational convention is

to let s_0 denote the signal that produces the larger output at sampling time T_0 ; that is, we assume

$$\hat{s}_0(T_0) > \hat{s}_1(T_0)$$

In the threshold device, the decision statistic $Z(T_0)$ is compared with a threshold z , and the decision is that s_0 was transmitted if the statistic $Z(T_0)$ is greater than the threshold z . If $Z(T_0)$ is less than z , the decision is that the signal s_1 was transmitted. This system always makes the correct decision in the absence of noise provided that the threshold is in the range

$$\hat{s}_0(T_0) > z > \hat{s}_1(T_0)$$

The channel noise $X(t)$ is a wide-sense stationary random process with autocorrelation function given by

$$R_X(u) = E\{X(t)X(t + u)\}$$

where $E\{Y\}$ denotes the expected value of a random variable Y . Without loss of generality, the channel noise process can be assumed to have zero mean (the mean value can always be absorbed in the signal set). The autocorrelation function for the noise at the output of the filter is given by

$$R(u) = \int_{-\infty}^{\infty} R_X(v)f(u - v) dv$$

where the function f is defined by

$$f(t) = \int_{-\infty}^{\infty} h(v)h(v + t) dv$$

These two expressions are equivalent to $R = R_X * f$ and $f = h * g$, respectively, where the function g is defined by $g(t) = h(-t)$. That is, g is the time reverse of the impulse response of the filter. Since $f(t) = f(-t)$, $h(v + t)$ can be replaced by $h(v - t)$ in the definition of the function f , but the form given above is usually easier to evaluate.

Let $P_{e,k}$ denote the probability that an error is made by the receiver when the signal $s_k(t)$ is the input to the channel. The error probability $P_{e,k}$ can be thought of as a *conditional probability of error* given that $Y(t) = s_k(t) + X(t)$. That is, if H_k denotes the event that the signal $s_k(t)$ is transmitted, the probability $P_{e,k}$ is the conditional probability of error given H_k .

When the signal $s_0(t)$ is transmitted, the receiver makes an error if $Z(T_0)$ is less than the threshold z .

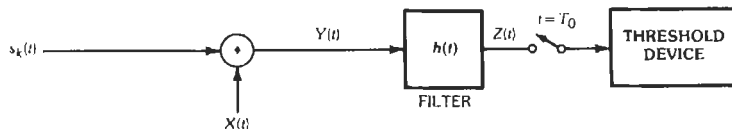


Fig. 9. Model of binary baseband data transmission system.

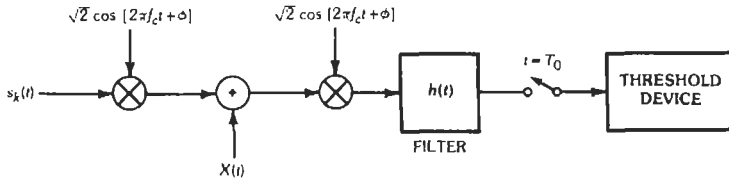


Fig. 10. Model of binary coherent RF communication system.

Thus,

$$P_{e,0} = P[Z(T_0) < z | H_0]$$

Given that $s_0(t)$ is the transmitted signal, $Z(T_0)$ is a Gaussian random variable that has mean $\hat{s}_0(T_0)$ and variance $\sigma^2 = R(0)$. Consequently,

$$P_{e,0} = Q[\{\hat{s}_0(T_0) - z\}/\sigma]$$

where the function Q is defined by

$$Q(x) = \int_x^\infty (2\pi)^{-1/2} \exp(-y^2/2) dy$$

The function $1 - Q(x)$ is the distribution function of a zero mean, unit variance, Gaussian random variable. A discussion of the function Q and its properties can be found in reference 11.

Similarly, if $s_1(t)$ is transmitted, $Z(T_0)$ is a Gaussian random variable with mean $\hat{s}_1(T_0)$ and variance $\sigma^2 = R(0)$, and an error occurs if $Z(T_0)$ is greater than the threshold z . Therefore, the error probability $P_{e,1}$ is

$$P_{e,1} = P[Z(T_0) > z | H_1] \\ = Q[\{z - \hat{s}_1(T_0)\}/\sigma]$$

Optimum Threshold for General Gaussian Noise—The choice of the threshold depends on the performance measure under consideration. One measure of performance that is often used is the maximum of the error probabilities,

$$P_{e,m} = \max\{P_{e,0}, P_{e,1}\}$$

The *minimax criterion* for selecting the threshold is to choose the threshold z to give the smallest possible value of $P_{e,m}$. The value of z that minimizes $P_{e,m}$ is given by

$$z_m = [\hat{s}_0(T_0) + \hat{s}_1(T_0)]/2$$

and this is referred to as the *minimax threshold*. Notice that the minimax threshold depends on the signal set $\{s_0, s_1\}$, the impulse response of the filter, and the sampling time T_0 .

It is convenient to introduce the signal

$$s_d(t) = \{s_0(t) - s_1(t)\}/2$$

and let $\hat{s}_d(t)$ denote the output of the filter when $s_d(t)$ is the input (i.e., $\hat{s}_d = s_d * h$). If the minimax threshold is employed, the error probabilities $P_{e,0}$ and $P_{e,1}$ are given by

$$P_{e,0} = P_{e,1} = Q[\hat{s}_d(T_0)/\sigma]$$

Whenever these two probabilities are equal, we denote the common value by P_e .

As an example, consider an antipodal signal set (discussed earlier in this chapter). Since $s_0(t) = -s_1(t)$, the signal $s_d(t)$ is given by $s_d(t) = s_0(t)$. The error probability for the minimax threshold is

$$P_e = Q[\hat{s}_0(T_0)/\sigma]$$

There are other criteria for selecting the threshold z , but they typically require additional information about the data sequence that is to be transmitted. For instance, if the data sequence is modeled as a random process with a known probability distribution, the average probability of error can be employed as a criterion. If p_k denotes the *a priori probability* that the binary digit k is transmitted, the *average probability of error* is

$$P_{e,A} = P_{e,0}p_0 + P_{e,1}p_1$$

The threshold that minimizes $P_{e,A}$ is given by

$$z_A = z_m + 2 \sigma^2 [\hat{s}_d(T_0)]^{-1} \ln(p_1/p_0)$$

If the data bits 0 and 1 are equally probable (i.e., $p_0 = p_1$), the thresholds z_A and z_m are the same. Otherwise, the second term on the right-hand side of the expression for z_A is a nonzero term that biases the decision in favor of the data bit with the largest a priori probability.

For applications in which unequal costs are associated with the two types of errors and the a priori probabilities p_i are known, it may be desirable to select the threshold to minimize the average cost. This is the well-known *Bayes criterion*, which is discussed in references 4, 9, and 10.

The Matched Filter—A white Gaussian noise process $X(t)$ has spectral density $S_X(f) = N_0/2$ for all f . This is just the Fourier transform of a delta function that has area $N_0/2$. If the channel noise process has this flat spectral density, the autocorrelation function of the noise at the output of the filter is given by

$$R(u) = (N_0/2)f(u)$$

where $f = h * g$, the convolution of the impulse response and its time-reverse. Letting $u = 0$, we find

$$\sigma^2 = R(0) = N_0 \|h\|^2/2$$

The norm of the function h is defined by

$$\|h\| = \left\{ \int_{-\infty}^{\infty} h^2(t) dt \right\}^{1/2}$$

Consider first the minimax criterion. The optimum threshold depends on the filter impulse response $h(t)$, and for each choice of the filter impulse response, the corresponding minimax threshold gives the smallest possible value of $P_{e,m}$. Let P_e denote this minimum value of $P_{e,m}$ (recall that $P_e = P_{e,0} = P_{e,1}$ if the minimax threshold is used).

The value of P_e depends on the impulse response of the filter. For a system with an *additive white Gaussian noise channel*, the minimum possible value of P_e is achieved by the impulse response

$$h_M(t) = s_d(T_0 - t)$$

Multiplication of the filter impulse response by a positive constant c does not change the error probability if the corresponding minimax threshold is used, so $h_M(t)$ can be any positive constant multiple of $s_d(T_0 - t)$. A filter with this impulse response is called a *matched filter*.

Given a binary signal set $\{s_0, s_1\}$, let E_i denote the energy in the signal s_i , and define E_A by

$$E_A = [E_0 + E_1]/2$$

The signal correlation r' is the inner product

$$r' = (s_0, s_1) = \int_{-\infty}^{\infty} s_0(t) s_1(t) dt$$

and $r = r'/E_A$ is the correlation coefficient for the given signal set. The value of $|r|$ is never greater than one; moreover, $r = -1$ only for antipodal signals, and $r = +1$ only if $s_0(t) = s_1(t)$.

The norm of signal s_d is related to the energy by

$$\|s_d\|^2 = E_A(1 - r)/2$$

The matched filter impulse response $h_M(t) = c s_d(T_0 - t)$ has norm

$$\|h_M\|^2 = c^2 \|s_d\|^2$$

If this filter is employed in the receiver, the minimax threshold is $z_m = c(E_0 - E_1)/4$, and the signal $\hat{s}_d(T_0)$ is given by

$$\hat{s}_d(T_0) = \|s_d\| \|h_m\| = c \|s_d\|^2$$

It follows that the error probability for a receiver with a matched filter and minimax threshold is

$$P_e = Q\{[E_A(1 - r)/N_0]^{1/2}\}$$

This is the smallest error probability that can be achieved (if a single bit is transmitted) with *binary* baseband data transmission or *binary* coherent RF communications over an additive white Gaussian noise channel.

Two important points should be made concerning the above expression for the minimum error probability. First, the result does not depend on the sampling time T_0 . This is because the matched filter automatically compensates for the sampling time by incorporating a delay in its impulse response. As a consequence of this, the sampling time can be selected for convenience of implementation. In particular, if the signal set $\{s_0, s_1\}$ is *time limited* to the interval $[0, T]$, the matched filter is causal for any choice of T_0 not less than T . For such a signal set, the customary choice for the sampling time is $T_0 = T$.

The second point is that the minimum probability of error depends on the three parameters E_A , N_0 , and r . The energy E_A can be increased only by increasing the power in the received signal, and the noise density can be decreased only by lowering the noise level in the receiver. Both of these require increased cost and complexity in the communication system. However, the parameter r can be decreased by proper signal design. As previously mentioned, r can be no smaller than -1 , and $r = -1$ if and only if the signals are antipodal. Antipodal signals have equal energy: $E_A = E_0 = E_1$. As a consequence, the subscripts on E can be omitted, and the error probability for antipodal signaling, an additive white Gaussian noise channel, and a matched-filter receiver is given by

$$P_e = Q\{[2E/N_0]^{1/2}\}$$

where E is the energy per data bit.

The Correlation Receiver—The matched-filter receiver is shown in Fig. 11. The decision statistic $Z(T_0)$ for this receiver can be expressed in terms of the channel output $Y(t)$ by

$$Z(T_0) = \int_{-\infty}^{\infty} Y(T_0 - u) h_M(u) du$$

The impulse response of the matched filter is defined in terms of the binary signal set $\{s_0, s_1\}$ by

$$h_M(t) = \{s_0(T_0 - t) - s_1(T_0 - t)\}/2$$

In terms of the signal

$$s_d(t) = \{s_0(t) - s_1(t)\}/2$$

the matched filter impulse response is

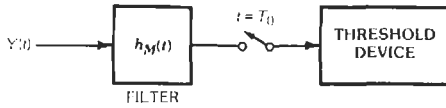


Fig. 11. Matched filter receiver.

$$h_M(t) = s_d(T_0 - t)$$

The statistic $Z(T_0)$ can therefore be expressed as

$$Z(T_0) = \int_{-\infty}^{\infty} Y(T_0 - u) s_d(T_0 - u) du$$

If the signals s_0 and s_1 are time limited to $[0, T]$ and the sampling time is $T_0 = T$, this expression reduces to

$$Z(T) = \int_0^T Y(t) s_d(t) dt$$

Thus, $Z(T)$ can be obtained by multiplying the channel output $Y(t)$ by the signal $s_d(t)$ and integrating from 0 to T , as shown in Fig. 12. The receiver that processes the channel output in this manner is known as the *correlation receiver*. The correlation receiver gives the same error probability as the matched-filter receiver.

Coherent Receivers for Binary and Quaternary RF Signals

For binary RF communications with signals of the form

$$s(t) = \sqrt{2} s_k(t) \cos[2\pi f_c t + \phi]$$

the matched-filter receiver for an additive white Gaussian noise channel is the receiver portion of the system shown in Fig. 10 with $h(t) = h_M(t)$. As in baseband systems, the correlation receiver can be substituted for the matched filter in coherent RF communication systems. The correlation receiver for coherent RF communications is shown in the diagram of Fig. 13.

Notice from Fig. 13 that the receiver must know the carrier phase ϕ . In practice, it suffices to have a good estimate of this phase, and there are many different kinds of tracking loops that can be employed in the receiver to provide such an estimate. A presentation of phase-tracking loops and related synchronization devices may be found in reference 5.

Binary phase shift keying (BPSK) employs binary RF signals with $s_d(t) = A/\sqrt{2}$ for $0 \leq t < T$ (see subsection

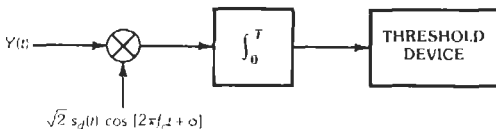


Fig. 13. Coherent correlation receiver for RF signals.

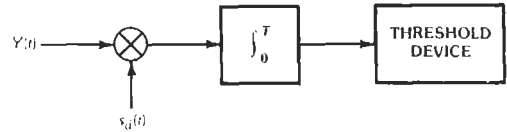


Fig. 12. Correlation receiver.

on BPSK). The optimum correlation receiver for coherent reception of BPSK is shown in Fig. 14. This is a special case of the receiver of Fig. 13. The square of the BPSK signal contains an unmodulated carrier at frequency $2f_c$ and phase 2ϕ , and this double-frequency carrier can be used to provide an estimate of the phase ϕ (this is the squaring loop*). This estimate is employed as the phase reference for coherent demodulation. An alternative method for coherent demodulation is to have the detection filter in the loop such as with a Costas loop or an I - Q loop.*

The coherent receivers for QPSK, OQPSK, and MSK are of the form shown in Fig. 15. These receivers consist of two branches, and each branch is a correlation receiver for an appropriate binary signal. The coherent references for the two branches differ in phase by $\pi/2$ radians, and in general there is a timing offset T' between the in-phase and quadrature branches. If the data rate is $R = 1/T$ bits per second in each case, then $T' = T$ for OQPSK and MSK, but $T' = 0$ for QPSK. The waveform $v(t)$ is the rectangular pulse $p_{2T}(t)$ for QPSK and OQPSK, and it is the sine pulse $\sin(\pi t/2T)$ for MSK.

Noncoherent Receivers

The signals employed in noncoherent communications are binary and M -ary signals of the type described in the subsections on FSK and MFSK. The signal set consists of M equal-energy signals of the form

$$s_k(t) = A a_k(t) \cos[2 f_c t + \theta_k(t) + \phi_k]$$

where $a_k(t)$ and $\theta_k(t)$ are baseband signals with bandwidths much smaller than the carrier frequency f_c .

For noncoherent demodulation, it is assumed that the receiver has no knowledge of the phase ϕ_k and that the phases for different signals are unrelated (e.g., ϕ_k and ϕ_j are statistically independent for $k \neq j$). Moreover, for some applications, a time-varying phase shift is introduced by the channel, so the phase of the received signal may change from one interval to the next.

* Reference 5.

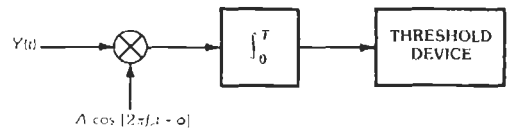


Fig. 14. Coherent correlation receiver for BPSK.

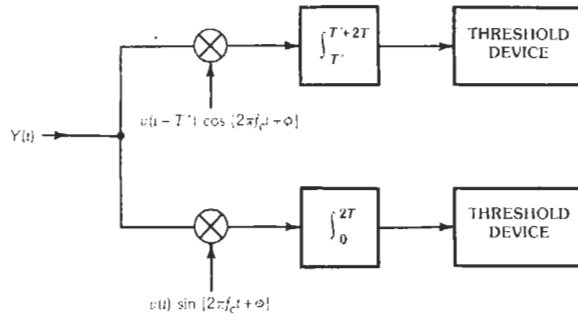


Fig. 15. Coherent receiver for quaternary RF signals.

The key component of the optimum noncoherent receiver for the signal set given above and an additive white Gaussian noise channel is the *noncoherent correlation detector (NCD)*, which is also referred to as an envelope correlation detector.* A block diagram of the noncoherent correlation detector for the signal $s_k(t)$ is shown in Fig. 16. This detector consists of in-phase and quadrature correlation receivers followed by squaring devices and a summing device. Thus the output of the noncoherent correlation detector is the square-law combination of the outputs of two correlation receivers whose phase references differ by $\pi/2$ radians. For convenience, let NCD/ k denote the noncoherent correlation detector for the k th signal in the set.

Let $P_{e,k}$ denote the probability that an error is made by the receiver when the signal $s_k(t)$ is sent. Two performance measures of interest are the maximum probability of error

$$P_{e,m} = \max\{P_{e,k}; k = 0, 1, \dots, M - 1\}$$

and the average probability of error

$$P_{e,A} = P_{e,0}p_0 + P_{e,1}p_1 + \dots + P_{e,M-1}p_{M-1}$$

In the expression for the average probability of error, p_k denotes the a priori probability that the signal $s_k(t)$ is sent. The receiver that gives the minimum possible

value of $P_{e,m}$ also gives the minimum possible value of $P_{e,A}$ for equally probable signals (i.e., $p_k = 1/M$ for each k). The block diagram for this optimum receiver is shown in Fig. 17. The decision of the receiver is based on a comparison of the outputs of the noncoherent correlation detectors: the receiver decides the i th signal was sent if

$$Z_i = \max\{Z_k; k = 0, 1, \dots, M - 1\}$$

If the signals have unequal a priori probabilities, some additional processing of the outputs of the noncoherent correlation detectors is necessary in order to give the minimum average probability of error.* This is also necessary for signals with unequal energy.

Since MFSK is a special case of M -ary signaling, an optimum receiver for MFSK can be obtained from Figs. 16 and 17. The signals for MFSK are

$$s_k(t) = A \cos[2\pi f_k t + \phi_k]$$

so the noncoherent correlation detector for MFSK uses reference signals $\cos(2\pi f_k t)$ and $\sin(2\pi f_k t)$ as shown in Fig. 18. An optimum noncoherent receiver for MFSK is obtained by employing this noncoherent correlation detector in the system of Fig. 17. An optimum noncoherent receiver for binary FSK is this system with $M = 2$.

* Reference 10.

* Reference 10.

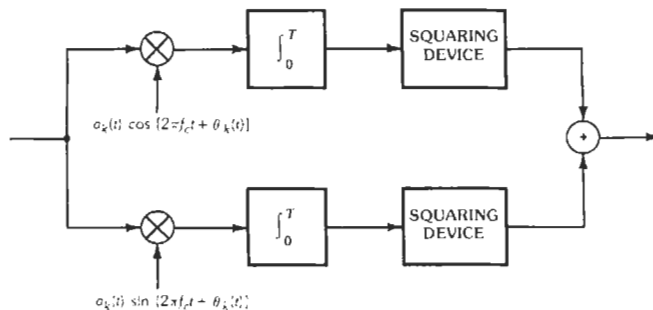


Fig. 16. Noncoherent correlation detector (NCD/ k).

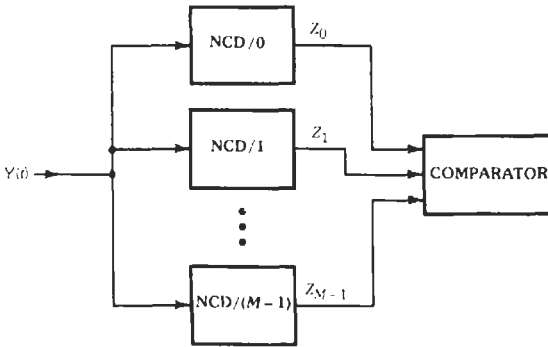


Fig. 17. Diagram of noncoherent receiver for M -ary communications.

If the channel noise is a band-limited white Gaussian noise process and the noise bandwidth is small compared to the tone frequencies f_k , there is an alternative implementation of the optimum noncoherent receiver that utilizes an *envelope detector*. In order to describe the response of the envelope detector to narrow-band inputs, define the signal

$$x_k(t) = v(t) \cos(2\pi f_k t) + w(t) \sin(2\pi f_k t)$$

in terms of baseband signals $v(t)$ and $w(t)$ that are narrow-band with respect to the frequency f_k . If $x_k(t)$ is the input to an envelope detector, the output is

$$e(t) = [v^2(t) + w^2(t)]^{1/2}$$

The signal $e(t)$ is called the *envelope* of the narrow-band signal $x_k(t)$.

Let $h_k(t)$ be the impulse response of a filter matched to the k th tone of the MFSK signal set; that is, let

$$h_k(t) = \cos(2\pi f_k t)$$

In contrast to the matched filters of the preceding subsection, this filter is *noncoherently matched* to the signal

$$s_k(t) = A \cos[2\pi f_k t + \phi_k]$$

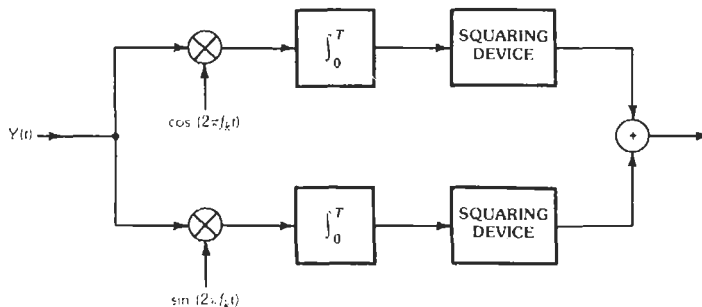


Fig. 18. Noncoherent correlation detector (NCD/ k) for MFSK signal set.

since no phase information is utilized in $h_k(t)$. The output of the system of Fig. 19 at time T is equivalent to the output of the noncoherent correlation detector. Thus the matched filter/envelope detector can be substituted for the noncoherent correlation detector in the system of Fig. 17 to obtain an alternative form of the optimum receiver.

ERROR PROBABILITIES FOR DIGITAL COMMUNICATION SYSTEMS

In the preceding section, the optimal receivers are given for the most commonly used digital modulation techniques. Although there are a large number of different modulation techniques, many of them give the same performance. Throughout this section, performance will be measured by the *bit error probability* P_b as a function of the *bit energy to noise density ratio* E_b/N_0 . Issues such as bandwidth and complexity of implementation will not be considered (see the subsection titled "Spectra of PSK and MSK Signals" for some comments on bandwidth of PSK and MSK signals).

Most binary communication systems of practical interest fall into two categories: coherent detection of antipodal signals and noncoherent detection of orthogonal signals. Moreover, certain quaternary signal sets such as QPSK, OQPSK, and MSK have the same bit error probability as binary antipodal signaling. For M -ary modulation, the primary case of interest is orthogonal signaling with noncoherent demodulation.

Coherent Systems With Additive White Gaussian Noise Channels

As discussed in the subsection titled "Receivers for Binary Baseband Data Transmission," binary coherent RF communication systems can be analyzed in terms of their low-pass equivalent systems. Thus, the results of that subsection give the error probabilities for optimum coherent demodulation of binary RF signaling over additive white Gaussian noise channels. If the channel noise has spectral density $N_0/2$, the error probability for

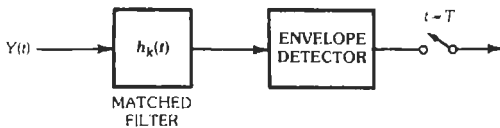


Fig. 19. Noncoherent matched filter/envelope detector.

binary signaling (baseband or RF) with an optimum receiver is

$$P_b = Q\{[E_A(1 - r)/N_0]^{1/2}\}$$

where,

r is the correlation coefficient for the binary signal set,

E_A is the average of the energies E_0 and E_1 for the two signals (see subsection cited above for formal definitions).

The binary signal sets of greatest interest are the *antipodal* ($r = -1$) and *orthogonal* ($r = 0$) equal-energy signals. If the two signals have equal energy, $E_A = E_0 = E_1$, and thus the *energy per bit* is E_A regardless of which bit is transmitted. For equal energy signals (binary or M -ary), it is customary to specify the performance in terms of the energy per bit, which is denoted by E_b .

It follows from the above comments that the bit error probability for optimum coherent detection of binary antipodal signals is

$$P_b = Q[(2E_b/N_0)^{1/2}]$$

For most applications, it is convenient to give the bit energy to noise density ratio in decibels. For this purpose, we define

$$(E_b/N_0)_{dB} = 10 \log_{10}(E_b/N_0)$$

A graph of P_b versus $(E_b/N_0)_{dB}$ for optimum demodulation and antipodal signaling is given in Fig. 20. Note that since binary PSK is just a special case of antipodal signaling, this curve applies to coherent detection of BPSK. Coherent detection of binary ASK with antipodal signals is another special case.

Many of the quaternary communication systems have the same bit error probability as binary antipodal signaling with an optimum receiver. The bit error probability for optimum coherent demodulation of QPSK, OQPSK, or MSK is given by

$$P_b = Q[(2E_b/N_0)^{1/2}]$$

exactly the same as for optimum coherent demodulation of BPSK. Consequently, the performance curve for QPSK, OQPSK, and MSK is the curve in Fig. 20 that corresponds to binary antipodal signaling.

Because the error probability $Q[(2E_b/N_0)^{1/2}]$ arises so frequently in the analysis of the performance of

coherent communication systems, a short table of its values is given in Table 1. The table lists $(E_b/N_0)_{dB}$, E_b/N_0 , and the error probability $Q[(2E_b/N_0)^{1/2}]$.

If the phase reference is not perfect, such as when it is estimated by some kind of tracking loop, the error probability for BPSK and related modulation schemes will be larger than given above. The *phase error* ϕ_e is the difference between the reference phase and the phase of the received signal. If ϕ_e is constant for the duration of the data bit, the bit error probability for BPSK is

$$P_b = Q[(2E_b/N_0)^{1/2} \cos(\phi_e)]$$

Coherent receivers for M -ary orthogonal signaling are of limited interest in applications. This is primarily because of the relative inefficiency of orthogonal signaling in a coherent system. For example, binary FSK requires twice the signal energy of BPSK to give the same bit error probability in a coherent system. This is true of any binary orthogonal signal set, since $r = 0$ implies

$$P_b = Q[(E_b/N_0)^{1/2}]$$

If a coherent receiver is to be employed, it makes sense to use a modulation scheme that is more efficient than orthogonal signaling.

The error probabilities for coherent demodulation of M -ary orthogonal signals are given in reference 5. The error probabilities for noncoherent demodulation of M -ary orthogonal signals are given in the next subsection. For large values of M , the performance of the coherent system is not much better than the performance of the noncoherent system (especially for low error probabilities). This is another reason there is little interest in implementing the more complex coherent receiver if M -ary orthogonal modulation is employed. For further discussion of this and related topics see reference 5.

Noncoherent Systems With Additive White Gaussian Noise Channels

The performance curve for binary FSK can be obtained by considering binary, orthogonal, equal-energy signals. For such signals, the bit error probability for optimum noncoherent detection is

$$P_b = 0.5 \exp(-E_b/2N_0)$$

which is also shown in Fig. 20.

Closely related to binary orthogonal signaling, although not always viewed as such, is *differential phase-shift keying (DPSK)*. In DPSK, binary digits are transmitted at a rate of $1/T$ bits per second by phase transitions at the ends of successive T -second intervals. A phase change of 180 degrees can represent a binary 1,

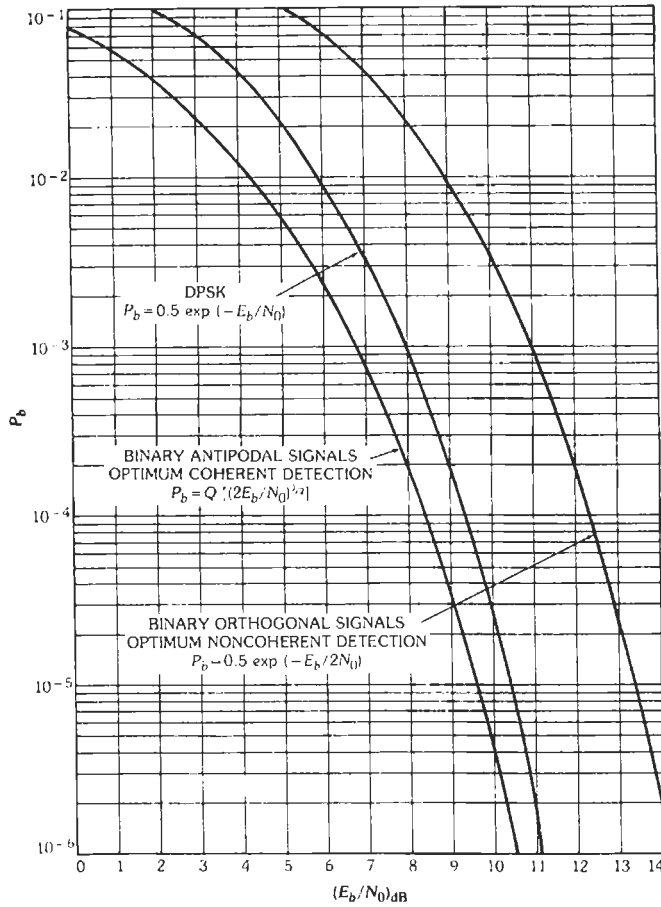


Fig. 20. Bit error probabilities.

and the absence of a phase change can represent a binary 0. This is simply differentially encoded PSK, and it can be demodulated either coherently or noncoherently.

The optimum noncoherent demodulator for DPSK is based on a noncoherent correlation detector (Fig. 16) in which the integration time is $2T$ seconds rather than T seconds. The optimality follows from the fact that

TABLE 1. VALUES OF ERROR PROBABILITY $Q[(2E_b/N_0)^{1/2}]$

$(E_b/N_0)_{dB}$	E_b/N_0	$Q[(2E_b/N_0)^{1/2}]$
5.0	3.16	5.95×10^{-3}
5.5	3.55	3.86×10^{-3}
6.0	3.98	2.39×10^{-3}
6.5	4.47	1.40×10^{-3}
7.0	5.01	7.73×10^{-4}
7.5	5.62	3.99×10^{-4}
8.0	6.31	1.91×10^{-4}
8.5	7.08	8.40×10^{-5}
9.0	7.94	3.36×10^{-5}
9.5	8.91	1.21×10^{-5}
10.0	10.00	3.88×10^{-6}
10.5	11.22	1.09×10^{-6}
11.0	12.59	2.62×10^{-7}
11.5	14.13	5.34×10^{-8}
12.0	15.85	9.03×10^{-9}

DPSK is a form of orthogonal signaling on the interval $[0, 2T]$: the signal corresponding to a 180-degree phase change is orthogonal to the signal corresponding to no phase change. Since the integration time is $2T$, the effective signal energy is $2E_b$. For optimum noncoherent demodulation of orthogonal signals of energy $2E_b$, the error probability is

$$P_b = 0.5 \exp(-E_b/N_0)$$

which is shown in Fig. 20.

The noncoherent correlation detector structure can be simplified to provide the alternative implementation of the optimal noncoherent demodulator for DPSK that is shown in Fig. 21. For this receiver the decision is made as follows. If

$$W_c Z_c + W_s Z_s > 0$$

the decision is that 0 was sent, and if

$$W_c Z_c + W_s Z_s < 0$$

the decision is that 1 was sent.*

A suboptimal receiver that is commonly suggested for demodulation of DPSK involves delay and multiplication of the received continuous-time signal, as shown in Fig. 22. The delay-and-multiply receiver does not produce the same decision statistic as the receiver shown in Fig. 21, and its probability of error is larger than that obtained by the receiver of Fig. 21. For very large values of the signal-to-noise ratio, the performance difference between the two receivers is small.

* See reference 7 for details.

Additional information on DPSK reception is given in references 6, 7, and 8.

Next we consider M -ary signaling with equal-energy, orthogonal signals. That is, the signal set $\{s_k: k = 0, 1, \dots, M - 1\}$ is such that $\|s_k\|^2 = E_s$ for all k , and $(s_k, s_j) = 0$ whenever $k \neq j$. The energy in the signal s_k is called the *energy per symbol* (denoted by E_s). A symbol error occurs if the receiver decides that s_j was sent when in fact s_k was sent and $k \neq j$. For an additive white Gaussian noise channel and an optimum noncoherent receiver, the symbol error probability P_s is given by

$$P_s = M^{-1} \sum_{n=2}^M \binom{M}{n} (-1)^n \exp\{-[(n-1)/n]E_s/N_0\}$$

This can be evaluated on a digital computer for values of M up to 32, although double precision computation may be required for values of M larger than about 16. The symbol error probability can be computed for larger values of M by numerical integration of an alternative expression for P_s ; for example, see equation (10-15) on p. 489 of reference 5.

If M -ary modulation is employed to transmit binary data, M should be a power of 2. If $M = 2^m$, each M -ary symbol can represent a unique sequence of m binary digits. The energy per bit is $E_b = E_s/m$. A symbol error can cause as few as one or as many as m bit errors. On the average, the number of bit errors per symbol error is $mM/2(M - 1)$ for M -ary orthogonal signaling with optimum noncoherent demodulation. Thus the bit error probability P_b is related to the symbol error probability by $P_b = MP_s/2(M - 1)$. Making the appropriate substitutions in the above expression for P_s , we find that the bit error probability is given in terms of the energy per bit by

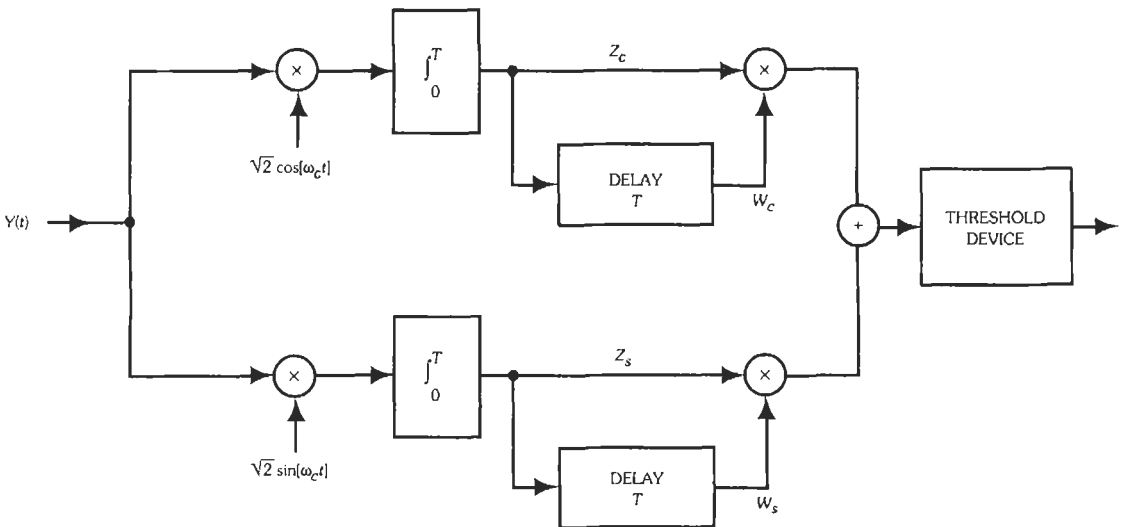


Fig. 21. Optimal noncoherent demodulator for DPSK.

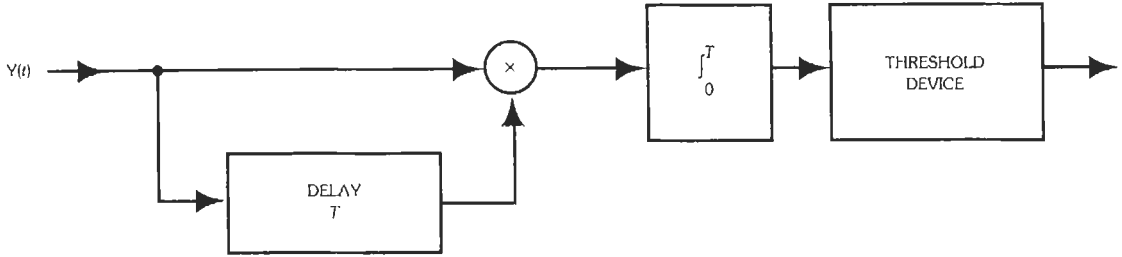


Fig. 22. The delay-and-multiply receiver—a suboptimal receiver for DPSK.

$$P_b = (2M - 2)^{-1}$$

$$\sum_{n=2}^M \binom{M}{n} (-1)^n \exp\{-[(n - 1)/n]mE_b/N_0\}$$

A table of values of the symbol error probability is given in Table 2 for $M = 2, 8,$ and 32 . The bit error probabilities can be obtained from this table by multiplying the symbol error probabilities by $M/2(M - 1)$. Curves of the bit error probability versus the bit energy to noise density ratio are shown in Fig. 23 for $M = 2, 4, 8, 16,$ and 32 . It can be shown that $P_b \rightarrow 0$ as $M \rightarrow \infty$ provided that E_b/N_0 is greater than $\ln 2$ (where \ln denotes the natural logarithm).

Noncoherent Systems With Nonselective Fading Channels

In some communication channels, there are several propagation paths between the transmitter and the receiver. The signal components arriving at the receiver from different paths add destructively or constructively depending on the difference in the propagation times and on any phase shifts introduced by the medium (e.g., reflections from objects). Thus the signal strength depends on the relative phasing of the components of the received signal, and this gives rise to the condition known as *fading*.

If the differential propagation times for the paths are small compared with the data symbol duration, and the

path strengths, propagation times, and phase shifts are nearly constant for the duration of the symbol, then there is no significant dispersion of the signal. In this case, the fading is referred to as *slow nonselective fading*. If the differential propagation times are large compared with the data symbol duration, signal dispersion will occur. Such channels are known as *time-dispersive* or *frequency-selective* fading channels. These channels may produce significant levels of intersymbol interference. A more complete discussion of fading and its effects on communication systems is given in reference 8, and the derivations of the results given in this section can be found in references 12, 13, and 14.

For *slow nonselective Rician fading*, the received signal is the sum of a nonfaded version of the transmitted signal and a slow nonselective Rayleigh faded version. The difference in propagation times for the two components is small enough that the channel is nonselective. The nonfaded component of the received signal is called the *specular component*, and the Rayleigh faded component is sometimes called the *random* or *scatter component*. The specular component may result from a direct path between the transmitter and receiver, and the faded component may arise from a large number of reflections.

Consider the transmission of a unit-amplitude RF signal over this Rician fading channel. The received signal is the sum of a deterministic component of amplitude A and a random component of amplitude V . The quantities A and V represent the gains associated

TABLE 2. SYMBOL ERROR PROBABILITIES FOR OPTIMUM NONCOHERENT DEMODULATION OF M -ARY ORTHOGONAL SIGNALING

$(E_b/N_0)_{dB}$	Symbol Error Probability, P_s		
	$M = 2$	$M = 8$	$M = 32$
3.0	1.84×10^{-1}	1.05×10^{-1}	5.03×10^{-2}
4.0	1.42×10^{-1}	5.38×10^{-2}	1.71×10^{-2}
5.0	1.03×10^{-1}	2.25×10^{-2}	4.03×10^{-3}
6.0	6.83×10^{-2}	7.25×10^{-3}	6.03×10^{-4}
7.0	4.08×10^{-2}	1.67×10^{-3}	5.08×10^{-5}
8.0	2.13×10^{-2}	2.53×10^{-4}	2.10×10^{-6}
9.0	9.42×10^{-3}	2.27×10^{-5}	3.64×10^{-8}
10.0	3.37×10^{-3}	1.06×10^{-6}	2.15×10^{-10}

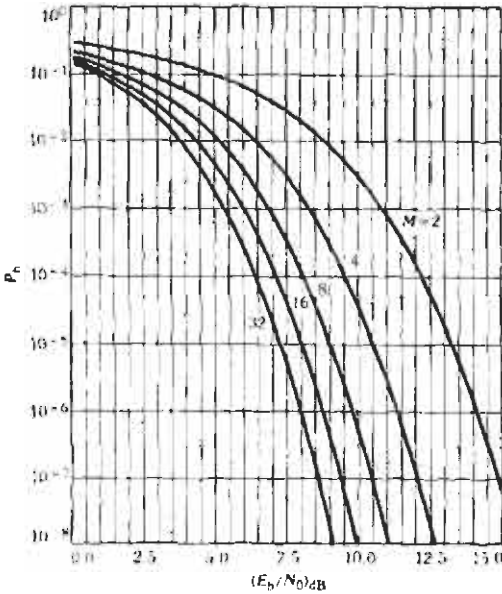


Fig. 23. Bit error probabilities for M -ary orthogonal signals and optimum noncoherent demodulation.

with the specular path and the scatter path, respectively. The random variable V has the Rayleigh distribution, and its probability density function is

$$f_V(v) = (2v/m_f) \exp\{-v^2/m_f\}$$

for $v > 0$, where m_f denotes the second moment or mean-square value of the amplitude of the faded component (i.e., $m_f = E\{V^2\}$). The phase between the two components is random and uniformly distributed on the interval $[0, 2\pi]$. The amplitude of the sum of the specular and faded components is a random variable Y which has a Rician distribution. The probability density function for the Rician distribution is

$$f_Y(y) \approx (2y/m_f) \exp\{(y^2 + A^2)/m_f\} I_0(2yA/m_f)$$

for $y > 0$, where I_0 is the zero-order modified Bessel function. The second moment for this distribution is $E\{Y^2\} = A^2 + m_f$. Notice that for $A = 0$ (no specular component), the Rician density reduces to the Rayleigh density, because $I_0(0) = 1$.

Let E' denote the transmitted energy per symbol in an M -ary noncoherent communication system. The energy per symbol in the received signal is $Y^2 E'$. The average energy per symbol in the received signal is

$$E_s = E\{Y^2\} E' = (A^2 + m_f) E'$$

Let $e_1 = A^2 E' / N_0$ and $e_2 = m_f E' / N_0$ be the symbol energy to noise density ratios for the specular and faded components, respectively. The conditional probability of symbol error given $Y \approx y$ is (see preceding subsection)

$$P_s(y) = M^{-1}$$

$$\sum_{n=2}^M \binom{M}{n} (-1)^n \exp\{-[(n-1)/n]y^2 E' / N_0\}$$

The average probability of symbol error is

$$P_s = \int_0^\infty P_s(y) f_Y(y) dy$$

$$= M^{-1} \sum_{n=2}^M \binom{M}{n} (-1)^n G[(n-1)/n, e_1, e_2]$$

where the function G is defined by

$$G(c, e_1, e_2) = (1 + ce_2)^{-1} \exp\{-ce_1/(1 + ce_2)\}$$

The bit error probability is obtained by multiplying the symbol error probability P_s by $M/2(M-1)$.

Two special cases are of interest. First, for the Rayleigh fading channel, $A = 0$ and $e_1 = 0$. It follows that $e_2 = E_s / N_0$. Moreover,

$$G(c, 0, e_2) = (1 + ce_2)^{-1}$$

so that the error probability for noncoherent reception of M -ary orthogonal signals with nonselective Rayleigh fading is given by

$$P_s = M^{-1} \sum_{n=2}^M \binom{M}{n} (-1)^n \{1 + [(n-1)/n] E_s / N_0\}^{-1}$$

The second important special case is the nonfading channel. If the received signal has no faded component, then $m_f = 0$, $e_2 = 0$, and $e_1 = E_s / N_0$. Since

$$G(c, e_1, 0) = \exp\{-ce_1\}$$

the error probability is

$$P_s = M^{-1} \sum_{n=2}^M \binom{M}{n} (-1)^n \exp\{-[(n-1)/n] E_s / N_0\}$$

as given in the preceding subsection.

For binary orthogonal signaling, each symbol represents one bit, so $P_s = P_b$. The resulting expression for the bit error probability in a Rician fading channel is

$$P_b = 0.5 G(0.5, e_1, e_2) = (2 + e_2)^{-1} \exp\{-e_1/(2 + e_2)\}$$

For Rayleigh fading, $e_1 = 0$ and $e_2 = E_s / N_0$. Since $E_b = E_s$ for binary signals, the bit error probability is

$$P_b = [2 + (E_b / N_0)]^{-1}$$

It is customary to present numerical data for the error probability in a Rician fading channel in terms of the

ratio of the power in the two components of the received signal. Let $g^2 = m_1/A^2$, so that $e_2 = g^2 e_1$. The signal-to-noise ratios e_1 and e_2 are related to E_s/N_0 by

$$e_1 = [1/(1 + g^2)] E_s/N_0$$

and

$$e_2 = g^2 e_1 = [g^2/(1 + g^2)] E_s/N_0$$

The Rician channel can be specified by giving either e_1 and e_2 or g^2 and E_s/N_0 . In the latter specification, E_s/N_0 is a measure of the average signal-to-noise ratio at the receiver, and g^2 specifies how the signal power is distributed between its two components. Specifically, $1/(1 + g^2)$ is the fraction due to the specular component, and $g^2/(1 + g^2)$ is the fraction due to the faded component. A substitution for e_1 and e_2 in the expressions above gives expressions for the error probability in terms of g^2 and E_s/N_0 .

Bit error probabilities for binary orthogonal modulation, noncoherent demodulation, and slow nonselective Rician fading are shown in Fig. 24 for five values of the parameter g^2 . The curve for $g^2 = 0$ is the same as for the additive white Gaussian noise channel (Fig. 20), and the curve for $g^2 = 10$ is very nearly the same as the error probability for Rayleigh fading. Table 3 lists values of the bit error probability for a Rayleigh fading channel, a Rician fading channel with $g^2 = 0.1$, and a nonfading channel (Rician with $g^2 = 0$).

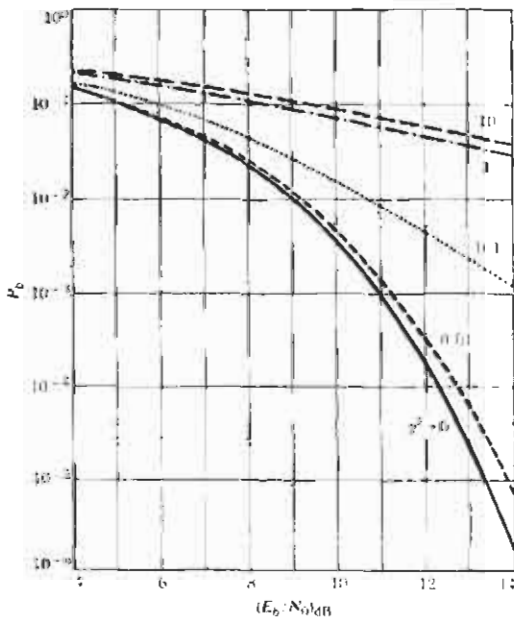


Fig. 24. Bit error probability for Rician fading.

SPREAD-SPECTRUM COMMUNICATIONS

A common trade-off in a communication system is performance versus bandwidth. That is, various aspects of communications performance can be improved at the expense of increased RF bandwidth. One example of this is the reduction of error probability in MFSK systems as the number of tones increases. More generally, the use of error-correcting coding gives lower error probabilities but requires a larger bandwidth than for uncoded signaling with the same information rate (see Chapter 25, Information Theory and Coding).

Spread-spectrum modulation is a very general method for improving other measures of system performance through the use of wideband signals. A properly designed spread-spectrum communication system can operate reliably in the presence of various types of radio-frequency interference (RFI) including multipath interference, multiple-access interference, and hostile jamming. In addition, when the signal power is spread over a large bandwidth, the signal has very small average power in any narrow-band slot. This means that the spread-spectrum system can share a frequency band with several narrow-band systems. It also means that it is difficult to detect the presence of the spread-spectrum signal by use of narrow-band equipment. The basic modulation formats and receiver structures are described below. Further details on the design and performance of spread-spectrum systems are given in references 24 through 40.

Direct-Sequence Spread-Spectrum Communications

The type of spread spectrum known as *binary direct-sequence spread spectrum* employs signals of the form

$$s(t) = A a_s(t) a_d(t) \cos[2\pi f_c t + \phi]$$

where,

$a_d(t)$ is a binary baseband data signal,

$a_s(t)$ is a baseband *spectral-spreading signal* with a bandwidth that is large compared to the data rate.

These binary baseband signals are of the form described in the section on baseband signal sets. Although only binary direct-sequence modulation is described here, there is a natural generalization to *quaternary direct sequence*, which is described in reference 24.

The data signal $a_d(t)$ consists of a sequence of positive and negative rectangular pulses; that is, it can be expressed as

$$a_d(t) = \sum_{n=-\infty}^{\infty} b_n p_T(t - nT)$$

TABLE 3. BIT ERROR PROBABILITIES FOR BINARY ORTHOGONAL SIGNALING, SLOW NONSELECTIVE FADING, AND NONCOHERENT DEMODULATION

$(E_b/N_0)_{dB}$	Rayleigh	Bit Error Probability $g^2 = 0.1$	No Fading
4.0	2.22×10^{-1}	1.61×10^{-1}	1.42×10^{-1}
6.0	1.67×10^{-1}	9.15×10^{-2}	6.83×10^{-2}
8.0	1.20×10^{-1}	4.18×10^{-2}	2.13×10^{-2}
10.0	8.33×10^{-2}	1.51×10^{-2}	3.37×10^{-3}
12.0	5.60×10^{-2}	4.41×10^{-3}	1.81×10^{-4}
14.0	3.69×10^{-2}	1.13×10^{-3}	1.76×10^{-6}
20.0	9.80×10^{-3}	2.48×10^{-5}	9.64×10^{-23}

where,

$p_T(t)$ is the rectangular pulse duration T ,
 b_n is either +1 or -1, depending on the data bit to be sent in the n th interval.

Similarly, the spectral-spreading signal can be written as

$$a_s(t) = \sum_{j=-\infty}^{\infty} a_j v(t - jT_c)$$

where $v(t)$ is the *chip waveform*, a time-limited pulse of duration T_c . The parameter T_c is called the *chip duration*, and the sequence

$$(a_j) = \dots, a_{-1}, a_0, a_1, a_2, \dots$$

is called the *signature sequence*. For reasonable choices of signature sequence and chip waveform, $1/T_c$ is a rough estimate of the bandwidth of the spread-spectrum signal. For most direct-sequence systems, $T_c \ll T$, so the bandwidth of the spread-spectrum signal is much larger than that of the data signal.

Usually, the data pulse duration T is an integer multiple of the chip duration. If $T = NT_c$ for some integer N , there are N chips per data pulse, and the bandwidth of the spread-spectrum signal is roughly N times the data rate. Usually the signature sequence is periodic, and the period p is equal to N or else $p \gg N$, depending on the nature of the application.

The optimum receiver for binary direct-sequence spread-spectrum modulation and an additive white Gaussian noise channel (i.e., thermal noise is the only interference) is the correlation receiver shown in Fig. 25 (cf. Fig. 14). The received signal $Y(t)$ is multiplied by

$$A a_s(t) \cos[2\pi f_c t + \phi]$$

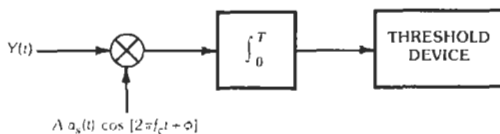


Fig. 25. Coherent correlation receiver for binary direct-sequence spread spectrum.

and the product is integrated over an appropriate interval of length T (this is illustrated for the interval $[0, T]$ in Fig. 25). Since the signals corresponding to $b_n = +1$ and $b_n = -1$ are antipodal, the bit error probability for binary direct-sequence spread-spectrum modulation and optimum coherent demodulation is

$$P_b = Q[(2E_b/N_0)^{1/2}]$$

the same as for binary BPSK with optimum coherent demodulation.

The motivation for considering direct-sequence spread spectrum is not performance improvement for channels in which the primary interference is thermal noise. Rather, it is the fact that the spectral-spreading signal can be selected to give the transmitted signal certain desirable properties and enable the receiver to discriminate against forms of interference other than thermal noise. For example, if the spread-spectrum signal is corrupted by additive RFI, and if we ignore the thermal noise for the moment, the received signal can be written as $Y(t) = s(t) + x(t)$, where $x(t)$ is an additive interference signal. The effects of this interference can be reduced by choosing the spectral-spreading signal to minimize

$$\int_0^T x(t) a_s(t) \cos[2\pi f_c t + \phi] dt$$

One application in which additive RFI is encountered is in a direct-sequence spread-spectrum *multiple-access* communication system.* In such a system, signals from the various transmitters have the same carrier frequency f_c . The spread-spectrum signals from other transmitters appear as additive RFI to a receiver that attempts to demodulate the signal $s(t)$. Let $s'(t)$ denote one of these other signals, and denote its spectral-spreading signal by $a'_s(t)$ and its signature sequence by (a'_j) . Since the signal $s'(t)$ is from a different transmitter than $s(t)$, it arrives at the receiver with a different propagation delay and phase angle. If u denotes the propagation delay and

* Reference 24.

ϕ' denotes the phase for $s'(t)$, the quantity to be minimized is

$$\int_0^T a_s'(t - u) a_s(t) dt \cos(\phi - \phi')$$

Since this should be small for all values of the phase angles, it is necessary to minimize the magnitude of

$$\int_0^T a_s'(t - u) a_s(t) dt$$

for all values of u . The goal is to select the signature sequences to accomplish this minimization (see references 24 and 28).

Similar considerations arise for channels with specular multipath, narrow-band interference, hostile jamming, and other forms of RFI. For complete discussions of the problems of signature sequence selection see references 24, 27, and 28. The problems of obtaining phase and timing references for direct-sequence spread-spectrum communication systems are discussed in reference 26. The performance of direct-sequence spread spectrum with various forms of RFI is considered in references 24, 25, and 36 through 38.

Frequency-Hop Spread-Spectrum Communications

Signals with very large RF bandwidths can be generated by a method known as *frequency hopping* in which the carrier frequency of a digital communication signal is changed, or "hopped," over a wide range of frequencies. If the digital communication signal is

$$c(t) = A a(t) \cos[2\pi f_c t + \theta(t) + \phi]$$

the resulting *frequency-hopped* signal is

$$s(t) = A a(t) \cos[2\pi f(t)t + \theta(t) + \phi]$$

The function $f(t)$, which describes the carrier frequency as a function of time, is called the (frequency) *hopping pattern*. The hopping pattern is generated by applying a random or pseudorandom sequence of inputs to a frequency synthesizer. Typically, the available RF bandwidth is partitioned into q nonoverlapping frequency intervals called *slots*, and the q different frequencies generated by the frequency hopper are the center frequencies for these slots.

A frequency-hop (FH) spread-spectrum signal with hopping rate R_h hops per second is a signal that has the form of $s(t)$ above for which the frequency $f(t)$ can change every $1/R_h$ seconds. The frequency is constant on intervals of length $T_h = 1/R_h$. The parameter T_h is called the *frequency dwell time* or *hop interval*. In contrast to direct-sequence spread-spectrum signals, which occupy the full RF bandwidth at all times, FH

spread-spectrum signals occupy only a small fraction of the RF bandwidth during a given hop interval.

Fast FH spread-spectrum systems have hopping rates that are larger than the data rate, and so the duration of a data pulse is larger than the hop interval. Since the transmission of a data pulse utilizes more than one of the q frequency slots, frequency diversity is obtained with fast FH spread-spectrum signaling. *Slow FH spread-spectrum systems* have hopping rates that are smaller than the data rate, and thus the hop interval is greater than the data symbol duration.

The total RF bandwidth of a slow FH spread-spectrum signal is approximately q times the bandwidth of the digital communication signal $c(t)$; it is virtually independent of the hopping rate. The total bandwidth of a fast FH spread-spectrum signal depends on the number of frequency slots and the hopping rate, but it does not depend very much on the data rate.

The multiple-access capability of frequency-hop spread spectrum is due to the fact that each signal occupies only the fraction $1/q$ of the bandwidth during each hop interval. Even for totally asynchronous operation of a large number of transmitters,* the hopping patterns can be designed such that the probability of interference between the signals from any two given transmitters during a given hop interval is no more than $2/q$. Since errors occur with high probability whenever the signals interfere, some form of error-control coding is necessary for typical multiple-access systems. In fact, error-control coding is virtually a requirement for frequency-hop spread-spectrum communication in the presence of any form of partial-band or pulsed interference. Convolutional codes and Reed-Solomon block codes (see Chapter 25) appear to be the most suitable error-correcting codes for use in frequency-hop spread-spectrum communication systems.

REFERENCES

General References for Digital Communications

1. Benedetto, S., Biglieri, E., and Castellani, V. *Digital Transmission Theory*. Englewood Cliffs, NJ: Prentice-Hall, 1987.
2. Blahut, R. E. *Digital Transmission of Information*. Reading, MA: Addison-Wesley, 1990.
3. Golomb, S. W., ed. *Digital Communications With Space Applications*. Englewood Cliffs, NJ: Prentice-Hall, Inc., 1964.
4. Helstrom, C. W. *Statistical Theory of Signal Detection*. 2d ed. New York: Pergamon Press, Inc., 1968.
5. Lindsey, W. C., and Simon, M. K. *Telecommunications Systems Engineering*. Englewood Cliffs, NJ: Prentice-Hall, Inc., 1973.

* Reference 29.

6. Proakis, J. G. *Digital Communications*. 2nd ed. New York: McGraw-Hill, 1989.
7. Pursley, M. B. *Introduction to Digital Communication Systems*. Reading, MA: Addison-Wesley, scheduled for 1993.
8. Stein, S. Part III of *Communication Systems and Techniques*. New York: McGraw-Hill Book Co., 1966.
9. Van Trees, H. L. *Detection, Estimation, and Modulation Theory*, Part I. New York: John Wiley & Sons, Inc., 1968.
10. Weber, C. L. *Elements of Detection and Signal Design*. New York: McGraw-Hill Book Co., 1968.
11. Wozencraft, J. M., and Jacobs, I. M. *Principles of Communication Engineering*. New York: John Wiley & Sons, Inc., 1965.
20. de Buda, R. "Coherent Demodulation of Frequency-Shift Keying With Low Deviation Ratio." *IEEE Transactions on Communications*, Vol. COM-20, June 1972, pp. 429-435.
21. Gronemeyer, S. A., and McBride, A. L. "MSK and Offset QPSK Modulation." *IEEE Transactions on Communications*, Vol. COM-24, August 1976, pp. 809-820.
22. Pasupathy, S. "Minimum Shift Keying: A Spectrally Efficient Modulation." *IEEE Communications Magazine*, July 1979, pp. 14-22.
23. Ungerboeck, G. "Channel coding with multilevel/phase signals." *IEEE Transactions on Information Theory*, Vol. 28, January 1982, pp. 55-67.

Noncoherent Communication Over Nonselective Fading Channels

12. Turin, G. L. "Error Probabilities for Binary Symmetric Ideal Reception Through Nonselective Slow Fading and Noise." *Proceedings of the IRE*, Vol. 46, September 1958, pp. 1603-1619.
13. Lindsey, W. C. "Error Probabilities for Rician Fading Multichannel Reception of Binary and N-ary Signals." *IEEE Transactions on Information Theory*, October 1964, pp. 339-350 (reprinted in reference 14).
14. Brayer, K., ed. *Data Communications Via Fading Channels*. New York: IEEE Press, 1975.

Bandwidth-Efficient Modulation (QPSK, OQPSK, MSK, CPFSK, and Trellis-Coded Modulation)

15. Amoroso, F. "The Bandwidth of Digital Data Signals." *IEEE Communications Magazine*, November 1980, pp. 13-24.
16. Amoroso, F., and Kivett, J. A. "Simplified MSK Signaling Technique." *IEEE Transactions on Communications*, Vol. COM-25, April 1977, pp. 433-441.
17. Anderson, J. B., Aulin, T., and Sundberg, C.-E. *Digital Phase Modulation*. New York: Plenum, 1986.
18. Biglieri, E., Divsalar, D., McLane, P. J., and Simon, M. K. *Introduction to Trellis-Coded Modulation with Applications*. New York: Macmillan, 1991.
19. Forney, G. D., Jr., Gallager, R. G., Lang, G. R., Longstaff, F. M., and Qureshi, S. U. "Efficient modulation for band-limited channels." *IEEE Journal on Selected Areas in Communications*, Vol. 2, September 1984, pp. 632-647.
24. Pursley, M. B. "Spread-Spectrum Multiple-Access Communications." In *Multi-User Communication Systems*. G. Longo, ed. Vienna and New York: Springer-Verlag, 1981, pp. 139-199.
25. Pickholtz, R. L., Schilling, D. L., and Milstein, L. B. "Theory of Spread-Spectrum Communications—A Tutorial." *IEEE Transactions on Communications*, Vol. COM-30, May 1982, pp. 855-884.
26. Holmes, J. K. *Coherent Spread Spectrum Systems*. New York: John Wiley & Sons, Inc., 1982.
27. MacWilliams, F. J., and Sloane, N. J. A. "Pseudo-Random Sequences and Arrays." *Proceedings of the IEEE*, Vol. 64, December 1976, pp. 1715-1729.
28. Sarwate, D. V., and Pursley, M. B., "Cross-correlation Properties of Pseudorandom and Related Sequences." *Proceedings of the IEEE*, Vol. 68, May 1980, pp. 593-619.
29. Geraniotis, E. A., and Pursley, M. B. "Error Probabilities for Slow-Frequency-Hopped Spread-Spectrum Multiple-Access Communications Over Fading Channels." *IEEE Transactions on Communications*, Vol. COM-30, May 1982, pp. 996-1009.
30. Stark, W. E. "Coding for frequency-hopped spread-spectrum communication with partial-band interference," Parts I and II. *IEEE Transactions on Communications*, Vol. 33, October 1985, pp. 1036-1057.
31. Pursley, M. B. "Frequency-hop transmission for satellite packet switching and terrestrial packet radio networks." *IEEE Transactions on Information Theory*, Vol. 32, September 1986, pp. 652-667.
32. Pursley, M. B. "The role of spread spectrum in packet radio networks." *Proceedings of the IEEE*, Vol. 75, January 1987, pp. 116-134.
33. Simon, M. K., Omura, J. K., Scholtz, R. A., and

- Levitt, B. K. *Spread Spectrum Communications*, Vol. I-III. Rockville, MD: Computer Science Press, 1985.
34. Ziemer, R. E., and Peterson, R. L. *Digital Communications and Spread Spectrum Systems*. New York: Macmillan, 1985.
35. Cooper, G. R., and McGillem, C. D. *Modern Communications and Spread Spectrum*. New York: McGraw-Hill, 1986.
36. Special Issue on Spread-Spectrum Communications, *IEEE Transactions on Communications*, Vol. COM-25, August 1977.
37. Special Issue on Mobile Spread-Spectrum Communications, *IEEE Transactions on Vehicular Technology*, Vol: VT-30, February 1981.
38. Special Issue on Spread-Spectrum Communications, *IEEE Transactions on Communications*, Vol. COM-30, May 1982.
39. "Spread Spectrum Communications I." *IEEE Journal on Selected Areas in Communications*, Vol. 8, No. 4, May 1990.
40. "Spread Spectrum Communications II." *IEEE Journal on Selected Areas in Communications*, Vol. 8, No. 5, June 1990.

25

Information Theory and Coding

Richard E. Blahut

Coding for Noiseless Channels	25-4
Capacity of Discrete Noiseless Channels	
State Diagrams and Trellises	
Source Compaction Codes	25-7
Source Models	
The Entropy Function	
Source Encoding	
Fixed-Length Block Codes	
Variable-Length Block Codes	
Variable-Length Tree Codes	
Universal Codes	
Coding for Discrete Noisy Channels	25-11
Mutual Information	
Channel Capacity	
Error-Control Codes	
Block Codes	
Convolutional Codes	
Continuous Channels and Sources	25-17
The Sampling Theorem	
Differential Entropy	
Entropy Power	
Capacity of a Continuous Channel	
The Additive Gaussian Noise Channel	
Waveform Channels	
Bit Energy and Bit Error Rate	
Signaling Without Bandwidth Constraints	
Signaling With a Bandwidth Constraint	

Decision Theory and Estimation Theory	25-22
Hypothesis Testing	
Estimation Theory	
The Matched-Filter Estimator	
Maximum Entropy and Minimum Discrimination	
Spectral Estimation	
Source Compression Codes	25-25
The Distortion-Rate Function	
Multiterminal Information Networks	25-26
Two-Way Channels	
Broadcast Channels	
Degraded Diversity Systems	
Remote Compaction of Dependent Data	
The Communication Games	25-27
The Jammer Saddle Point	
Spectrum Spreading	
Cryptography	

Information theory is a discipline centered around a common mathematical approach to the study of the collection and manipulation of information. It studies the theoretical basis of such activities as observation, measurement, data compression, data storage, communication, estimation, decision making, and pattern recognition. Many complex and expensive systems are built for automating or expanding these operations. Information theory attempts to guide the development of such systems based on a study of the possibilities and limitations inherent in mathematics and probability theory.

The communication problem, represented in Fig. 1, is the archetypal problem of information theory; much of the underlying structure and semantics of information theory is suggested by the communication channel. A source of information is to be connected to a user of information by a channel. A communication system is provided to prepare the source output for the channel and to prepare the channel output for the user. It consists of a device between the source and the channel called the encoder/modulator, and another device between the channel and the user called the demodulator/decoder.

It is conventional to partition the major functions of a modern communication system in the manner shown in

the block diagram of Fig. 2. Data from the data source are first processed by a source encoder, whose purpose is to represent the source data more compactly. A block of source data is represented by a sequence of symbols, usually binary, called the "source codeword." The data then are processed by the channel encoder, which transforms the symbols from a sequence of source codewords into another sequence of symbols called the "channel codeword." The channel codeword is a new, longer sequence that has more redundancy than the source codeword. Each symbol in the channel codeword might be represented by a bit or a group of bits. Next, the modulator represents each symbol of the channel codeword by its corresponding analog signal from a finite set of possible analog symbols. The sequence of analog symbols, called a "waveform," is transmitted through the channel. Because the channel is subject to various types of noise, distortion, and interference, the channel output differs from the channel input. The demodulator converts the noisy received sequence of analog symbols, possibly mutually interfering, into a sequence of discrete symbols of the channel codeword alphabet (based on a best estimate of the transmitted signal). The demodulated sequence of symbols is called the "received word." Because of channel noise, the demodulator sometimes makes errors, so the symbols



Fig. 1. The communication problem.

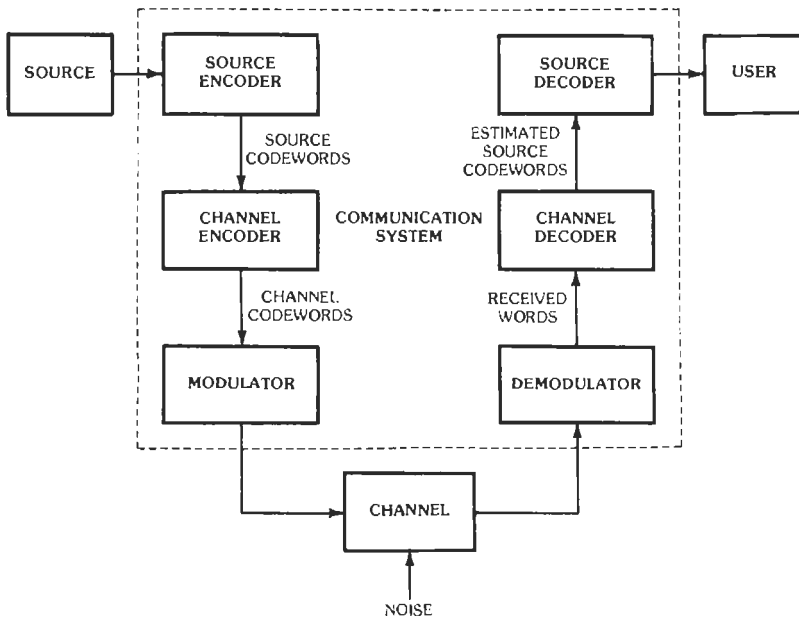


Fig. 2. Block diagram of a digital communication system.

of the received word do not always match those of the channel codeword.

The function of the channel decoder is to use the redundancy in a channel codeword to correct the errors in the received word, and then to produce an estimate of the source codeword from it. If all errors are corrected, the estimated source codeword matches the original source codeword. The source decoder performs the inverse operation of the source encoder and delivers its output to the user. The source encoder and decoder are studied under the terms "data compaction" and "data compression." The channel encoder and decoder are commonly split into two functions: implementing error control to negate the effects of channel noise, and preparing the sequence of transmitted symbols to be compatible with channel constraints. These are studied under the terms "error control codes" and "constrained channel codes." The modulator and the demodulator are studied under the term "modulation theory."

Although we use the terminology of communication theory, the model is general and applies to a great variety of situations. One can interpret many other information-handling systems, such as mass storage systems, in terms of this model. It only is necessary to identify the boundaries between the boxes. This is arbitrary and depends on the goals of a particular analysis. Usually, the source, channel, and user are identified with those parts of the system that are fixed, and the encoder/decoder and modulator/demodulator are identified with those parts of the system that are subject to design. Therefore, in different circumstances, the identification of these functions may be different.

The operation of the encoders and decoders is to map strings of symbols from one alphabet into strings of symbols from a second alphabet. The two alphabets are often the same, but they need not be. A *block code* breaks the input data stream into blocks of fixed length k and encodes each block into a codeword of fixed length n ; these are concatenated to form the output data stream. A *variable-to-fixed-length block code* breaks the input data stream into blocks of variable length and encodes the blocks into codewords of fixed length n that are concatenated to form the output data stream. A *fixed-to-variable-length block code* breaks the input data stream into blocks of fixed length k and encodes these into codewords of variable length that are concatenated to form the output data stream. A *tree code* breaks the input data stream into frames of length k_0 that are encoded into codeword frames of length n_0 with the encoding map depending on the previous m input data frames. The codeword frames are concatenated to form the output data stream. A tree code is called a *variable-to-fixed-length tree code* or a *fixed-to-variable-length tree code* when the input or output frames, respectively, are of variable length. A tree code with a finite encoding memory of m frames is called a *sliding block code* if the encoding operation is time invariant, and it is called a *convolutional code* if the encoding operation is both linear and time invariant.

Introductory information-theory textbooks intended for engineers are listed as references 1 through 4 at the end of this chapter. Other books and papers devoted to special topics will be cited in the appropriate section only if the topic is not treated within the general textbooks of the field.

CODING FOR NOISELESS CHANNELS

A discrete channel is a system by means of which an arbitrarily long sequence of symbols, each chosen from a finite set of I symbols $\{a_0, \dots, a_{I-1}\}$, can be transmitted from one point to another. The transmission of symbol a_i requires a certain time duration, t_i seconds, which is not necessarily the same for all i . A noiseless channel is one in which the output is completely determined by the input—errors do not occur. It is not always true that all possible sequences of symbols from the set $\{a_i\}$ can be transmitted through the channel. Some channels, called constrained channels, forbid certain sequences of symbols from being transmitted.

Teletypewriters and telegraphy are two simple examples of discrete channels that are historically important. In the teletypewriter case, there are 32 symbols, each of the same duration, and any sequence of the 32 symbols is allowed. Each symbol can be used to represent five bits of information. If the system transmits r symbols per second, it is natural to say that the channel has a capacity of $5r$ bits per second. This does not mean that the teletypewriter channel will always be transmitting information at this rate. Whether or not the actual rate reaches this maximum possible rate depends on how the source of information is connected to the channel.

For the telegraphy channel, convention has fixed the symbols as a dot, a dash, a letter space, and a word space. We formalize these symbols as follows: (1) a dot, consisting of line closed for one unit of time and then line open for one unit of time; (2) a dash, consisting of three time units of closure and one unit open; (3) a letter space, consisting of three time units of line open; (4) a word space, consisting of six time units of line open. We also impose the restrictions on allowable sequences that no space may directly follow another space. This we take as the formal definition of the telegraphy channel.

The Morse code is one system of encoding information for this channel. However, one may properly question the efficiency of the Morse code. Is there a limit on the information that can be conveyed through the telegraphy channel, and does the Morse code achieve this limit? These questions are answered by information theory.

Capacity of Discrete Noiseless Channels

The capacity, C (in units of bits/second), of a discrete noiseless channel is defined by

$$C = \lim_{T \rightarrow \infty} (1/T) \log_2 N(T)$$

where $N(T)$ is the number of allowed sequences of symbols of duration T .

The limit in the definition will exist and be finite in most cases of interest. From the definition, it is clear that about 2^{CT} different messages can be transmitted through the channel in T seconds for large enough T . We say that the channel can transmit C bits per second. It is easily seen that in the teletypewriter example this definition of capacity reduces to the previous result of $5r$ bits per second.

The evaluation of capacity is more difficult if the symbols are of different length, as in the telegraphy channel, or if certain sequences are forbidden, also as in the telegraphy channel. Suppose first that all sequences of the symbols a_0, \dots, a_{l-1} are allowed and these symbols have durations t_0, \dots, t_{l-1} that are integer multiples of one time unit. Let $N(t)$ represent the number of sequences of duration t . We can set up a recursive equation to find $N(t)$. A sequence of length t can be produced from a sequence of length $t - t_0$ by appending symbol a_0 , or it can be produced from a sequence of length $t - t_1$ by appending symbol a_1 , and so on. That is

$$N(t) = N(t - t_0) + N(t - t_1) + \dots + N(t - t_{l-1})$$

According to a well-known result in the study of finite-difference equations, $N(t)$ is asymptotic for large t to λ^t , where λ is the largest real root of the characteristic equation:

$$x^{-t_0} + x^{-t_1} + \dots + x^{-t_{l-1}} = 1$$

Then the capacity, C , is equal to $\log_2 \lambda$.

Even when there are restrictions on allowed sequences, we may still be able to write down a difference equation by inspection. In the case of the telegraphy channel:

$$N(t) = N(t - 2) + N(t - 4) + N(t - 5) + N(t - 7) + N(t - 8) + N(t - 10)$$

as we see by counting sequences of symbols according to the last or next-to-last symbol occurring. Hence, C equals $\log_2 \lambda$ where λ is the largest real zero of the polynomial $x^{-10} + x^{-8} + x^{-7} + x^{-5} + x^{-4} + x^{-2} = 1$. Solving this, we find $C = 0.539$ bit per unit of time. This is the maximum rate at which information can be conveyed by the telegraphy channel. One may devise many codes whose rates are close to the capacity of the channel, but never greater.

The Morse code shown in Table 1 is a widely used code for transmitting written text over the telegraphy channel. The Morse code combines both source encoding, described in the next section in terms of the *entropy* of the source, and channel encoding, described above.

The Morse code does not exploit the Markov structure of natural language. We will evaluate the code for use with a memoryless source whose 27 output letters (including the space) occur with the same probabilities as in English text.

The entropy of this memoryless model of English is 4.03 bits per letter, and the telegraphy channel capacity is 0.539 bit per unit time. Hence, an optimum code uses an average of 7.48 units of signaling time per source output letter. The Morse code in Table 1 uses an average of 9.296 units of signaling time per source output letter, which is 124 percent of the time needed by the optimum code. This establishes the amount by which the Morse code could be improved for use with the memoryless model of English text. Of course, an optimum code may be too complex for an operator to learn. The Morse code is an excellent compromise between performance and simplicity.

State Diagrams and Trellises

A constrained channel is one that does not accept input sequences containing any of a certain collection of forbidden subsequences. Such a channel can be described by a state diagram. Each state, s_0, s_1, \dots, s_{m-1} , corresponds to a recent past history of channel inputs. Loosely, we say that the channel is in one of these states, but we mean that its past history is described by that state. For each state, only certain symbols from the set a_0, \dots, a_{l-1} can be transmitted next. When one of these has been transmitted, the state changes to a new state depending both on the old state and on the particular symbol transmitted.

Sometimes the state diagram is augmented by a time axis so that one can see how the channel state changes with time. It is then called a *trellis*.

The telegraphy channel gives a simple illustration of a state diagram. There are two states. The state specifies whether or not a space was the last symbol transmitted. In state s_0 , only a dot or a dash can be sent next, and the state always changes. In state s_1 , any symbol can be transmitted, and the channel changes state if a space is sent; otherwise it remains in the same state. The conditions are indicated in the state diagram shown in Fig. 3.

The trellis diagram for the telegraphy channel is shown in Fig. 4. The Morse code is a collection of some of the paths through this trellis. An optimum code would use all of the paths, but the mapping from the set of source output sequences to the set of trellis paths might be a complicated one.

The digital magnetic-recording channel is another important example of a constrained channel. The channel input alphabet is binary, and a typical set of channel constraints, called run-length constraints, are that at least r zeros follow every one and that no more than s zeros occur in sequence. These constraints are imposed to minimize intersymbol interference and to facilitate clock recovery at the receiver.

In another kind of application, certain bit sequences

TABLE 1. THE MORSE CODE

Character	Per Letter Probability P_i	International Morse	Letter Duration l_i
Space	0.1859		6
A	0.0642	.-	9
B	0.0127	...-	13
C	0.0218	---.	15
D	0.0317	-. .	11
E	0.1031	.	5
F	0.0208	...-	13
G	0.0152	---.	13
H	0.0467	11
I	0.0575	..	7
J	0.0008-	17
K	0.0049	-. -	13
L	0.0321	13
M	0.0198	--	11
N	0.0574	-. .	9
O	0.0632	---	15
P	0.0152-	15
Q	0.0008	---.	17
R	0.0484	-. .	11
S	0.0514	9
T	0.0796	-	7
U	0.0228	...-	11
V	0.0083-	13
W	0.0175-	13
X	0.0013	---.	15
Y	0.0164	---.	17
Z	0.0005	----	15

$\sum_i p_i \log p_i = 4.03$ Bits $\sum_i p_i l_i = 9.296$ Time Units

are used for higher-level protocols or punctuation in the communication system and so must be forbidden in the data. A code must be used to translate from the set of unconstrained binary sequences into the set of sequences that satisfy the constraints.

Fig. 5 is a state diagram for a channel with $r = 2$ and $s = 7$ run-length constraints. Fig. 6 is the trellis diagram for the same channel; the trellis is drawn with state s_2 as the initial state. The characteristic equation

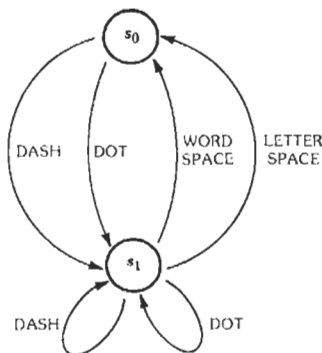


Fig. 3. State diagram for the telegraphy channel.

for this state diagram is found to be $x^8 - x^5 - x^4 - x^3 - x^2 - x - 1 = 0$. The capacity, C , of the channel, equal to the base-two logarithm of the largest zero of the characteristic polynomial, is 0.518 bit/bit. Codes exist that record 0.518 bit of information for each channel bit recorded, but no codes can do better. Practical codes in use fall a little short of the channel capacity. The Franaszek code* of Fig. 7 is a variable-to-variable-length code whose rate is 0.5. Hence, it achieves 96 percent of the channel capacity and yet is quite simple. A code that comes closer to channel capacity may be quite complex.

By inspection, one can see that any codeword of the Franaszek code followed by any other codeword satisfies the channel run-length constraints, and the codewords are uniquely decodable without punctuation. This last condition, known as the prefix condition, is due to the fact that no codeword looks like the beginning of any other codeword. This is an important requirement of a variable-length block code since the decoder must be able to break into codewords any concatenated string that is not punctuated.

Codes more elaborate than the Franaszek code are known, but they cannot be understood without develop-

* Reference 5.

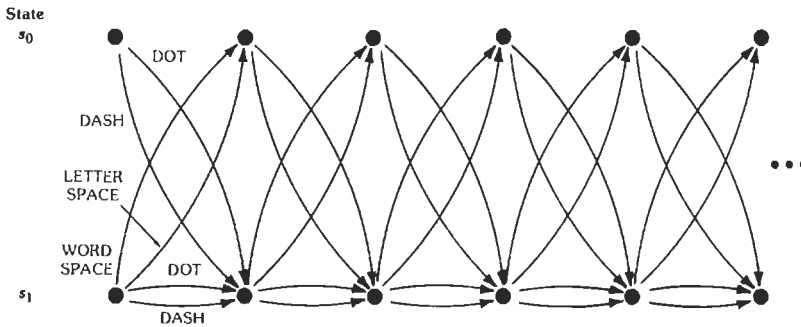


Fig. 4. Trellis diagram for the telegraphy channel.

ing more theory. The Adler-Hassner finite-state tree codes* encode a continuous bit sequence into a continuous bit sequence without word boundaries. Fig. 8 shows the encoder and decoder for an Adler-Hassner code for the $r = 1, s = 3$ run-length-limited channel. This channel has a channel capacity of 0.5515 bit/input. The rate one-half Adler-Hassner code implemented in Fig. 8 achieves 91 percent of the channel capacity. Binary symbols are clocked from the encoder into the channel and from the channel into the decoder at twice the source rate. Another Adler-Hassner code, say a rate $1/20$ code, could achieve a larger fraction of a channel capacity, but the implementation would be much more complex.

* Reference 6.

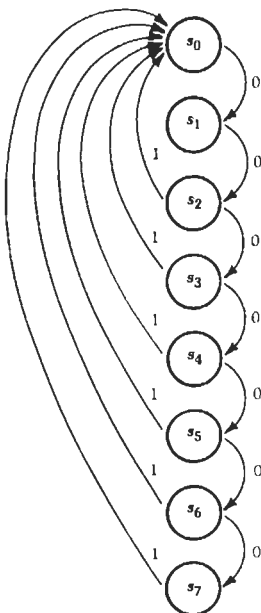


Fig. 5. State diagram of $r = 2, s = 7$ run-length-limited channel.

SOURCE COMPACTION CODES

Source compaction codes are used to represent the output of a data source more efficiently. Data sources such as facsimile, voice, digital recording, data tables, or word text can produce many millions or billions of bits. However, in their natural form, the data from these sources can be highly redundant. Practical data compaction codes are now available that can reduce considerably the number of bits needed to encode many such sources.

Source Models

An information source produces messages by generating a sequence of letters from a fixed alphabet of permitted symbols called the source alphabet. The alphabet may be finite, in which case the source is called a discrete source, or the source alphabet may be continuous, such as the set of real numbers, in which case the source is called a continuous source. A source might also put out continuous functions on the time axis. Sampling techniques are available to make this source into a time-discrete source.

The output of a discrete information source is a random sequence of symbols from a finite alphabet containing J symbols given by $\{a_0, a_1, \dots, a_{J-1}\}$. The sequence is produced according to some probability rule. The sources that are the easiest to study mathematically are those known as ergodic sources. An *ergodic source* is one in which every sequence of symbols produced by the source is the same in statistical properties. If it is observed long enough, such a source will produce, with probability approaching one, a sequence of symbols that is "typical." In simple terms, this means that a sufficiently long sequence from the source will almost always produce the same frequency of occurrence of symbols and symbol combinations.

An information source is said to be *memoryless* if successive symbols generated by the source are statistically independent. That is, a source is memoryless if each symbol is selected without influence from all

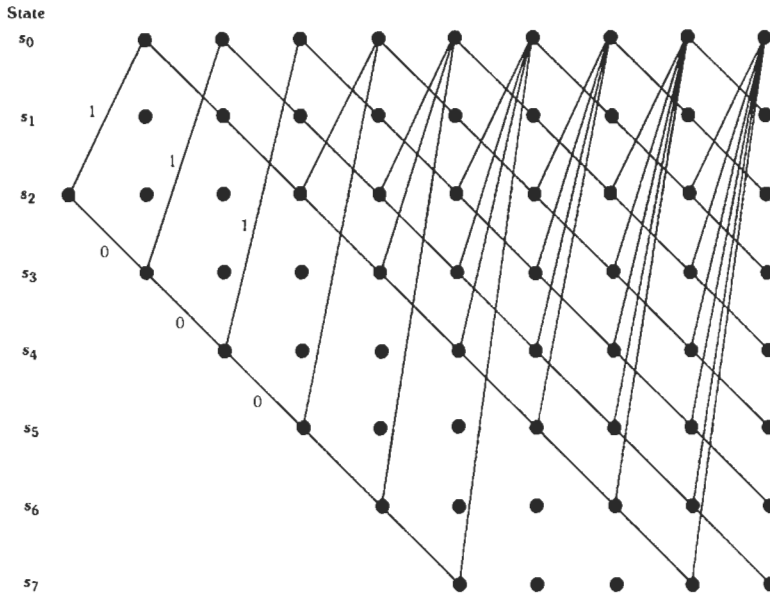


Fig. 6. Trellis diagram of $r = 2, s = 7$ run-length-limited channel.

previous symbols. If previously selected symbols influence the selection of a symbol, then the source is said to possess memory. If the selection of a symbol is influenced only by the immediately preceding symbol, the source is known as a *Markov source*. If the selection is influenced by the m previously selected symbols, the source is called an *m th-order Markov source*.

The Entropy Function

If a source output a_j occurs with probability $p(a_j)$, then the amount of information associated with the known occurrence of the event is defined to be

$$I(a_j) = -\log p(a_j)$$

In the definition, when logarithms are to the base 2, the information is measured in units of *bits*. When logarithms are taken to the base e , the information is measured in units of *nats* (a shortened form of natural units). Usually, the nat is more convenient for theoretic

cal developments, and the bit is more convenient for the final result. The conversion factor is: one nat equals 1.443 bits.

If the probabilities of selecting the source symbols are $p(a_0), p(a_1), \dots, p(a_{J-1})$, respectively, then the information generated each time a symbol is selected is $-\log_2 p(a_j)$ bits. Since the symbol a_j will, on average, be selected $np(a_j)$ times in a total of n selections, the average amount of information obtained from n selections is

$$np(a_0) \log_2 p(a_0)^{-1} + \dots + np(a_{J-1}) \log_2 p(a_{J-1})^{-1}$$

bits. Divide by n to obtain the average amount of information per symbol selection. This is known as the *average information*, the *uncertainty*, or the *entropy*, $H(p)$.

$$H(p) = -\sum_{j=0}^{J-1} p(a_j) \log_2 p(a_j) \text{ bits/symbol}$$

or

$$H(p) = -\sum_{j=0}^{J-1} p_j \log_2 p_j$$

where p_j is an abbreviation of $p(a_j)$.

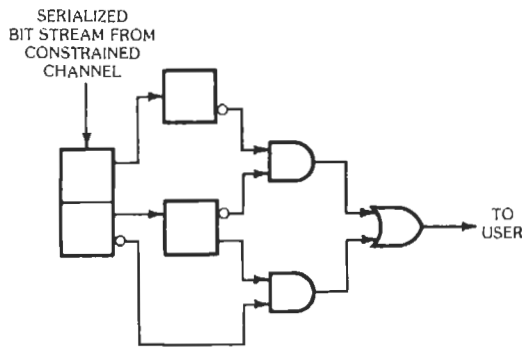
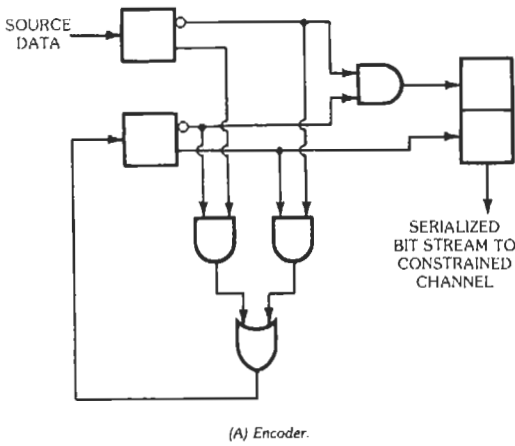
The entropy associated with n selections from the statistically independent set is equal to n times the entropy per single selection. The term entropy is used

INFORMATION WORDS	CODEWORDS
11	0100
10	1000
000	000100
010	001000
011	100100
0010	00001000
0011	00100100

Example:

1101010... → 01000010001000...

Fig. 7. The $r = 2, s = 7$ Franaszek rate one-half code.



- NOTES:
1. Every 1 followed by a 0.
 2. At most three consecutive 0's.
 3. Code rate equal to one half.
 4. Channel symbol rate twice source symbol rate.

Fig. 8. Encoder and decoder for an Adler-Hassner code.

since the function is the same as that derived in statistical mechanics for the thermodynamic quantity entropy.

For a first-order Markov source, the memory can be expressed by a conditional probability distribution. This conditional probability may be written as $p(a_j|a'_k)$ where a_j is any of the possible source symbols and a'_k is the source symbol emitted at the previous time instant.

The entropy is defined as an average

$$H(p) = - \sum_{k=0}^{J-1} p_k \sum_{j=0}^{J-1} p_{j|k} \log_2 p_{j|k}$$

where $p_{j|k}$ is an abbreviation of $p(a_j|a_k)$. This can also be rewritten

$$H(p) = - \sum_{k=0}^{J-1} \sum_{j=0}^{J-1} p_{jk} \log_2 p_{j|k}$$

where $p_{jk} = p_k p_{j|k}$.

The entropy function can be justified by its role in source compaction coding theorems, by its role in combinatorics, and by its intuitive properties as a measure of uncertainty. These properties are as follows:

1. The entropy function is continuous in p . This is a reasonable property, since it says that small changes in the probability distribution only make small changes in the uncertainty.
2. $H \geq 0$ and $H = 0$ if and only if all the p_j are zero except one which is unity. This a reasonable property of a measure of uncertainty, since if the outcome of a selection is sure, there is no uncertainty.
3. For a given n , $H(p) \leq \log n$ and is equal to $\log n$ if and only if the p_j are equal to $1/n$. This is a reasonable property which says that the more symbols available for selection without preference, the larger the uncertainty.

A binary source is one that contains only two symbols, a_0 and a_1 ; their probabilities are given by p and $1 - p$, respectively. The binary entropy function

$$H_b(p) = -p \log_2 p - (1 - p) \log_2 (1 - p)$$

is shown in Fig. 9.

The output of a binary source is a binary digit. The distinction between the bit used as a measure of information and the bit used as a binary output symbol

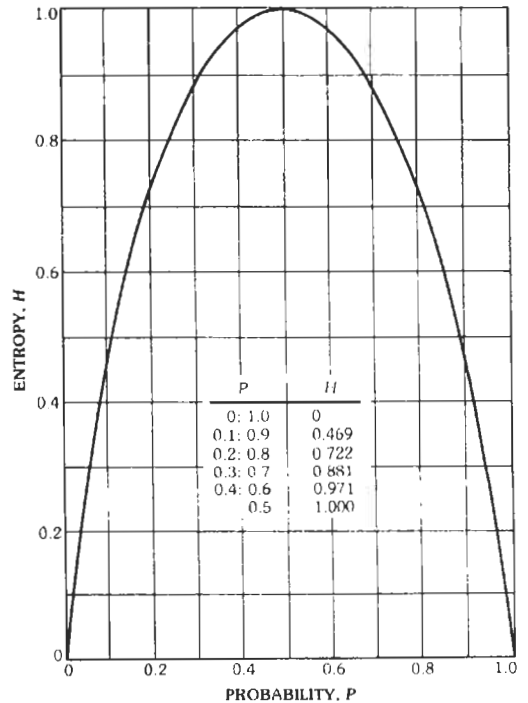


Fig. 9. The binary entropy function.

should be carefully noted. Fig. 9 shows that on average the amount of information provided by a binary source is always equal to or less than 1 bit/bit (one bit of information per data bit). The binary source provides one bit of information for each selected symbol only when the two symbols are equiprobable.

Source Encoding

The entropy function can be interpreted as the average amount of information necessary to specify which symbol has been produced by the source. If a source selects n symbols, where n is a very large number, then with high probability, it will select a sequence from the set of 2^{nH} different typical sequences each having a probability of occurrence of about $(1/2)^{nH}$. This is a direct physical interpretation of H . It means that, theoretically, a very long typical sequence of n q -ary symbols selected by the source can be encoded and retransmitted using only nH binary digits, each carrying one bit of information.

Usually, the binary codeword will be embedded in some longer binary text. Sometimes punctuation symbols are used and are not counted as part of the codeword. However, the more enlightened view is to incorporate the punctuation symbols into the code alphabet and devise an even more compact code. Then the code must be a self-punctuating code. Symbols to punctuate the code become implicit to the codeword.

Fixed-Length Block Codes

A fixed-length block code needs no punctuation. Compaction of a source with J symbols encodes n source output symbols into a codeword of length k bits, where n and k are fixed. There are J^n possible source output blocks the length n , so the blocklength of a binary code must be $\log_2 J^n$ bits if every possible output block is to be encoded. This requires $\log_2 J$ codeword bits per source output symbol. In general, this is much greater than the source entropy rate.

A fixed-length block code can encode at the entropy rate but then must allow the possibility of an error. A source compaction block code with rate a little above the entropy provides only $n(H(p) + \epsilon)$ bits, so $2^{n(H(p) + \epsilon)}$ codewords are assigned to the "typical" source output blocks—those blocks of greatest probability. By picking n large, the probability of a nontypical—hence noncodable—source output block can be made many orders of magnitude smaller than the probability of failure of the equipment. A fixed-length block code for source compaction must rely on this kind of reliability argument in order to encode at a rate near the entropy.

Variable-Length Block Codes

A variable-length block code encodes n source output symbols into k channel symbols, where either n or k is variable. The Huffman codes constitute the most com-

mon example. For them, k is variable and depends on the particular block of n symbols observed, while n is fixed.

The variable-length codewords of a Huffman code are constructed from knowledge of the probabilities of each of the J^n source output blocks of length n . We will construct a simple example using $n = 1$. The construction, illustrated in Fig. 10, proceeds as follows.

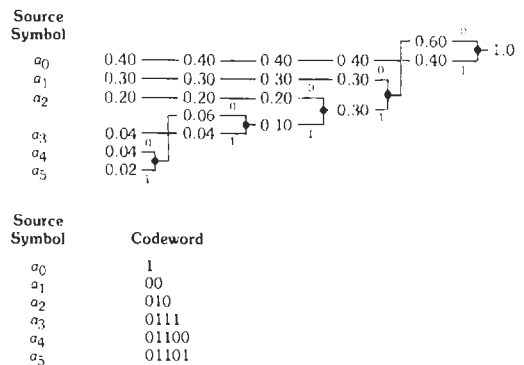
For convenience, the symbols are listed in order of decreasing probability. The two symbols of lowest probability are merged into a single symbol whose probability is the sum of the two constituent probabilities. At each step, this same procedure is repeated, merging two symbols into one symbol and adding their probabilities. The process stops when one symbol remains. Finally, read the tree from right to left, labeling the branches that leave each node either with a zero or with a one. The codeword is the string of labels from the rightmost node back to the original symbol.

A Huffman code can be formed just as easily for blocks of length n . Simply replace the source symbols in the construction with blocks of source symbols, and replace the symbol probabilities with probabilities of blocks. By choosing a large enough n —usually an n quite small will do—the average codeword length can be made as close to the source entropy as desired, and so a nearly optimal Huffman code can be constructed.

A Huffman code is a prefix code, sometimes said to be self-punctuating. It needs no explicit punctuation. Successive n blocks of source symbols can be encoded one after the other and concatenated. By decoding the codewords in a first-in, first-out fashion, the source symbols are uniquely recovered.

Variable-Length Tree Codes

A variable-length tree code for source compaction encodes indefinitely long strings of source output symbols into indefinitely long strings of code symbols. The encoding operation has a sliding structure whereby, as a



Average codeword length = 2.06 binary digits/symbol. The entropy of the source = 1.999 bits/symbol.

Fig. 10. Example of Huffman encoding.

few source symbols enter the encoder, a few codeword bits leave the encoder. The relationship between the number of symbols entering the encoder and the number of bits leaving it is variable, depending on the particular source symbols to be encoded.

One of the earliest variable-length tree codes for source encoding is the Elias code. Because of its simplicity, it is an excellent tutorial example. We describe an Elias code for a binary source with binary source alphabet $\{a_0, a_1\}$ and probability distribution $p = \{0.7, 0.3\}$. The entropy of this source is 0.88 bit. Imagine that the semi-infinite sequence of source output bits is a binary representation of a real number, r , in the interval $[0, 1]$. Refer to Fig. 11. The first symbol tells whether r is in the interval $[0, 0.7]$ or in the interval $[0.7, 1]$. This interval is itself subdivided in the same proportion, and the next source bit selects one of these two intervals. At each iteration, the process repeats, and, in this way, the process continues indefinitely.

The codeword, on the other hand, is a conventional binary representation of point r . As soon as enough source symbols are received to determine whether r is in the interval $[0, 0.5)$ or the interval $[0.5, 1]$, a codeword bit can be transmitted. The encoding is variable length since the number of source symbols needed to produce one codeword bit is random.

The decoder can begin its task after receiving only a few codeword bits. For example, if the binary sequence starts with 011 . . . , the point represented must lie between 0.375 and 0.50; hence the first symbol from the source must be a_0 . If the binary sequence starts with 0110, the point represented must lie between 0.375 and

0.4375; hence the first three symbols from the source must be $a_0a_0a_1$.

The Elias code is not practical because of precision problems. Errors in early calculations, no matter how small, will eventually cause encoding and decoding errors. A practical variable-length tree code, due to Pasco and Rissanen,* is available, but the description is more difficult.

Universal Codes

A source compaction code usually requires a probabilistic model of the source in order to attain its best performance. The encoder will encode to near the entropy of the source only if it is given a satisfactory model of the source. Some source compaction codes, called *universal source codes*, implicitly construct their own model of the source as they go along (assuming that past source sequences are representative of future sequences) and so encode near the entropy of whatever source they are given.

The Lempel-Ziv code is a universal variable-to-fixed-length source-compaction block code that does not require an externally constructed source model. The technique, illustrated in Fig. 12, is to break up the string of source output symbols into substrings of variable length but not longer than a largest allowed length. The encoder has a buffered copy of the raw data that has already been encoded. At each iteration, it searches for a prefix of the data yet waiting to be encoded within the data already encoded. The longest prefix of which a copy can be found among the symbols recently encoded, together with one more innovation symbol, becomes the next substring.

A substring is encoded in three parts: a binary-encoded pointer telling where a copy of the substring begins within recently encoded data; a binary-encoded number giving the length of the substring; and the value of the innovation symbol. By including an innovation symbol, the encoder cannot degenerate into an unencodable situation.

The Lempel-Ziv decoder must keep a copy of recently decoded data in which it looks to recreate subsequent substrings that it must decode. The decoder contains a finite-length buffer equal in length to the encoder finite-length buffer. In decoding a substring, the decoder buffer contains the same symbols that the encoder buffer contained when that substring was encoded. Hence, the decoder can use the pointer and the substring length to reconstruct the next substring.

CODING FOR DISCRETE NOISY CHANNELS

The deep-space Gaussian noise channel and the noisy binary channel are two important examples of noisy channels. The first example is a continuous-time, continuous-amplitude channel that is usually made into a

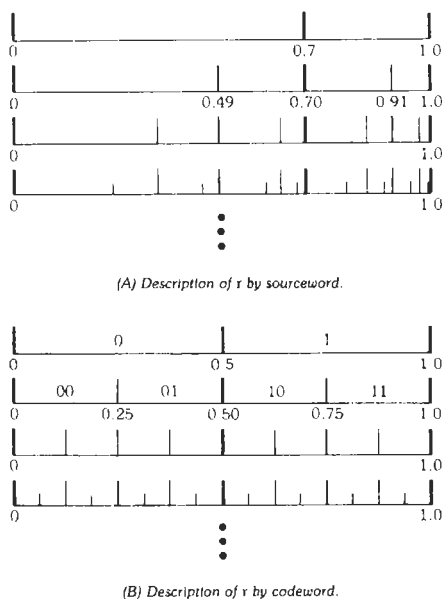


Fig. 11. Construction of an Elias code.

* Reference 5.

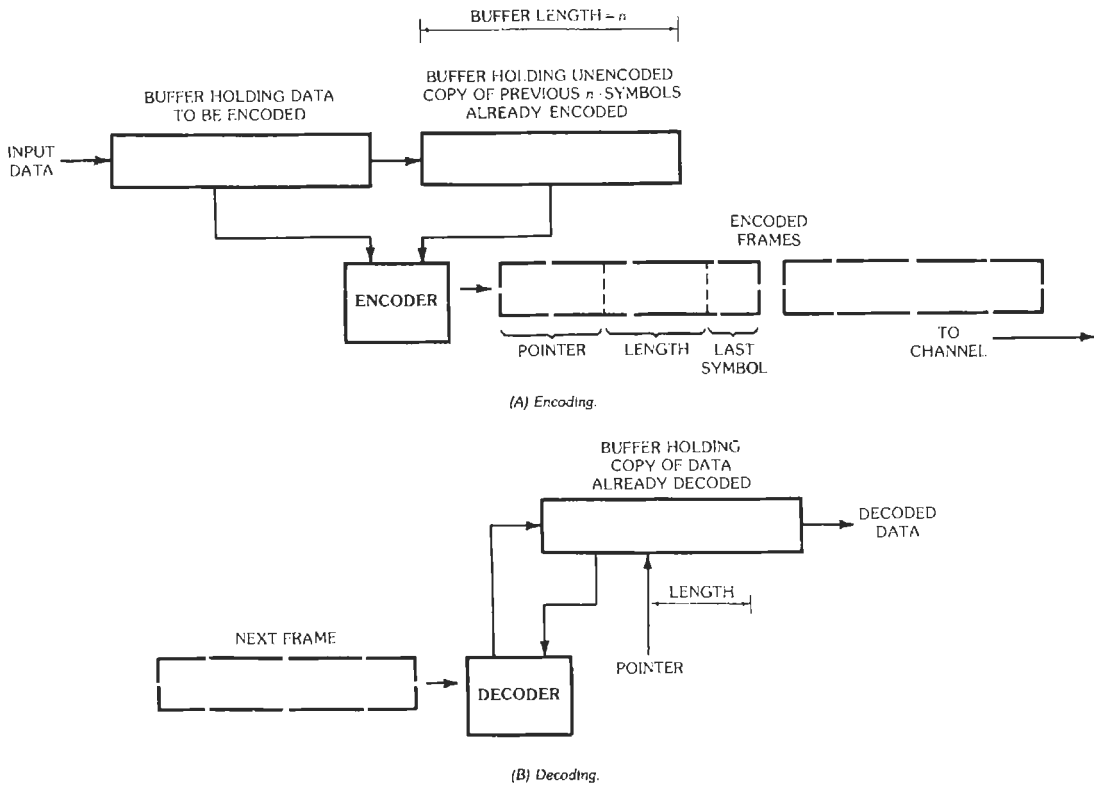


Fig. 12. Lempel-Ziv encoder and decoder.

discrete channel by the modulator/demodulator; the second is a discrete channel.

A *discrete channel* is a system in which a sequence of letters chosen from a finite set of symbols $\{a_0 \dots a_{I-1}\}$ can be transmitted from one point to another. A *noisy channel* is one for which the output symbol is not completely determined by the input symbol; only some probability distribution on the set of output symbols is determined by the input symbol. If the probability distribution is independent of previous inputs or outputs from the channel, the channel is called *memoryless*. Let $\{b_0, \dots, b_{J-1}\}$ be the set of channel output symbols; possibly $J \neq I$. Let $Q(b_j|a_i)$, abbreviated Q_{ji} , be the probability that symbol b_j is received given that a_i was sent. It is called the *channel transition probability*, and the J by I matrix $\{Q_{ji}\}$ is called the *transition matrix* of the channel. Some simple channels are shown in Fig. 13.

Information can be sent reliably through a discrete noisy channel by the use of elaborate crosschecking techniques known as error control codes. A noisy channel (such as a magnetic tape) may also have constraints. For such a channel, common practice is to treat the translation code that satisfies the constraints separately from the error control, and then use the two codes in series as shown in Fig. 14. Single codes that combine both functions neatly have not yet been

discovered. An inner encoder/decoder matches its input symbols to the channel. When the channel makes an error, the inner decoder makes an error. The outer encoder/decoder has no constraints on the input sequence.

Mutual Information

Let $p(a_i)$, abbreviated p_i , be the probability that input symbol a_i is sent through the channel. The probability, $q(b_j)$, abbreviated q_j , of symbol b_j being the channel output is

$$q_j = \sum_{i=0}^{I-1} p_i Q_{ji}$$

Given output b_j , the probability, $P(a_i|b_j)$, abbreviated P_{ij} , that symbol a_i was transmitted is given by the Bayes rule

$$P_{ij} = \frac{p_i Q_{ji}}{\sum_i p_i Q_{ji}}$$

The I by J matrix $\{P_{ij}\}$ is called the *backward transition matrix* of the channel.

The average information per symbol carried by the input random variable, X , which takes value a_i with a probability $p(a_i)$, is the entropy.

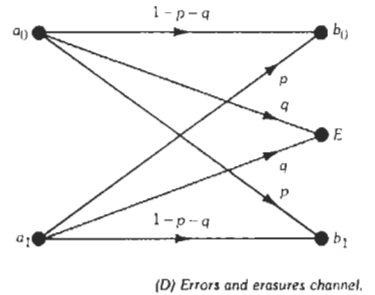
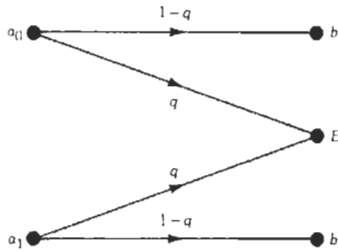
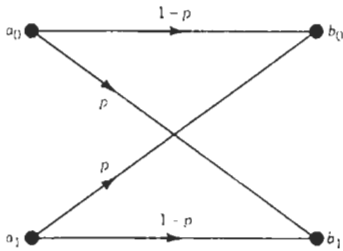
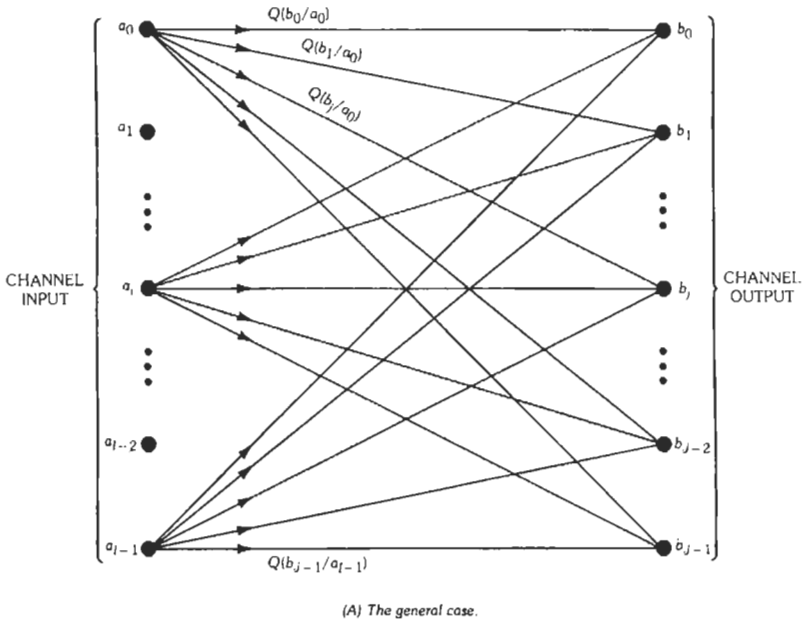


Fig. 13. Some discrete memoryless channels.

$$H(X) = - \sum_{i=0}^{I-1} p_i \log_2 p_i$$

Conditional on b_j being received, the average information per symbol carried by input random variable X is

$$H(X|b_j) = - \sum_{i=0}^{I-1} P_{ij} \log_2 P_{ij}$$

The conditional input entropy, or *equivocation*, is then defined as an expectation

$$H(X|Y) = - \sum_{j=0}^{J-1} q_j \sum_{i=0}^{I-1} P_{ij} \log P_{ij}$$

bits/symbol. On average, this is the remaining uncertainty in the channel input after the channel output has been observed.

The difference between the average uncertainty in channel input X before and after channel output Y is received is called the *average mutual information* of the channel.

$$I(X; Y) = H(X) - H(X|Y)$$

$$= - \sum_{i=0}^{I-1} p_i \log p_i + \sum_{j=0}^{J-1} \sum_{i=0}^{I-1} q_j P_{ij} \log P_{ij}$$

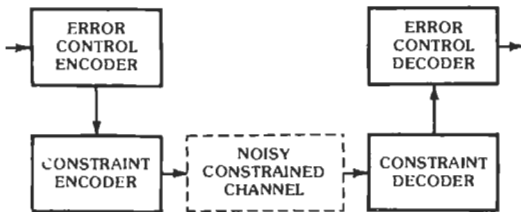


Fig. 14. Channel coding for noisy constrained channels.

By the Bayes formula, this becomes

$$I(X; Y) = \sum_{i=0}^{I-1} \sum_{j=0}^{J-1} p_i Q_{ji} \log \frac{Q_{ji}}{\sum_i p_i Q_{ji}}$$

The average mutual information is a measure of the average amount of information about the channel input that can be received by the user by observing the symbol at the output of the channel.

The average mutual information has a number of important and satisfying properties:

1. Average mutual information is nonnegative and is strictly positive unless the channel output is independent of the channel input.
2. If $I = J$ and the channel is noiseless (i.e., $Q_{ji} = 1$ if $j = i$), then $I(X; Y) = H(X)$, and the average mutual information between the output and the input of the channel is equal to the average information into the channel.
3. Average mutual information is symmetric in X and Y . That is $I(X; Y) = I(Y; X)$ where

$$I(Y; X) = \sum_{j=0}^{J-1} \sum_{i=0}^{I-1} q_j P_{ij} \log \frac{P_{ij}}{\sum_j q_j P_{ij}}$$

Channel Capacity

The *channel capacity*, C , is defined to be the maximum rate at which information can be transmitted through a channel. The *fundamental theorem of information theory* says that at any rate below channel capacity, an error control code can be designed whose probability of error is arbitrarily small. Intuitively, in a well-designed message, an isolated channel input symbol a_i should occur with a probability p_i such that the average mutual information is maximized.

For the noisy channel, with transition matrix Q , the channel capacity, C , is expressed mathematically by

$$C = \max_p \sum_{i=0}^{I-1} \sum_{j=0}^{J-1} p_i Q_{ji} \log \frac{Q_{ji}}{\sum_i p_i Q_{ji}}$$

where the maximum is over all probability distributions on the input alphabet. For simple channels, the capacity can be evaluated by finding the maximum analytically. For more difficult channels, efficient computational algorithms exist.

The capacity of the binary symmetric channel with transition matrix

$$\begin{bmatrix} 1 - p & p \\ p & 1 - p \end{bmatrix}$$

is

$$C = 1 + p \log_2 p + (1 - p) \log_2 (1 - p)$$

bits per input symbol. The capacity of the binary erasure channel with erasure probability q is

$$C = 1 - q$$

bits per input symbol. This is somewhat surprising, since it says that the erasures cause a loss in capacity exactly equal to the fraction of symbols erased despite the fact that the encoder does not know which nq symbols will be erased. The capacities of these two simple channels are shown in Fig. 15.

Error-Control Codes

Powerful codes are now available for use with q -ary symmetric channels whenever q is a prime or a prime power. Codes in use can achieve extremely small probability of symbol error, but the known codes are not good enough to achieve code rate near the channel capacity. Error control codes are routinely used for many purposes: to reduce required link signal-to-noise ratio; to protect against various forms of pulsed and other interference; to protect against hostile interference in military systems; to protect against noise and media defects in magnetic recording systems; and so forth.

Error-control codes in common use work with a code alphabet of size 2^m where m is a positive integer. For m equal to eight, the code alphabet consists of the set of all eight-bit bytes. It is common to use byte-organized codes because errors in practical systems tend to occur in bursts and because the hardware is simpler.

Error-control codes can be used for error detection or error correction. When they are used for error detection, a request for retransmission (called *automatic repeat request (ARQ)*) is fed back to the transmitter. When they are used for error correction (called *forward error correction (FEC)*), the link operates without interruption, and the error control is transparent to other levels of the system. Forward error correction can operate with noisy links and low transmitted power, whereas ARQ cannot; ARQ will play a reduced role in future systems.

English-language textbooks on error-control codes include Berlekamp (reference 8), Blahut (reference 9), Clark and Cain (reference 10), Lin and Costello (reference 11), and Peterson and Weldon (reference 12).

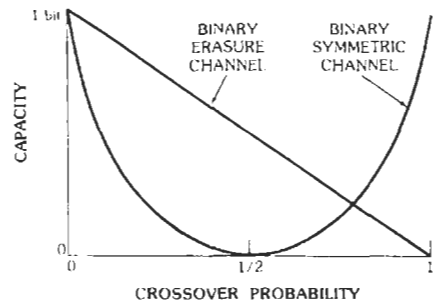


Fig. 15. Capacity of binary symmetric and binary erasure channels.

Block Codes

An (n, k) block code for error control encodes a block of k information symbols into a block of n codeword symbols. Each symbol is an m -bit byte. The rate, R , of the code is equal to k/n .

The rules for constructing good block codes for error control make use of the arithmetic of Galois fields. A Galois field with 2^m elements, denoted $GF(2^m)$, is an arithmetic system containing the operations of addition, subtraction, multiplication, and division defined in such a way that most arithmetic and algebraic procedures are valid. The arithmetic operations themselves are unconventional but have the enormous advantage that there is no overflow nor round-off error. Since the error control code is used to process bit packages but not to do real computations, it does not matter that the arithmetic rules are unconventional. Chart 1 shows addition and multiplication tables for several simple Galois fields. Notice that $GF(2)$ and $GF(3)$ are modulo 2 and modulo 3 arithmetic, respectively, but $GF(4)$ is *not* modulo 4 arithmetic. Modulo 4 arithmetic cannot form a field because $2 \cdot 1 = 2 \cdot 3 \pmod{4}$, so division by 2 does not behave properly. Large fields such as $GF(256)$ are very important in practice, but the multiplication tables are too large to show here.

Chart 2 shows the set of codewords of the Hamming $(7, 4)$ code, a binary code—symbols in $GF(2)$ —that can correct a single bit error. This simple code has sixteen codewords. Chart 2 also shows some of the codewords of the Reed-Solomon $(7, 5)$ code, an octal code—symbols are in $GF(8)$ —that can correct a single octal symbol in error. This code has 8^5 , or 32 768, codewords. This second example is actually a very small code although the number of codewords is already

too great to enumerate. This is why it is important to use the computational structure of a Galois field to construct the encoders and decoders. Reed-Solomon codes as large as a $(256, 224)$ code over $GF(2^8)$ are now quite practical. This code consists of 224 information symbols followed by 32 parity symbols and can correct 16 symbol errors; a symbol error is a symbol that is wrong in any possible way. A symbol in $GF(2^8)$ may be used to represent eight bits by the user, yet on the channel may be represented by a 256-ary symbol.

The most important block codes in practice are those known as the *BCH codes* and the *Reed-Solomon codes*. They are important because efficient decoding algorithms exist for them. These decoders are based either on the *Berlekamp-Massey algorithm* or on the *Euclidean algorithm*, and they have a hardware cost proportional to nt where n is the blocklength and t is the number of errors to be corrected.

Convolutional Codes

A convolutional code encodes a stream of information symbols into a stream of codeword symbols. The duration of the stream is so long that it is effectively infinite and does not enter into the design of the encoder and decoder. An information sequence is shifted into the encoder beginning at time zero and continuing indefinitely into the future. The stream of incoming information symbols is broken into segments of k_0 symbols called information frames, which often may be as short as one symbol in practice. The encoder can store m frames. During each frame time, a new information frame is shifted into the encoder, and the oldest information frame is shifted out and discarded. At the end of any frame time, the encoder has stored the most recent m frames, a total of mk_0 information symbols. At the beginning of a frame, from the new incoming information frame and the m previously stored information frames, the encoder computes a single codeword frame of length n_0 symbols. This codeword frame is shifted out of the encoder as the next information frame is shifted in. Hence, the channel must transmit n_0 codeword symbols for each k_0 information symbols. The rate, R , of the convolutional code is defined as $R = k_0/n_0$.

The *constraint length*, v , of a convolutional code is defined as the number of memory stages in a minimum encoder. The complexity of a decoder is often q^v . Binary convolutional codes used in practice may have a constraint length in the range of about seven to forty. Fig. 16 shows an encoder for a much simpler convolutional code, one of constraint length 2 with $k_0 = 1$ and $n_0 = 2$. A trellis description of the code is shown in Fig. 17. The convolutional code is the set of all semi-infinite binary words that may be read off any path through the trellis. A "zero" information bit entering the encoder at any node is encoded into the upward path out of that node, and a "one" information bit is encoded into the lower path. This convolutional code is able to correct any number of error events each containing two bit

CHART 1. EXAMPLES OF FINITE FIELDS
 $GF(2)$

+		0	1	·		0	1
0	0	0	1	0	0	0	0
1	1	0	1	1	0	1	1

$GF(3)$

+		0	1	2	·		0	1	2
0	0	0	1	2	0	0	0	0	0
1	1	2	0	1	1	0	1	2	1
2	2	0	1	2	2	0	2	1	2

$GF(4)$

+		0	1	2	3	·		0	1	2	3
0	0	0	1	2	3	0	0	0	0	0	0
1	1	0	3	2	1	1	0	1	2	3	1
2	2	3	0	1	2	2	0	2	3	1	2
3	3	2	1	0	3	3	0	3	1	2	3

CHART 2. EXAMPLES OF BLOCK CODES

Hamming (7, 4) Code							Reed-Solomon (7, 5) Code						
Information				Parity			Information				Parity		
0	0	0	0	0	0	0	0	0	0	0	0	0	
0	0	0	1	0	1	1	0	0	0	0	1	6	3
0	0	1	0	1	1	0	0	0	0	0	2	7	6
0	0	1	1	1	0	1	0	0	0	0	3	1	5
0	1	0	0	1	1	1				⋮			
0	1	0	1	1	0	0	0	0	0	1	0	1	1
0	1	1	0	0	0	1	0	0	0	1	1	7	2
0	1	1	1	0	1	0	0	0	0	1	2	6	7
1	0	0	0	1	0	1	0	0	0	1	3	0	4
1	0	0	1	1	1	0				⋮			
1	0	1	0	0	1	1	0	0	0	7	0	7	7
1	0	1	1	0	0	0	0	0	0	7	1	1	4
1	1	0	0	0	1	0	0	0	0	7	2	0	1
1	1	0	1	0	0	1	0	0	0	7	3	6	2
1	1	1	0	1	0	0				⋮			
1	1	1	1	1	1	1	0	0	1	0	0	7	3
							0	0	1	0	1	1	0
							0	0	1	0	2	0	5
							0	0	1	0	3	6	6
										⋮			

errors provided the error events are spaced far enough apart for the decoder to clear one error event before it sees the next.

Convolutional codes in use have a larger constraint length. The most important decoders for convolutional codes either use the *Viterbi algorithm* or use *sequential*

decoding. The Viterbi algorithm has a hardware cost proportional to 2^v and so is practical only for small v . Sequential decoding, using either the *Fano algorithm* or the *stack algorithm*, has a hardware cost that grows only slowly with v , but has a random decoding delay and will occasionally overflow any finite-sized decoding buffer.

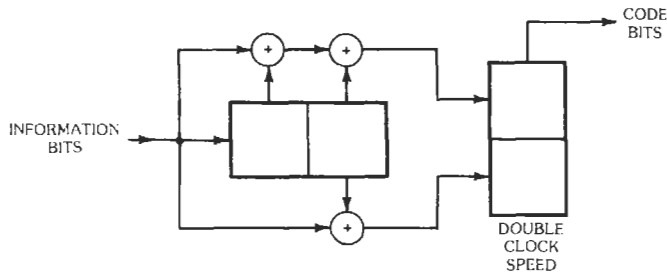


Fig. 16. Encoder for a convolutional code with constraint length 2.

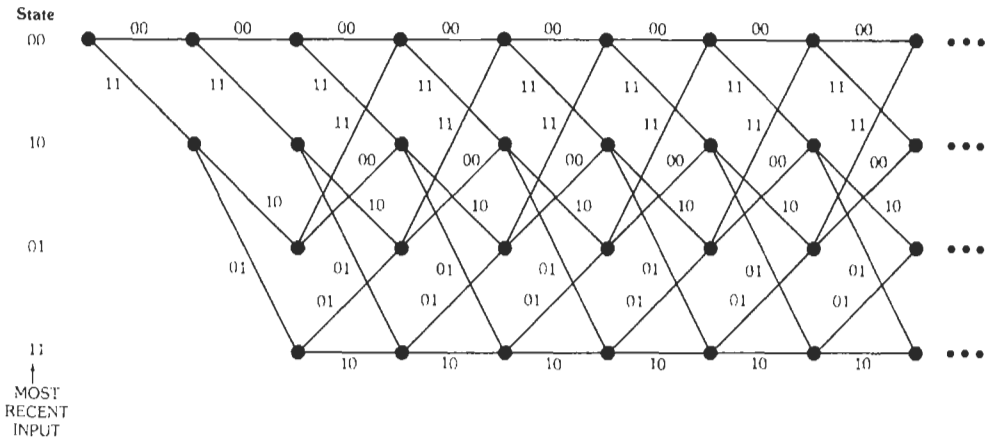


Fig. 17. Trellis diagram for a convolutional code with constraint length 2.

CONTINUOUS CHANNELS AND SOURCES

A discrete information source generates information at a finite rate; the entropy rate, which measures the information generated, is finite. A continuous information source can assume any one of an infinite number of amplitude values and so requires an infinite number of binary digits for its exact specification. The entropy rate is infinite. An immediate consequence of this is that in order to transmit the output of a continuous information source and recover it exactly, a channel of infinite capacity is required. Since, in practice, every continuous channel is perturbed by noise and therefore has a finite capacity, it is not possible to transmit the output of a continuous source over a channel and recover it exactly; there is always some distortion.

We distinguish among a discrete source or channel, a continuous source or channel defined on discrete time instants, and a waveform source or channel defined on a continuum of time points.

A waveform channel can be made into a discrete channel by choosing a finite set of modulation waveforms for the channel. The relationship between waveform design and information theory is discussed in Wozencraft and Jacobs.* A waveform source can be made into a continuous source by sampling, provided the bandwidth is finite, and into a discrete source by sampling and quantization provided a small amount of distortion can be introduced.

The Sampling Theorem

The sampling theorem is an important aid in the design and analysis of communication systems involving the use of continuous time functions of finite

bandwidth. The theorem states that, if a function of time, $f(t)$, contains no frequencies of W hertz or higher, then it is completely determined by giving the value of the function at a series of points spaced $(2W)^{-1}$ seconds apart. The sampling rate of $2W$ samples per second is called the *Nyquist rate*.

If $f(t)$ contains no frequencies of W hertz or higher, then it can be recovered from its samples by the Nyquist-Shannon interpolation formula:

$$f(t) = \sum_{n=-\infty}^{+\infty} f(n/2W) \{[\sin \pi(2Wt-n)]/\pi(2Wt-n)\}$$

The sampling theorem makes no mention of the time origin of the samples; it is only the spacing of the samples that matters.

If function $f(t)$ is negligible in magnitude outside a time interval T and has negligible energy at frequencies higher than W hertz, it can be specified by $2TW$ ordinates. If a Gaussian noise process with rectangular spectrum is sampled at the Nyquist rate, the samples are independent.

Differential Entropy

The differential entropy of a continuous random variable, X , with probability density function $p(x)$ is defined as

$$H(X) = - \int_{-\infty}^{+\infty} p(x) \log_2 p(x) dx$$

The differential entropy is *not* the limiting case of the entropy; the entropy of a continuous distribution is infinite. The differential entropy is not invariant under coordinate transformations. By itself, it has no fundamental physical meaning, but it occurs often enough to have a name.

* Reference 13.

The Gaussian probability distribution

$$p(x) = [1/(\sqrt{2\pi} \sigma)] \exp(-x^2/2\sigma^2)$$

plays a major role in information problems. Of all probability distributions with variance σ^2 , the Gaussian distribution has the largest differential entropy, given by

$$H(X) = \frac{1}{2} \log_2 (2\pi e\sigma^2)$$

The differential entropy of a vector of length n whose components are independent identically distributed Gaussian random variables is

$$H(X) = \frac{1}{2} n \log_2 (2\pi e\sigma^2)$$

Consequently, if a Gaussian noise process of ideal rectangular spectrum is sampled at the Nyquist rate of $2W$ samples per second, we say the differential entropy rate of the source is $W \log_2 (2\pi e\sigma^2)$. A Gaussian noise process whose spectrum is constant for frequencies below W hertz and is zero for larger frequencies is called *band-limited white noise*.

Entropy Power

The *entropy power* of a random signal is defined to be the power of band-limited white noise having the same differential entropy rate and bandwidth as the original noise.

If a random signal has an entropy H (bits), the power of band-limited white noise having the same entropy rate is given by

$$N_e = (1/2\pi e)2^{2H}$$

The power N_e is the entropy power of the random signal.

It should be noted that since white noise has the maximum differential entropy for a given power, the entropy power of any random signal is less than or equal to its actual power.

Capacity of a Continuous Channel

The average mutual information between two continuous random variables X and Y with joint probability density function $P(x, y)$ is

$$\begin{aligned} I(X; Y) &= H(X) - H(X|Y) = H(Y) - H(Y|X) \\ &= H(X) + H(Y) - H(X, Y) \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} P(x, y) \log_2 [P(x, y)/p(x)p(y)] dy dx \\ &= \int_{-\infty}^{\infty} p(x) \int_{-\infty}^{\infty} Q(y|x) \end{aligned}$$

$$\log_2 \frac{Q(y|x)}{\int_{-\infty}^{\infty} p(x) Q(y|x) dx} dy dx$$

bits. Even though the differential entropy terms appearing here are not invariant under a coordinate transformation, their difference is. The mutual information has a fundamental significance.

A continuous channel usually has some constraint on the input probability distribution; an average power constraint or a peak power constraint is the most common. Therefore, $p(x)$ must satisfy one or more constraint equations of the form

$$\int_{-\infty}^{\infty} e(x)p(x)dx \leq S$$

where $e(x)$ is some nonnegative function. The average power constraint is

$$\int_{-\infty}^{\infty} x^2 p(x) dx \leq S$$

The capacity of the continuous channel is the maximum value of the average mutual information

$$C = \max_{p(x)} I(X; Y)$$

where the maximum is taken over the set of probability distributions that satisfy the constraints.

The Additive Gaussian Noise Channel

Calculation of the capacity can be a difficult task requiring a computer, but the important case of an additive Gaussian noise channel subject to an average power constraint can be solved analytically.

The channel output, Y , is given by $Y = X + Z$ where X is the channel input and Z is the noise. Since X and Z are independent, $H(Y|X) = H(Z)$, and $I(X; Y) = H(Y) - H(Z)$. For Gaussian noise, $H(Z) = \frac{1}{2} \log_2 2\pi eN$ where N is the noise power. Then

$$C = \max_{p(x)} 2W[H(Y) - \frac{1}{2} \log_2 2\pi eN]$$

Factor $2W$ converts the units of capacity to bits/second. Output Y has variance $S + N$ where S is the average power constraint. Entropy $H(Y)$ is largest if Y is Gaussian, which will be true if X is Gaussian of variance S . Then

$$\begin{aligned} C &= W \log_2 2\pi e(S + N) - W \log_2 2\pi eN \\ &= W \log_2 (1 + S/N) \end{aligned}$$

This is Shannon's formula for the capacity of the additive Gaussian noise channel with an average power

constraint. The channel input signal that achieves channel capacity has an amplitude distribution described by the Gaussian function.

If a channel has additive but non-Gaussian noise and a power constraint, then the capacity may be difficult to calculate exactly. However, upper and lower bounds on the capacity are easy to obtain. At average transmitted power S , the capacity, C , in bits/second, is bounded by the inequalities

$$W \log_2 [(S + N)/N_e] \geq C \geq W \log_2 [(S + N)/N]$$

where,

- W is the bandwidth,
- N is the average noise power,
- N_e is the entropy power of the noise.

Waveform Channels

An additive Gaussian noise waveform channel with a general transfer function is shown in Fig. 18. The capacity of the waveform channel under Gaussian noise is equal to the capacity one would obtain if the transfer function were approximated by many thin ideal rectangular transfer functions of different amplitudes. The capacity is given by the so-called "water-pouring" formulas, written parametrically in terms of θ .

$$C(S, \theta) = \frac{1}{2} \int_{-\infty}^{\infty} \max[0, \log \frac{\theta}{N(f)|H(f)|^2}] df$$

$$S(\theta) = \int_{-\infty}^{\infty} \max[0, \theta - N(f)|H(f)|^2] df$$

The reason these are called the water-pouring formulas can be understood from Fig. 19. The input information is imagined as "poured" into a vessel whose shape is defined by $H(f)$ and $N(f)$. This produces the optimum spectral shape of input waveform $S(f) = \max[0, \theta - N(f)|H(f)|^2]$ that achieves the channel

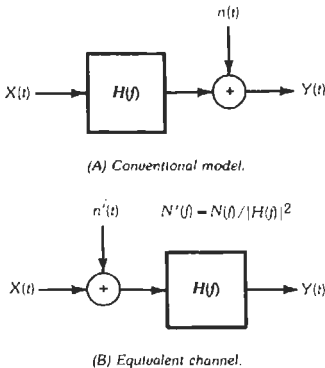


Fig. 18. Gaussian-noise waveform channel.

capacity. The important lesson given by the water-pouring principle is that optimum waveforms put most of their energy in the spectral region where the channel is good and little or no energy in the spectral region where the channel is poor. The optimum strategy is exactly the opposite of an often-used equalization strategy that "boosts" the skirts of the channel by adding extra gain there.

Bit Energy and Bit Error Rate

The performance of a digital communication system is measured by the probability of bit error, also called the *bit error rate* (BER). On an additive Gaussian noise channel, the bit error rate can always be reduced by increasing transmitted power, but it is by the performance at low transmitted power that one judges the quality of a digital communication system. The better of two systems, otherwise the same, is the one that can achieve a desired bit error rate with the lower transmitted power.

Given a message, $s(t)$, of duration T containing K information bits, the bit energy, E_b , is given by

$$E_b = E_m/K$$

where

$$E_m = \int_0^T s(t)^2 dt$$

is the message energy.

Bit energy E_b is calculated from the message energy and the number of information bits at the input to the encoder/modulator. At the input to the channel, one may find a message structure in which he perceives a larger number of bits. The extra symbols may be parity symbols for error control, or symbols for frame synchronization or channel protocol. These symbols do not represent transmitted information, and their energy must be amortized over information bits. Only information bits are used in calculating E_b .

For an infinite-length message of rate R information bits/second, E_b is defined by

$$E_b = S/R$$

where S is the message average power.

In addition to the message energy, the receiver also sees a white noise signal of one-sided spectral density N_o watts/hertz. Only the ratio E_m/N_o or E_b/N_o affects the bit error rate because the reception of the signal cannot be affected if both the signal and the noise are doubled. Signaling schemes are compared by comparing their respective graphs of BER versus required E_b/N_o .

It is possible to make precise statements about values of E_b/N_o for which good waveforms exist; these are a

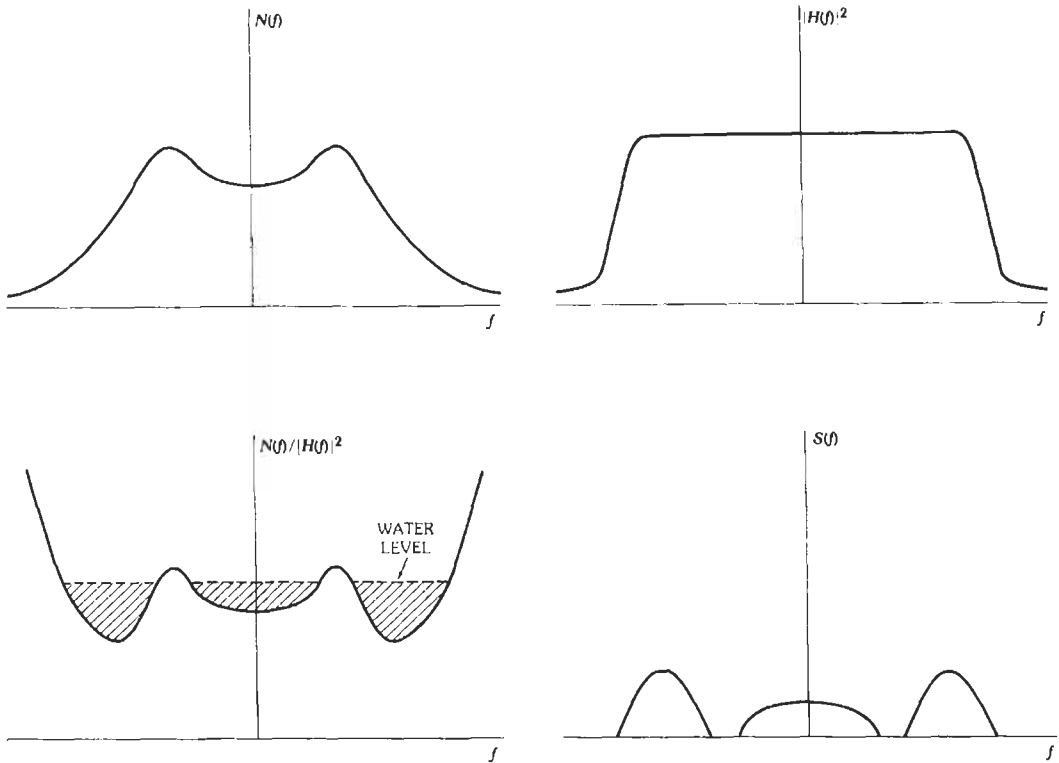


Fig. 19. Water pouring.

consequence of the channel capacity formula for the ideal rectangular bandpass channel in additive Gaussian noise. Let the signal power be $S = E_b R$ and the noise power be $N = N_o W$. Then

$$C/W = \log_2 (1 + RE_b/N_o W)$$

Define the spectral bit rate, r (measured in bits per second per hertz), by

$$r = R/W$$

The spectral bit rate, r , and E_b/N_o are the two most important figures of merit of a digital communication system.

Since the rate, R , is less than but can be made arbitrarily close to the capacity, C , the capacity formula becomes

$$E_b/N_o > (2^r - 1)/r$$

but E_b/N_o can be arbitrarily close to the bound by designing a sufficiently sophisticated digital communication system. This inequality, shown in Fig. 20, tells us that increasing the bit rate per unit bandwidth increases the required energy per bit. This is the basis of

the energy/bandwidth trade of digital communication theory where increasing bandwidth at a fixed information rate can reduce power requirements.

Every communication system can be described by a point lying below the curve of Fig. 20. Any communication system that attempts to operate above the curve will lose enough data through errors so that its actual data rate will lie below the curve. By the fundamental theorem of information theory, for any point below the curve one can design a communication system that has as small a bit error rate as one desires. The history of digital communications can be described in part as a series of attempts to move ever closer to this limiting curve with systems that have very low bit error rate. Such systems employ both modem techniques and error-control techniques.

If bandwidth W is a plentiful resource but energy is scarce, then one should let W go to infinity, or r to zero. Then we have

$$E_b/N_o \geq \log_2 2 = 0.69$$

This is a fundamental limit. Ratio E_b/N_o is never less than -1.6 dB, and by a sufficiently expensive system one can communicate with any E_b/N_o larger than -1.6 dB.

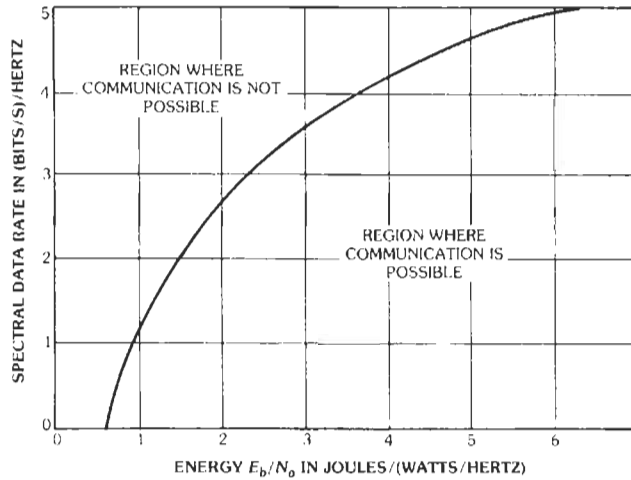


Fig. 20. Capacity of baseband Gaussian noise channel.

Signaling Without Bandwidth Constraints

If the bandwidth is much larger than the data rate, then in principle one can signal arbitrarily closely to channel capacity by using an M -ary signaling alphabet. An M -ary signaling alphabet is a collection of M sufficiently distinct waveforms. Usually M is chosen equal to 2^k for some k . The modulator maps each k -bit word from the channel encoder into one of the 2^k waveforms in the signaling alphabet. The demodulator compares the received signal in each signaling time interval with each of the 2^k possible transmitted waveforms and chooses the most likely.

There are a great many sets of M -ary signal alphabets that can be used. Those with the best performance known are a type known as *simplex waveforms*, but orthogonal waveforms are almost as good and are usually used in practice. The advantage of M -ary signaling can be understood as the resultant of two opposing forces. Since the modulator puts k bits into one waveform, it has k times as much energy as it would for a binary waveform. Since the demodulator must make a 2^k -way decision rather than a binary decision, it needs more energy in the waveform to preserve a small error probability. However, the additional energy needed grows with k more slowly than the additional energy available. The net effect is a system that can operate at a very low E_b/N_0 ; however, the waveforms occupy a bandwidth much greater than the bit rate. Fig. 21 shows the performance of a 2^k -ary family of orthogonal waveforms. The figure shows that E_b/N_0 of -1.6 dB is attainable for large k . Practical systems are usually limited to a k of about six. In the figure, E_b is the energy per data bit delivered; because some of these data bits will be incorrect, the energy per information bit is larger.

Signaling With a Bandwidth Constraint

When bandwidth is expensive, one no longer attempts to transmit with very small E_b/N_0 . Now it becomes important to transmit with large spectral bit

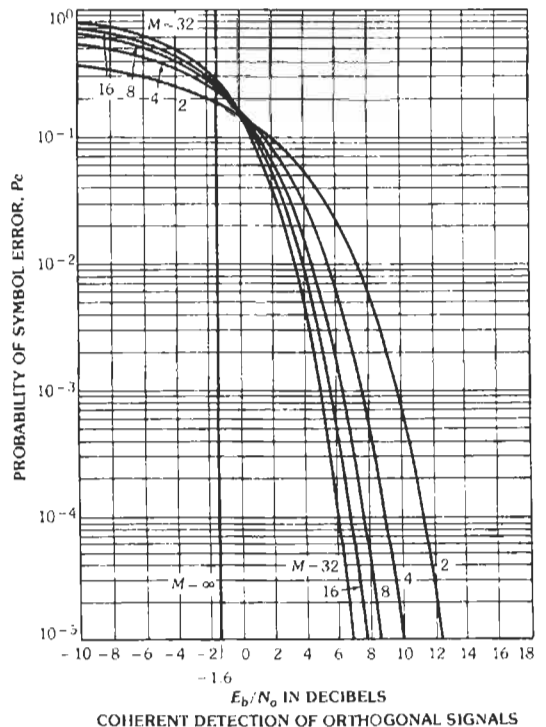


Fig. 21. Performance of M -ary signaling.

rate, r (bits/second)/hertz. Good performance can be achieved by careful combination of modulation and error control.

The Ungerboeck convolutional codes are a class of convolutional codes that provide good combinations of spectral bit rate and E_b/N_o . The rate 2/3 codes are designed to work with a set of eight fixed complex numbers, called a signaling constellation, which represent phase/amplitude modulation patterns. Two such signaling constellations are shown in Fig. 22. An encoder for the constraint-length-four Ungerboeck code for the PSK signaling constellation is shown in Fig. 23. This code can be used as a plug-in replacement for the popular uncoded four-phase PSK modulator. The information rate is still two bits per symbol. There is no change in the channel symbol rate, so the coded system has the same bandwidth as the uncoded system and transmits the same number of information bits per symbol. Hence, the user of the system is unaware of the presence of the code. However, the system now can run at a lower E_b/N_o or signal-to-noise ratio; the constraint-length-four code has a coding gain of 5.7 decibels.

DECISION THEORY AND ESTIMATION THEORY

Decision theory is concerned with the problem of deciding between a set of hypotheses when given a collection of imperfect measurements. Estimation theory is concerned with the problem of selecting the best value of a parameter from a continuum of possible values when given a collection of imperfect measurements. Estimation theory is also concerned with the selection of a best waveform from a collection of waveforms when given a collection of imperfect measurements.

General discussion of decision theory and estimation theory can be found in Davenport and Root* and Van Trees.†

Hypothesis Testing

The simplest problem of decision theory, called hypothesis testing or detection, is to decide between two mutually exclusive hypotheses. In the case of a radar,

* Reference 14.
 † Reference 15.

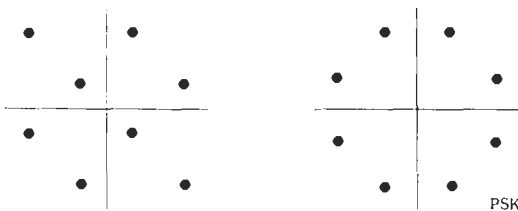


Fig. 22. Two 8-ary signaling constellations.

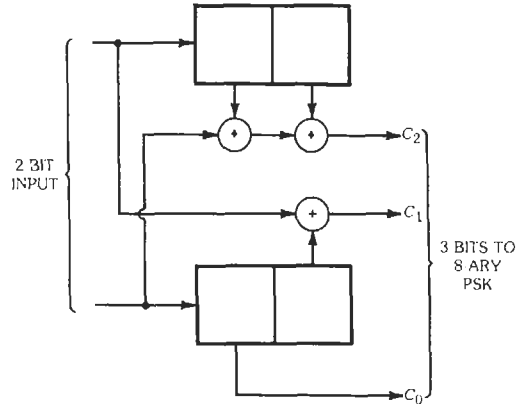


Fig. 23. Diagram of encoder for a constraint-length-four Ungerboeck code.

the hypotheses are *target present* or *target absent*. In the case of signal acquisition, the hypotheses are *signal present* or *signal absent*. In the general case, the hypotheses are called the *null hypothesis*, H_0 , and the *alternate hypothesis*, H_1 . The problem is to decide which hypothesis is correct by collecting data and processing them. The data are randomly distributed with a probability distribution that depends on the true hypothesis. The set of data may actually be quite extensive, but it is enough to think of it as a simple measurement whose outcome can only be an element of a finite set of K elements called the measurement space, and indexed by k . The theory applies to arbitrary data sets simply by replacing the scalar measurement by an appropriate vector measurement.

Associated with each hypothesis is a probability distribution on the measurement space. If H_0 is true, then q_{0k} gives the probability that k will be the measurement outcome; if H_1 is true, then q_{1k} gives this probability. A simple measurement consists of an observation of a realization of the random variable. A measurement is observed, and the problem is to decide whether hypothesis H_0 or hypothesis H_1 is true. A hypothesis testing rule is a partition of the measurement space into two disjoint sets U_0 and $U_1 = U_0^c$. If the measurement is an element of U_0 , we decide that H_0 is true; if it is an element of U_1 , we decide that H_1 is true. Each hypothesis testing rule can be described as such a partition.

Accepting hypothesis H_0 when H_1 actually is true is called a type-one error, and the probability of this event is denoted by α . Accepting hypothesis H_1 when H_0 actually is true is called a type-two error, and the probability of this event is denoted by β . Obviously,

$$\alpha = \sum_{k \in U_0} q_{1k} \quad \beta = \sum_{k \in U_1} q_{0k}$$

A method for finding the optimum decision regions is given by the Neyman-Pearson theorem. This theorem

expresses the decision regions in terms of a parameter, T , called the threshold, and a function called the log-likelihood ratio given by

$$l(k) = \log (q_{0k}/q_{1k})$$

The Neyman-Pearson theorem says that the sets parameterized by T ,

$$U_0(T) = \{k|l(k) \geq T\}$$

$$U_1(T) = \{k|l(k) < T\}$$

are an optimum family of decision rules in the sense that no decision rule can have both type-one error and type-two error better than any of these rules. The type-one and type-two error probabilities can be traded by varying the threshold, T .

The log-likelihood ratio for a block $v = (k_1, \dots, k_n)$ of independent identically distributed measurements of length n can be written as a sum of per-letter log-likelihood ratios

$$l(v) = \log [q_0(v)/q_1(v)] = \sum_{i=1}^n \log [q_{0k_i}/q_{1k_i}]$$

because the probability distributions on v are products of the single letter probability distributions

$$q_0(v) = \prod_{i=1}^n q_{0k_i} \quad q_1(v) = \prod_{i=1}^n q_{1k_i}$$

The Neyman-Pearson decision regions for the block measurement are expressed in terms of $l(v)$ by the inequalities $l(v) < T$ and $l(v) \geq T$.

The same description applies even when the measurements are not independent or identically distributed. However, the log-likelihood ratio, $l(v)$, then cannot be written as the sum of per-letter log-likelihood ratios.

The expected value of the log-likelihood ratio with respect to q_0 is the function known as the *discrimination*.

$$L(q_0, q_1) = \sum_{k=0}^{K-1} q_{0k} \log (q_{0k}/q_{1k})$$

The discrimination is closely related to the entropy and the average mutual information. It is useful for forming bounds on the probability of error.

Estimation Theory

Estimation theory is concerned with the problem of finding a best value for an unknown parameter when only imperfect data are available. The estimation of the unknown parameter depends on measurements that are random variables. The quality of the estimate is limited by the quality of the measurements. These limitations can be expressed by means of information-theoretic

bounds; one such bound is known as the Cramer-Rao inequality.

The simplest estimation problem involves an unknown real parameter, θ , to be determined, and a random variable, X . Understanding this problem leads to understanding of the more general problems: the estimation of several parameters based on the observation of a finite number of random variables; the estimation of several parameters based on observing a sample waveform of a stochastic process; and the estimation of a function on an interval based on observing a sample waveform of a stochastic process.

The random variable has a continuous probability distribution, $q(x|\theta)$, conditional on the parameter θ . The unknown θ must be estimated based upon an observation of X . The estimate of θ given the measurement X is a function $\hat{\theta}(x)$. The estimate, $\hat{\theta}$, is a random variable, since it is a function of the random variable X . Estimation theory studies various criteria for making good estimates, that is, for selecting the function $\hat{\theta}(x)$.

The quality of an estimator is judged by its bias

$$\bar{\theta} = E[\hat{\theta}(x)]$$

and by its variance

$$\sigma_{\hat{\theta}}^2 = E\{[\hat{\theta}(x) - \bar{\theta}]^2\}$$

Intuitively, one hopes to choose the estimator, $\hat{\theta}(x)$, so that its bias is zero and the variance is as small as possible. Such an estimator is called a *minimum-variance unbiased estimator*. Sometimes, however, the minimum-variance unbiased estimator does not exist, or another estimator has some advantage, and so another estimator is used.

The Cramer-Rao bound is a lower bound on the variance of any estimator. For an unbiased estimator of a single parameter θ , the most common form of the bound is

$$\sigma_{\hat{\theta}}^2 \geq \{E\{(\partial/\partial\theta) \log q(x|\theta)\}^2\}^{-1}$$

The same bound holds even if X is a vector of measurements. The Cramer-Rao bound can also be expressed in terms of the discrimination

$$\sigma_{\hat{\theta}}^2 \geq \left\{ \lim_{\theta' \rightarrow \theta} \frac{L(q_{\theta}; q_{\theta'})}{(\theta - \theta')^2} \right\}^{-1}$$

where,

$$L(q_{\theta}; q_{\theta'}) = \int q(x|\theta) \log [q(x|\theta)/q(x|\theta')] dx$$

The Matched-Filter Estimator

Suppose that a known waveform, $s(t)$, has spectrum $S(f)$ and that additive Gaussian noise has spectral

density $N(f)$. The received noisy signal with an unknown time of arrival τ is

$$v(t) = s(t - \tau) + n(t)$$

The Cramer-Rao bound says that any estimator of the time of arrival satisfies

$$\sigma_\tau^2 \geq \left[\int_{-\infty}^{\infty} \frac{|2\pi f S(f)|^2}{N(f)} df \right]^{-1}$$

The optimal estimator, if the noise is white and the signal-to-noise ratio is sufficiently high, is the matched-filter estimator. It achieves the Cramer-Rao bound. The matched-filter estimator passes $v(t)$ through a filter with impulse response $s^*(-t)$ (with some fixed delay to make it realizable) and estimates τ from the peak of the filter output. If the signal-to-noise ratio is not sufficiently high, then the optimal estimator is not known.

Similar remarks apply to the problem of estimating the frequency offset, ϕ , of an otherwise known band-pass waveform

$$v(t) = s(t) \cos [2\pi(f_0 + \phi)t + \theta] + n_R(t) \cos 2\pi f_0 t - n_I(t) \sin 2\pi f_0 t$$

where $s(t)$ has spectrum $S(f)$ and is a known and finite energy waveform whose bandwidth is much less than the known carrier frequency, f_0 . A small offset in frequency ϕ is to be estimated, and $n_R(t)$, $n_I(t)$ are independent covariance stationary Gaussian processes with identical power density spectrum $N(f)$, which is known. The Cramer-Rao bound says that any estimator of frequency ϕ has an error variance satisfying

$$\sigma_\phi^2 \geq \left[\int_{-\infty}^{\infty} \frac{|S'(f)|^2}{N(f)} df \right]^{-1}$$

where $S'(f)$ is the derivative of $S(f)$. The matched-filter estimator of frequency consists of a bank of matched filters. Each filter has impulse response of the form $s^*(-t) \cos 2\pi(f_0 + \Delta)t$ for Δ in a range of interest. The Δ of the filter with the largest output provides the estimate of ϕ . If the noise is white, and the signal-to-noise ratio is sufficiently high, the matched-filter estimator achieves the Cramer-Rao bound, and so is optimal. This matched-filter estimator of frequency offset often can be approximated by a phase-locked loop. If the signal-to-noise ratio is not sufficiently high, the optimal estimator of frequency is not known.

Maximum Entropy and Minimum Discrimination

The *Jaynes maximum entropy principle* of data reduction says that when reducing a set of data into the form of an underlying model, one should be maximally noncommittal with respect to missing data. If one must

estimate a probability distribution subject to constraints on it, then one should choose the probability distribution of maximum entropy consistent with the constraints. For example, if $q(x)$ is a probability density function of a single variable, x , and the mean and variance of $q(x)$ are known and otherwise $q(x)$ is unknown, then one should estimate that $q(x)$ is a Gaussian probability density function with the given mean and variance.

The *Kullback principle of minimum discrimination* is a more general principle that applies when one is given a distribution, $p(x)$, that is a prior estimate of $q(x)$ and a set of constraints on $q(x)$. The principle states that, of the distributions that satisfy the constraints, one should choose as the new estimate that $q(x)$ that minimizes the discrimination

$$L(q,p) = \int q(x) \log [q(x)/p(x)] dx$$

If the probability distribution is discrete and $p(x)$ is a uniform distribution, then this reduces to the maximum entropy principle.

Spectral Estimation

Spectral estimation is the problem of estimating the power spectrum of a stochastic process given partial data, usually only a finite number of samples of the autocorrelation function of limited accuracy. Most spectral estimation methods in use ignore the noise on the autocorrelation samples and treat only the problem of dealing with the missing autocorrelation samples.

The simplest and most popular spectral estimation procedure is to choose as the estimated spectrum the Fourier transform of the known values of the autocorrelation function, possibly tapering the known values of the autocorrelation function or padding them with some zero components. The Fourier transform techniques implicitly augment the measured components of the autocorrelation function with artificial components: either a periodic continuation of the known values or some zero components.

The *Jaynes maximum entropy principle* applied to spectral estimation yields an alternative procedure known as *maximum-entropy spectral estimation*. If the correlation samples are given at a set of equispaced lags, then the maximum entropy estimate of the stochastic process turns out to be

$$v_i = - \sum_{j=1}^L h_j v_{i-j} + n_i$$

which is the output of an *autoregressive filter* excited by discrete-time white noise. The filter coefficients are chosen so that the correlation coefficients of this process agree with the known correlation coefficients. Hence the filter coefficients are the solution of the matrix equation

$$\sum_{j=1}^L h_j R_{i-j} = -R_i$$

where the correlation coefficients, R_i , for $i = -L, \dots, L$, are known and R_{-i} equals R_i . An efficient algorithm for solving this set of equations is the *Levinson algorithm*.

SOURCE COMPRESSION CODES

The average amount of information required to describe a source output symbol is equal to the entropy of the source. Sometimes it is not convenient or practical to retain all this information. It then is no longer possible to maintain an exact reproduction of the source. An analog source has infinite entropy, so distortion must always be present when the source output is passed through a channel of finite capacity.

Data compression is the practice of intentionally reducing the information content of a data record. This should be done in such a way that the least-distorted reproduction is obtained. Information theory finds the performance of the optimum compression of a random source of data; a naive encoder will have greater distortion.

The output of a discrete information source is a random sequence of symbols from a finite alphabet containing J symbols given by $\{a_0, a_1, \dots, a_{J-1}\}$. A memoryless source produces the j th letter with probability p_j where p_j is strictly positive. The source output is to be reproduced in terms of a second alphabet called the reproducing alphabet, often identical to the source alphabet, but not always so. For example, the reproducing alphabet might consist of the union of the source alphabet and a single new element denoting "data erased."

A distortion matrix is a J by K matrix with nonnegative elements ρ_{jk} for $j = 0, \dots, J - 1$ and $k = 0, \dots, K - 1$ that specifies the distortion associated with reproducing the j th source letter by the k th reproducing letter. Without loss of generality, it can be assumed that for each source letter a_j there is at least one reproducing letter b_k such that the resulting distortion ρ_{jk} equals zero. Usually the distortion in a block is defined as the arithmetic average of the distortions of each letter of the block. This is called a per-letter fidelity criterion.

An important distortion matrix is the probability-of-error distortion matrix. For this case the alphabets are identical. For example, take $J = K = 4$ and

$$\rho = \begin{bmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{bmatrix}$$

This distortion matrix says that each error is counted as one unit of distortion. A different data-compression

problem with the same source alphabet and reproducing alphabet is obtained if one takes the distortion matrix

$$\rho = \begin{bmatrix} 0 & 1 & 2 & 1 \\ 1 & 0 & 1 & 2 \\ 2 & 1 & 0 & 1 \\ 1 & 2 & 1 & 0 \end{bmatrix}$$

This distortion matrix says that, modulo four, an error of two units is counted as two units of distortion.

A source compression block code of blocklength n and size M is a set consisting of M sequences of reproducing letters, each sequence of length n . The source compression code is used as follows. Each source output block of length n is mapped to that one of the M codewords that results in the least distortion.

The entropy of the output of the data compressor is less than that of the original source and therefore can be encoded into a smaller number of bits.

The Distortion-Rate Function

Data compression is a deterministic process. The same block of source symbols always produces the same block of reproducing letters. Nevertheless, if attention is restricted to a single source output symbol without knowledge of the previous or subsequent symbols, then the reproducing letter is not predetermined. The letter, b_k , into which source letter a_j is encoded becomes a random variable even though the block encoding is deterministic. This random variable can be described by a transition matrix, Q_{kj} . Heuristically, we think of Q_{kj} as describing an artificial channel that approximates the data compression. Each time the source produces letter a_j , it is reproduced by letter b_k with probability Q_{kj} .

To obtain the greatest possible compression, it seems that this conditional probability should, on average, result in the smallest possible mutual information between the source and the reproduction provided that the average distortion is less than the allowable average distortion. This heuristic discussion motivates the following definition.

The distortion-rate function, $D(R)$, is given by

$$D(R) = \min \sum_j \sum_k p_j Q_{kj} \rho_{jk}$$

where the minimum is over all probability transition matrices Q connecting the source alphabet and the reproducing alphabet that satisfy $I(p;Q) \leq R$. The definition is justified by the source compression theorems of information theory. Intuitively, if rate R bits per source letter is specified, then any compression must provide a distortion of at least $D(R)$. Conversely, compression to a level arbitrarily close to $D(R)$ is possible by appropriate selection of the compression scheme.

The distortion-rate function can be evaluated analytically for simple sources and distortion matrices. Fig. 24

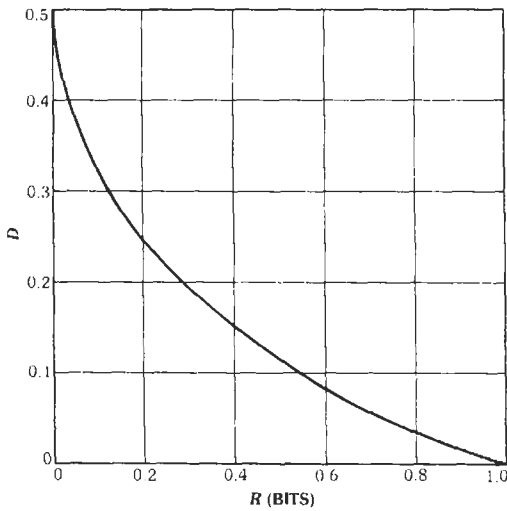


Fig. 24. Distortion-rate function for a binary symmetric source.

shows $D(R)$ for the simplest example, a binary memoryless source with equiprobable zeros and ones, and a probability-of-error distortion matrix. This figure shows that if only one-half bit per bit is used to represent the output of a binary equiprobable source, then at least 11 percent of the reproduced bits must be in error. Fig. 25 shows $D(R)$ for Gaussian sources with mean-square distortion measure. This figure shows that

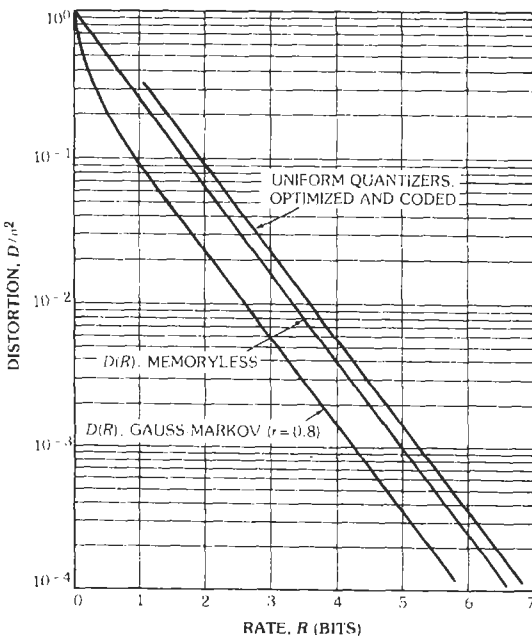


Fig. 25. $D(R)$ for Gaussian sources.

the conventional equispaced quantizer is not optimal, but the performance loss is quite small. Source compression codes are discussed further in Berger.*

MULTITERMINAL INFORMATION NETWORKS

Information networks involving several sources or several channels can have conflicting requirements imposed by the several users. An efficient system tries to satisfy several goals at once, perhaps to interlock several independent messages into one efficient waveform, perhaps to share a multiple-access channel with several users, or perhaps to break a message into two distorted replicas that together contain enough information to reconstruct the original message without distortion. Problems of these kinds are not efficiently implemented in today's communication technology. Multiterminal information theory studies how to solve such problems efficiently. Most such problems are only partially solved. Further discussion of multiterminal information theory can be found in the survey articles by van der Meulen,† Berger,§ and El Gamal and Cover.**

Two-Way Channels

A two-way channel has two terminals, as shown in Fig. 26. Each terminal attempts to get a message across to the other terminal through the two-way channel, but the transmission in one direction interferes with the transmission in the other direction. The problem is to design the terminals to achieve high data rates in both directions simultaneously.

Many practical channels are intrinsically two-way channels, but the designer elects to break them into two one-way channels using a technique such as time division or frequency division. This is a simple and workable solution, but one wishes to know if there is a performance penalty. What is the two-way capacity region, and what sort of communication scheme is optimal? This problem, one of the most difficult ones in multiterminal information theory, has been open for twenty years. Some highly simplified models have been solved or partially solved. These solutions show that time division and frequency division are not optimal.

Broadcast Channels

A broadcast channel consists of a single transmitter and multiple receivers as shown in Fig. 27. If the same message is to be sent to all receivers, then the task of designing the communication waveform is no different than when there is a single receiver. However, if a

* Reference 16.
 † Reference 17.
 § Reference 18.
 ** Reference 19.

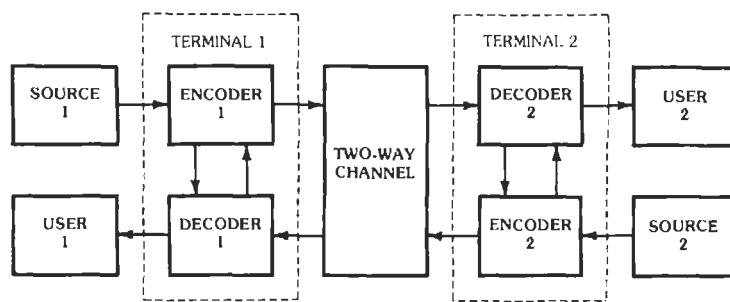


Fig. 26. The two-way channel.

different message is to be sent to each receiver, the problem is more complex. Conventionally, the channel is divided according to the needs of the receivers, using time-division or frequency-division signaling techniques. Optimum techniques are more tightly interlocked and do better. Time-division or frequency-division multiaccess signaling cannot achieve the capacity region of a broadcast channel.

Degraded Diversity Systems

A *diversity system* sends the same information to a user twice through two channels so that if one channel is broken, the message will still arrive. A *degraded diversity system* is more subtle. It sends half the message through each of two channels but in such a way that either half suffices to reconstruct a degraded copy of the message. A high-fidelity reproduction is obtained if both channels are intact; a low-fidelity reproduction is obtained if only one channel is intact. The simplest model of a degraded diversity system is shown in Fig. 28. The theory shows that the data can be encoded so that either side channel can reconstruct the data from the half-rate code with a probability of error of 20.7 percent, and no better encoding exists.

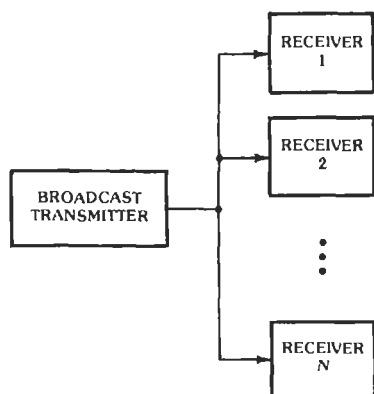


Fig. 27. Broadcast channel.

Remote Compaction of Dependent Data

A fundamental multiterminal source coding problem is shown in Fig. 29. Data sources X and Y are dependent. The data are to be sent to a common site at rates R_x and R_y , but the encoders do not communicate with each other. The admissible rate region is shown in Fig. 30. Encoder/decoders can be built with rate pairs corresponding to any point in the upper-right region. It is quite remarkable that this region would not be any larger even if the two encoders could each see the output of both sources.

THE COMMUNICATION GAMES

The communication problem is given a new dimension of complexity and becomes a two-person game when an adversary is introduced. The theory that is then developed has its roots in both information theory and game theory. The transmitter and receiver together constitute one player, or team, and the adversary is the second player. The adversary may have a variety of purposes. It may be his goal to interrupt communication, in which case he is called a jammer or a spoofer; it may be his goal to determine the specific message transmitted, in which case he is called a wiretapper or eavesdropper; it may be his goal simply to detect the existence of communication, in which case he is called a monitor; or it may be his goal to determine the location of the transmitter, in which case he is called a locator.

The Jammer Saddle Point

The jamming game is a game played by the jammer and the transmitter. The jammer selects the noise or interference properties of the channel. That is, the jammer selects a probability transition matrix, Q , connecting the input and output alphabets. The transmitter designs a communication waveform without knowledge of which channel noise was selected by the jammer. Variations are possible wherein the jammer

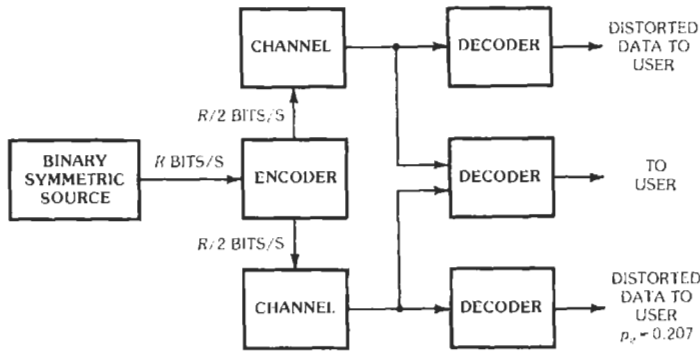


Fig. 28. A degraded diversity system.

has partial knowledge about the transmitted signal, but these variations will not be considered here. The jammer and the transmitter each must operate within certain constraints. The most important case is one in which the jammer and transmitter are each limited in average power, and this is the case we discuss.

The set of channels that the jammer can use is restricted in that he has a limited average transmitter power. Let Q_N denote the set of channel noise distributions that satisfy power constraints on the jammer,

$$Q_N = \{Q: \int z^2 Q(z) dz \leq N\}$$

and let P_S denote the set of probability distributions on the channel input that satisfy average power constraints on the transmitter.

$$P_S = \{p: \int x^2 p(x) dx \leq S\}$$

The jammed capacity is given by:

$$C(S, N) = \max_{p \in P_S} \min_{Q \in Q_N} I(p; Q) \\ = \min_{Q \in Q_N} \max_{p \in P_S} I(p; Q)$$

The first expression gives the capacity of the channel if the jammer must first choose the channel noise statistics before the transmitted waveform is designed. The second expression gives the capacity if the transmitted

waveform is designed before the jammer chooses the channel noise. The well-known minimax theorem of game theory says that the two expressions are equal.

For the discrete-time, continuous-amplitude channel with transmitter and jammer both limited in average power, the minimax jammer strategy is to transmit Gaussian noise, and the minimax transmitter strategy is to encode assuming Gaussian noise. From the point of view of the jammer, the jammed capacity is the data rate to which the jammer can hold the transmitter. From the point of view of the transmitter, the jammed capacity is the data rate that the transmitter can achieve in the presence of the jammer.

The most important application is the band-limited waveform channel in which the transmitter and the jammer are each limited in average power. The jammer saddle point is

$$C(S, N) = W \log(1 + S/N)$$

This expression is linear in the bandwidth, W . If the signal power, S , and the jammer power, N , are fixed, the capacity can be increased by increasing W . This is the basis for the practice of spectrum spreading discussed below. If the spectral power density, N_o (watts/hertz), is fixed rather than N , then

$$C(S, N_o) = W \log(1 + S/WN_o)$$

The capacity is bounded as W increases. Spectrum

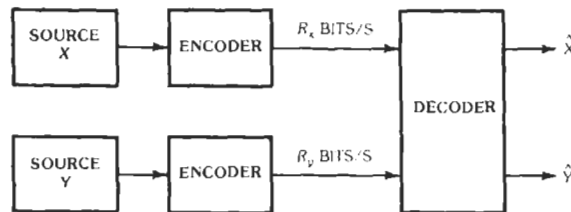


Fig. 29. Remote compaction of dependent data.

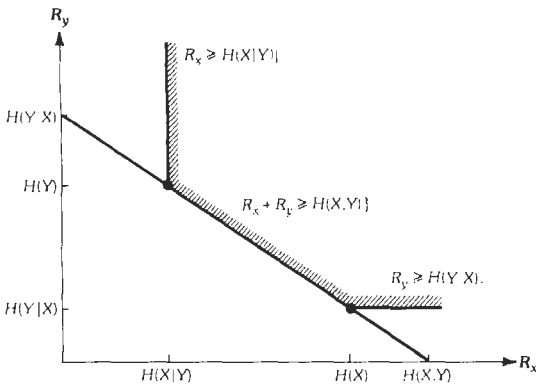


Fig. 30. The Slepian-Wolf admissible rate region.

spreading is useful only when the jammer's total power is limited. It is worthless if his spectral power density is limited instead. Spectrum spreading does not improve performance against broad-band white noise.

Spectrum Spreading

A *spread-spectrum* system is one that takes a signal of small bandwidth, say γW , where γ is much smaller than one, and converts it into a signal of high bandwidth, W . This is done as shown in Fig. 31 by modulating the information waveform with a wideband waveform that is not known to the jammer. The wideband waveform may be noise-like (hence called pseudo-noise, or PN), or a time-varying carrier frequency (called frequency hopping). The wideband signal is transmitted through the channel, received, and then compressed back into the original low-bandwidth signal by demodulating the information waveform from the wideband waveform. This scheme requires that the transmitter and receiver share knowledge of the bandspreading waveform through a secure channel. If the jammer knows nothing

about the wideband waveform, his "best" strategy is to use wideband noise of bandwidth W . After the signal spectrum is compressed, the jammer signal will still have bandwidth W , but the signal will have bandwidth γW . A bandpass filter will reduce the jammer noise power by the factor γ . The reciprocal of the factor γ is called the processing gain of the spread-spectrum system.

A spread-spectrum system does not make optimum use of the wide bandwidth. This can be seen by looking at the channel capacity. The channel capacity of the available wideband channel of bandwidth W and jammer power N is

$$C = W \log (1 + S/N)$$

However, the channel capacity of the spread-spectrum channel of compressed bandwidth γW and processed jammer power γN is

$$C_{SS}(\gamma) = \gamma W \log (1 + S/\gamma N)$$

whereas if only the narrow bandwidth, γW , is used, there is no reduction of jammer power, and the channel capacity is

$$C_{NB}(\gamma) = \gamma W \log (1 + S/N)$$

If the channel must transmit R bits per second, the wideband channel requires a signal-to-noise ratio of at least

$$S/N = 2^{(R/W)} - 1$$

while the spread-spectrum channel requires a signal-to-noise ratio of

$$S/N = \gamma (2^{(R/\gamma W)} - 1)$$

which can be considerably larger than what is needed in the wide bandwidth but is smaller than

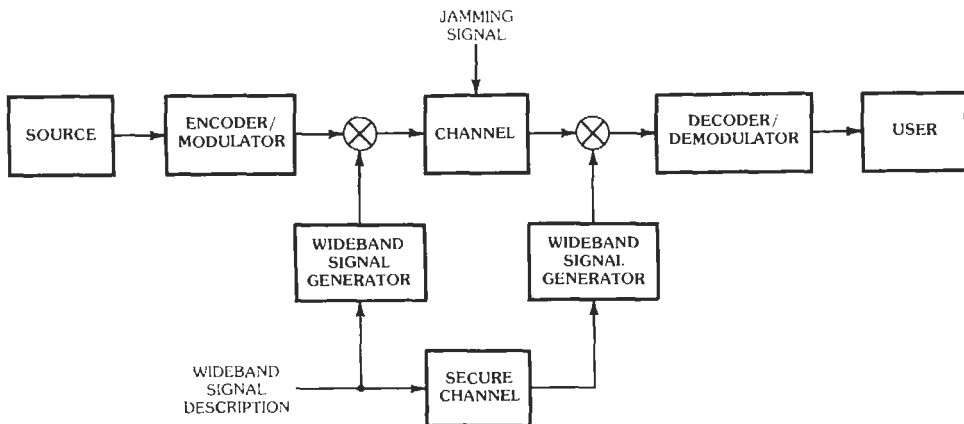


Fig. 31. Spread-spectrum signaling.

$$S/N = 2^{(R/\gamma W)} - 1$$

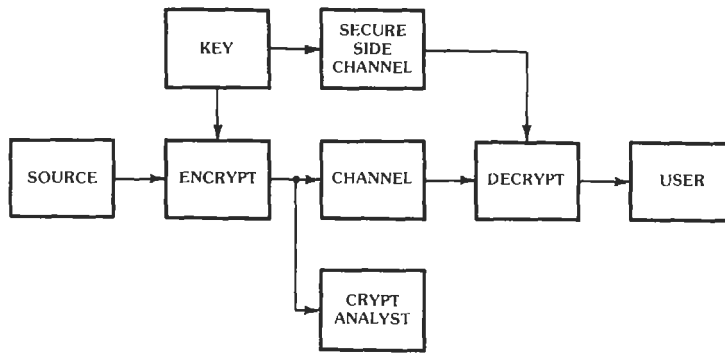
which is required if there is no spectrum spreading at all. Spread-spectrum by itself improves performance, but much better improvement is theoretically possible using the wide bandwidth.

The use of M -ary orthogonal signaling, spread-spectrum, and error-control codes in combination will make better use of the wide bandwidth. In addition, error-control codes are necessary in a spread-spectrum system to protect against partial-time or partial-band jamming tactics. A well-designed spread-spectrum system will use an error-control code and a modulation

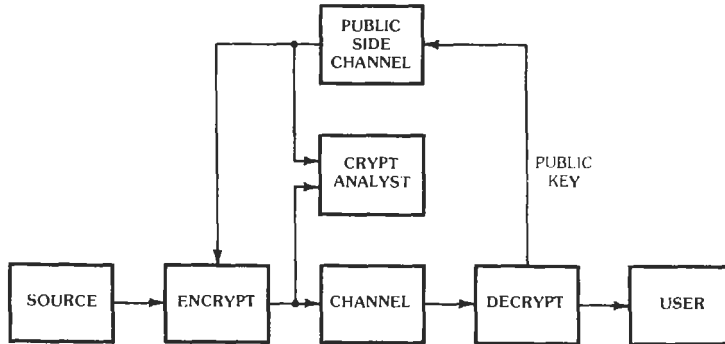
waveform that operate at a small E_b/N_0 and then spread the spectrum of this waveform to fill out the available frequency band.

Cryptography

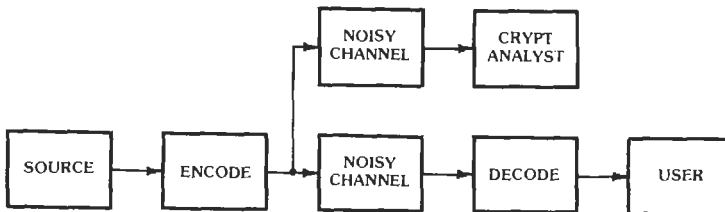
Three major kinds of cryptosystem are shown in Fig. 32. The first two kinds require some type of key; the third kind relies on noise in the channel to hide the message. While at first it may seem that the wire-tap channel (Fig. 32C) is not a true cryptosystem, in fact it would be misleading to ignore it. All cryptosystems have internal compartments where the message or the



(A) Conventional cryptosystem.



(B) Public-key cryptosystem.



(C) Wire-tap channel.

Fig. 32. Types of cryptosystem.

key exists in the clear. The designer relies on the fact that any portion that leaks out will be buried in the ambient noise. Hence the wire-tap channel is present in every cryptosystem. A secure side channel for transmitting a key is of this type.

A conventional cryptosystem (Fig. 32A) is secure from a direct attack by the cryptanalyst if the key is nearly as long as the message. A so-called "one-time pad" has one random bit in the key for each bit of the message, and it is not reused. The one-time pad can be penetrated only by an attack on the secure side channel. Since it is difficult to transmit and store long keys, most cryptosystems have short keys. Then there is redundancy in the encrypted message that the cryptanalyst can exploit.

A public-key cryptosystem (Fig. 32B) differs from a conventional cryptosystem in that it relies on one-to-one functions, $f(x)$, whose inverses, $f^{-1}(x)$, are very difficult—even computationally intractable—to deduce. These are sometimes called one-way functions. Hence, the user randomly makes up an $f^{-1}(x)$, finds $f(x)$, and publicly announces that he can decrypt messages that are encrypted with $f(x)$. The cryptanalyst always knows the encoding key, but is unable to convert this into a decoding key.

REFERENCES

1. Gallager, R. G. *Information Theory and Reliable Communication*. New York: John Wiley & Sons, Inc., 1968.
2. McEliece, R. J. *The Theory of Information and Coding*. Reading, Mass.: Addison-Wesley Publishing Co., Inc., 1977.
3. Ash, R. B. *Information Theory*. Interscience, 1965.
4. Viterbi, A. J., and Omura, J. K. *Principles of Digital Communication and Coding*. New York: McGraw-Hill Book Co., 1979.
5. Franzaszek, P. A. "Sequence-State Coding for Digital Transmission." *Bell System Tech. J.*, 1968, pp. 113-157.
6. Adler, R. L., Coppersmith, D., and Hassner, M. "Algorithms for Sliding Block Codes." *IEEE Transactions on Information Theory*, Vol. IT-29, 1983; pp. 5-22.
7. Rissanen, J., and Langdon, G. G. Jr. "Universal Modeling and Coding." *IEEE Transactions on Information Theory*, Vol. IT-27, 1981, pp. 12-23.
8. Berlekamp, E. R. *Algebraic Coding Theory*. New York: McGraw-Hill Book Co., 1968.
9. Blahut, R. E. *Theory and Practice of Error-Control Codes*. Reading, Mass.: Addison-Wesley Publishing Co., Inc., 1983.
10. Clark, G. C., and Cain, J. Bibb. *Error-Correction Coding for Digital Communications*. New York: Plenum Press, 1981.
11. Lin, S., and Costello, D. J. Jr. *Introduction to Error Correction Codes.*, Englewood Cliffs, N.J.: Prentice-Hall, Inc., 1983.
12. Peterson, W. W., and Weldon, E. J. Jr. *Error Correcting Codes*, 2nd ed. Cambridge, Mass.: The M. I. T. Press, 1971.
13. Wozencraft, J. M., and Jacobs, I. M. *Principles of Communication Engineering*. New York: John Wiley & Sons, Inc., 1965.
14. Davenport, W. B. Jr., and Root, W. L. *Random Signals and Noise*. New York: McGraw-Hill Book Co., 1958.
15. Van Trees, H. L. *Detection, Estimation, and Modulation Theory*, Part 1. New York: John Wiley & Sons, Inc., 1968.
16. Berger, T. *Rate Distortion Theory: A Mathematical Basis for Data Compression*. Englewood Cliffs, N.J.: Prentice-Hall, Inc., 1971.
17. Van der Meulen, E. C. "A Survey of Multi-Way Channels in Information Theory: 1961-1976." *IEEE Transactions on Information Theory*, Vol. IT-23, 1977, pp. 1-37.
18. Berger, T. "Multiterminal Source Coding," in *The Information Theory Approach to Communications*, G. Longo, ed. New York: Springer-Verlag, 1977.
19. El Gamal, A., and Cover, T. M. "Multiple User Information Theory." *Proceedings IEEE*, Vol. 68, 1980, pp. 1066-1083.
20. Special Issue on Spread Spectrum, *IEEE Transactions on Communication Theory*, Vol. COM-30, May, 1982.

26

Computer Communications Networks

Fouad A. Tobagi and Mario Gerla

- The Structure of Computer Networks 26-2
 - Switching Techniques
 - Types of Computer Communications Networks
- Network Functions 26-7
- The Physical Layer 26-9
- The Data Link Control Layer 26-13
 - Asynchronous DLC Protocols
 - Synchronous DLC Protocols
- Multiaccess Link Control 26-18
 - Fixed Assignment Techniques
 - Random Access Techniques
 - Centrally Controlled Demand Assignment
 - Demand Assignment With Distributed Control
- The Network Layer 26-23
 - The Network Services
 - Routing
 - Congestion Control
 - CCITT Recommendation X.25
- Higher-Level Protocols 26-29
 - The Transport and Session Layers
 - The Presentation Layer
- High Speed Local Area Networks 26-30
 - A Taxonomy of Local Area Networks
 - Low- and Medium-Speed LANs
 - High-Speed LANs
 - Supercomputer Networks
 - Ultrahigh-Speed Networks
- High Speed Wide Area Networks 26-36
 - Bandwidth and Traffic Allocation in ATM Networks
 - Congestion Control in ATM

About two decades ago, a great challenge to many computer system designers was to enhance the processing power of computers and make them available to a large number of users on a time-sharing basis. The design of fast arithmetic and control units as well as the design of complex operating systems were identified among the key tasks needed to accomplish this objective. As time-sharing progressed, it was soon realized that such large resources would not be effectively utilized unless the problem of connecting remote user terminals to the central computing facility was adequately solved, thus allowing a large population of users to share the facility. With this problem, the field of computer communications came into existence. The early communications systems were merely *terminal access networks*. The next stage in the evolution consisted of the creation of the so-called distributed *resource-sharing networks*. The goal here is to interconnect computers and their users at various geographically distributed sites in order to allow the sharing of hardware and software resources developed at all sites by all users connected to the network. Today the term *computer communications* has a much broader meaning. Not only does it refer to communications among computers and their users, but it also refers to all kinds of communications applications (among humans and among machines) that make use of the computer as a tool. Examples of such applications are voice communications, electronic mail, facsimile, image transfer, process control, etc. The goals are also somewhat more diverse. Instead of communications per se, the driving force may very well be that of cost or reliability of computing power. High reliability may be achieved by having alternative sources of computing made available via a communications network. Lower cost may be achieved via distributed computing architectures that are simpler to design, cheaper to build, and easier to maintain, and of which communication is an intrinsic part. Thus, while in the seventies one might have characterized the computer communications field as having been in its research and development phase, the eighties and beyond will be marked by the wide use of this relatively new technology for a large repertoire of applications.

In the following, the structure of computer communications networks is examined and their building blocks identified. First, the various types of networks in existence and the switching techniques in use are described. Then the general functions of a network needed in establishing a communications path among remote users are described, and the organization of these functions into a standard layered architecture called the *Open Systems Interconnection reference model* is discussed. Using this reference model as a guide, we then examine in more detail the functions in each of the layers and highlight standards whenever applicable. Finally, we introduce a class of network architectures which has emerged in the past few years, namely, the class of high-speed networks. The OSI model also applies, with some modifications, to this

class, but the requirements, characteristics, and applications of high-speed networks are so different from those of conventional networks that a separate treatment at the end of this chapter is justified. Two sections are dedicated to this subject, namely, high-speed local area networks and high-speed wide area networks.

THE STRUCTURE OF COMPUTER NETWORKS

A computer-communication network typically comprises a collection of computing resources called *hosts*, a collection of users, some of which are associated with the hosts, and a so-called *communication subnet* that connects them (Fig. 1.) The communication subnet consists of two basic components: the communication channels and the switching elements (or switching nodes). Depending on (1) the physical medium used for the communication channels, (2) the subnet topology according to which the communication channels and the switching elements are interconnected to form a network, and (3) the switching technique used in providing a physical path among two or more communicating parties, several computer network types may be identified.

Switching Techniques

There are two basic types of switching techniques: *circuit switching* and *message switching*. In circuit switching, a total path of connected lines is set up from the origin to the destination at the time the call is made, and the path remains allocated to the source-destination pair (whether used or not) until it is released by the communicating parties. The switches, called circuit switches (or office exchange in telephone jargon), have no capability of storing or manipulating users' data on their way to the destination. The circuit is set up by a special signaling message that finds its way through the network, seizing channels in the path as it proceeds. Once the path is established, a return signal informs the source to begin transmission. Direct transmission of data from source to destination can then take place without any intervention on the part of the subnet.

In message switching, the transmission unit is a well defined block of data called a *message*. In addition to the text to be transmitted, a message comprises a *header* and a *checksum*. The header contains information regarding the source and destination addresses as well as other control information; the *checksum* is used for error control purposes. The switching element is a computer referred to as a *message processor*, with processing and storage capabilities. Messages travel independently and asynchronously, finding their own way from source to destination. First the message is transmitted from the host to the message processor to which it is attached. Once the message is entirely received, the message processor examines its header, and accordingly decides on the next outgoing channel on which to transmit it. If this selected channel is busy,

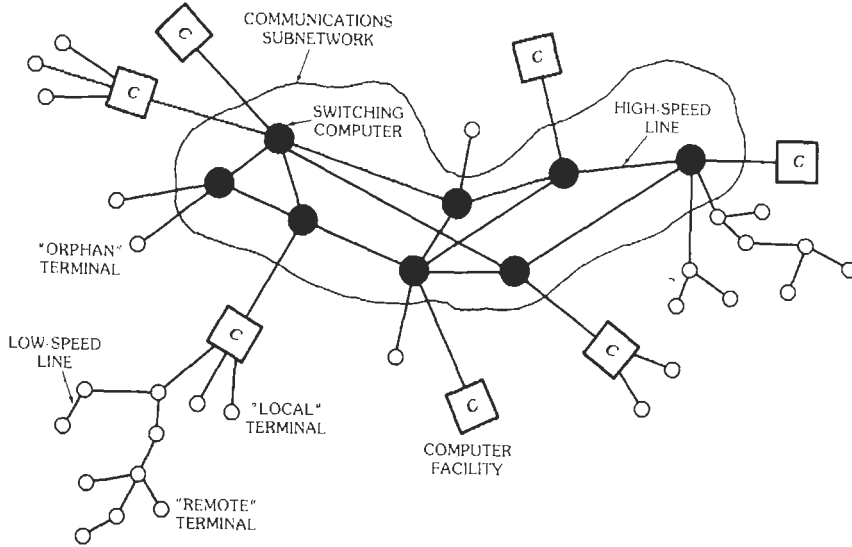


Fig. 1. The structure of a computer-communication network. (From Leonard Kleinrock, *Computer Applications, Vol. II of Queuing Systems*, © 1976 by John Wiley & Sons, Inc., New York.)

the message waits in a queue until the channel becomes free, at which time transmission begins. At the next message processor, the message is again received, stored, examined, and transmitted on some outgoing channel, and the same process continues until the message is delivered to its destination. This transmission technique is also referred to as the *store-and-forward transmission* technique.

A variation of message switching is *packet switching*. Here the message is broken up into several pieces of a given maximum length, called *packets*. As with message switching, each packet contains a header and a checksum. Packets are transmitted independently in a store-and-forward manner.

With circuit switching, there is always an initial connection cost incurred in setting up the circuit. It is cost-effective only in those situations where once the circuit is set up there is a guaranteed steady flow of information transfer to amortize the initial cost. This is certainly the case with voice communication in the traditional way, and indeed circuit switching is the technique used in the telephone system. Communication among computers, however, is characterized as *bursty*. Burstiness is a result of the high degree of randomness encountered in the message-generation process and the message size, and of the low delay constraint required by the user. The users and devices require the communication resources relatively infrequently; but when they do, they require a relatively rapid response. If a fixed dedicated end-to-end circuit were to be set up connecting the end users, then one must assign enough transmission bandwidth to the circuit in order to meet the delay constraint with the consequence that the resulting channel utilization is low. If the circuit of high bandwidth were set up and

released at each message transmission request, then the set-up time would be large compared to the transmission time of the message, resulting again in low channel utilization. Therefore, for bursty users (which can also be characterized by high peak-to-average data rate requirements), store-and-forward transmission techniques offer a more cost-effective solution, since a message occupies a particular communications link only for the duration of its transmission on that link; the rest of the time it is stored at some intermediate message switch and the link is available for other transmissions. Thus the main advantage of store-and-forward transmission over circuit switching is that the communication bandwidth is dynamically allocated, and the allocation is done on the fine basis of a particular link in the network and a particular message (for a particular source-destination pair).

Packet switching achieves the benefits discussed so far and offers added features. It provides the full advantage of the dynamic allocation of the bandwidth, even when messages are long. Indeed, with packet switching, many packets of the same message may be in transmission simultaneously over consecutive links of a path from source to destination, thus achieving a "pipelining" effect and reducing considerably the overall transmission delay of the message as compared to message switching. It tends to require smaller storage allocation at the intermediate switches. It also has better error characteristics and leads to more efficient error recovery procedures, as it deals with smaller entities. Needless to say, packet switching presents design problems of its own, such as the need to reorder packets of a given message that may arrive at the destination node out of sequence.

Fig. 2 illustrates the three switching modes for an

example of a communication subnet involving four nodes and three transmission links. The figure shows the advantage of packet switching over message switching and of message switching over circuit switching. Clearly, this comparison involves a number of trade-offs, and it is not hard to imagine situations where the conclusion may be reversed. This depends on a number of factors, among others the number of hops from source to destination, the length of the message, the amount of overhead incurred in the header and checksum of each message and packet, the circuit set-up delay, etc. For example, fast circuit switching accom-

plished with solid-state implementation of the switches can achieve a set-up time on the order of a fraction of a second (a few milliseconds per switch), which renders this technique a viable solution. In most cases, however, the picture shown in Fig. 2 reflects the real situation.

The advantages of store-and-forward transmission are obtained at the expense of higher processing and storage capabilities at all switches. Overall, however, the present economics are in favor of store-and-forward. Roberts[†] has shown that since the early 60s, the

[†] Reference 1.

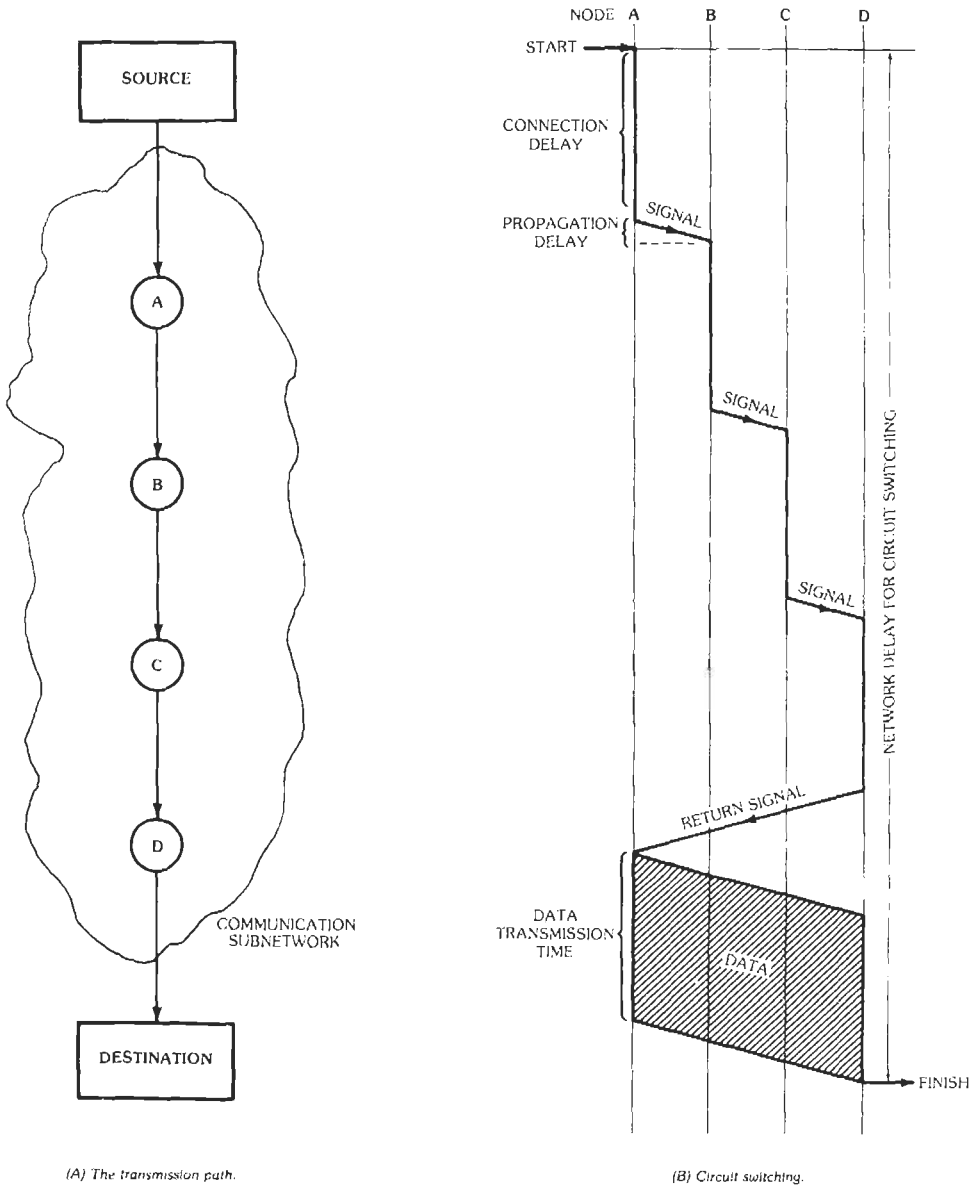


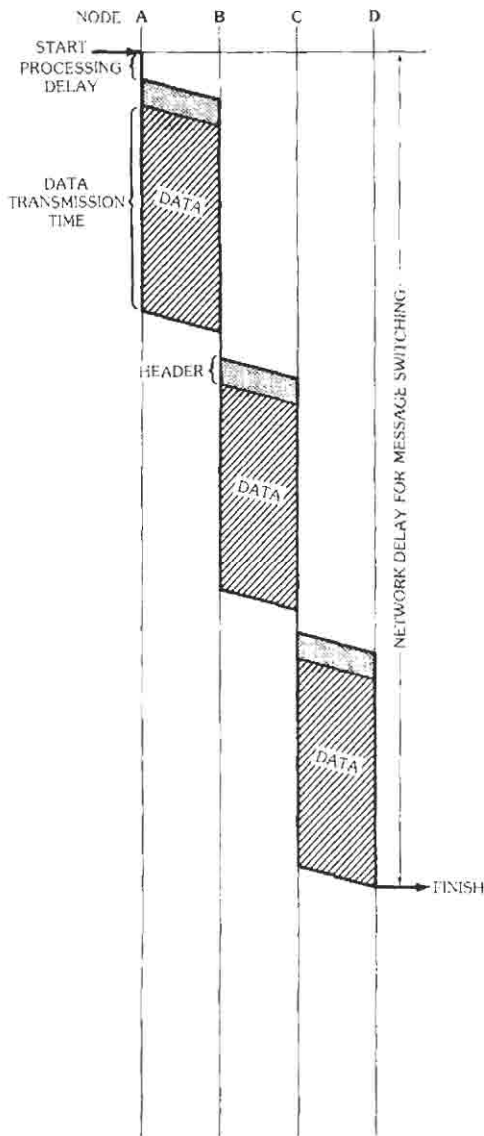
Fig. 2. Comparison of network delay for circuit, message, and packet switching. (From Leonard Kleinrock, Computer Applications,

incremental cost of computing needed to send one megabit of data through a nationwide network has been decreasing at a much faster rate than the incremental cost of the communication land lines in a national net, with a crossover point having occurred about 1969. Fig. 3 summarizes Roberts' findings concerning these cost trends (the reader is referred to the portion of the figure corresponding to the range up to 1974).

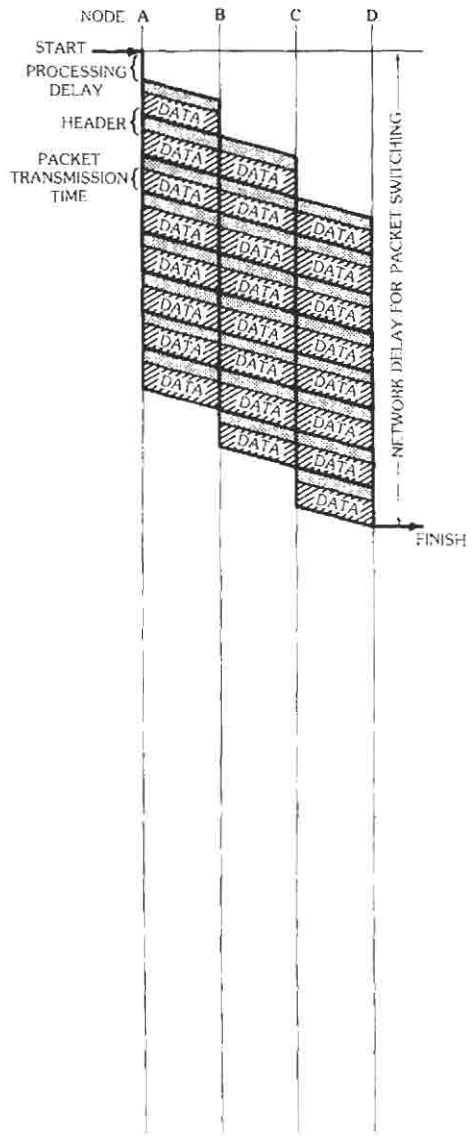
In the development that follows, the three terms message switching, packet switching, and store-and-forward are used interchangeably unless the distinction is made explicitly.

Types of Computer Communications Networks

A distinction is first made between *circuit-switched networks* and *packet communication networks*. Clearly, circuit-switched networks are those that use circuit switching. When used for local communication, these networks have a star topology with a single switch in the center, and they may often be privately owned (e.g., CBX). Long-haul or nationwide networks, on the other hand, have a hierarchical mesh topology: subscribers are connected to local central offices via two-wire



(C) Message switching



(D) Packet switching

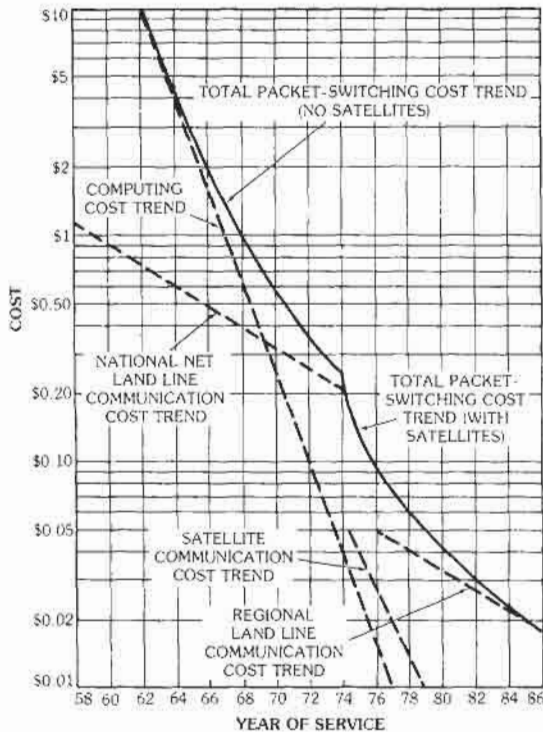


Fig. 3. Incremental cost for sending 1 kilopacket (1 megabit) through a nationwide network. (From Leonard Kleinrock, *Computer Applications, Vol. II of Queueing Systems*, © 1976 by John Wiley & Sons, Inc., New York.)

connections to form a star local network; local offices are connected to toll offices via toll connecting trunks; toll offices are interconnected via very high bandwidth intertoll trunks and other intermediate switching offices to form a highly redundant mesh topology.

Packet communication networks are all based on packet switching and may be of three types. The first is referred to as the *point-to-point store-and-forward* type. In this type, packet switches are interconnected by point-to-point data channels to form a mesh topology. Each channel is used only by the two switches adjacent to it, one in each direction; thus there exists no contention. The channels are usually full duplex, allowing transmission to take place in both directions simultaneously. Examples of such networks are the ARPANET, the Cigale Network, TELENET, TYMNET, DATAPAC, TRANSPAC, EURONET, etc.† This type is usually for large geographical scope, although it is also used for terminal access and local area communication; in the latter case, the network typically has a star topology with a packet switch at its center.

The second type of packet communication network is the *multiaccess-broadcast* type. This type of network

consists of a *single* transmission medium that is shared by all subscribers; they access the medium according to some multiaccess scheme, each access being for the duration of a single packet. Any user's transmission is heard by all other users, hence the broadcast attribute. The single-hop broadcast nature of these networks achieves full connectivity at a very small cost: each subscriber is connected to the common channel through an interface that listens to all transmissions and accepts packets addressed to it. Examples of this type of network can be found in both long-haul and local-area communication systems. One example is that of a *satellite channel* (Fig. 4): a satellite transponder in a geostationary orbit above the earth can receive signals from any earth station in its coverage pattern and can transpond these signals to all such earth stations (unless the satellite uses spot beams). Broadband satellite channels offer a cost-effective alternative to store-and-forward networks for long-haul communication, as can be easily seen from Roberts' projections in Fig. 3. Another example is that of *ground packet radio* systems where the radio medium is used for terminal access and local-area communication among a large population of terminals, possibly mobile (Fig. 5). The multiaccess and broadcast attributes are achieved when all users share a common radio frequency, employ omnidirectional antennas (thus facilitating communication among the mobile users), and are in line-of-sight and within range of each other. Yet another example is the so-called *broadcast bus* system in which all users are connected to a single cable via simple passive taps (Fig. 6). These systems prove to be ideal for short-distance

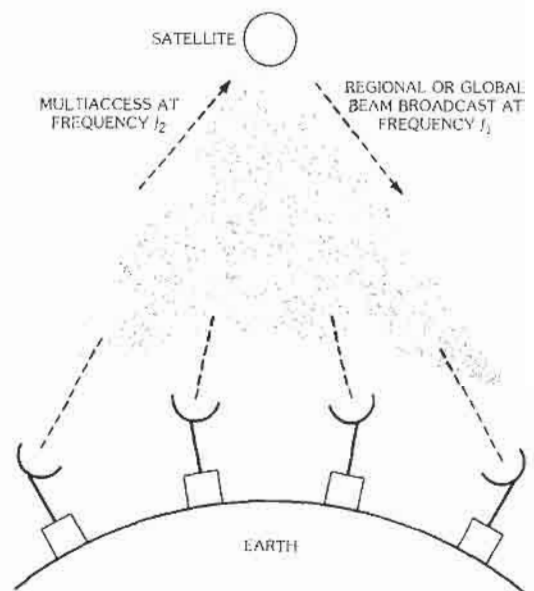


Fig. 4. A satellite channel. (From Leonard Kleinrock, *Computer Applications, Vol. II of Queueing Systems*, © 1976 by John Wiley & Sons, Inc., New York.)

† References 2 and 3.

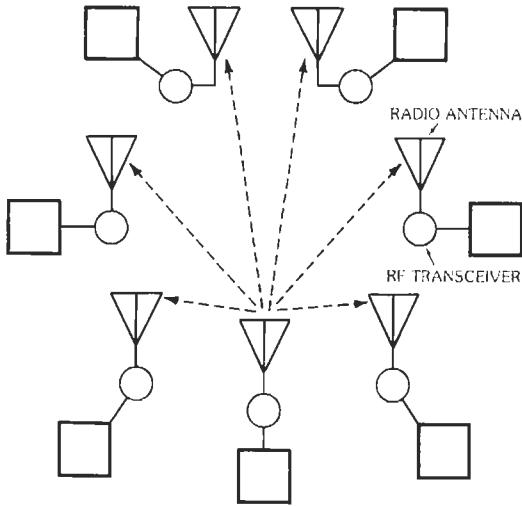


Fig. 5. A single-hop fully connected ground packet radio network.

local-area communication within a building (up to a few kilometers) involving a large and often variable number of inexpensive devices requiring interconnection. The simple topology and simple and inexpensive connection interfaces of broadcast bus systems provide great flexibility in accommodating the high degree of variability in such environments while achieving a desired level of reliability. A final example that falls into this category is the *ring* network (Fig. 7). Each node is connected only to two neighbors via point-to-point unidirectional links. Messages (or packets) are circulated around the ring (hence achieving the broadcast feature), incurring a delay of only one or more bit times within each of the intermediate nodes. The multiaccess schemes available to operate multiaccess-broadcast networks are discussed in more detail below in the section on multiaccess link control.

A third type of packet network can be identified. It is the (multihop) *store-and-forward multiaccess/broadcast* type, which combines the features exhibited in the two types mentioned above. The best (and perhaps only) example is DARPA's packet radio network (PRNET). (See Fig. 8.) Its goal is to provide direct communication by a ground radio network among mobile users over wide geographical areas. This requires store-and-forward switches called repeaters to become integral components of the system. Here too, for easy communication among mobile users and for rapid deployment in military applications, all users employ omnidirectional antennas and share a high-speed radio channel.

Today all these types of networks have progressed considerably, and each plays an important role of its own. Communication systems that are efficient and easy to use involve several types of networks coexisting and integrated to form a single system, the constituents of

which are totally transparent to the user. The interconnection of several networks forms what is referred to as an internetworking environment.

NETWORK FUNCTIONS

The basic function of a computer network is to make it possible for geographically remote end users to communicate. "End users" means any entities that require to communicate; these are application processes that may reside in computers, terminals, or any interface between a piece of equipment (or a human) and the communication network itself. Communication should be accomplished efficiently, by making efficient use of the communications resources, and in a way that is responsive to the specific needs of the end users. It should also be made possible in spite of various types of errors that may occur in the transmission process, and in spite of the differences that exist between the end users with respect to the data formats used, the data rates supported, the patterns of intermittency, etc.

As this global function is a very complex one, it is best accomplished by dividing it into a sequence of more elementary functions that are organized in some fashion. The precise definition and structure of this aggregate set of functions is called a network *architecture*. Most typically, network architectures follow a

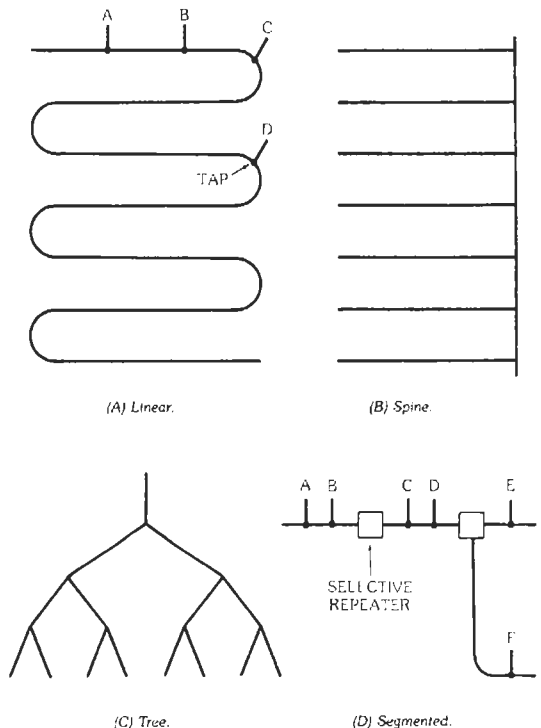


Fig. 6. Cable topologies. (From Andrew S. Tanenbaum, Computer Networks, © 1981. Used by permission of Prentice-Hall, Inc., Englewood Cliffs, N.J.)

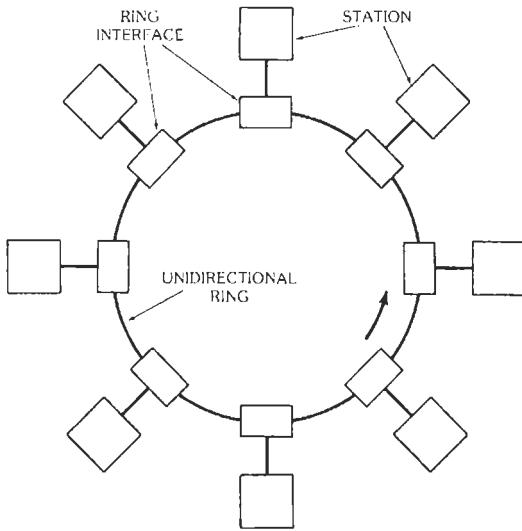


Fig. 7. A ring network. (From Andrew S. Tanenbaum, Computer Networks, © 1981. Used by permission of Prentice-Hall, Inc., Englewood Cliffs, N.J.)

linearly hierarchical model; i.e., the functions are organized into a linear succession of so-called *layers*. Such architectures are referred to as *layered architectures*. At the lowest level of the hierarchy, for example, are those functions that concern themselves with providing a physical connection between the end users and with the transmission of raw bits (i.e., regardless of their meaning) across the physical transmission resources. These functions constitute the physical layer. At the highest level, we find those portions of the application processes residing at the end users' machines that are the originators and destinations of the communication requests, and which directly serve the end user as far as the global communication function is concerned. The other intermediate layers comprise all other functions ranging from detecting and correcting errors to achieving efficient utilization of the resources; performing routing functions; preventing congestion; regulating the flow of data; preventing unfairness; supporting the patterns of intermittency and specific service requirements of the end users; and finally resolving the differences among the users pertaining to formats, character codes, device control procedures, and the like. In such a hierarchical structure, each layer is considered to be wrapping the lower layers and isolating them from the higher layers. Each layer offers a well specified service that it provides to the adjacent layer above it. In supplying that service, the layer obtains service from the adjacent layer just below it and performs its own functions. The rules for performing the functions at a given layer, once specified, are called *protocols*. The interactions between adjacent layers (at their boundaries) are called *interfaces*. The rules for achieving interactions and communication between ad-

acent layers, once specified, are called *interface protocols*. Also essential to a layered architecture is the concept of *peer protocols* or *peer interaction*. The processes implementing protocols within a layer at one machine (host or intermediate mode) communicate only with the processes implementing protocols within the same layer at another machine. Compatibility between protocols of the same level at (virtually) adjacent machines must be guaranteed. This layered approach to network functions for their description and design is now universal, and most (if not all) existing networks follow a layered architecture.

Layered architectures offer a number of advantages. The first is the prevention of ad hoc methods of network design, which, due to the intricacies involved, would certainly lead to a proliferation of designs difficult to understand, implement, or interact with. This is particularly beneficial when systems designed by different organizations or manufacturers are to be interconnected. The second advantage is that of modularity: If one specifies the interfaces between adjacent layers in a general way that is independent of the particular protocols designated at the interacting layers, then it becomes possible to modify the protocols at some layer without affecting the rest of the system. This is particularly beneficial in order to take advantage of newly emerging technologies, and in order to experiment with new protocols without the need for a major redesign of the entire system. For these reasons, the International

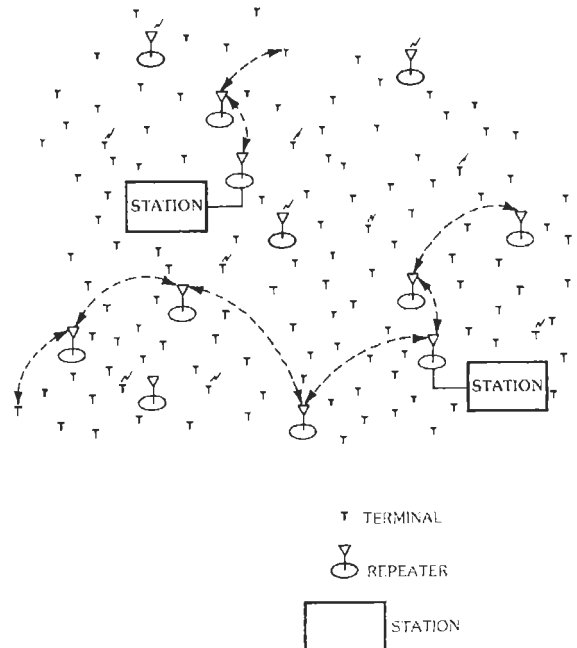


Fig. 8. A multihop packet radio network.

Organization for Standardization† (ISO) and other standards organizations were urged to come up with standards for “Open Systems Interconnections” (OSI). The objective is that by conforming to those international standards, a system will be capable of interacting with all other systems obeying the same standards throughout the world. The standard architecture adopted by ISO, also the best known one, is a layered architecture with seven layers. It is only an architecture, in the sense that it only defines the services to be performed by a layer for the next higher layer, independent of how these services are performed. Accordingly, it is often referred to as ISO’s OSI reference model. It is the first stage toward a complete standardization of network functions, although very few protocols and interfaces have been standardized so far.

The number of layers decided upon is the result of long debates and of a number of principles followed in the process. Some of these principles are:

1. To collect “similar” functions into the same layer and separate “manifestly different” functions into separate layers.
2. To create a layer of those functions that are easily localized and that may be totally redesigned in the future when taking advantage of new technological advances, without the need to change the services or the interfaces with the adjacent layers.
3. To create boundaries at those service points where the interactions across the boundaries are minimized, those which past experience has demonstrated to be successful, and those where it may be useful at some point in time to standardize the corresponding interface.

Another principle calls for a limitation on the number of layers so as to keep the engineering task of describing them and integrating them a simple one, but it allows further subgrouping of functions within a layer so as to form sublayers that may be bypassed if the corresponding services are not needed.

The seven layers of the ISO OSI reference model (Fig. 9) are: (1) the physical layer, (2) the data link layer, (3) the network layer, (4) the transport layer, (5) the session layer, (6) the presentation layer, and (7) the application layer. Specific protocols and their description are considered in subsequent sections.

While the ISO reference model is used here as a guide to describe network functions and protocols, it is to be noted that there are a few existing network architectures, either designed for experimental networks, such as the ARPANET, or supplied by computer

manufacturers, such as the Digital Equipment Corporation DECNET and the IBM Systems Network Architecture (SNA). These architectures are layered but do not correspond exactly to the ISO reference model. See Table 1. For more information on these architectures, the reader is referred to references 2 and 3.

THE PHYSICAL LAYER

The physical layer is concerned first of all with the transparent transmission of a bit stream (regardless of its meaning) across physical communication resources. In local networks, the physical medium may be a twisted pair, a coaxial cable, an optical fiber, or radio, and it may be privately owned. For long-haul links, the medium may be either copper wires or optical fibers as in terrestrial links, or radio as in satellite links; it is supplied by a common carrier. The methods used for the transmission of bits across the physical medium depend on the type of medium used. The description of these techniques is out of the scope of this chapter, but suffice it to mention here that bit transmission is done either in analog form by means of *modems* or in digital form by means of *line drivers*, and it involves determining such attributes as data encoding, timing, voltage levels, data rates, type of operation (half-duplex or full-duplex, synchronous or asynchronous), etc.

The second concern of the physical layer is to provide a physical interface between the end-user machine—which may be a terminal, a computer, or any other data-processing box, and which is referred to by the telecommunications administrations as the *Data Terminal Equipment* (DTE)—and the termination point of the communications circuit—that is the modem, line driver, etc., referred to as the *Data Circuit-Terminating Equipment* (DCE). (See Fig. 10.) The definition of such an interface requires the determination of four important characteristics: the mechanical, electrical, functional, and procedural characteristics.

The mechanical aspects pertain to the point of demarcation, which most typically consists of a pluggable connector. They include specifications of the connector, its latching and mounting arrangements, its location with respect to the DCE, etc. Fig. 11 illustrates the various connectors known, along with their ISO identification numbers. The number of pins per connector ranges from 9 to 37.

The electrical aspects pertain to the electrical characteristics of the generators and receivers. They specify such parameters as the range of the signal voltage level, rise-time characteristics of the generator, data signaling rates as a function of the interconnecting cable distances, generator and receiver impedances, etc. They also include specifications regarding the reference with respect to which signal levels are measured; two cases exist: the “unbalanced” case where a single common return lead is used (or perhaps one common return for each direction), and the “balanced” case where each

† The International Organization for Standardization (ISO) is a voluntary nontreaty group, the membership of which includes the principal standardization body of each represented nation. The US member body is the American National Standards Institute (ANSI). ANSI is a nonprofit, nongovernmental organization. It serves as the national clearing house and coordinating activity for voluntary standards in the US.

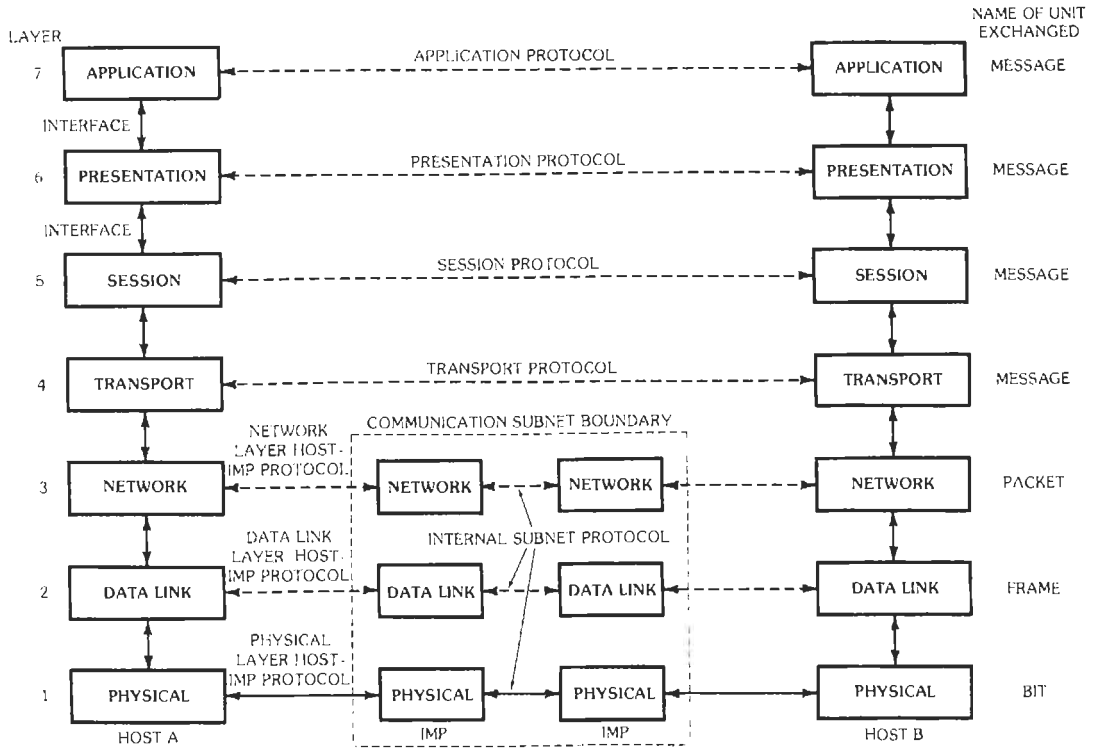


Fig. 9. The ISO OSI reference model. (From Andrew S. Tanenbaum, Computer Networks, © 1981. Used by permission of Prentice-Hall, Inc., Englewood Cliffs, N.J.)

interchange circuit uses a pair of wires creating a differential signal. A balanced configuration is more desired because it permits longer distances and higher data rates.

The functional characteristics specify the number of interchange circuits used and their assigned functions.

The functions are classified into four categories: data, control, timing, and grounds. Timing circuits are used for bit (and sometimes byte) synchronization, whereas control circuits are used for the exchange of status information, commands, and responses. Some interfaces employ only one function per interchange circuit

TABLE I. APPROXIMATE CORRESPONDENCES AMONG THE VARIOUS NETWORK ARCHITECTURES†

Layer	ISO	ARPANET	SNA	DECNET
7	Application	User	End user	Application
6	Presentation	Telnet, FTP	NAU services	
5	Session	(None)	Data flow control	(None)
4	Transport	Host-host	Transmission control	Network services
3	Network	Source to destination IMP	Path control	Transport
2	Data link	IMP-IMP	Data link control	Data link control
1	Physical	Physical	Physical	Physical

† From Andrew S. Tanenbaum, Computer Networks, © 1981. Used by permission of Prentice-Hall, Inc., Englewood Cliffs, N.J.

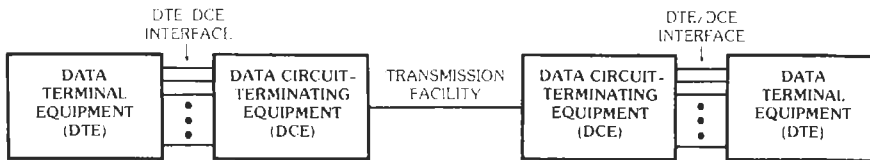


Fig. 10. DTE/DCE interface. (From Paul E. Green, Jr., *Computer Network Architectures and Protocols*, © 1982. Used by permission of Plenum Publishing Corp., New York.)

(e.g., CCITT V.24†), whereas others multiplex DCE control functions and data over a single "data" interchange circuit in each direction, resulting in a considerably more compact interface (e.g., CCITT X.24).

The procedural characteristics consist of the set of procedures for using the interchange circuits, in order to achieve the transmission of bit streams as well as to provide maintenance test loops.

Three examples of widely used physical interface protocols are considered here to illustrate the concepts and physical characteristics discussed above. These are EIA RS-232-C, § EIA RS-449, and CCITT X.21. They are all serial data DTE/DCE interfaces. The dominant interface in use today between terminals and modems is RS-232-C. The connector is the 25-pin ISO 2110 (Fig. 11). Its electrical characteristics are compatible with CCITT recommendation V.28. This specifies a single-ended generator that produces a 5 to 15 volt signal with respect to signal ground, negative for binary 1 and positive for binary 0; a single common return lead is used for all interchange circuits (thus unbalanced); the generator rise time is fast (1 ms to cross ± 3 V); the data signaling rates are limited to below 20 Kbits/s; the cable distances are limited to within 15 m. The functional characteristics of RS-232-C are compatible with CCITT recommendation V.24. Recommendation V.24 defines 43 interchange circuits, one function per circuit; RS-232-C uses 21 of the 43 interchange circuits; not all 21 circuits are needed in every application (for example, timing circuits are omitted for asynchronous applications). The procedural characteristics of RS-232-C are many and complex. They contain procedures equivalent to those in recommendation V.24 describing the interrelationships between interchange circuits, as well as those in recommendation X.20bis and X.21bis, which

specify asynchronous and synchronous operation on a public data network for DTEs designed to interface with V-series asynchronous and synchronous modems, respectively. (As indicated earlier in a footnote, the V-series recommendations correspond to data communications over telephone facilities.) For a detailed description of the procedural characteristics of RS-232-C, see reference 3.

In 1973, EIA began work on a new interface with improved performance over RS-232-C (longer cable distance and higher maximum data rates) and with additional interface functions such as maintenance loopback testing. The result was RS-449, published in November 1977. The new interface was carefully designed so as to allow interoperability with existing RS-232-C equipments. The RS-449 interface uses connectors of the same family as that used for RS-232-C: it uses a 37-pin connector (ISO 4902) for the basic interface, and a separate 9-pin connector if a "secondary channel" operation is in use (Fig. 11). Furthermore, a careful pin assignment plan is chosen to minimize cross talk in multipair cables and to facilitate the design of an adapter to RS-232-C. As in RS-232-C, the functional characteristics of RS-449 are compatible with V.24, one function per interchange circuit. There are 10 new circuits with respect to RS-232-C, among which we note the following: send common (SC) and receive common (RC), which are the common return leads for all unbalanced interchange circuits employing one wire in the direction toward the DCE and the DTE, respectively; terminal in service (IS), which indicates to the DCE whether the DTE is operational or not; local loopback (LL), remote loopback (RL), and test mode (TM), which are used in checking and testing; select standby (SS) and standby indicator (SB), which are used when standby facilities are in use in order to facilitate the rapid restoration of service when a failure has occurred. The electrical characteristics are defined in RS-423-A (compatible with V.10/X.26) and RS-422-A (compatible with V.11/X.27). Standard RS-423-A, referred to as the "new unbalanced" electrical characteristics, specifies a single-ended generator that produces 4 to 6 volts with respect to the common return; a single common return is used for each direction; data signaling rates are up to 3 Kbits/s over cable distances of 1000 meters and can be as high as 300 Kbits/s for distances up to 10 meters. Standard RS-422-A, referred to as the "new balanced" characteristics, specifies a balanced generator producing a 2 to 6 volt differential signal; each interchange circuit employs

† The International Telegraph and Telephone Consultative Committee (CCITT) is a committee of the International Telecommunication Union (ITU), a specialized agency of the United Nations Organization. The CCITT work on data communications is focused in two study groups. CCITT Study Group XVII is responsible for data communications over telephone facilities. Its work is contained in V-series recommendations. CCITT Study Group VII is responsible for data communications over data networks. Its work is contained in X-series recommendations.

§ The Electronic Industries Association (EIA) is a trade association that represents manufacturers in the US electronics industry. EIA standards on data communications are published in the RS-series.

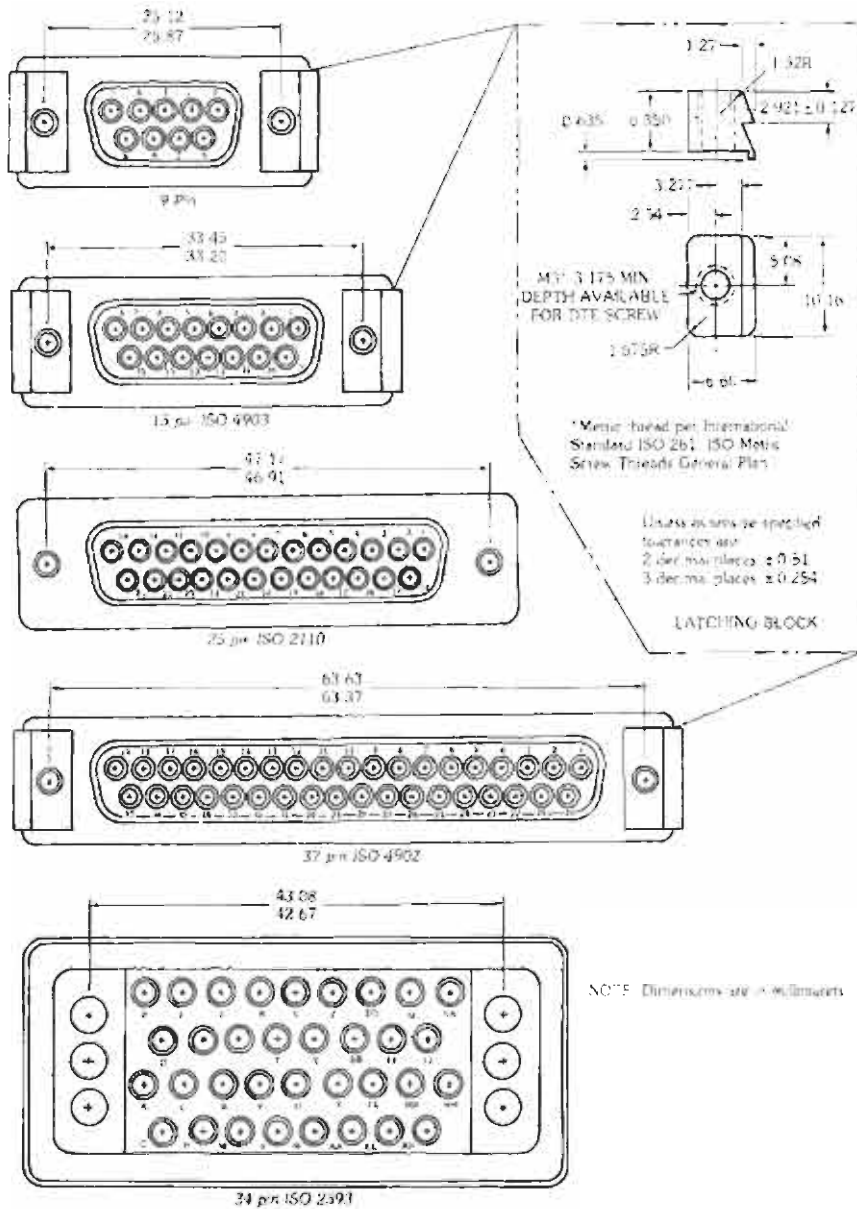


Fig. 11. DCE connectors. (From Paul E. Green, Jr., Computer Network Architectures and Protocols, © 1982. Used by permission of Plenum Publishing Corp., New York.)

two wires; data rates are up to 100 Kbits/s over cable distances of 1000 meters and can be as high as 10 Mbits/s over 10 meters. Standard RS-449 uses the new unbalanced characteristics, RS-423-A, for all interchange circuits when the data rate is below 20 Kbits/s. This permits direct interoperability with RS-232-C since RS-423-A (V.10/X.26) is interoperable with both V.28 (RS-232-C) and RS-422-A (V.11/X.27). At data rates above 20 Kbits/s and up to 2 Mbits/s, the new balanced FIA RS-422-A must be used on 10 specific

interchange circuits (send data, receive data, send timing, . . .) employing two wires per circuit, while RS-423-A is used for all other circuits.

As for the procedural characteristics of RS-449, all those procedures defined in RS-232-C carry over to RS-449. The newly added test and standby functions are conceived on the basis of action-reaction and thus are simple. For example, after the DTE turns on the local loopback (I.L), it waits until the DCE responds with test mode on. The DTE can then proceed with test

data transmission on the send data circuit, expecting to receive them back on the receive data circuit. Deactivation follows a similar procedure.

The physical characteristics for a general-purpose DTE/DCE interface for synchronous digital transmission on public data networks are defined in CCITT recommendation X.21. The DCEs are linked by means of real digital circuits, possibly through circuit-switching equipment. When circuit-switched services are provided, X.21 also specifies the call control procedures, i.e., the protocols by which to establish (and later to disconnect) a physical connection between two DTEs. Since the establishment of a circuit clearly requires knowledge of the address of the remote DTE (or its corresponding DCE) and the understanding of the various call progress signals, and since these requirements are supposed to be met only at the network layer (as discussed below), there is a general consensus that the call control procedure span the first three layers with the resulting service being offered by the network layer to the transport layer. Thus, these issues are dealt with in the network layer. The benefit of such a viewpoint is that of consistency in the network architecture, whereby the classification of functions is the same regardless of whether the data network is of the circuit-switched type or the packet-switched type. Indeed, in packet networks, the network layer takes care of opening a "connection," called the "virtual circuit," with the remote DTE before the flow of packets takes place.

Simplicity and enhanced performance were prime objectives in the design of X.21. It permits interface operation over distances considerably greater than that available with V.28, and for the synchronous data rates specified by X.1—that is, 600, 2400, 4800, 9600, and 48 000 bits/s. Accordingly, the new balanced electrical characteristics (Recommendation X.27, introduced above as EIA RS-422-A) are specified for the DCE side of the interface. To allow flexibility in DTE design at the data rates 600, 2400, 4800, and 9600 bits/s, the DTE is permitted to use either the new balanced or the new unbalanced (X.26 or RS-423-A) electrical characteristics. The mechanical interface for X.21 is the 15-pin ISO 4903 connector. There is a careful assignment of interchange circuits to connector pin numbers, a pair of pins for each interchange circuit. The number of interchange circuits is reduced to 5 (Fig. 12). A transmit (T) circuit and a receive (R) circuit are used to convey both user data and network control information, depending on the state of the control (C) circuit and the indication (I) circuit; bit timing is continuously provided by a signal element timing (S) circuit; a sixth interchange circuit that provides byte timing information is optional; a signal ground circuit is also provided.

Among all the procedures specified in X.21, those associated with the quiescent phase are agreed upon to be within the physical level. Two quiescent signals are defined for the DCE: *DCE not ready*, indicating that no service is available, and *DCE ready*, indicating that the DCE (network) is ready to enter the operational phase.

Three quiescent signals are defined for the DTE: *DTE uncontrolled not ready*, indicating that the DTE is unable to enter operational phases because of abnormal conditions, *DTE controlled not ready*, indicating that the DTE is operational but is temporarily unable to enter operational phases, and *DTE ready*, indicating that the DTE is ready to enter operational phases. Fig. 13 shows the various combinations of quiescent states of the X.21 interface and the possible link transitions between these states.

THE DATA LINK CONTROL LAYER

The data link control layer offers the capability for error-free conveyance of messages between machines that are physically connected by a communication channel. It consists of two sublayers. The upper sublayer deals with such issues as synchronization, error control, and link management and is present in all cases. The lower sublayer deals with the multiaccess link control and is present only when there is a shared channel that provides an any-to-any topological connectivity between stations. The focus in this section is on the upper sublayer. Multiaccess link control is described in the following section.

Asynchronous DLC Protocols

The DLC protocols used in early teleprinters are of the *asynchronous* type; the amount of data transmitted at a time is one character (5 to 8 bits—Fig. 14); when

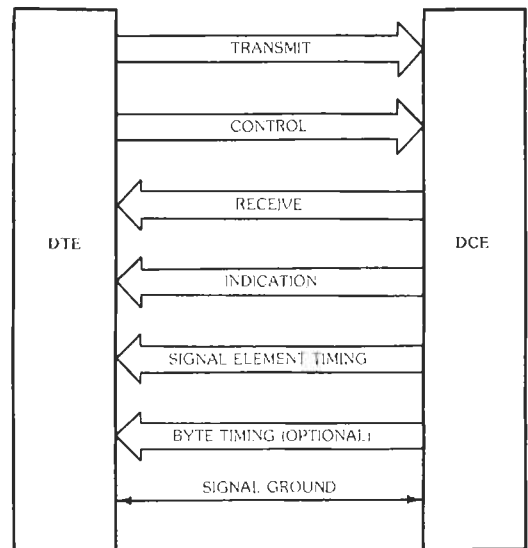
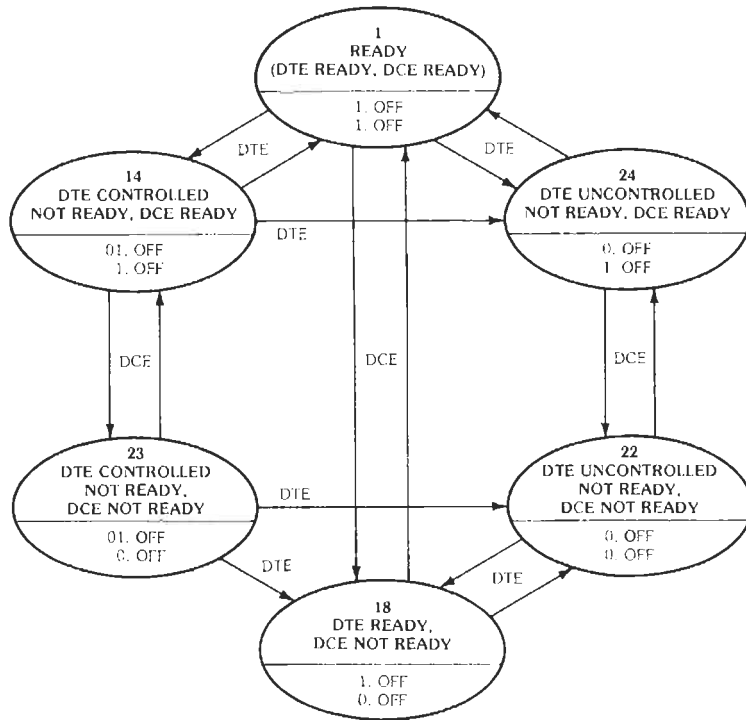
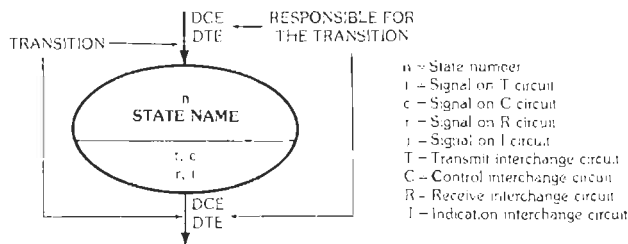


Fig. 12. CCITT Recommendation X.21 DTE/DCE interface. (From Paul E. Green, Jr., *Computer Network Architectures and Protocols*, © 1982. Used by permission of Plenum Publishing Corp., New York.)



LEGEND: Each state is represented by an ellipse wherein the state name and number are indicated, together with the signals on the four interchange circuits which represent that state. Each state transition is represented by an arrow, and the equipment responsible for the transition (DTE, or DCE) is indicated beside that arrow.



0 and 1 Refers to steady binary conditions
 01 Refers to alternate binary 0 and binary 1 conditions
 OFF and ON Respectively, refer to continuous OFF (binary 1) and ON (binary 0) conditions.

Fig. 13. CCITT Recommendation X.21 quiescent states. (From Paul E. Green, Jr., Computer Network Architectures and Protocols, © 1982. Used by permission of Plenum Publishing Corp., New York.)

no transmission is taking place, the channel is in the idle state, which is represented by a continuous 1. When a character is to be transmitted, the terminal transmits a 0 (the start bit) followed by the character (of fixed length), and then followed by a 1 (the idle state) for a minimum period of time called the stop interval. The receiver determines the initial sampling from the 1 to 0 transition, using a clock sixteen times the transmission rate. The sampling period is known from the bit transmission rate. Bit synchronism is maintained only during the transmission time of a character. Anytime following the stop interval, another character can be transmitted.

The start-stop protocols are simple but present several disadvantages, among others the overhead of a start bit and a stop interval associated with each character, and the lack of any inherent link control capability.

Synchronous DLC Protocols

To improve the efficiency of a data link and provide link control capability, *synchronous* protocols were introduced. A block of several characters referred to as a *frame* is transmitted as a whole, with bit synchronism

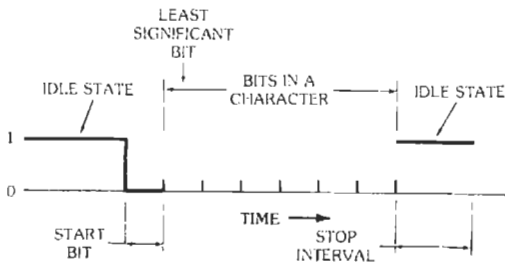


Fig. 14. Asynchronous DLC protocol.

maintained during the entire transmission time of the frame. These protocols provide the following functions: (1) frame synchronization (a mechanism needed to delimit the beginning and end of transmission blocks, and to acquire bit synchronization), (2) error control (mechanisms needed to detect errors, acknowledge correctly received blocks, request the retransmission of incorrectly received blocks, and control the sequence of blocks to identify lost and duplicate blocks), (3) link management (mechanisms needed to establish a data link over a communication facility that has been idle and, in the case of multipoint facilities, to identify sender and receiver), (4) flow control (mechanisms needed to regulate the flow of information), and (5) functions to permit the recovery from abnormal conditions such as illegal sequences, loss of response, etc. Two types of synchronous DLC protocols exist: character-oriented and bit-oriented.

Character-Oriented DLC Protocols—These protocols are suitable for two-way alternate operation on full-duplex or half-duplex multipoint, switched, or dedicated link configurations. These protocols started with IBM's binary synchronous communication (BSC) protocol in the late 1960s; following it, there has been a proliferation of different protocols by various manufacturers, which are optimized for specific implementations. Although of the same type (all are character-oriented), these are incompatible with each other. A standardization effort was undertaken by ANSI and ISO. Today, the most widely used protocols are IBM's BSC, ANSI X3.28 (1971, updated 1976), ISO's IS/745 (1975), and DEC's DDCMP. To illustrate this type of DLC protocol, ANSI 3.28 is used here.

In this protocol, the definitions of a message and a block are as follows: While a message is an ordered sequence of characters arranged to convey information from the originator user to the destination user (including a header portion), a block is a group of characters arranged for technical or logical reasons that is transmitted as a unit. A block may contain an entire message or a portion of a message.

Ten ASCII characters are designated as communication control characters. The first is SYN (synchronization character). It is used at the beginning of a data stream (following an idle time) to allow the receiver to

establish byte synchronization. A minimum of two SYN characters are usually used to guarantee proper operation. This character is also used following the data stream (during the idle time) to allow the receiver to maintain synchronization, if this is desired. Two characters, ENQ (enquiry) and EOT (end of transmission), are used to manage the link. Using ENQ, a station solicits a response (such as status or identification) from another user. In particular, this character is used by a master station in a nonswitched multipoint configuration to poll a secondary station, and thus to establish a link between it and the secondary station. The EOT character is used to indicate the end of transmission of one or more messages and to relinquish the data link. Four characters are used to delimit blocks that are being transmitted: SOH (start of heading) and STX (start of text) are used to designate beginning of a block containing message header information and message text, respectively; ETB is used to designate end of transmission block; and ETX (end of text) is used at the end of the last block of a message. Two characters, ACK and NAK, are used as affirmative response and negative response (e.g., a block received in error) to a sender, respectively. Since all these characters have a specific control meaning, they should not normally appear in the user's data. To remove such a restriction and allow control characters to be in the user's data stream (that is, to achieve transparency), a tenth character, DLE (data link escape) is added; for any control character to be interpreted as such, it must be preceded by DLE. If DLE itself is to be transmitted as data, then it is transmitted as the sequence DLEDLE.

For error detection, three techniques are available: the *vertical redundancy check* (VRC), the *longitudinal redundancy check* (LRC), and the *cyclic redundancy check* (CRC). The vertical redundancy check consists of appending one parity bit to each character; thus it can only detect an odd number of bit errors in a character. The longitudinal redundancy check consists of computing parity on the entire sequence of successive characters in a block, and transmitting the result as an additional check character; thus it is vulnerable to double bit errors in the row of characters. With CRC, the bit stream constituting a block or a packet is treated as a polynomial, $M(x)$, with the least significant bit considered to be the coefficient of x^0 . With k bits in the block, the polynomial is of degree $k - 1$. A generator polynomial, $G(x)$, of degree r is agreed upon by both sender and receiver. After appending r zero bits to the low-order end of the message to form the corresponding polynomial $x^r M(x)$, the latter is divided by $G(x)$ (polynomial division modulo 2), and the remainder is subtracted from $x^r M(x)$ (modulo 2). The result, $T(x)$, is the checksummed message. With the proper choice of $G(x)$, the CRC technique can safely detect single bit errors, double errors, any odd number of errors, and any burst error of length $\leq r$. For burst errors longer than $r + 1$, there is a nonzero probability of a bad message going undetected. Three international standards for a CRC generator exist. These are

$$\text{CRC-12} = x^{12} + x^{11} + x^3 + x^2 + x^1 + 1$$

$$\text{CRC-16} = x^{16} + x^{15} + x^2 + 1$$

$$\text{CRC-CCITT} = x^{16} + x^{12} + x^5 + 1a$$

Of the three techniques (VRC, LRC, and CRC), CRC offers the best protection against undetected errors.

Error-free transmission of messages is accomplished by the use of the ACK and NAK control characters in conjunction with an error-detection procedure and the retransmission of blocks with errors. Moreover, timers are provided at both the sender and the receiver in order to recover from abnormal conditions. The *response timer* (typically 2–3 seconds) at the sending station protects against invalid or missing response, and the receive timer (about 500 ms) protects against failure to receive or recognize ETB or ETX. With these definitions and features, ANSI 3.28 specifies a number of specific protocols, each suitable to a particular scenario (or situation). For example, “subcategory 2.4” is for two-way alternate transmission on a nonswitched multipoint configuration with centralized operation; “subcategory D1” is for message-independent blocking, with cyclic checking, alternating acknowledgements, and transparent heading and text.

Bit-Oriented DLC Protocols—A number of deficiencies exist in character-oriented DLC protocols. First of all, such protocols are best suited for the two-way alternate transmission mode (i.e., half-duplex transmission) and thus make inefficient use of full-duplex lines. Secondly, error checking is done only on text, leaving control sequences unprotected, and thus leading to complex recovery procedures. Thirdly, a single data link function is performed with each transmission (e.g., either data transmission, or acknowledgement, or polling command, etc.), which leads to a larger number of turnarounds and to an unsatisfactory ratio of data transfer exchange to control exchange. Finally, character-oriented DLC protocols are too rigid and do not allow for easy expansion.

The new bit-oriented DLC protocols have been designed to overcome all the above limitations. Examples are ANSI's Advanced Data Communication Control Protocols (ADCCP), ISO's High-Level Data Link Control (HDLC), and IBM's Synchronous Data Link Control (SDLC). The ADCCP and HDLC protocols are compatible. In this discussion, ADCCP is used to illustrate bit-oriented DLC protocols.

Bit-oriented DLC protocols have been designed to support all possible environments: point-to-point and multipoint configurations using two-way alternate or two-way simultaneous operation over switched and nonswitched transmission lines, with both terrestrial and satellite connections. Three data transfer modes have been defined: the normal response mode (NRM) and the asynchronous response mode (ARM) for use in point-to-point and multipoint configurations, and the asynchronous balanced mode (ABM) for use in point-

to-point configurations. In NRM and ARM, a so-called *primary station* controls the operation of the data link, while one or more *secondary stations* act as subservient to the primary. The primary issues commands and receives responses; the secondary station receives commands and issues responses in accordance with the command received and the mode of operation. In NRM, a secondary station initiates transmission only as a result of receiving explicit permission to do so. This mode is suitable for polled multipoint operation between a central location and a number of outlying stations. In ARM, a secondary station can initiate transmission without waiting for an explicit permission. It is suitable for a single primary and a *single* activated secondary wishing to transmit freely to one another without the overhead of polling. The ABM mode provides a *balanced* type of data transfer between two stations referred to as *logically equal* or *combined* stations. Each station operates as a primary for its data transfer, and is thus capable of initializing the link, activating the other combined station, and logically disconnecting the link. Typically, each combined station may be a host computer, an intelligent network node (e.g., a packet-switching node), or a highly intelligent terminal that has the capability to control the data link.

In bit-oriented DLC protocol, the basic transmission unit is called a *frame*. In order to achieve enhanced capabilities over character-oriented DLC protocols, a well defined frame format as shown in Fig. 15 is used, including an address field, a control field, and a frame check sequence (FCS).

The flag sequence (F) is the unique eight-bit pattern 01111110 used to delimit the start and end of each frame, and to fill idle time between frames during the transmission of multiple frames. In order to prohibit the occurrence of the flag in the address, control, information, and FCS fields, and thus achieve transparency, a technique called bit-stuffing is used. The transmitter inserts a 0 bit following any five contiguous 1 bits encountered in the above mentioned fields; the receiver deletes the 0 bit following five contiguous 1 bits following a 0 bit anywhere before receiving the closing flag sequence.

The address field in a command frame identifies the station (secondary or combined) that is to receive the command; in a response frame, it identifies the station (secondary or combined) that is sending the frame. Thus in NRM and ARM, the address always identifies the secondary station; in ABM, it identifies the response-generating portion of a combined station. There are two mutually exclusive address-field options: either a single octet accommodating up to 256 stations, or multiple octets recursively extendable by having the first bit of each octet be a 1 bit if the octet is the last one in the address or a 0 bit otherwise. Broadcast addressing is also possible by using the all-zero address.

The control field identifies the function and purpose of the frame. Three different formats exist, defining the three types of frames: information transfer type (I

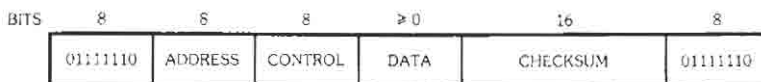


Fig. 15. Frame format for bit-oriented protocols. (From Andrew S. Tanenbaum, Computer Networks, © 1981. Used by permission of Prentice-Hall, Inc., Englewood Cliffs, N.J.)

frame), supervisory type (S frame), and unnumbered type (U frame). See Fig. 16. The P/F bit in the control field provides for a check-pointing mechanism allowing a response frame to be logically associated with the appropriate initiating command frame. In NRM, the primary station sets the P/F bit to 1 to poll a secondary station. The secondary station sets P/F to 1 in the last frame sent in response to a poll. In ARM and ABM, receipt of a frame with the P/F bit set to 1 causes the responding station to set the P/F bit to 1 in the next appropriate frame.

The I frames are used to transfer user information across a data link. A so-called *window mechanism* is used to guarantee sequential error-free delivery of the information without degrading the performance of the link. Successive frames are numbered sequentially (modulo some number 2^n), where the numbers range from 0 to the maximum, $2^n - 1$. The frame sequence number is transmitted in the control field $N(S)$. The sender maintains a list of consecutive sequence numbers called the *sending window* corresponding to frames it is permitted to send. Likewise, the receiver maintains a list of consecutive sequence numbers called the *receiving window* for frames it is permitted to receive. Most acknowledgements in these DLC protocols are handled in a *piggyback* fashion. A receiver that has correctly received all frames numbered up to some number, k , will send in the $N(R)$ field of the next I frame the number $N(R) = k + 1$; i.e., the sequence number of the next frame the receiver is expecting. In case an I frame is not ready for transmission, the receiver will send an S frame with $N(R) = k + 1$. Upon receiving an acknowledgement, the sender advances the lower edge of its sending window to $N(R)$ and the upper edge to $N(R) + W_S$. The sequence numbers of frames already transmitted that are within the sender window represent frames sent but as yet not acknowledged. The window size, W_S , represents the maximum number of such frames. It is determined as a function of the frame transmission time and the end-to-end propagation delay so as to "fill

the pipe" and achieve high throughput. This feature is particularly important when dealing with wideband satellite channels where the ratio of round-trip delay to frame transmission time is high. The receive window, W_R , is determined as a function of the buffer space available at the receiver. The sending window, W_S , and receive window, W_R , need not be the same. The window mechanism protects against the loss of a frame and the loss of acknowledgement, detects and filters out duplicates, and guarantees that the network level receives the frames in sequential order. It is important, however, that the condition $W_S + W_R \leq 2^n$ be satisfied. In the ADCCP and HDLC standards, two sequence ranges are available: the unextended case, which is modulo 8, and the extended case, which is modulo 128.

There are four supervisory frames. Receive ready (RR) is used to indicate that the receiver is ready to receive I frames; receive not ready (RNR) is used to indicate a temporary busy condition. Both RR and RNR S frames acknowledge I frames already received. Reject (REJ) is used to request retransmission of all I frames starting with $N(R)$; selective reject (SREJ) is used to request the retransmission of a single designated I frame, $N(R)$.

There are 32 U-frame commands and/or responses that can be defined, although fewer than 32 are in fact determined in the current standards. They do not contain a sequence number and thus are not accounted for in the normal flow of frames. They are mainly used in the control of the data link. Examples are: the set mode commands (SNRM, SARM, SABM, etc.), disconnect command (DISC), an unnumbered poll command, an unnumbered information (UI) command/response, an unnumbered acknowledgement (UA), and a frame reject (FRMR). (The frame reject is used to indicate that a received frame was in error and cannot be corrected by retransmission, such as, for example, an information field greater than the maximum established length, receipt of a control field that is either invalid or not implemented, etc.)

CONTROL FIELD FOR	CONTROL FIELD BITS							
	1	2	3	4	5	6	7	8
INFORMATION TRANSFER COMMAND/ RESPONSE (I FRAME)	0	N(S)			P/F	N(R)		
SUPERVISORY COMMANDS/ RESPONSES (S FRAME)	1	0	S	S	P/F	N(R)		
UNNUMBERED COMMANDS/ RESPONSES (U FRAME)	1	1	M	M	P/F	M	M	M

Fig. 16. Control field formats in ADCCP. (From Paul E. Green, Jr., Computer Network Architectures and Protocols, © 1982. Used by permission of Plenum Publishing Corp., New York.)

The FCS consists of the 16-bit checksum obtained from the CRC error detection technique using the CRC-CCITT generator. An optional 32-bit FCS is also available if a higher degree of error detection is needed. The corresponding generator polynomial is $x^{32} + x^{26} + x^{23} + x^{22} + x^{16} + x^{12} + x^{11} + x^{10} + x^8 + x^7 + x^5 + x^4 + x^2 + x + 1$. Just as with character-oriented DLC protocols, classes of procedures are defined in order to provide organization and direction for the application of the bit-oriented data link control procedures. Certain commands belong to all the classes, while others are viewed as being optional. Examples are: UNC, 3, 4—unbalanced NRM with SREJ and UI capabilities; BAC, 2, 8—balanced ABM with REJ and the restriction on the use of I frames as commands only.

MULTIACCESS LINK CONTROL

Multiaccess protocols arise whenever a communication channel is shared by many independent contending users. Two major factors contribute to a multiaccess situation: the need to share an expensive communication channel in order to achieve its efficient utilization, and the need to provide a high degree of connectivity for communication among independent subscribers. The main issue of concern is how to control access to the common channel to allocate the available bandwidth efficiently. The solutions to this problem form the set of protocols known as *multiaccess protocols*. These protocols and their performance differ according to the environment and the specific requirements to be satisfied. In satellite channels, there is an inherent long propagation delay of approximately 0.25 s, a delay that is usually long compared to the transmission time of a packet given the bandwidth usually available. In ground radio environments, the opposite is true: the transmission radius is relatively small due to a limited power supply at the radio units and a heavy signal attenuation, and the data rate is relatively low, so that the propagation delay is relatively short compared to the transmission time of a packet. Finally, local area environments are characterized by a large and often variable number of inexpensive devices; they call for networks with simple topologies and inexpensive connection interfaces; the propagation delay is short but the bandwidth can be high; the propagation delay may or may not be larger than the packet transmission time.

Multiaccess schemes are evaluated according to various criteria: bandwidth utilization, message delay, the ability to support traffic of different types simultaneously (with different priorities, variable message lengths, and differing delay constraints), and robustness (defined here as the insensitivity to errors resulting in misinformation).

Multiaccess protocols differ by the static or dynamic nature of the bandwidth allocation algorithm, the centralized or distributed nature of the decision-making process, and the degree of adaptivity of the algorithm to changing needs. Accordingly, these protocols can be grouped into several classes described hereafter.

Fixed Assignment Techniques

Fixed assignment techniques consist of allocating the channel to the users, independently of their activity, by partitioning the time-bandwidth space into slots that are assigned in a static predetermined fashion. These techniques take two common forms: *orthogonal*, such as frequency division multiple access (FDMA) or synchronous time division multiple access (TDMA), and "*quasi-orthogonal*," such as code division multiple access (CDMA).

The FDMA technique consists of assigning to each user a fraction of the bandwidth and confining its access to the allocated subband. Orthogonality is achieved in the frequency domain. The TDMA technique consists of assigning fixed predetermined channel time slots to each user; the user has access to the entire channel bandwidth, but only during its allocated slots. Here, signaling waveforms are orthogonal in time. The TDMA approach is more complex to implement than FDMA, but an important advantage is the connectivity that results from the fact that all receivers listen to the same channel while senders transmit on the same common channel at different times. Accordingly, with TDMA, many network realizations, both in ground and satellite environments, are easier to accomplish.

The CDMA technique allows overlap in transmission both in the frequency and time coordinates. It achieves orthogonality by the use of different signaling codes in conjunction with matched filters (or equivalently, correlation detection) at the intended receivers. Each user is assigned a particular code sequence, which is modulated on the carrier with the digital data modulated on top of that. Two common forms exist: the frequency-hopped SSMA and the phase-coded SSMA. In the former, the frequency is periodically changing according to some known pattern; in the latter, the carrier is phase modulated by the digital data sequence and the code sequence. Multiple orthogonal codes are obtained at the expense of increased bandwidth requirements (in order to spread the waveforms). With CDMA, there is also a lack of flexibility in interconnecting all users (unless, of course, matched filters corresponding to all codes are provided at all receivers). However, CDMA has the advantage of allowing the coexistence of several systems in the same band, as long as different codes are used for different systems.

Random Access Techniques

Since data traffic in computer communication is characterized as bursty, a more advantageous approach than fixed assignment is to provide a single sharable high-speed channel to the large number of users. The strong law of large numbers then guarantees that, with a very high probability, the demand at any instant will be approximately equal to the sum of the average demands of that population. When dealing with shared channels in a packet-switched mode, one must be prepared to resolve conflicts that arise when more than one demand

is placed on the channel. For example, in packet-switched channels, whenever a portion of the transmission of one user overlaps with the transmission of another user, then the two collide and "destroy" each other (unless a code division multiple-access scheme is used). The existence of some positive acknowledgement scheme permits the transmitter to determine if a transmission is successful or not.

ALOHA—Pure ALOHA permits a user to transmit any time it desires. If a user transmits a packet, and within some appropriate time-out period following its transmission it receives an acknowledgement from the destination, then it knows that no conflict occurred. Otherwise, it assumes that a collision occurred and it must retransmit. To avoid continuously repeated conflicts, the retransmission delay is randomized across the transmitting devices, thus spreading the retry packets over time. A slotted version, referred to as *slotted ALOHA*, is obtained by dividing time into slots of duration equal to the transmission time of a single packet (assuming constant-length packets). Each user is required to synchronize the start of transmission of its packets to coincide with the slot boundary. When two packets conflict, they will overlap completely rather than partially, providing an increase in channel efficiency over pure ALOHA. Due to conflicts and idle channel time, the maximum channel efficiency available with ALOHA is less than 100 percent, 18 percent for pure ALOHA and 36 percent for slotted ALOHA. Note that, although the maximum achievable channel utilization is low, the ALOHA schemes are superior to fixed assignment schemes when there is a large population of bursty users and low packet delay is of the essence. (See reference 3.)

Carrier Sense Multiple Access (CSMA)—In the CSMA technique, an attempt is made to avoid collisions by listening to the carrier due to transmission from another user before transmitting, and inhibiting transmission if the channel is sensed busy. This is advantageous when the propagation delay between any source-destination pair is small compared to the packet transmission time. Many CSMA protocols exist; they differ according to the action that a terminal takes to transmit a packet after sensing the channel. In all cases, however, when a terminal learns that its transmission has incurred a collision, it reschedules the transmission of the packet according to a randomly distributed delay. At this new point in time, the transmitter senses the channel again and repeats the algorithm dictated by the protocol. In *nonpersistent CSMA*, a ready terminal senses the channel and operates as follows. If the channel is sensed idle, the terminal transmits the packet. If the channel is sensed busy, then the terminal schedules the retransmission of the packet to some later time according to the retransmission delay distribution. At this new point in time, it senses the channel and repeats the algorithm described.

In the *1-persistent CSMA* protocol, a ready terminal senses the channel and operates as follows. If the

channel is sensed idle, it transmits the packet with probability one. If the channel is sensed busy, it waits until the channel goes idle and then immediately transmits the packet with probability one. In *1-persistent CSMA*, whenever two or more terminals become ready during a packet transmission period, they wait for the channel to become idle (at the end of that transmission), and then they all transmit with probability one. A conflict will also occur with probability one. Randomizing the starting time of transmission of packets accumulating at the end of a transmission period reduces interference and improves performance. The *p-persistent* scheme involves including an additional parameter p , the probability that a ready packet persists ($1 - p$ being the probability of delaying transmission by τ seconds, where τ is the maximum propagation delay among all pairs). Parameter p is chosen to reduce the level of interference while keeping the idle periods between any two consecutive nonoverlapped transmissions as small as possible.

The CSMA technique has been applied to ground radio (e.g., PRNET), and to local area communications (e.g., ETHERNET). In ETHERNET, CSMA is used on a tapped coaxial cable to which all the communicating devices are connected. On the coaxial cable, in addition to sensing carrier, it is possible for the transceivers to detect collisions. This is achieved by having each transmitting device compare the bit stream it is transmitting to the bit stream it sees on the channel. When transmitting users detect interference among several transmissions (including their own), they abort the transmission of colliding packets. This variation of CSMA is referred to as *carrier sense multiple access with collision detection (CSMA-CD)*.

The performance of CSMA is heavily dependent on the ratio, a , of propagation delay to packet transmission time. The maximum throughput of a CSMA protocol degrades significantly as a gets larger. For a ratio $a = 0.01$, nonpersistent CSMA achieves a channel utilization equal to 0.815, a significant improvement over the ALOHA schemes.

While until recently most of the concepts described in this section had been realized in experimental systems (namely, the ALOHA System, PRNET, and Xerox's experimental ETHERNET), it is important to note that today many contention systems of the ETHERNET type are available on the market. Examples are the Hyperchannel and the Hyperbus of Network Systems Corporation, Z-Net of Zilog, Omninet of Corvus, and ETHERNET itself. The latter has been announced as a product made available jointly by Xerox Corporation, Digital Equipment Corporation, and INTEL. Complete specifications of the data link and physical link protocols have been issued and constituted the basis of a standard for the IEEE Computer Society Project 802 on the standardization of local networks. A key feature that distinguishes this product from other already available systems is the LSI implementation of many of the data link and physical link protocols. The LSI implementation of network protocols clearly marks

a trend in the evolution of computer networking, a trend that is indicative of the existence of a wide market and the need to provide reasonably priced components.

Busy-Tone Multiple Access (BTMA)—In ground radio environments, it is possible for two terminals to be within range of the intended receiver, but out of range of each other or separated by some physical obstacle opaque to radio signals. The existence of hidden terminals in such an environment significantly degrades the performance of CSMA. The hidden-terminal problem can be eliminated by frequency dividing the available bandwidth into two separate channels, a busy-tone channel and a message channel, thus giving rise to *busy-tone multiple access* (BTMA). As long as a node senses carrier on the message channel, it transmits a (sine wave) busy-tone signal on the busy-tone channel. It is by sensing carrier on the busy-tone channel that nodes determine the state of the message channel. The action that a node takes pertaining to the transmission of the packet is again prescribed by the particular protocol being used, similar to those described for CSMA.

Capture—Capture in narrow-band channels can be defined as the ability of a receiver to receive a packet successfully (with nonzero probability) although the packet is partially or totally overlapped by another packet transmission. Capture is mainly due to a discrepancy in receive power between two signals that allows the receiver to receive the stronger correctly; both distance and transmit power contribute to this discrepancy. Clearly, capture improves the overall network performance, and, by means of adaptive transmit power control, it allows one to achieve either fairness to all users or intentional discrimination.

Spread-Spectrum Multiple Access (SSMA)—The SSMA technique is considered here to be a CDMA scheme. In one form of SSMA for packet radio, all transmitters employ the same code. Security, coexistence with other systems, and the ability to counteract the effects of multipath and capture are key benefits of SSMA. Contrary to the case of nonspread systems, the effect of interference in SSMA is minimized by the "capture effect," defined as the ability of the receiver to "lock on" one packet while all other overlapping packets appear as noise. The receiver locks on a packet by correctly receiving the preamble appended to the front of the transmitted packet. As long as the preambles of different packets do not overlap in time, and the signal strength of the late packets is not too high, capture of the earliest packet occurs with a high probability. In essence, SSMA allows a packet to be captured at the receiver, while CSMA allows a user to capture the channel. It is possible to use CSMA in conjunction with SSMA, but the channel sensing is more difficult. This mode will have the benefit of keeping away all users within hearing distance of the transmitter and thus help keep the capture effect and antijamming capability of the system at the desired level.

Centrally Controlled Demand Assignment

Demand assignment techniques require that explicit information regarding the need for the communication resource be exchanged. These techniques may be either centralized, whereby a central scheduler performs the assignment, or distributed, whereby all stations take active part in the assignment. Centrally controlled techniques are addressed in the present subsection.

Circuit Oriented Systems—In circuit oriented systems, the bandwidth is divided into FDMA or TDMA subchannels that are assigned on demand. The satellite SPADE system, for example, has a pool of FDMA subchannels that are allocated on request. It uses one subchannel operated in a TDMA fashion with one slot per frame permanently assigned to each user to handle the requests and releases of FDMA circuits. Intelsat's MAT-1 system uses the TDMA approach. The TDMA subchannels are periodically reallocated to meet the varying needs of earth stations.

The Advanced Mobile Phone Service (AMPS), introduced by Bell Laboratories, is another example of a centrally controlled FDMA system. The uniqueness of this system, however, lies in an efficient management of the spectrum based on space division multiple access (SDMA). That is, each subchannel in the pool of FDMA channels is allocated to different users in separate geographical areas, thus considerably increasing the spectrum utilization. To accomplish space division, the AMPS system has a cellular structure and uses a centralized handoff procedure (executed by a central office) that reroutes the telephone connections to other available subchannels as the mobile users move from one cell to another.

Polling Systems—In packet oriented systems, polling consists of having a central controller send polling messages to the terminals, one by one, asking the polled terminal to transmit. If the polled terminal has something to transmit, it goes ahead; if not, a negative reply (or absence of reply) is received by the controller, which then polls the next terminal in sequence. Polling requires this constant exchange of control messages between the controller and the terminals and is efficient only if (1) the round-trip propagation delay is small, (2) the overhead due to polling messages is low, and (3) the user population is not a large bursty one.

Adaptive Polling or Probing—The primary limitation of polling in lightly loaded systems is the high overhead incurred in determining which of the terminals have messages. A modified polling technique called *probing*, based on a tree searching algorithm, helps decrease this overhead. This technique assumes that the central controller can broadcast signals to all terminals. First the controller interrogates all terminals, asking if any of them has a message to transmit, and repeats this question until some terminals respond by putting a signal on the line. When a response is

received, the central station divides the population into subsets (according to some tree structure) and repeats the question to each of the subsets. The process is continued until the terminals having messages are identified. When a single terminal is interrogated, it transmits its message. This probing technique can be made adaptive by having the controller start a cycle by probing groups of smaller size as the probability of terminals having messages to transmit increases.

Split-Channel Reservation Multiple Access (SRMA)—An attractive alternative to polling is the use of explicit reservation techniques. In dynamic reservation systems, it is the terminal that makes a request for service on some channel whenever it has a message to transmit. The central scheduler manages a queue of requests and informs the terminal of its allocated time. In SRMA, the available bandwidth is divided into two channels, one used to transmit control information and the second used for the data messages themselves. The request channel is operated in a random access mode (ALOHA or CSMA). Upon correct reception of the request packet, the scheduling station computes the time at which the backlog on the message channel will empty and transmits back to the terminal an answer packet containing the address of the terminal and the time at which it can start transmission.

Demand Assignment With Distributed Control

There are two reasons why distributed control is desirable. The first is reliability: with distributed control the system is not dependent on the proper operation of a central scheduler. The second is improved *performance*, especially when dealing with systems with long propagation delays, such as those using satellite channels. The basic element underlying all distributed algorithms is the need to exchange control information among the users, either explicitly or implicitly. Using this information, all users then execute independently the same algorithm, with the result that there is some coordination in their actions.

Reservation-ALOHA—Reservation-ALOHA for a satellite channel is based on a slotted time axis where the slots are organized into frames of equal size. The duration of a frame must be greater than the satellite propagation delay. A user that has successfully accessed a slot in a frame is guaranteed access to the same slot in the succeeding frame, and this continues until the user stops using it. “Unused” slots, however, are free to be accessed by all users in a slotted ALOHA contention mode. A slot in a frame is an unused slot if in the *preceding* frame it either was idle or contained a collision. Users need simply to maintain a history of the usage of each slot for just one frame duration. Since no request is explicitly issued by the user, this scheme has been referred to as an *implicit reservation* scheme. Clearly, Reservation-ALOHA is effective only if the

users generate stream type traffic or long multipacket messages. Its performance will degrade significantly with single packet messages, since every time a packet is successful the corresponding slot in the following frame is likely to remain empty.

A First-in First-out (FIFO) Reservation Scheme—In this scheme, reservations are made explicitly. Time division is used to provide a reservation sub-channel. The channel time is slotted as before, but every so often a slot is divided into V small slots that are used for the transmission of reservation packets (as well as possibly acknowledgements and small data packets); these packets contend on the V small slots in a slotted ALOHA modc. All other slots are data slots and are used on a reservation basis, free of conflict. To execute the reservation mechanism properly, each station must maintain information on the number of outstanding reservations (the “queue in the sky”) and the slots at which its own reservations begin. These are determined by the FIFO discipline based on the successful reservations received. To maintain synchronization of control information at the proper time and to acquire the correct count of packets in the queue if out-of-sync conditions do occur, each station sends information regarding the status of its queue in its data packet. This information is also used by new stations that need to join the queue. The robustness of this system is achieved by a proper encoding of the reservation packets to increase the probability of their correct reception at *all* stations.

A Round-Robin (RR) Reservation Scheme—The basis of this scheme is fixed TDMA assignment, but with the major difference that “unused” slots are assigned to the active stations on a round-robin basis. This is accomplished by organizing packet slots into equal-size frames of duration greater than the propagation delay and such that the number of slots in a frame is larger than the number of stations. One slot in each frame is permanently assigned to each station. To allow other stations to know the current state (used or unused) of its own slot, each station is required to transmit information regarding its own queue of packets piggy-backed in the data packet header (transmitted in the previous frame). A zero count indicates that the slot in question is free. All stations maintain a table of the queue lengths of all stations, allowing them to allocate among themselves unassigned slots in the current frame. Round-robin or other scheduling disciplines can be used. A station recovers its slot by deliberately causing a conflict in that slot, which other users detect.

Distributed Tree Retransmission Algorithms—Tree algorithms are based on the observation that a contention among several active sources is completely resolved if and only if all the sources are somehow subdivided into groups such that each group contains at most one active source. (See probing in the section on centrally controlled demand assignment.) Each source corresponds to a leaf on a *binary* tree. The channel time

axis is slotted, and the slots are grouped into pairs. Each slot in a pair corresponds to one of the two subtrees of the node being visited. Starting with the root node of the tree, we let all terminals in each of the two subtrees of the root transmit in their corresponding slots. If any of the two slots contains a collision, then the algorithm proceeds to the root of the subtree corresponding to the collision and repeats itself. This continues until all the leaves are separated into sets such that each of them contains at most one packet. This is known to all users, as the outcome of the channel is either a successful transmission or an idle slot. Collisions caused by the left subtree (first slot of a pair) are resolved prior to resolving collisions in the right subtree. This scheme provides a maximum throughput of 0.347 packets/slot. Clearly, a binary tree is not always optimum. If, each time a return to the root node is made, the tree is reconfigured according to the current traffic conditions, it can be shown that the optimum tree is binary everywhere except for the root node, whose optimum degree depends on traffic conditions.

The preceding four schemes have been proposed for satellite channels. All assumed fixed-size slots, and thus can be implemented in systems that have been built for synchronous TDMA. If used in systems with small propagation delay, such as ground radio, then they will perform significantly better. Due to the inherent small propagation delay in ground radio and local environments, other access modes with distributed control are also possible if all devices are in line-of-sight and within range of each other. A description of these follows.

Minislotted Alternating Priorities (MSAP)—

The MSAP technique is a "carrier-sense" version of polling with distributed control. The time axis is slotted with the slot size again equal to the maximum propagation delay (and referred to hereafter as a minislot). All users are synchronized and may start transmission only at the beginning of a minislot. Users are considered to be ordered from 1 to M . When a packet transmission ends, the alternating priorities (AP) rule assigns the channel to the same user that transmitted the last packet (say user i) if it is still busy; otherwise the channel is assigned to the next user in sequence (i.e., user $[i(\text{mod } M) + 1]$). The latter (and all other users) detects the end of transmission of user i by sensing the absence of carrier over one minislot. At this new point in time, either user $[i(\text{mod } M) + 1]$ starts transmission of a packet (which will be detected by all other users) or it is idle, in which case a minislot is lost and control of the channel is handed to the next user in sequence. The overhead at each poll in this scheme is one minislot. Scheduling rules other than AP are also possible, such as round-robin and random order.

The Assigned-Slot Listen-Before-Transmission Protocol—Time is minislotted and divided into frames, each containing an equal number of minislots (say L). To each minislot of a frame is assigned a given subset of M/L users. A user with a packet ready for

transmission in a frame can sense the channel only in its assigned minislot. If the channel is sensed idle, transmission takes place; if not, the packet is rescheduled for transmission in a future frame. Parameter M/L is adjusted according to the load placed on the channel. For high throughput, $M/L = 1$ is found to be optimum, and the scheme becomes a conflict-free one that approaches MSAP. For very low throughput, $M/L = M$ (i.e., $L = 1$) is found to be optimum; this corresponds to pure CSMA. In between the two extreme cases, intermediate values of M/L are optimum.

The URN Scheme—The time axis is divided into packet slots, and all users are synchronized. Assuming that all users know the exact number, n , of busy users, the scheme consists of giving full access right (i.e., the right to transmit with probability 1) to some number, k , of users. A successful transmission will result if there is exactly one busy user among these k . The probability of such an event is maximized when $k = [M/n]$, where $[M/n]$ denotes the integer part of M/n . Assume the system is lightly loaded (for instance $n = 1$). A large number of users are given access right (in the example $n = 1$, the number is $k = M$), but only a few and hopefully only one will make use of it (in the example $n = 1$, a successful transmission takes place). As the load increases, k decreases and the access right is gradually restricted. For the extreme case of $n = M$, $k = 1$ and the scheme converges to TDMA. One possible scheme for estimating n with good accuracy is to include a single reservation minislot at the beginning of each data slot. An idle user that turns busy sends a standard reservation message of few bits. All users are able to detect the following three events: no new busy users, one new busy user, and more than one new busy user (termed an erasure). As it is impossible with this minimal overhead to estimate the exact number of new busy users when the latter is greater than one, errors in estimation result; however analysis and simulation have shown that this error is negligible and, furthermore, that the scheme is insensitive to small perturbations in n . This last statement is even more important with respect to the robustness of the scheme, since it means that all users need not have exactly the same estimate for n . As for coordinating the selection of the k users, an effective mechanism is the use of synchronized pseudorandom generators at all users, which allow them to draw the same k pseudorandom numbers. Another mechanism, referred to as a round-robin slot sharing window mechanism, consists of having a window of size k move over the population space. When a collision occurs, the window stops and decreases in size. When there is no collision, the tail of the window is advanced to the head of the previous window, and the size is again set to k as determined by n .

Distributed Control Algorithms in Local Area Networks—In addition to the random-access schemes described previously, all above algorithms are also applicable to local-area (broadcast) bus networks, as

these exhibit the required characteristics of small propagation delay and full connectivity. But in local-area communication, a slightly different topology, namely the *ring* (or loop), has also been widely considered. As described previously in connection with Fig. 7, in the ring topology messages are not broadcast but rather passed from node to node along unidirectional links, until they return to the originating node. A simple scheme suitable for a ring consists of passing the access right sequentially from node to node around the ring. (Note that in a ring, the physical locations of the nodes define a natural ordering among them.) One implementation of this scheme is exemplified by the Distributed Computing System's network where an 8-bit *control token* is passed sequentially around the ring. Any node with a ready message may, upon receiving the control token, remove the token from the ring, send the message, and then pass on the control token. Another implementation consists of providing a number of *message slots* that are continuously transmitted around the ring. A message slot may be empty or full; a node with a ready message waits to see an empty slot pass by, marks it as full, and uses it to send its message. A still different strategy is known as the *register insertion* technique. Here, a message to be transmitted is first loaded into a shift register. If the ring is idle, the shift register is just transmitted. If not, the register is inserted into the network loop at the next point separating two adjacent messages; the message to be sent is shifted out onto the ring while an incoming message is shifted into the register. The shift register can be removed from the network loop when the transmitted message has returned to it.

Priority-Oriented Demand Assignment (PODA)—In the context of a satellite channel, PODA has been proposed as the ultimate scheme that attempts to incorporate all the properties and advantages seen in many of the previous schemes. It has provision for both implicit and explicit reservations, thus accommodating both stream and packet-type traffic. It may also integrate the use of both centralized and distributed control techniques, thus achieving a high level of robustness. Channel time is divided into two basic subframes, an information subframe and a control subframe. The information subframe contains scheduled packets and packet streams that also contain, piggybacked, control information such as reservations and acknowledgements. The control subframe is used exclusively to send reservations that cannot be sent in the information subframe in a timely manner. In order to achieve integration of centralized and distributed assignments, the information subframe is further divided into two sections, one for each type. Access to the control subframe (which is divided into slots accommodating fixed-size control packets) can take any form that is suitable to the environment. It can be by *fixed assignment* (TDMA) if the number of stations is small (giving rise to the so-called FPODA), or by *contention* as in ALOHA if the stations have a low duty cycle (giving

rise to CPODA), or a combination of both. The boundary between the control subframe and the information subframe is not fixed, but varies with the demand placed on the channel. As in the FIFO and RR reservation schemes, distributed control is achieved by having all stations involved in this type of control keep track of their queue length information. Priority scheduling can thus be achieved. For stream traffic, a reservation is made only once and is retained by each station in a stream queue. Centralized assignment may be used when delay is not the crucial element. This scheme has been proposed in the context of a satellite channel but may be applied to other environments as well.

THE NETWORK LAYER

The network layer consists of those functions that control the transportation of data from source-host to destination-host. It serves directly transport entities residing at the network hosts, relieving them from any concern about network issues such as switching, routing, and congestion control. The network-layer functions are implemented at all switching nodes of a network. This layer is the highest one that resides at the switching nodes. It makes use of the data-link layer to accomplish the error-free transmission of data over individual links.

The Network Services

The nature of the services provided by the network layer varies considerably depending on the switching technique used in the communication subnet and the transport-layer requirements pertaining to the delivery of data from one end to the other.

A network-layer protocol for a circuit-switched network is found in CCITT Recommendation X.21 (1972). As stated earlier, X.21 is a general-purpose interface between DTE and DCE for synchronous operation on public data networks. When circuit-switched services are provided, X.21 includes a data-link-layer function and a network-layer function needed for call establishment. The data-link layer is character oriented and includes only the minimum elements necessary for basic operation, namely character synchronization (using two or more SYN characters) and error detection (using odd parity). The network layer clearly defines the procedures used in processing the various phases of call requests, incoming calls, facility requests, call progress, and call clearing. The reason these procedures belong to the network layer as opposed to the physical layer (although the end result is a direct physical connection between the two end hosts) is that the network layer receives the remote DTE (or DCE) address from the transport layer and makes use of the data-link layer in processing the calls.

With packet switched networks, two types of network services exist: the *datagram service* and the *virtual*

circuit (VC) service. In the VC model, the network layer provides the transport layer with a "perfect" connection: no errors, no duplicates, and all packets are delivered in order. In the datagram model, the network layer accepts messages from the transport layer and simply makes the best effort to deliver them independently (and not necessarily in order). The implementation of the datagram mode is simple; it merely consists of a routing algorithm that attempts delivery of the messages to their destination. For VC, in addition to routing the messages, error control and sequencing must be implemented at the end nodes. In theory, packets may travel on different routes and arrive in any order. Resequencing at the destination node would then be required. In practice, the implementation of a VC in the subnet is by establishing a route between the source and destination end nodes at connection time, and by continuously using the same route for all packets belonging to the same virtual circuit. The implementation requires the packet to carry a virtual circuit number, and each node to contain a table with an entry for each VC traversing it, relating incoming packets from an adjacent neighbor with a VC number to an outgoing link and a VC number on that link. Forwarding packets to the destination node is then straightforward. All nodes use the first-come-first-served service discipline, thus preserving the order of packets for each VC. End-to-end reliability and sequencing are achieved by means of a window mechanism, which also provides flow control. It is to be noted, however, that the network layer does not achieve complete end-host to end-host reliability, since it is subject to node and link failures; depending on the environment and on the application, higher level functions (at the transport layer) must exist to guarantee that reliability.

Routing

In multihop store-and-forward networks, the network layer includes a routing algorithm that is responsible for deciding on which output link a packet should be transmitted. Although one primary objective is that each packet reaches its destination, there are several other objectives that are also very important, such as to minimize packet transit times, to avoid congestion and deadlocks, to maximize the network throughput, etc. It is also desirable that the algorithm be simple, robust, stable, and fair to all users. There is a broad spectrum of routing algorithms. They vary according to various attributes pertaining to the nature of decision-making (centralized or distributed), the degree of adaptivity, the frequency of updates, etc. Several routing techniques are described below.

Directory Routing—In directory routing, each node maintains a table with one row for each destination. The row gives one or several outgoing links together with relative weights assigned to them. Upon receipt of a packet with a given destination address, the node simply performs a table look-up and chooses one

of the alternatives, using the relative weights as probabilities. The selection of routes and their weights may be based on the number of hops. If the source-destination traffic requirement of the network is stationary, then it is possible to use routes that minimize the average message delay in the network. For a given destination, the routes from an intermediate node to that destination are entirely determined by the tables and independent of the source. This is not restrictive, since if node j is on the optimal path from source i to destination k , then the optimal path from j to k should follow the same route. This is known as the optimality principle. As a result, the set of routes from all sources to a given destination form a tree with the destination as a root; such a tree is called the *sink tree* for that destination.

Hierarchical Routing—If the size of the network is large, then hierarchical routing is used. The network is partitioned into regions; node addresses are hierarchical and contain a region number and a node number within the region. In the table at each node, there is an entry for each destination in the region in which the node is located, and an entry for each of the other regions. Hierarchical routing decreases the overhead incurred in terms of storage and processing requirement. If the network is very large, a hierarchy with more than two levels may be needed.

Static Versus Dynamic Routing—Static routing refers to the case in which the table content is fixed. Static routing is adequate if the topology and traffic conditions do not change much. Dynamic routing refers to the case in which table contents change as the network condition changes. This is also referred to as *adaptive routing*. For example, if routing tables are based on minimizing message delay, then the routes are modified as the traffic pattern changes. There is a wide range of adaptivity depending on the frequency of changes, the type and amount of information used, and the means for implementing the changes. For example, static routing may be used, but changed only when there is a failure. To construct the best routing tables in the nodes at all times, information is needed about the instantaneous state of the network and its traffic. Unfortunately, it is not possible for the nodes to have complete and up-to-date information about the entire network. To provide it would also constitute too great an overhead. Several practical alternatives are presented below.

Centralized Routing—In centralized routing, a node is designated as the *routing control center (RCC)*. Each node periodically sends status information to the RCC. The RCC thus acquires global information, based upon which it computes optimal routes. New routing tables are periodically distributed to the nodes in the network. While centralized routing may achieve global optimal and relieve the node from the task of routing computation, it has some drawbacks: the information collected at the RCC may be old due to the delay in the

network; the communication overhead incurred in collecting status information and distributing routing information may be substantial; the reliability of the entire network rests on the proper operation of the control center.

Isolated Routing—Isolated routing is the most extreme case of decentralized routing. In isolated routing, each node makes its own routing decision based on information it has at hand. The *hot potato* technique consists of passing the packet on as quickly as possible, by sending (or queueing) it on the outgoing link with the shortest queue. Variations of this shortest-queue routing are obtained by applying various biases. For example, the link selected may be determined by a combination of the weights assigned to the static alternate routes and the queue size at each. Another isolated technique is the *backward learning* technique. It consists of having a node attempt to estimate the number of hops (or delay) of a route going from it to some destination (starting with an outgoing link) by measuring the number of hops (or delay) incurred by packets arriving from that destination on that route (i.e., on that outgoing link). To implement this technique, using the delay measure for example, each packet is time-stamped when it sets off on its journey, and from this time-stamp each node compiles a table of information about delays. One main problem with that implementation is that the delays measured are incurred by packets traveling in the direction opposite to that of concern.

Delta Routing (or Hybrid Routing)—The delta routing algorithm consists of using both central and local decisions. Using information periodically sent to it by the nodes, the RCC computes the k best paths for each pair of nodes, where only paths that differ in the initial link are considered. The RCC then sends to each node all equivalent paths (i.e., those with cost or delay differing by less than some number δ) for each of its possible destinations. In routing a packet, the node may choose any of the equivalent paths either at random, or by choosing the line with the smallest current cost (or delay). By adjustment of k and δ , the scheme can be made more or less centralized. Transpac, the French public packet switching network, uses delta routing.

Distributed Algorithms—In this class of algorithms, the nodes exchange information about delays by sending control messages to one another. To keep the overhead low, this information is exchanged only among adjacent nodes. Each node communicates to its neighbors its estimate of the minimum delay to every other node of the network. When receiving such estimates from its neighbors, a node adds to them its own delay to reach each of the neighbors, and selects the best outgoing link for each destination. Information exchange may take place either periodically at regular intervals or asynchronously, such as when the estimates change by more than some amount. The old ARPANET routing algorithm was of this type.

Session Routing and Logical Circuit Routing—In the routing algorithms discussed above, the routing decision is made for different packets independently. In session routing, the route is chosen when a session is established. All packets in the session go on the same path. In logical circuit routing, the route is chosen by means of a route set-up packet when the virtual circuit is established. The route setup packet finds its way to the destination using any of the schemes described above. All packets belonging to a virtual circuit are then transmitted on the same route. Packets must only carry the virtual circuit number. Routing in the subnet is implemented as described in the subsection headed "The Network Services." Different logical circuits between the same pair of hosts may take different routes.

Broadcast Routing—Broadcast routing refers to those techniques by which to deliver a packet originating at some source to all possible destinations in the network. Of course, this can be accomplished by sending multiple copies of the same packet, one for each destination. However, more efficient techniques exist. In the *multidestination routing* scheme, a packet is issued with a list of destinations. At an intermediate node, a copy of the packet is sent out on an outgoing line if the latter is the best route for some destination (i.e., on the sink tree for that destination); the copy will then contain the list of destinations that are to use that line. After some number of hops, each copy will contain a single address and is treated as a normal packet. A more efficient technique makes use of the sink tree associated with the source of broadcast. Assuming knowledge of the spanning sink tree in both directions, each node will broadcast the packet on all links belonging to the spanning tree except the one on which it arrived. Normally, the spanning sink tree is known only in the direction toward the sink. In that case, an approximation of the above algorithm is as follows. Copies arriving at some node on a link belonging to the spanning sink tree are repeated on all links except the one on which they came. Copies arriving on all other links are discarded. This algorithm is called *reverse path forwarding*.

Congestion Control

Congestion in an uncontrolled network is inevitable due to the fact that all resources (line capacities, buffer space, and processing capability) are limited. It is also often due to the protocols in use, such as, for example, the need to retransmit packets that are in error, the need for sequencing, etc. Some types of congestion may be relieved by the routing element that, if made dynamic, would attempt to route traffic on underutilized paths. But unfortunately, routing is not sufficient to prevent congestion altogether; it merely helps reduce it or delay it. Other flow control procedures are needed to prevent congestion. Many of these functions are present in the network layer and thus are presented here. There are,

however, other flow control procedures that reside at other layers of the hierarchy of protocols. Fig. 17 illustrates the various levels at which flow control is exercised. As will be clear from the following discussion, flow control helps prevent loss of efficiency, deadlocks, and unfair allocation of the resources.

Hop-Level Flow Control—Channel-Queue-Limit Flow Control: The *direct store-and-forward deadlock* occurs in nonlossy networks when all buffers in some node A are destined to an adjacent node B and vice versa. Channel-queue-limit flow-control techniques are used to prevent such deadlocks. They consist of partitioning the buffers available at a node among the various outgoing channels. Complete partitioning is such a scheme that does not allow any sharing of buffer among the channels. Sharing with maximum queues, with minimum queues, or with maximum and minimum queues are variants. In addition to preventing deadlocks, these techniques achieve more efficient and fairer operation of the network by avoiding hogging.

Structured Buffer Pool Flow Control: Given that a channel-queue-limit flow control is implemented, thus preventing the direct store-and-forward deadlock, another type of deadlock, called the *indirect store-and-forward deadlock*, might still occur. Consider a closed chain of consecutively adjacent nodes $N_0, N_1, \dots, N_{n-1}, N_0$. Assume all buffers on the outgoing queue from N_i to N_{i-1} are occupied by packets destined to N_{i+2} , for $i = 0, 1, \dots, n-1$ (addition on indices is modulo n). This situation clearly leads to a deadlock. A solution to this problem is given by the structured buffer pool technique. Packets arriving at an outgoing link queue are classified according to the number of hops they have already traveled. With N nodes in the network, the maximum number of hops is $N - 1$. The buffers available for that queue are allocated to the various classes (either in a fixed fashion or adjustable as needed). A packet of class k can use all buffers available to classes 1 through k . It has been shown that such a technique avoids completely cycles such as those described above that lead to the indirect S/F deadlock.

Hop-Level Flow Control With Virtual Circuits: The above techniques are used with both datagram and virtual circuit modes. In virtual circuit modes, it is also possible to exercise selective flow control on individual

virtual circuit streams. An example consists of setting a maximum limit, M , on the number of packets for each VC in transit at each intermediate node. Limit M may be fixed or dynamically adjusted. This scheme is effective in slowing down VCs that directly feed into congested areas; by backpressure the control is exercised at the source, which in response reduces its input. In the absence of this control, congestion is bound to spread to other areas of the network and affects other sources originally not responsible for the congestion.

Entry-to-Exit Flow Control—This level of flow control is primarily aimed at preventing buffer congestion at the exit node that may be due to conditions of overload of local lines connecting the exit node to the host. If the exit node must reassemble the messages before handing them to the host, then the problem is further complicated by reassembly and resequence deadlocks. An example illustrating a reassembly buffer deadlock is shown in Fig. 18. In this figure, $A_i, i = 1, 2, 3, 4$ are packets belonging to the same message. Similarly, $B_j, j = 1, 2, 3, 4$ and $C_k, k = 1, 2, 3$ represent packets of two other messages. Node 3 is the exit node. It is currently assembling message $A_1A_2A_3A_4$. Packet A_2 is missing. The deadlock occurs because packet A_2 cannot make it through node 2, where the buffers are full. An example illustrating resequence deadlocks is given in Fig. 19. Here A, B, \dots, K are consecutive packets. The exit node (node 3) must deliver packets to host 1 in sequence. The deadlock takes place because node 3 cannot deliver any of its packets, and packet A cannot reach node 3 via node 2.

The most common solution to reassembly deadlocks is to reserve a reassembly buffer for each message entering the network. The solution to the sequence deadlock is simply to discard out-of-sequence messages at the destination. The ARPANET applies this technique. A logical pipe is assumed to exist for all messages from some host A to some host B . Each pipe is individually flow controlled by a window mechanism of size $W = 8$. Messages arriving out-of-range are discarded. A request for next message (RFNM) is issued by the destination node to the source node to permit the source node to issue the next message. Buffer allocation for message reassembly at the destination is

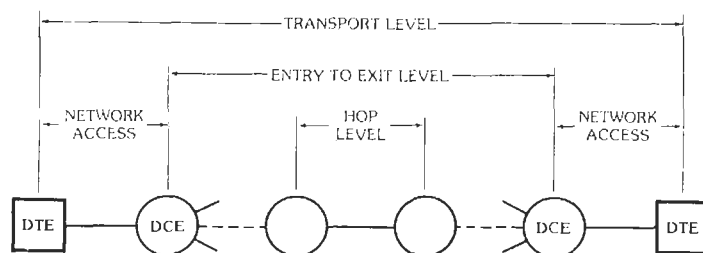


Fig. 17. Flow control levels. (From Paul E. Green, Jr., *Computer Network Architectures and Protocols*, © 1982. Used by permission of Plenum Publishing Corp., New York.)

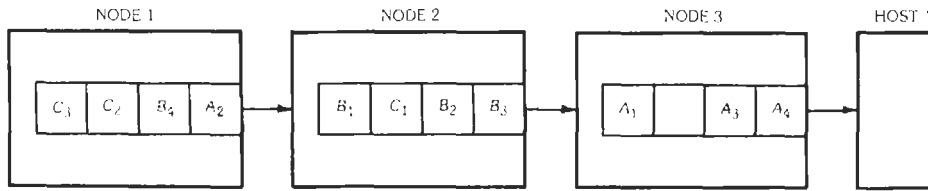


Fig. 18. Reassembly buffer deadlock. (From Paul E. Green, Jr., Computer Network Architectures and Protocols, © 1982. Used by permission of Plenum Publishing Corp., New York.)

guaranteed prior to message transmission by the source node. In GMDNET, a network based on virtual circuits, entry-to-exit flow control is exercised individually on each VC. A window mechanism is used for each VC with the window size large enough to attain efficient utilization. However, the window size is variable: it is reduced if the destination is slow in accepting packets; it is increased otherwise.

Network Access Flow Control—Network access flow control is aimed at throttling the external input to the network with a view toward preventing overall internal network congestion. This control is based on measurements of internal network congestion, where the measures collected may be local (e.g., buffer occupancy at the entry node), global (e.g., total number of empty buffers available in the network), or selective (e.g., some measure of congestion on the path leading to the destination). Three techniques are described hereafter.

The Isarithmic Scheme: This scheme consists of setting an upper limit on the number of packets that can circulate in the network at any one time. It is based on the concept of *permits*. A number of permits are present in the network. Each packet offered to the network must secure a permit at the entry node before being admitted. Once a permit is so secured, it remains in use until the corresponding packet is delivered to its destination and hence exits from the network. The permit becoming free is added to the pool of permits at the node where it is freed. This scheme as described functions adequately if the traffic pattern is uniform and balanced. Otherwise, one must avoid having a large number of permits accumulate in certain parts of the network at the expense of other parts. This is accomplished by setting limits on the number of permits at each node. The

question of how to distribute the permits in the network then arises. The performance of the network under this scheme is sensitive to the permit distribution algorithm in use. A few such algorithms have been investigated analytically and by simulation.

The Input Buffer Limit Scheme: At each node, a distinction is made between (external) input traffic and transit traffic. The input is throttled based on buffer occupancy data at the entry node. This scheme rests on the observation that congestion at the entry node is a good indication of global congestion due to the back-pressure effect propagating internal network congestion back to the source nodes. Different versions exist. The structured buffer pool technique discussed previously falls into this category. Input traffic (which is in the 0-hop class) is discarded when all class-zero buffers are occupied. The number of class-zero buffers has been shown to have a great effect on the network performance under heavy load. Another version defines a limit, N_i , smaller than the total number of buffers, N , and new packets are limited to N_i while transit packets have no limitation imposed upon them. The performance here is sensitive to the ratio N_i/N . Another version yet consists of discarding an input packet if the total number of packets in the entry node exceeds a given threshold.

Choke Packet Scheme: In this scheme, a link is said to be congested if its utilization over some history window is beyond some threshold. A path is said to be congested if any of its links is congested. Along with routing information, congestion information is propagated in the network. When a packet is received by a node for some destination whose path is congested, then (1) if the packet is new external input, it is dropped; (2) otherwise, it is forwarded and a choke packet is sent back to the source, slowing the source. The path to the destination is unblocked gradually as no choke packets are received during a specified period of time.

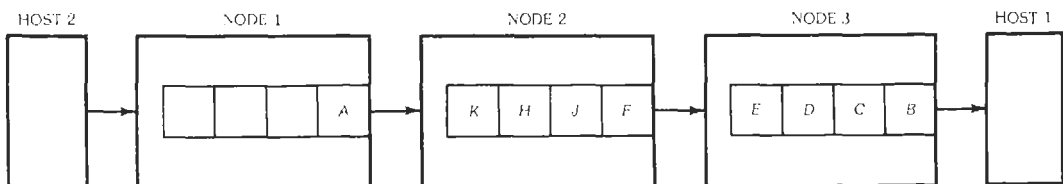


Fig. 19. Resequencing deadlock. (From Paul E. Green, Jr., Computer Network Architectures and Protocols, © 1982. Used by permission of Plenum Publishing Corp., New York.)

CCITT Recommendation X.25

Recommendation X.25 is a widely accepted standard entitled "Interface Between a DTE and a DCE for Terminals Operating in the Packet Mode on Public Data Networks." It is a set of three peer protocols: a physical level, which is the same as in Recommendation X.21 (discussed in the section "The Physical Layer"); a link control level, which is essentially HDLC; and a network layer, which is discussed in the present subsection.

Recommendation X.25 is based on the concept of virtual circuits. These virtual circuits can be either temporary or permanent. A temporary virtual circuit referred to as a virtual call requires three phases: setting up the call, the data-transfer phase, and disconnecting the call. In setting up a virtual call, a logical channel number is assigned to the call. With a permanent virtual call, the numbers are assigned when the customer leases the facility from the common carrier. In such a case, the two phases of setting up and disconnecting the circuit do not exist.

It is important to note that Recommendation X.25 describes only the DTE/DCE interface. It is not concerned about how the network operates and how virtual circuits within the network are implemented. Recommendation X.25 defines the services rendered to the transport layer of the ISO model that resides in the user's equipment. Thus X.25 is known as a network access protocol standard. In many cases, the DTE is a terminal that handles characters rather than packets and that uses the start-stop asynchronous DLC protocol. In such cases, common carriers operating X.25 networks provide interface machines to which these terminals are connected. These interface machines, called PAD (Packet Assembly/Disassembly), are extensions to X.25 but not part of it.

The network layer of X.25 defines explicitly the formats of packets and the operational phases of a virtual circuit connection. Two types of packets exist: data packets and control packets. A data packet has a variable-length information field to carry user data. A control packet can be of different types, many of which are discussed below.

When a DTE sets up a virtual call, it selects a free logical channel number from the set allocated to it. Such a number consists of two parts: a logical channel group number and a logical channel number. A DTE has available to it up to 15 logical channel groups (addressed with 4 bits) and up to 255 logical channels within each group (addressed with 8 bits). Using a CALL REQUEST packet, the DTE sends the selected logical number along with its own address, the address of the destination, and other facility codes to the DCE to which it is attached. (For efficiency purposes, the CALL REQUEST packet may also carry user data up to 16 bytes.) The request to establish a connection arrives at the remote DTE from its attached DCE as an INCOMING CALL packet. The called DTE decides whether to accept or refuse the call. If it accepts it, then a CALL ACCEPTED packet is sent back to the DCE.

The indication of acceptance of the call is achieved via a CALL CONNECTED packet, which is sent by the local DCE to its calling DTE. The disconnection of a call follows similar procedures. It uses the control packets CLEAR REQUEST, CLEAR CONFIRMATION, and CLEAR INDICATION. Once a connection is established, the flow of data packets can take place. Data packets are sequenced, and the flow is regulated according to a window mechanism as described in the section "The Data Link Control Layer." As with ADCCP and HDLC, the two sequence ranges of 8 and 128 are available. The control packets RECEIVE READY and RECEIVE NOT READY, similar to those described for ADCCP and HDLC, are used in a supervisory mode to acknowledge data packets. Other packets such as INTERRUPT, RESET, RESTART, and their confirmations are used for further control and to recover from certain minor and major problems. The INTERRUPT packet contains only one byte of user data and has no sequence number. It is transmitted as quickly as possible, preempting the queues of normal data packets. It is delivered (out of sequence) even when data packets are not being accepted. It is normally used to convey the fact that a terminal user has hit the break key to stop a flow of data from a distant computer. An INTERRUPT packet is acknowledged by the local DCE with an INTERRUPT CONFIRMATION packet. The RESET commands are used to reinitialize the window parameters corresponding to a particular virtual circuit to zero. Both DTEs and DCEs can initiate a reset. Reset packets contain up to two bytes of data explaining the reasons for the reset. These may be: the remote DTE is out of order, the network is congested or one of its nodes had failed, etc. Following a reset, the DTE has no knowledge about the status of outstanding (unacknowledged) packets. Recovery must be accomplished by higher levels. The *restart* condition provides a mechanism to recover from major failures. It has the effect of clearing all the virtual calls that a DTE has connected and resetting the permanent virtual calls, bringing the user/network interface to the state it was in when the service was initiated.

Since X.25 is an interface between a DTE and the local DCE to which it is attached, the receipt of an acknowledgement by a DTE from the DCE means that the latter has received the packet, but does not imply that the remote DTE has received it. To provide some level of end-to-end acknowledgement, a bit called the delivery confirmation bit, or *D-bit*, has been provided in the packet format. The DTE sets *D* to 1 to request confirmation of delivery to the remote DTE. An acknowledgement for this packet by the DCE would then have to be a guarantee of its delivery to the remote DTE.

In addition to the standard features of X.25, there are also some optional features that may be requested by a user when a call is set up. Examples are: reverse charging, closed user group (where members of a closed group can communicate only with one another), packet retransmission (where a DTE can ask its DCE to retransmit one or several data packets, and this is

accomplished by the use of the REJECT command as in ADCCP), etc.

HIGHER-LEVEL PROTOCOLS

This section is concerned with layers four through seven of the ISO model. The brevity of this section is due to three factors: (1) many of the concepts and techniques needed at these higher levels, such as sequencing, acknowledging, and flow control, are similar to those used in lower layers and have already been described and discussed in previous sections; (2) many of the functions residing at these layers have not yet been well defined, and there are no internationally accepted standards yet; and (3) some of the functions (especially concerning the presentation and application layers) are not strictly speaking communication functions, and thus are considered out of the scope of this chapter.

The Transport and Session Layers

Internally to a host computer, application programs are represented by processes. Communication between such processes across a network is accomplished by means of "reliable connections" called *sessions* between specific pairs of processes, and it is the role of the session layer to establish and maintain such connections. It is possible to imagine that each such connection is achieved as a virtual circuit that the network layer provides. In that case, the session layer is trivial, the transport layer is nonexistent, and most of the work is done by the network layer. However, there are many problems with that design. The first problem arises when the network layer provides only a datagram service: the session layer will then have to be constructed appropriately so as to overcome the deficiencies of the network. In such a design, the session layer will also have to interface to the networks differently, depending on their type. It will have to be concerned with managing the buffers at the host where the processes reside, guaranteeing end-to-end reliability between the end hosts (which the network does not guarantee even when it is providing the virtual circuit service, since crashes, resets, and restarts can occur), and controlling the flow. As the protocols for many of these functions again differ depending on the network, it is best to render the session layer independent from network issues by the introduction of the transport layer. Thus, while the transport layer provides truly reliable host-to-host communication, the session layer connects two processes together in a session without worrying about any of the implementation details of the actual network. It merely requests from the transport layer its services. The functions of a session layer are thus considerably reduced. At session setup time, it establishes the conventions necessary for the connection, conventions relating to the data-transfer mode (half-duplex versus full-duplex), the character codes, the flow control

window sizes, how to recover from failures at the transport layer, etc. It also allows processes to refer to destinations using symbolic names by doing itself the mapping onto transport addresses.

The transport layer within a host defines a set of *transport addresses* or *sockets* through which connections are established. A transport address consists of a *network* number, a *host* number, and a *port* number assigned by the host. A CCITT numbering plan exists which uses 14-digit numbers. The first three identify the country; the fourth identifies the network within the country; the last ten digits are assigned by each network operator, some to indicate hosts, the remaining ones to indicate ports. There is a distinct difference between a transport address and a network layer address (which defines a virtual circuit), and the mapping need not be one-to-one. The transport layer uses network addresses to create transport connections. If the mapping is not one-to-one, then *multiplexing* is said to be in effect. In *upward multiplexing*, several transport connections are multiplexed onto the same network virtual circuit connection. In *downward multiplexing*, a single transport connection uses several virtual circuits. The choice of one type of multiplexing or the other is a function of the charging procedure used in the public network, the traffic volume exercised by a transport connection, and the window size of virtual circuits within the network.

The transport layer functions to establish and close connections as well as to control the flow across each connection. The techniques are similar to those described for lower levels. However, due to the ultimate end-to-end reliability that is to be guaranteed, several problems need to be solved. First of all, it is important to guarantee that old packets that for some reason have been delayed in the subnet but whose sequence numbers fall within the current receive window of the destination are not delivered to the destination. This is done by limiting the amount of time that a packet can exist in the network and by using a sequence space so large (e.g., 32 bits) that no packet can live for a complete cycle. For this method to be successful, however, it is important to have each connection initialize its sequence numbers to a value higher than the previous end sequence number by some margin so as to guarantee it is higher than any existing packet number. The second problem to deal with is to guarantee synchronization between the source host and destination host with regard to their respective initial sequence numbers. The solution to this problem is the *three-way handshake* mechanism. It consists of the following. The sender sends a CALL REQUEST packet with initial sequence number *S*. The destination responds with a CALL ACCEPTED with a sequence number *A* along with the source sequence number, *S*. Upon receipt of this packet, the source sends a packet with *S* and *A*. Thus, only when source and destination have confirmation of the combined sequence numbers would they consider the connection to have been established. The control of flow over a connection between source and destination is achieved by using a window mechanism similar to that of lower levels. The

window sizes must be chosen appropriately. To be effective, it is important in some cases to distinguish acknowledgements from the destination confirming correct reception of packets from so-called *credits*, which are sent by the destination to inform the sender that receive buffers are available. Acknowledgement of a packet allows the sender to free its own buffer, while receipt of a credit allows the sender to advance its send window.

The Presentation Layer

The presentation layer performs transformation on the data generated by or destined to application programs before they are sent to or after they arrive from the session layer, respectively. The objectives and nature of such transformations are diverse. In order to save on communication bandwidth, the presentation layer compresses text and reduces the amount to be transmitted across the network. In order to guarantee the security of data, the presentation layer encrypts it at the source and decrypts at the destination. Many compression and encryption techniques are available. They are out of the scope of this chapter, and thus are not described here.

Another objective of the presentation layer is to resolve the differences that exist among the various equipments connected to the network. Advances of this nature have been made in two areas: (1) resolving the differences that exist among the various types of terminals in the networks leading to the so-called *virtual-terminal protocols*, and (2) resolving the differences that exist among the various hosts so that files stored in one host can be transferred to another host. This has led to the so-called *file transfer protocols*.

Terminals are basically of three types: *scroll-mode* terminals, which have no intelligence and generally use the start-stop protocol; *page-mode* terminals, which handle about 25 lines of 80 characters each at a time and have some local editing capabilities; and *form-mode* terminals, which are micro-based, have more intelligence, and handle specific forms that the operator fills out using local editing facilities. Besides the differences that exist among the different types of terminals, many differences exist within each class with regard to such attributes as the character set, line length, carriage-return and tab rules, cursor addressing, blinking, local editing capabilities, etc. Virtual-terminal protocols attempt to hide the differences by mapping real terminals onto a hypothetical network virtual terminal. For *scroll-mode* terminals, such a conversion is straightforward. Each time a new version of such a terminal type is to be supported, the corresponding conversion rules are added to the presentation layer at the host computer. For the more sophisticated types, a *data structure model* has proven useful and effective. Each end of the session has a data structure representing the state of the virtual terminal. Every time the application program changes the data structure on its machine, the presentation layer sends a message to the remote

machine telling it how to change its data structure, and vice versa.

While the terminal conversion problem has been adequately solved, file transfer has proven to be much more difficult, and no general solutions exist. In particular, the concept of a network standard file format is not attractive due to the great difficulties involved in converting a file from one machine format to another. There are, however, several ad hoc file-transfer protocols particular to specific networks, e.g., ARPANET's FTP (reference 2).

HIGH SPEED LOCAL AREA NETWORKS

The advances in technology and the pull of new, more demanding user applications are stimulating the development of *Local Area Networks (LANs)* at speeds much higher than those available just a few years ago. High speed LANs in the 100 Mbps range are now commercially available, and prototypes operating at speeds of 1 Gbps are being tested. In this section, we trace the evolution of the LAN technology from low to high speed, identify the technology push and the market pull, survey the novel high speed LAN products that are now coming to market, and examine some of the most recent LAN proposals and prototypes designed to operate at speeds beyond the Gbps range.

The first LAN products were introduced in the late '70s in offices, campuses, and production plants in order to interconnect with each other and to a growing number of terminals, workstations, and peripherals. Typical applications were word processing, interactive computing, data-base queries, and file transfers—a predominantly data-oriented environment with peak data transfer speeds on any user session well below 1 Mbps. Thus, LAN speeds from 1 to 10 Mbps were adequate for these applications.

In the mid '80s, the original LAN environment underwent an important series of transformations which forced LAN designers to look ahead for new network technologies. First, the proliferation of LANs on campuses and industrial plants led to the need to interconnect them—via other LANs capable of carrying aggregate traffic several orders of magnitude higher than the traffic found in an individual LAN. Thus, speeds of up to 100 Mbps and geographical coverage of several miles are now required. In addition to the data applications, new applications and services start now to emerge: voice, video, imaging. For example, in an office we desire to combine the telephone and the data terminal in the same multimedia workstation, and to carry both voice and data on the same cable. In a campus environment, we wish to distribute video lectures and to support video conferencing. In a hospital, we wish to permit the radiologist to retrieve X rays from an image data base.

New high-speed LAN requirements emerge also in supercomputer labs. An important application here is "scientific visualization," i.e., the ability to visualize,

through a high-definition graphics screen, a physical phenomenon (e.g. the fluid dynamic turbulence of a jet engine) which is simulated in the supercomputer. Currently, most of the visualization set-ups are supported by high-speed dedicated point-to-point connections between the supercomputer and the graphics workstation. Data rates of up to 100 Mbps are observed on such connections. Specialized, high-speed "supercomputer" LANs are now proposed to reduce the cost and increase the flexibility of such interconnections, by sharing a common network. Many other applications besides visualization will benefit from the supercomputer LAN; examples are memory-to-memory transfers, distributed processing and load sharing among supercomputers, etc.

Advances in light-wave technology have greatly influenced the development of high-speed LANs. Performance and cost are the principal advantages of light-wave technology, as compared to copper-based or radio-frequency technologies. The improvement in light-wave technology has been rapid and steady in the areas of optical fibers, light sources, light detectors, couplers, amplifiers, connectors, switches, and signal processors. Advanced research in other areas of light-wave technology—such as coherent processing, WDM, and photonic switching—is also yielding new results that can be exploited in constructing even higher-performance and more-cost-effective networks.

Whereas today's high-speed LANs are based on simple light-wave technology, the LANs of the future will rely on more advanced light-wave principles to achieve higher performance and increased functionality. High-speed LANs such as FDDI and DQDB employ direct-detection, LED-based transmission over multimode optical fibers, which can easily accommodate data rates in the range of about 100 Mbps. Using LDs to transmit over single-mode optical fibers, higher-speed networks can reach speeds in the range of about 1 Gbps. Research proposals for future LANs and MANs are often postulated on the emergence of WDM, which would allow the easy implementation of multichannel networks. Such proposals therefore depend on the successful development of the coherent techniques that multichannel transmission requires, especially in the areas of modulation, heterodyne and homodyne detection, and the tunability of lasers and detectors. More advanced research proposals might even assume that all-optical transmission is to be used, in which optical signals are switched by purely photonic means, thereby eliminating undesirable electro-optical conversions within the network. Such advanced proposals thus rely on technological progress in areas such as photonic switching and passive optical amplification.

Furthermore, current LANs based on light-wave technology use simple topologies to distribute optical signals essentially in a point-to-point manner. These LANs regenerate the optical signal by means of electro-optical conversions at every network interface, which eases considerably the difficulty of meeting the link budget. Link distances are usually constrained to only a

few kilometers, and network interfaces are typically connected to or isolated from the optical fiber through an optical bypass switch. However, such simple point-to-point topologies limit the designer's flexibility and make the implementation of multichannel networks problematic. By making effective use of optical couplers and broadband optical amplifiers, one can design LANs and MANs with nonpoint-to-point topologies, e.g. trees and stars. Tree-Net, for instance, can be used as a physical tree topology for distributing optical signals to a large number of users over a wide geographical area, using only passive optical couplers and optical amplification at the root of the tree. The commercialization of a broadband optical amplifier would be necessary before tree-like multichannel networks become a reality. All-optical networks could be constructed as arbitrary mesh-connected topologies using linear combiners/dividers or wavelength-selective switches. These devices are not currently available, but active research into their design is under way.

Technological progress in light-wave transmission must be accompanied by progress in digital electronics. High-speed data transmission relies on fast electronics that can provide serial data to a transmitter at network rates. The speed of digital electronic circuits must therefore grow as rapidly as the speed of light-wave transmission systems. Given that the bandwidth of a single-mode optical fiber is in the Tbps range, it is clear that WDM must be used to allow a serial electronic interface to tap into this enormous bandwidth. In order to transmit at Gbps speeds, electronic transceivers must be clocked at GHz rates, which suggests the use of high-performance integrated circuitry such as GaAs or ECL. Since the processing of data-link-layer protocols is normally hardware-assisted in high-speed LANs, it is essential to implement state machines as integrated circuits that can keep up with the transmission system. Some LANs, such as HPPI, employ complex, large switching fabrics to operate at very high speeds, which also implies the need for technological advances in integrated-circuit technology.

A Taxonomy of Local Area Networks

In the past LANs have been classified according to several different criteria (topology, access protocol, transmission medium, etc). For the purposes of this section, which is focused on high-speed LANs, the most appropriate criterion for classification is data rate. Thus, we define the following four classes:

- (1) Low- and Medium-Speed LANs: Data rate up to 10 to 20 Mbps. This class includes the great majority of conventional LANs (Ethernet, Token ring, etc.), as well as the recently introduced wireless LANs.
- (2) High-Speed LANs: Data rates ranging from 50 to 150 Mbps. Typical entries are FDDI and DQDB. The most common medium is the fiber.

- (3) Supercomputer LANs: Data rate in the order of 800 Mbps. Typical entries are UltraNet and Nectar. Due to the limited geographical range, the most common transmission medium is the cable (32 or 64 parallel connections).
- (4) Ultragigabit LANs: Data rates well above the Gbps range. Although no commercial offerings exist in this area, several experimental prototypes are being developed (e.g., ATT Multihop Network). The medium is the single-mode fiber. The technology of choice is Wavelength Division Multiplexing (WDM).

Low- and Medium-Speed LANs

The appearance of Ethernet in the late 70s marks the birth of first generation LANs. Within the next few years, the Token ring and Token bus LAN architectures also become well established, and are all eventually standardized by the IEEE 802 committee.† These first generation LANs have comparable characteristics in terms of speed (around 5–10 Mbps), geographical range (around 1000 feet), and station connectivity (a few dozen). They offer, however, markedly different features so as to gain a competitive advantage in different applications. For example, Ethernet prevails in office and campus environments because of its low cost, good fault tolerance, and easy maintenance. Token ring and Token bus offer guaranteed delay and therefore better real-time traffic support. In particular, the Token bus was chosen as the MAC layer for MAP (Manufacturing Automation Protocol) because of its suitability for carrying the real-time control messages common in manufacturing plants.

A detailed survey of first generation protocols is clearly beyond the scope of this section; the interested reader can consult the excellent survey by Tobagi in reference 5. Rather, we would like to comment on the ability to scale up the speed of these LANs by one or two orders of magnitude, from 10 Mbps to the 100–1000 Mbps range. At the same time, we explore the possibility of stretching the geographical coverage from the 1-km range to the 10-km range. These speed and range upgrades are of interest because of the broadening of LAN applications in the early 80s.

If we denote by τ the end-to-end propagation delay on the bus and by P the frame transmission time, we find that Ethernet performance degrades dramatically when the ratio $a = \tau/P$ becomes 1. For channel speed = 100 Mbps, bus length = 10 miles, and frame length = 1000 bits, channel utilization becomes 0.03. The reason for this dramatic drop is that it takes more time (on the order of τ) to ensure that the transmission will not suffer a collision than to transmit the frame itself (P seconds).

The Token bus does not fare much better. The time to pass the token from one station to the next in the polling

sequence can be as large as τ plus the token processing time. Assuming that a station transmits a single frame per token visit, the channel utilization is given by the following expression:

$$\text{Token bus utilization} \leq P/(P + \tau) = 1/(1 + \tau/P)$$

Thus, again, utilization decreases as the product of bandwidth times length ($B \times L$) increases.

In contrast with Ethernet and Token bus, the Token ring does not have any intrinsic limitations imposed by the hardware. In principle, the ring can be maintained fully utilized even at very high $B \times L$ products if stations attach their frames to the token one after the other forming a "train." The higher $B \times L$, the larger is the number of frames simultaneously outstanding on the ring. However, the IEEE 802.5 standard implementation requires that only one station be allowed to maintain one (or more) frames outstanding on the ring at any one time.‡ This requirement is dictated by the particular implementation of priorities in IEEE 802.5. Namely, stations announce their priority in the header of a passing frame. Thus, the station currently holding the token must wait for the header of its frame(s) to come back before it can reissue the token at the highest announced priority. Assuming that each station transmits one single frame per token acquisition, and letting τ' be the ring latency (propagation + shift register delay), the maximum ring utilization is $\leq 1/(1 + \tau'/P)$. Again, we find that performance degrades as $B \times L$ increases.

The interested reader can verify that in general first generation LAN architectures do not scale well with speed and distance.§ This is not to say that such LANs are obsolete; they are simply designed to satisfy requirements of the 10-Mbps, 1000-foot type.

High-Speed LANs

In the early '80s, two separate developments contributed to the start of a new generation of LANs (and MANs): the rapid advancement of fiber-optic technology, and the emergence of very-high-speed local and metropolitan applications. In the former area were the fabrication of low-cost, low-attenuation fibers and of high quality optical transmitters, detectors, and passive components. In the latter area came the proliferation of private LANs and the need to interconnect them across metropolitan regions through a high-speed backbone. We must thus introduce a new version of "stretched" LAN, called the Metropolitan Area Network (MAN), which shares many common characteristics with the LAN, and therefore will be covered along with the LANs in this survey.

We also note growing interest in videoconferencing

† Reference 2.

‡ Reference 2.

§ Reference 5.

and in the integration of voice, data, and video communications in a single network. The integration is motivated both by the desire to reduce network cost (by sharing facilities) and by the need to support new applications (e.g., multimedia communications).

To satisfy these requirements, a new family of LANs is required, with data rates in the order of 100 Mbps and above, with geographical coverage on the order of 10 km, and with the ability to support a large number of stations (in the hundreds). Also, good fault recovery properties are required, since metropolitan networks span areas outside of the customer and network operator control, and are therefore more vulnerable to link failures.

The most natural way to meet the new challenge is to start from first generation LAN architectures, replacing the cable with fiber and modifying the protocols so as to remove the $B \times L$ limitations observed in previous schemes. This is exactly what happens in the case of the ring.

Rings—FDDI (Fiber Distributed Data Interface) is a 100 Mbps, multimode fiber ring which implements a "multitoken" protocol (i.e. several stations can simultaneously maintain outstanding frames on the ring). See Fig. 20. The multitoken variant removes the inefficiencies of the 802.5 token ring protocol and makes it possible to operate the ring at high speed over large distances (up to 200 km). Priorities can, of course, no longer be handled with header reservations. Thus, new priority mechanisms based on maximum token rotation time are implemented. Fault tolerance to link failures is achieved by using a double-counter rotating ring topology.

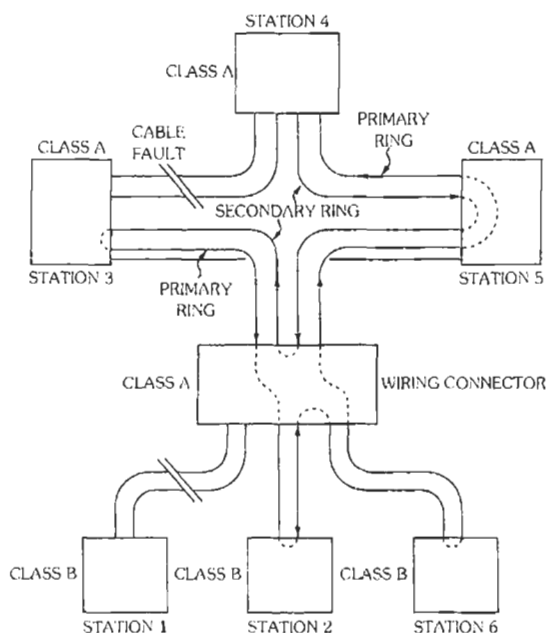


Fig. 20. FDDI topology.

gy. Isochronous traffic support is provided by a variant of FDDI (called FDDI-II). In FDDI-II the ring is slotted. Each slot corresponds to a channel that can be reserved by isochronous traffic. The asynchronous traffic is still managed by the token protocol through the use of unreserved slots.

In FDDI and FDDI-II, each frame is removed by the origin, after a full rotation around the ring. Efficiency can be improved (by as much as 4 times on average) if the destination removes the frame, thus allowing spatial reuse of portions of the ring. A recently proposed buffer insertion ring, the Meta Ring, does in fact offer destination removal and spatial reuse, upgrading the FDDI speed to 800 Mbps (using a 2-counter rotating ring).

Linear Buses—The linear bus topology was also investigated for possible fiber-optic upgrade. It was found that the fiber bus cannot be used effectively in a bidirectional mode (such as in Ethernet), because of the optical coupler implementation. Thus, two candidate configurations emerged: twin parallel unidirectional buses, and a single, unidirectional folded bus.‡ Since CSMA-CD and token bus protocols do not scale up to high speeds, a new access protocol, called "implicit" token (or "ordered access") was introduced. Basically, the implicit token is a burst of carrier issued by a designated end station. Backlogged stations append their frames in sequence after the token, forming a train. Upon completion of a cycle, the end station starts the next cycle. Good efficiency at high $B \times L$ speeds is achieved by virtue of the fact that no time is wasted in passing the token from one station to the next (in contrast to what was happening in IEEE 802.4). In fact, stations are automatically "polled" in sequence, based on their position on the bus. Furthermore, no token processing is required. One drawback of linear fiber buses, however, is fiber and coupler attenuation. This problem is corrected by regenerating and amplifying the signal at each interface. Fault tolerance to link failures is achieved by coiling the topology in a pseudo loop. Real-time (isochronous) traffic is carried in reserved slots, or in alternating cycles. A large number of implicit token fiber-bus architectures were proposed and implemented in the '80s. Among them were Express-Net, U-Net, and FasNet.‡

An alternative approach to multiple access in linear buses is reservation. This scheme differs drastically from the implicit token approach in that stations are given the opportunity to place reservations in a separate logical channel (implemented by using special fields in the data packet headers, or dedicated portions of a cycle). Thus, the order of packet transmissions is first come first served, as opposed to the round robin order in token schemes.

The best known implementation of a reservation strategy is DQDB (Distributed Queue Dual Bus), a 150-Mbps twin bus in which each station records and

‡ Reference 6.

maintains the reservations in a local queue. In ideal situations, all local queues are perfectly consistent with each other, thus yielding an exact FCFS schedule. The strong points of DQDB are simplicity and extreme efficiency of the distributed reservation mechanism. Reservations are piggybacked on data packets, thus avoiding cycle latency and providing channel efficiency close to unity for a single sending station. As in previous protocol proposals, the bus is slotted to provide isochronous traffic support. The bus is installed as a pseudoloop topology to permit recovery from single link failures (Fig. 21). DQDB has encountered rapid and broad acceptance, and has been chosen as the MAN standard by the IEEE 802.6 committee.

Trees—During the '80s, new LAN and MAN topologies such as the tree and the mesh also emerged. Trees are particularly desirable for the efficient coverage of metropolitan areas. A well known fiber-optics tree architecture is Hubnet. Hubnet implements a random-access protocol with capture (i.e., if there is a collision, one of the contentents prevails and captures the channel). Contention resolution is implemented at the nodes (which, therefore, include an active processing element).

Another tree architecture with random access and capture is CAMB (Collision Avoidance Multiple Broadcast). In CAMB, tree nodes have the ability to route frames on-the-fly within subtrees (instead of sending all frames to the root of the tree as in Hubnet). Subtree routing yields spatial reuse, and therefore better overall throughput, at the cost of additional processing in the nodes.

A general limitation of the above tree architectures is the inability to guarantee delays and, therefore, to support real-time traffic. It will be shown later that this limitation can be overcome by using tree topologies with passive nodes and token or reservation access schemes.

Grids—Spatial reuse is also the driving force behind mesh architectures. The best known mesh architecture is MSN (Manhattan Street Network), organized as a square grid. To avoid store and forward buffering and therefore improve efficiency at high speeds, a very simple routing scheme, deflection routing, is implemented at the nodes (i.e., grid crosspoints). Recall that in a grid each node has exactly two inputs and two

outputs, and consider the case of both inputs competing for the same output. One packet is sent to the correct output. The other, instead of being buffered, is sent to the "wrong" output, and is thus forced to take a "round-the-block" detour. This feature, albeit somewhat wasteful in terms of delivered throughput, makes it possible to avoid nodal buffering, thus speeding up the operation. The positive aspects of MSN are the good metropolitan coverage, the simplicity of protocols and nodal implementations, and the potentially unlimited throughput offered by a topology with multiple concurrently available paths. In this sense, MSN is a major step forward with respect to more traditional LAN and MAN architectures, where a single broadcast channel (bus, tree, or ring) inevitably becomes the bottleneck. The drawback of the mesh topology, however, is the inability to provide efficient broadcast and isochronous traffic support, properties in which the bus, tree, and ring topologies excel. It will be seen later how the path concurrency of MSN and the broadcast and isochronous support features of more conventional MANs can be combined in a single architecture, the WON (Wavelength division Optical Network).

Supercomputer Networks

Another important development in the area of high-speed local and campus-wide networking is the emergence of the networks especially designed to connect supercomputers, massive storage devices, frame buffers, and high-end graphics workstations. Typical applications are scientific visualization, load sharing among supercomputers, and, more generally, distribution of tasks among supercomputers and workstations. Very high data rates (in the order of Gbps) are required to support real-time applications and an interactive mode. A high-speed interface for such "supercomputer" networks has been proposed, namely HPPI (previously known as High Speed Channel), and is evaluated by an ANSI committee. We will mention two network architectures which have been developed to interface with HPPI and to satisfy the above needs. One is Ultranet, basically a datagram network, consisting of high performance ultrabus switches interconnected by intelligent, high throughput bridges. In contrast, the CP* network, developed by the Los Alamos National Laboratory, is a fast circuit switched network. It features crossbar switches and basically sets up a circuit for each message exchanged across the network. Both Ultranet and CP* operate at near Gbps speed. They have been designed for a specific application and therefore cannot be easily extended to satisfy the general requirements of a metropolitan network. On the other hand, a challenging question for network designers is whether the future LAN and MAN architectures will be sufficiently powerful and flexible to satisfy HPPI users.

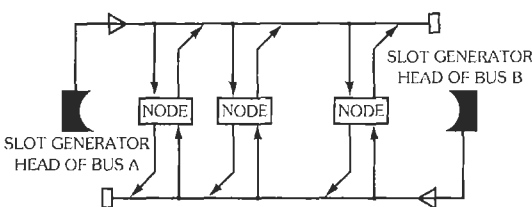


Fig. 21. DQDB topology.

Ultragigabit Networks

Many of the above architectures have been successfully demonstrated at speeds ranging between 100 and 200 Mbps, and some have been pushed up to Gbps speeds. These speeds are quite adequate for introductory LAN and MAN solutions. Traffic projections for the mid '90s, however, show that a single MAN will be required to carry simultaneously a large number of very high speed connections (e.g., scientific visualization, high-quality videoconferencing, multimedia communication, etc.) along with the support of a large population of more conventional users. The statistical aggregation of these traffic components in a large metropolitan area will probably require an overall network capacity in the order of hundreds of Gbps. Clearly, current LANs and MANs cannot deliver such performance. First, the presence of active processing elements along the network path poses a limitation on speed. In fact, intermediate station interfaces along the DQDB bus or the FDDI ring must process in real time, at line speed (as part of the access protocol) some of the incoming information (as opposed to buffering it and processing it later). For example, in DBDQ the Busy and Reservation bits must be read and set "on the fly" in the packet headers. Likewise in FDDI, the header of a packet must be inspected at channel speed while it transits through the shift register at each interface.

Another limitation common to most of the above architectures stems from the fact that loops and rings cannot be easily (and effectively) configured to cover a distributed metropolitan area. The common solution is to resort to a multilayer architecture with a backbone loop at the top of the hierarchy and several layers of distribution loops. The weak point of the multilayer configuration, however, is the bridge that connects the layers. The bridge is basically a store-and-forward processor which buffers packets, inspects the headers, and performs some simple flow control and traffic filtering functions. But, simple as these functions may be, they still pose a practical limit on the throughput that can be handled by a bridge, probably in the order of 100K packets/second (state of the art bridges today can forward at most 20K packets/second). Furthermore, it will be difficult to satisfy the tight delay constraints of voice and video connections across several bridges.

Manhattan Street Networks can potentially yield the desired throughput (by exploiting parallelism, spatial reuse, and locality) and can provide efficient metropolitan coverage. However, they cannot efficiently support broadcast and real time services. Also, they require fairly sophisticated controls in order to avoid unfairness and congestion.

Finally, none of the above architectures are well suited to exploiting WDM (Wavelength Division Multiplexing). WDM makes it possible to maintain multiple channels in a fiber and can, in principle, expand the effective data rate of a single fiber from the Gbps to the Tbps range. In order to implement WDM, an FDDI or

DQDB network requires a splitter, a filter, and a combiner at each interface, plus a wideband amplifier to make up for the power loss—a very expensive proposition!

The above considerations suggest that in order to achieve ultragigabit speeds, the network must be totally "transparent" to WDM, and thus must avoid any processing components (e.g., shift registers, busy-bit manipulators, bridges, etc.) along the communications path between source and destination. This has led to the development of a new generation of LANs and MANs, based on passive optical components and WDM techniques.

The key components of a wavelength division optical network are:

- Passive optical network, including couplers, splitters, combiners, and optical amplifiers. Electrically controllable couplers and switches are also required for some proposed architectures.
- Wavelength division gear, namely: tunable lasers, detectors (conventional and coherent), and filters. For some architectures, wavelength switches are also required.
- Network management, to control wavelength allocation and dynamically reconfigure the virtual topologies according to traffic patterns and user demands.

Several wavelength division proposals and prototypes have been reported in the literature. The most common architecture consists of a passive optic broadcast medium (star, tree, bus)—the so-called "mass of glass"—on which several wavelengths (a few dozen with conventional, coarse WDM; thousands with coherent, dense WDM) can be multiplexed. The multiplexing and switching are done at the periphery of the "mass of glass." Indeed, the main challenge is how to use this abundance of bandwidth efficiently (i.e., channel access control).

Two main approaches have been proposed. The first approach is the time/frequency division approach. A global, time-slotted reference frame is defined (as on a satellite multiple access channel). A source/destination pair must reserve in advance a slot in the frame for the transmission of each packet and must agree on the frequency (i.e., wavelength) to be used. This system can achieve high efficiency, but it involves high overhead for source/destination coordination and high cost (and technology risk) due to the quickly retunable transmitters and receivers.

The second approach assigns a limited number of fixed wavelengths to each user station (say two for transmission and two for reception). By means of these wavelengths, users can interconnect with each other with a "virtual" topology (as opposed to the "physical" topology—e.g., a star or a tree—corresponding to the fiber layout). Fig. 22 shows a virtual topology. In this case a "perfect shuffle" topology, mapped on a

physical tree. Transmissions between arbitrary source/destination pairs are accomplished via "multihop" (i.e., store and forward) through the user stations. The main advantages of the multihop solution are lower transmitter and receiver costs at each station and lower control overhead. On the negative side, however, broadcast, multicast, and real-time traffic support is not very efficient. In view of these limitations, hybrid solutions combining the benefits of the two approaches are now under investigation.

In contrast to the broadcast, fully transparent, mass-of-glass topology, other approaches propose the use of mesh type topologies and wavelength sensitive optical switches. Light paths can then be established in the same manner as connections in a conventional circuit switched network. An attractive feature of these networks is the spacial wavelength reuse.

HIGH SPEED WIDE AREA NETWORKS

Progress in fiber optics and high-speed switching technology and increasing user demands have recently stimulated the research and development of *Broadband Integrated Services Digital Networks (B-ISDNs)*. B-

ISDNs will carry several services (data, voice, video, graphics, etc.) on very high speed fiber optics link with user interfaces and interswitch trunks operating at 150 Mbps speeds.

For the implementation of B-ISDN, CCITT has targeted Asynchronous Transfer Mode (ATM) as the switching and multiplexing technology of choice. ATM networks use fixed-size data units called *cell* which are routed independently according to their destination address in a connection oriented environment. The advantage of using ATM resides in its flexibility in accommodating the variety of new and existing services requirements, and in the statistical gain obtained by multiplexing bursty sources.

ATM networks pose many research challenges ranging from the design of fast packet switching fabrics to the development of capacity planning tools. This is because, with the decreasing cost of fiber trunks, and the higher data rates required by the users, the network bottleneck has shifted from the communications channels to the switches. Consequently, while conventional packet networks were designed to minimize link costs, ATM networks are designed with the notion that the switch is the most expensive resource.

† Reference 7.

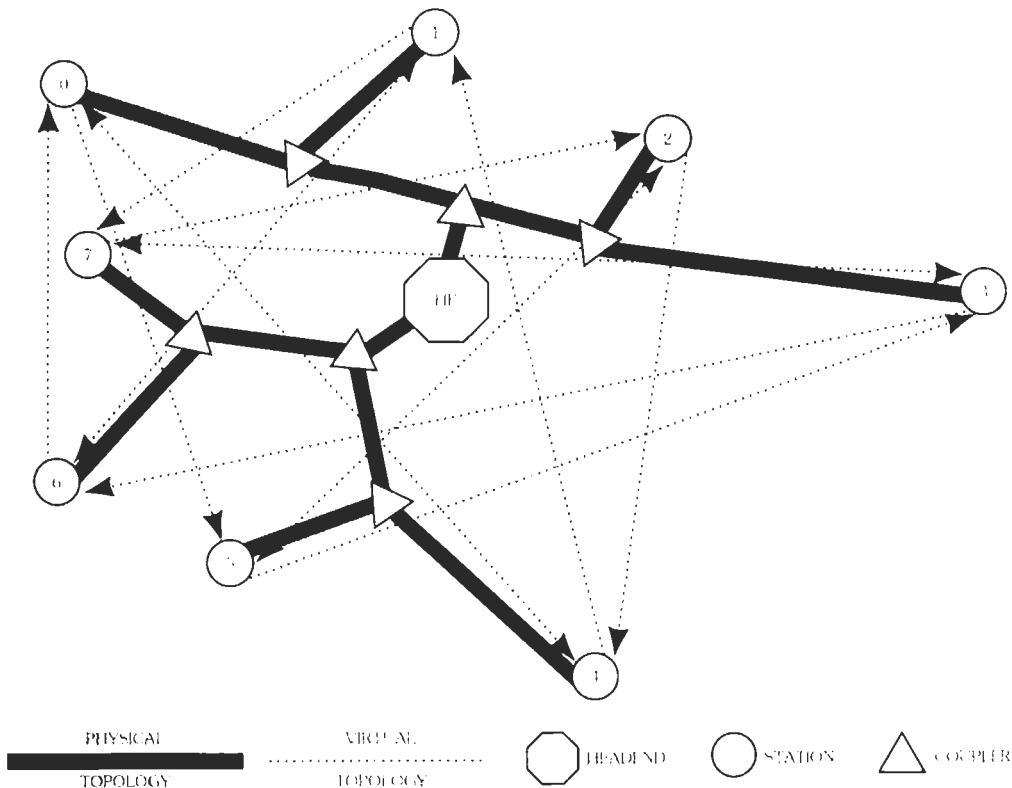


Fig. 22. Optical network (WDM) topology.

In order to enable the fast packet switch to process millions of packets per second, the internal network protocols have been drastically streamlined. In particular, almost all flow control features have been dropped from the link and packet level. The lack of internal flow control mechanisms combined with the inadequacy of buffers to handle sustained overloads at such high speeds and the strict response time requirements of voice and video traffic make the problem of flow control, congestion protection, and traffic management of B-ISDNs a formidable one.

In this section, we focus on two aspects of the ATM network which make its design and operation unique, namely bandwidth allocation and congestion control.

Bandwidth and Traffic Allocation in ATM Networks

One of the most critical issues in the design and operation of an ATM network is bandwidth and traffic management. This is because of two main reasons: (1) the bulk of ATM traffic will be of a real-time nature (voice, video, etc.), which is intolerant of cell loss and retransmission and which requires bandwidth and delay guarantees; and (2) the network level protocols have been streamlined (for processing efficiency) to the point that little or no flow control can be applied to a source after the user connection has been accepted. Thus, congestion must be *prevented* through careful allocation of bandwidth ahead of time.

In an ATM network, bandwidth can be allocated at the *physical level* and at the *transport* (i.e., ATM) level. At the physical level, the ATM topology can be dynamically reconfigured by adding/removing trunks between ATM switches. This allocation of bandwidth is made possible by the presence of Synchronous Transfer Mode (STM) Digital Cross Connects (DCS). We will refer to this allocation as STM allocation. At the ATM level, we can allocate bandwidth to individual Virtual Circuits (ATM-VC allocation) as well as to Virtual Paths (ATM-VP allocation).

ATM-VC allocation is the conventional type of allocation and the most commonly studied by researchers. The other two allocation schemes, however, are also very important and play a key role in the efficient use of bandwidth resources and the prevention of congestion.

STM Allocation—In order to explain how STM allocation works, we recall that ATM switches will be connected by bundles of 150 Mbps fiber optics trunks. These trunks are obtained from an underlying ‘pool’ of fiber facilities interconnected by circuit switches (DCS) operating in the STM mode. Thus, a trunk in the ATM network may consist of a chain of segments connected by DCS switches. In other words, the ATM network is ‘embedded’ into the STM network.

The topology design of the ‘embedded’ packet switched network differs from that of a traditional

packet switching (P/S) network. In traditional designs, we minimize total trunk cost subject to delay constraints. Here, the available facilities are given, and therefore the cost is fixed (at least for the short term). The problem thus becomes one of optimally configuring the P/S network topology and capacities, within the constraints set by the underlying C/S network. Furthermore, this reconfiguration can be carried out dynamically and can be ‘tuned’ to traffic fluctuations (i.e., STM allocation).

To illustrate the point, consider the network shown in Fig. 23. From the original (physical) topology, several logical packet topologies can be derived. The embedded topology of Fig. 24 is identical to the physical topology, whereas the topology in Fig. 25 has introduced a number of ‘express pipes’ between remote nodes. Express pipes reduce the number of intermediate hops along the path and thus reduce store-and-forward delay and nodal processing overhead. They also simplify the congestion control problem since, when a pipe becomes congested, the sources can be immediately stopped.

Express pipes also reduce the number of packet switch terminations. Note that the topology of Fig. 24 requires 192 packet switch terminations, whereas the topology of Fig. 25 (which has more express pipes) requires 152 terminations. Using a fully piped network with a direct pipe (of capacity 4, say) also between nodes B and D will further reduce the number of terminations to 144. This is a very important point, since current trends in transmission and processing costs indicate that terminations costs will soon dominate the cost of fiber trunks. Thus, adding more express pipes will lead to overall cost reductions.

There are, however, drawbacks in the implementation of fully piped (i.e., fully connected) logical topologies. For example, in networks with a large number of nodes the bandwidth may become too fragmented, and the advantage of statistically multiplexing several sessions on the same trunk may be lost. Thus, a good balance must be struck between express pipes and large trunks.

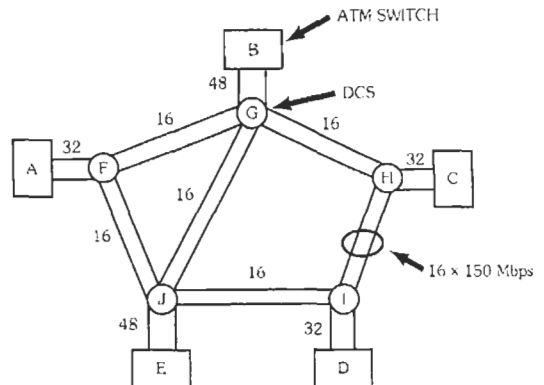


Fig. 23. Backbone topology.

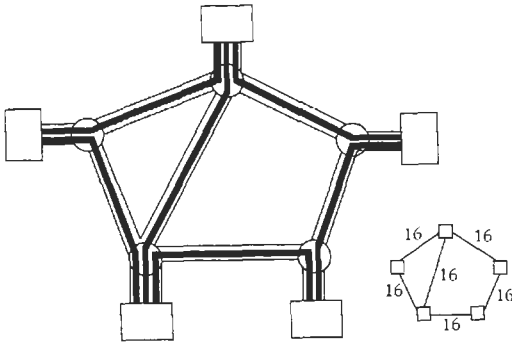


Fig. 24. Embedded topology A.

ATM-VP Allocation—We now move on to describe the *ATM-VP allocation* technique. We recall that a virtual path in ATM is basically a “virtual pipe” established between an origin/destination pair of ATM switches, for the purpose of multiplexing several VCs. The VP may traverse several transit nodes (also called ATM cross connects). The VC’s are carried through the transit nodes transparently; i.e., the VC header need not be inspected at such nodes. A VP may be set up at network initialization and is typically assigned a “peak” bandwidth. This is basically a “contract” between the VP users and the network: the users will not exceed the declared peak bandwidth (otherwise, cells will be dropped). Likewise, the network commits the peak bandwidth along the path (instead of implementing statistical allocation among paths). Thus, from the bandwidth-sharing point of view, there is little difference between ATM-VP and STM allocation: both schemes allocate the peak, with no statistical sharing among “pipes.” From the implementation standpoint, however, we recall that the VP is a “virtual” pipe, while the STM connection is a “physical” pipe. This implies that in VP the bandwidth can be allocated/deallocated much more dynamically than in STM.

The aggregation of several VCs into a VP pipe yields many advantages, including: (1) lower VC call set overhead; (2) lower VC bandwidth allocation and policing overhead (only the entry node gets involved); (3) connectionless traffic support; (4) better network congestion protection; and (5) better fault tolerance if redundant VPs are in use. On the negative side, the VP peak bandwidth allocation precludes some of the statistical sharing advantages available when bandwidth is allocated (statistically) to individual VCs.

Congestion Control in ATM

Recall that congestion control in conventional P/S nets is based on the principle of *detection* (of congestion) and *recovery*. When congestion builds up in the

network, causing buffers to overflow, some mechanisms (e.g., selective back pressure in VC nets, or “choke” packets in datagram nets) are triggered to relay this information to the traffic sources so that further inputs are slowed down or stopped.

Unfortunately, detection and recovery schemes are not well suited to ATM nets for the following reasons: (1) the protocols have been streamlined, thus precluding the implementation of sophisticated flow control mechanisms (e.g., rotating windows in X.25, “choke” packets, etc.); (2) voice and video traffic (the predominant component in ATM nets) cannot tolerate slowdowns or interruptions; and (3) at hundreds-of-Mbps trunk speeds and coast-to-coast round-trip delays, no buffer is large enough to absorb the temporary overload while the flow control mechanisms are trying to stop the traffic sources.

A better approach is to use *prevention* (instead of detection and recovery). For prevention, the key assumption is that traffic rates on user sessions can be predicted. This is generally true for voice and video connections. With this assumption in place, the network can perform bandwidth allocation and can verify during call setup that all the trunks along the path can accept the declared load. Alternatively, a background “bandwidth control” procedure can be used to keep track of the amount of bandwidth available (on the average) on various paths from a source to all destinations. Each node maintains a table with values of available bandwidth to each destination. When a connection request, with associated average bandwidth requirement, is submitted to the network, it is accepted only if a path with adequate bandwidth to satisfy such a request is found in the table. Otherwise, the request is rejected. Note that the bandwidth control implementation in the switch does not impact packet processing performance, since the bandwidth control algorithm is run in the background.

With the above preventive methods in place, congestion is, in principle, avoided since the network never

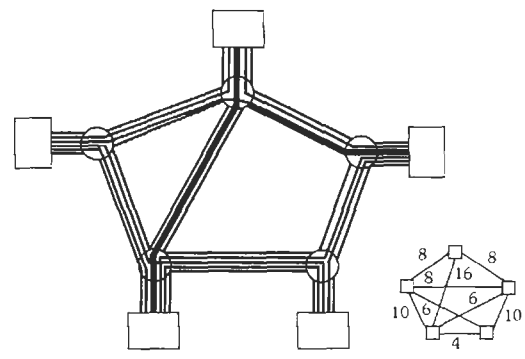


Fig. 25. Embedded topology B.

accepts more than it can carry. In practice, however, users may sometimes violate the initial bandwidth contract, either intentionally or, most commonly, because they estimated (and declared) their requirements incorrectly. In this case, congestion may still occur in the network, and it must be eliminated by using appropriate congestion control schemes.

There are several approaches to ATM congestion control. The first one is the *brute force* approach. With this approach, the network will indiscriminately drop packets when buffers are full. This alternative, however, is unacceptable for real-time traffic, and it may cause even more congestion (because of end-to-end retransmissions) in data traffic. Another possibility is *selective dropping* of least significant packets. With this approach, voice and video are encoded in such a way that less significant bits are carried in separate packets. When congestion builds up, the less significant packets are dropped, causing degradation (but no disruption) of service. This alternative, however, works only on voice/video and allows (at best) only up to 50% regulation of traffic volume. Furthermore, either scheme is unfair because it affects users indiscriminately, instead of selecting and penalizing those users who are abusing the network (i.e., who are transmitting at a higher than initially declared rate).

Another approach is to apply *input rate control*. Each user is monitored (or *policed*) for compliance with the initial bandwidth declaration. If the user exceeds the declared bandwidth, the excess traffic will be dropped directly at the source. A popular, extensively studied enforcement mechanism is the "leaky bucket," which is basically a queue with prenegotiated (at call setup time) service rate and buffer storage. If the user burst exceeds the buffer capacity, subsequent cells are dropped. The leaky bucket is definitely more fair than indiscriminate dropping. It suffers, however, from several inefficiencies. For one thing, the parameters cannot be easily set in such a way as to get an accurate enforcement. Furthermore, it is quite possible that the extra traffic (although in violation of the contract) could actually be tolerated by the network at that particular time (and therefore should be admitted because there is enough extra bandwidth available on the path).

An approach related to input rate control is Distributed Source Control (DSC). Each source is assigned a "smoothing interval" "T" and an allowance "W" of packets which can be transmitted in such interval. "T" and "W" are properly chosen so that the network has sufficient internal buffering to accommodate the corresponding input rate without risk of overflow. If a burst is generated by the source, it must be distributed over several smoothing intervals. If the burst overruns the

source buffers, it is dropped at the source before entering the network (similar to what was happening in a leaky bucket). This packet loss at the source may occur even if the network happens to have enough bandwidth to tolerate the burst. In fact, DSC resource allocation tends to be quite conservative.

These considerations lead to yet another approach which is based on *marking*, rather than discarding, abusive packets. More specifically, the policer sets the CLP (Cell Loss Priority) bit to one in the abusive packet header. The marked packet is forwarded thus in the network. If congestion occurs, marked packets are discarded first. A further refinement on the marking scheme is possible in conjunction with the previously mentioned bandwidth control scheme. If the policer knows how much residual bandwidth is available on the path, then it can accept (and mark) abusive traffic at a rate up to the value of the available bandwidth. Traffic beyond that value is discarded. This way, even the marked traffic has a very good chance of making it through.

REFERENCES

The material presented in this chapter is primarily based on the content of the books listed below. By consulting these, the reader will be able to obtain more detailed information about this field and trace the original contributions and contributors.

1. Kleinrock, Leonard. *Computer Applications*. Vol. II of *Queueing Systems*. New York: John Wiley & Sons, Inc., 1976.
2. Tanenbaum, Andrew S. *Computer Networks*, 2nd ed. Englewood Cliffs, N.J.: Prentice-Hall, Inc., 1988.
3. Green, Paul E. Jr., ed. *Computer Network Architectures and Protocols*. New York: Plenum Publishing Corp., 1982.
4. Martin, James. *Computer Networks and Distributed Processing: Software, Techniques, and Architecture*. Englewood Cliffs, N.J.: Prentice-Hall, Inc., 1981.
5. Tobagi, F. "Multiaccess Link Control," from Sunshine, C., ed., *Computer Network Architectures and Protocols*, 2nd ed. New York: Plenum Publishing Corp., 1989.
6. Gerla, M., Rodriguez, P., and Yeh, C. "Token Based Protocols for High Speed Optical Fiber Networks." *Journal of Lightwave Technology*, Vol. LT-3, No. 3, June 1985.
7. Minzer, S. E. "Broadband ISDN and asynchronous transfer mode (ATM)." *IEEE Communications Magazine*, Sept. 1989.

27

Satellite and Space Communications

Geoffrey Hyde and Pier L. Bargellini

Evolution of Communication Satellite Systems and
Their Role in Worldwide Communications 27-4

Satellite Orbits 27-6

Information Transmission in Space 27-7

The Choice of the Orbit 27-8

Solar Eclipses 27-9

Nongeostationary Orbits 27-11

Elevation and Azimuth Angles 27-11

The Choice of Frequency 27-12

Link Budgets 27-13

Spacecraft Architecture 27-14

Transponders 27-15

Overall Transmission System Considerations 27-19

Analog Transmission Systems	27-21
FDM/FM (Single-Carrier Case)	
FDM/FM/FDMA (Multicarrier Case)	
Digital Transmission Systems	27-23
Demand-Assigned Multiple Access (DAMA)	27-24
Time-Division Multiple Access (TDMA)	27-24
Satellite-Switched Time-Division Multiple Access (SS-TDMA)	27-27
Other Services	27-28
Networks	27-31
Spacecraft Antennas	27-31
Propagation	27-32
Faraday Rotation	
Ionospheric Scintillations	
Tropospheric Effects	
Clear-Sky Noise Temperature	
Hydrometeors	
Rain Attenuation	
Sky Noise Temperature with Rain	
Diversity	
Depolarization	
Earth Stations	27-40

Advances in rocketry and microwave engineering inspired early proposals for communications satellites. Experiments of the late 1950s and early 1960s culminated in the successful launch of *Early Bird* in 1965, by Comsat (the Communications Satellite Corporation), established by act of the US Congress for this purpose, from which time satellite communications can be dated.

In space systems, earth stations operate in conjunction with orbiting spacecraft that probe the space environment—the earth as observable from space, the moon, a planet, or any other celestial body. In satellite communications systems, two or more stations located on or near the earth communicate via satellites that serve as relay stations in space. In both instances, control and monitoring of the spacecraft require that telemetry and command links be added to the main function of the mission. Space systems include terrestrial missions (e.g., earth and/or sea surface observations of different kinds), weather satellites, and navigation satellites. Beyond the earth, space systems can be classified in terms of the mission range (i.e., cislunar, lunar, translunar, or planetary) as well as in terms of the specific nature of the observations to be carried out.

Communications satellite systems are classified in terms of their territorial coverages, e.g., global, regional, or national (domestic); in terms of the type of services offered, e.g., fixed, mobile, maritime, aeronautical, etc., or point-to-point, broadcasting, commercial, military, amateur, experimental, etc.; or in terms of their orbit, e.g., geostationary (GEO), medium earth orbit (MEO), or low earth orbit (LEO).

Regulatory bodies (FCC, ITU) often categorize capabilities by service. In satellite communications, there are three broad service categories: fixed, mobile, and broadcast. Fixed satellite services (FSS) cover links between satellites and fixed (nonmoving) earth stations. Mobile service covers satellite links to stations that may be in motion (mobile), including ships (maritime mobile [MMSS]), aircraft (aeronautical mobile [AMSS]), and land vehicles (land mobile [LMSS]). Broadcast services include TV (DBS-TV) and audio (DBSA). It should be noted that services are evolving, and while the FCC and ITU designations will change very slowly, the reality is changing rather more quickly, as will be discussed briefly in the section on Other Services below. Examples of this are DTH (direct to home), which covers a broader swath than just TV as interactive services emerge, and DARS (digital audio radio service), which covers nationwide satellite digital radio to automobiles (and possibly homes).

The ITU internationally and the FCC in the United States allocate frequency bands for the use of these services, and where the bands overlap, they designate which service has priority.

The environment of space affects the design of communications systems in several ways that make it dif-

ferent from the design of terrestrial systems. Major differences are:

- A. Space and satellite communications systems can and often do cover distances far exceeding those encountered on earth.
- B. As spacecraft power and allocated bandwidth are limited resources, trade-offs of space and earth segment design characteristics affect the overall system cost.
- C. As the conditions along the signal paths are much more time invariant in space than on earth, it is possible to design space-to-space communications systems with great precision. Space–earth and earth–space signal paths traverse the troposphere and ionosphere and are subject to the vagaries thereof, but these tend not to be so severe as those encountered on long terrestrial paths (see section on Propagation).

Three categories characterized by different environmental constraints can be identified, as follows.

Spacecraft-to-Spacecraft: In principle, the designer has maximum freedom in the choice of the operating frequency. The major difficulty resides in maintaining tracking between spacecraft.

Earth-to-Spacecraft (Up-Link): The choice of the operating frequency is primarily determined by the availability of spectral windows in the signal path. The window boundaries are dictated by absorption and dispersion phenomena in the ionosphere and troposphere, the attenuation and depolarization effects of hydrometeors (rain, snow, ice, etc.) and also by the spectral and spatial distributions of natural noise sources. On earth, the ubiquitous availability of electrical power and the relatively benign environment make it possible to use large amounts of transmitter power enhanced by high-gain, large antennas that must be precisely aimed at the spacecraft.

Spacecraft-to-Earth (Down-Link): Launch-vehicle limitations restrict the size of the spacecraft transmit and receive antennas, and the spacecraft receiver is affected by the background noise of the earth. The spacecraft transmitter is limited by the power available, which in turn is limited by the size of the solar cell arrays, again launch-vehicle limited. These limitations are increasingly less severe. Spacecraft originally (1963–1965) had dipole antennas and only a few hundred watts of dc electrical power from solar arrays. In the year 2000, spacecraft on orbit had 15-meter antennas and dc power up to 15 kW.

In earlier days, limited spacecraft transmitter power and antenna sizes (gain) were compensated for by electrically large earth-station receive antennas and low-noise receivers. Today, spacecraft eirps (equivalent isotropically radiated power) have reached the

point where handsets with antennas not unlike those used in terrestrial cell phones are coming into service. The trend toward multibeam narrow-beam spacecraft antennas has increased the requirements for precise spacecraft positioning, station-keeping, and antenna pointing with a consequent increase in the amount of onboard fuel required for the above-mentioned functions. The very stringent reliability requirements necessary to guarantee spacecraft design lifetime are achieved by painstaking selection of the equipment for all subsystems and provision for redundancy of critical components.

EVOLUTION OF COMMUNICATION SATELLITE SYSTEMS AND THEIR ROLE IN WORLDWIDE COMMUNICATIONS

From 1965 until around 1980, communication satellites were used chiefly to provide long-distance high-capacity communications links across the oceans, archipelagos, and large landmasses. The capacity and reliability of satellite circuits exceeded those of other communications systems such as HF, VHF, and UHF radio (ionospheric, troposcatter, and line-of-sight circuits), as well as coaxial cables. The majority of traffic was telephony, with data transmission growing at an accelerated rate. The relaying of TV programs across the oceans was made possible for the first time through the use of satellites. TV distribution and TV Satellite News Gathering (SNG) are major and growing niches. Geostationary satellites proved perfectly capable of interfacing with existing land-based communications systems as echo problems were circumvented through the use of echo cancellers. Analog transmission techniques initially adopted are still in use, albeit no longer dominant, while digital transmission now is the dominant form and continues its rapid expansion. In point-to-point links, satellites carrying "transparent" transponders have served at times as "cables in the sky." However, it became clear from the very beginning of satellite communications that the most valuable characteristic of satellites was their capacity of providing multipoint and multimedia communications links in star-configured networks (multiple access). Therefore, the efforts of systems planners and spacecraft designers were directed toward a more efficient use of the multiple-access characteristic.

Enhancement of the multiple-access capability implied the addition to the basic function of signal amplification and frequency-band changing found in "transparent transponders" the functions of switching (beam and/or transponder) and of signal processing of various kinds such as demodulation, remodulation, storage, buffering, etc. With the rise of digital-signal onboard processing through the past two decades, some satellites can be increasingly thought of as signal processing and routing nodes in the sky.

The evolution of spacecraft design was characterized by increases of size, mass, and power, made possible by the availability of more powerful launchers, in combination with substantial improvements in the design of the various onboard subsystems in both the communications and noncommunications areas. Within less than 25 years, communications capacity was increased by more than three orders of magnitude with only about a fifty-fold increase in spacecraft mass and power. This process continues, although at a lesser rate. As the power and mass available for the communications payload has increased, those payloads have grown in capability, and the overall system trade-offs have shifted toward placing more burdens on the spacecraft payload. Onboard processing, especially for the communications signals but also for the satellite bus housekeeping, continues to increase. The total radiated rf power has increased by orders of magnitude, and antenna diameters have reached 15 meters (46 feet), resulting in great increase in cirp. One of the most significant consequences of the above-mentioned evolution is the reduction in the size of earth-station aperture and transmitters. Such reductions and the opening of communications at the higher frequency bands allocated to satellite systems made it possible to install earth terminals directly at the customers' premises. The introduction of low-powered, very small antenna terminals (VSATS) allowed the development of numerous independent communications networks linking industries, banks, and other institutions. All of this resulted not only in great cost reductions but also in enabling users a degree of operational freedom hitherto unavailable through conventional communications facilities. Although the majority of the current VSAT systems operate with the peripheral terminals communicating with and/or through a central station (hub) on earth, it is envisaged that in the future the functions of the hub will be available on board more sophisticated satellites. In this manner, the double-hop situation encountered in present VSAT systems will be eliminated with consequent various advantages. The process of continued reduction in the size of the terminals has, in the late 1990s, led to the introduction of ultra-small antenna terminals (USATs) in the FSS, and in the year 2000, handsets not unlike and compatible with those used in terrestrial mobile communications, for example, the GSM.

The use of satellites has facilitated worldwide TV distribution for special events in real time and, in most cases, with delays as required by the different time zones, regularly scheduled programs with participants in two or more locales connected by satellite, as well as the SNG services noted above. In the United States, the nationwide distribution of TV programs to cable networks has become an industry of its own with little or no competition from other means of transmission, including fiber-optic cables. For rural areas not reached by conventional TV broadcast stations and not served by cable TV systems, higher levels of power radiated by satellites made possible direct reception of

TV from satellites by receive-only earth terminals (TVRO) equipped with parabolic antennas with diameters between 2 and 5 m. These developments have paved the way for direct TV broadcast from satellites operating in the K_u -band with powers between 100 and 250 watts and received with antennas down to 0.5 m in diameter. Direct TV broadcasting from satellites, earlier initiated in Europe and Japan, has finally become widespread in the United States. Other direct to home (DTH) satellite services are starting to appear, including interactive ones such as Internet service via satellite.

At this point in time, almost all voice and data transmission is digital. Analog TV, i.e., "old TV," is a decreasing sector as TV goes digital. For example, DBS is digital. The Internet is digital. Digital radio services such as DARS are emerging. Analog is dying out. This is in spite of the fact that voice, video, and radio are analog at the end users (throat, eye, and ear, respectively). One major reason is that digital methods provide source encoding and processing such that compression techniques can, have, and are being developed so that more and more voice circuits and TV channels are being transmitted through less and less bandwidth. Because bandwidth is a limited resource in satellite communications, this is highly desirable. The earliest digital voice service over satellite used 56 kbps to achieve "toll quality." Today (year 2000) one can achieve very near toll quality with 4 kbps using advanced coding techniques. Comparable advances are available for TV. The other reason is that the much greater eirps permit the use of digital modulation schemes that can provide more bits per Hertz as part of the systems trade-offs.

An area in which satellites are bound to maintain a significant role is mobile communications. After the establishment of the MARISAT system in 1976, the follow-on INMARSAT global system underwent a considerable expansion marked by the sequence of three generations of ever-increasing capacity providing ship-to-shore as well as air-to-ground communications services on a global scale. Satellite systems capable of serving earth vehicles such as trucks, trains, and automobiles have been implemented; they provide not only communications services but also additional important functions such as determination of vehicle location, monitoring of the status of critical shipment loads, etc. The 1990s saw the emergence of several satellite mobile systems and numerous satellite mobile communications providers. There has been a blurring of the lines among the satellite services for marine, aeronautical, and ground mobile users.

As noted above, satellite mobile service in the year 2000 has been extended to hand-held sets not much different in size than those used in terrestrial wireless service. The consequent issues of compatibility and complementarity with terrestrial-based mobile communications systems are being addressed. Indeed, at this time, handsets that can access a satellite system and the European GSM have been introduced. The emergence of "software radio" systems makes it likely

that the issues of compatibility will be overcome by chip sets in the handsets in the next few years. The so-called third generation of mobile services in Europe is being planned so that satellite-based services complement the terrestrial networks. More of such efforts are in the offing.

The emergence of terrestrial and submarine fiber-optic cables in the early and mid-1980s induced a shift in the use of satellite systems toward those applications for which they are either indispensable or better suited. Clearly, the two above-mentioned cases of mobile communications and direct broadcasting are best performed via satellites in many regions of the world. However, even for overland and transoceanic routes, satellites will remain competitive vis-à-vis optical cables not only in terms of economic factors but also especially in consideration of the complementary characteristics of cable and satellite systems. The need to have satellite circuits at readiness to restore interrupted cable circuits has been repeatedly demonstrated. Quick restoration of communications services via satellite during often long-term submarine cable outages has been frequent and effective with considerable, albeit unadvertised, benefits to the users. From a topology viewpoint, a cable that constitutes a link and a satellite that functions as a node are bound to provide, respectively, fixed trunk capabilities and nodal connectivity to serve many users, especially in the case of thin traffic routes. In addition, because of the portability of smaller satellite earth terminals, where temporary and/or emergency services are needed, communications services can be provided or restored in a much shorter time frame than can terrestrial services.

The integration of satellite and cable circuits, which already has proven beneficial, has become indispensable in the case of the Integrated Digital Services Network (ISDN) currently under development. The definitions of the characteristics of the ISDN by consultative bodies such as the ITU-T and the ITU-R lead to the conclusion that satellites need to be taken into consideration not just as another transmission medium, but also because of their unique capabilities of higher connectivity, flexibility, assignment of capacity to meet rapidly changing traffic needs, and readiness to confront emergencies. Compatibility and complementarity of fiber-optic cables and satellite systems, already well assured with advantages to users at the current ISDN transmission rates, must be pursued and extended throughout the development phases of the future higher-transmission-rates ISDN (B-ISDN) hierarchy. Communications satellites will maintain a very important role within the overall ISDN development scenario, possibly widening it in terms of network growth anticipation, planning, design, and operation.

Satellite-based services are beginning to emerge for the Internet. There is every evidence that this trend will continue. Although one can expect the bulk of Internet traffic to be carried over terrestrial facilities, for the same reasons that satellite mobile services will always be needed, so will satellite-based Internet services.

At the dawn of the 21st century, improving and changing technology are the enablers for expanding services and lower costs. The single greatest area of change is digital technology, both with regard to hardware and software. Microwave antenna technology and solar cell power generation efficiency improvements are also critical areas. As these technologies evolve, one can expect satellite communications to continue to thrive through the foreseeable future.

SATELLITE ORBITS

Satellites move around the earth as planets do around the sun. Kepler's laws apply:

First Law: Planetary orbits are elliptical with the sun at a focus.

Second Law: The radius vector from the sun to a planet sweeps equal areas in equal times.

Third Law: The ratio of the square of the period of revolution and the cube of the ellipse semimajor axis is the same for all planets.

These laws can be mathematically derived in terms of Newtonian mechanics and universal gravitation.

When the satellite mass is negligible with respect to the mass of the earth, the balance between gravitational and centrifugal force leads to the vector differential equation

$$d^2\vec{r} / dt^2 + \mu\vec{r}^3 = 0 \tag{Eq. 1}$$

in which \vec{r} is the radius vector from the center of the earth to the satellite and

$$\mu = GM = gR^2 = 3.99 \times 10^{14} \text{ m}^3/\text{s}^2$$

where

$$G = 6.67 \times 10^{-11} \text{ N m}^2/\text{kg}^2$$

$$M = 5.98 \times 10^{24} \text{ kg}$$

$$g = 9.81 \text{ m/s}^2$$

$$R = 6.38 \times 10^6 \text{ m, the radius of the earth}$$

As the centrifugal and gravitational forces are aligned and opposed, the satellite moves on a plane (orbit plane). Lengthy manipulations of Eq. 1 lead to Kepler's laws. The first law is verified by showing that the satellite moves along a conic on the orbit plane, the equation for which is

$$r = p / (1 + e \cos v) \tag{Eq. 2}$$

where

p is a parameter,

e is the eccentricity,

v is the central angle (true anomaly).

The value of the eccentricity determines the type of the conic.

$$e \begin{cases} = 0 & \text{circle} \\ < 1 & \text{ellipse} \\ = 1 & \text{parabola} \\ > 1 & \text{hyperbola} \end{cases}$$

With reference to Fig. 1, the following relationships, which hold for elliptical orbits, are frequently used:

$$x = r \cos v = a(\cos E - e)$$

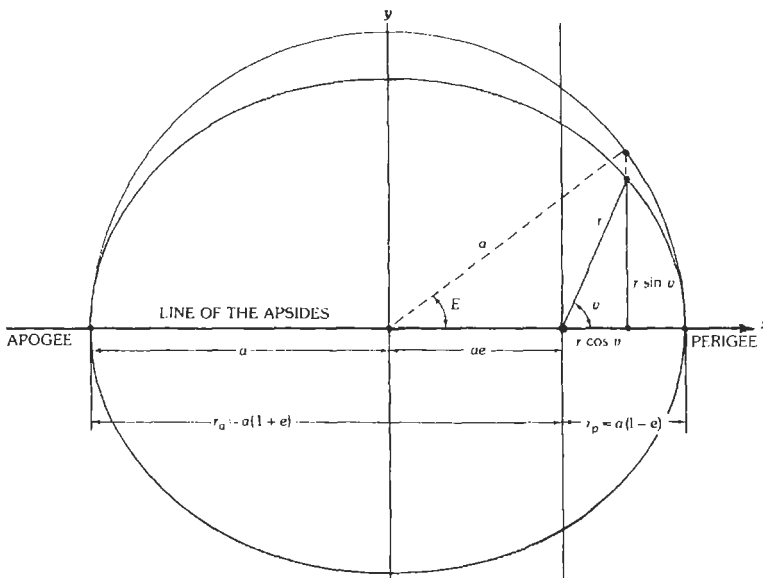


Fig. 1. Elliptical-orbit geometry.

$$y = r \sin v = a \sin E (1 - e^2)^{1/2}$$

Eccentricity:

$$e = c/a = (r_a - r_p)/(r_a + r_p)$$

Semimajor axis:

$$a = (r_a + r_p)/2$$

Apogee distance:

$$r_a = a + c = a(1 + e)$$

Perigee distance:

$$r_p = a - c = a(1 - e)$$

Locus parameter:

$$p = a(1 - e^2) = 2r_a r_p / (r_a + r_p)$$

Semiminor axis:

$$b = a(1 - e^2)^{1/2} = (r_a r_p)^{1/2}$$

Kepler's second and third laws can be derived from Eq. 1 through manipulations based on conservation of energy and angular momentum considerations whereby

$$v^2/2 - \mu/2 = \text{constant} \quad (\text{Eq. 3})$$

and

$$\vec{r} \times \vec{v} = \vec{h} \quad (\text{Eq. 4})$$

The orbit period, P , is then

$$P = 2\pi a^{3/2} / \mu^{1/2} \quad (\text{Eq. 5})$$

Angle E (eccentric anomaly) is the central angle measured from the x axis to the vertical projection of the satellite point over the circle of radius a . The true anomaly, v , and the eccentric anomaly, E , are related by any of the expressions:

$$\cos E = (e + \cos v)/(1 + e \cos v)$$

$$\cos v = (e - \cos E)/(e \cos E - 1)$$

$$\tan (v/2) = [(1 + e)/(1 - e)]^{1/2} \tan (E/2)$$

For the hypothetical case of uniform motion on the circle circumscribing the ellipse, where t = time, the fictitious angle M (mean anomaly) can be defined as

$$M = 2\pi t/P \quad (\text{Eq. 6})$$

Angles M and E are related by the equation

$$M = E - e \sin E \quad (\text{Eq. 7})$$

It can also be shown that

$$v^2 = 2\mu/r - \mu/a \quad (\text{Eq. 8})$$

and that the orbit velocity vector has two constant components respectively normal to the radius vector and the major axis of the ellipse.

To define the satellite position in inertial space, three additional angles are required. These are: the

inclination, i , of the orbit plane with respect to the equatorial plane of the earth; the right ascension (longitude), Ω , of the ascending node [i.e., the angle between the direction, γ , of the vernal equinox (point of Aries) and the intersection of the orbit and equatorial planes]; and finally the argument of perigee, ω , i.e., the angle measured in the orbit plane between the direction of the ascending node and that of the perigee.

The nonsphericity of the earth, the nonuniform distribution of its mass, as well as the effects of lunar and solar gravitation make the actual motion of satellites depart from that described so far, which is based on the simple two-body problem.

The gravitational potential of the earth, which is no longer simply proportional to the inverse of the distance, can be expressed as a series of spherical functions. Two motions known as *secular variations* result from the above-mentioned causes; these are nodal regression and rotation of the line of the apsides.

Nodal regression is a rotation of the orbit plane in the direction opposite to the satellite motion around the axis of rotation of the earth. The rate of this motion (in degrees per day) is

$$\dot{\Omega} = -[10/(1 - e^2)^2] (R/a)^{7/2} \cos i \quad (\text{Eq. 9})$$

Rotation of the line of the apsides is rotation of the ellipse major axis around the center of the earth on a fixed orbit plane. Apogee and perigee move then with respect to the earth, and the rate of motion of the argument of perigee is also a function of the orbit inclination angle, i , or:

$$\dot{\omega} = [5/(1 - e^2)^2] (R/a)^{7/2} (5 \cos^2 i - 1)^\circ/\text{day} \quad (\text{Eq. 10})$$

The critical value $i = \cos (1/\sqrt{5}) = 63^\circ 24'$ stops this motion. Smaller values of i make the ellipse rotate in the same sense as the satellite motion along it; higher values of i make the ellipse rotate in the opposite sense. The effect becomes smaller for higher orbits.

Gravitational forces due to the moon and the sun also influence satellite motion. For low-altitude orbits, the gravitational field of the earth is preponderant and solar and lunar effects can be neglected. For high-altitude orbits, gravitational disturbances of the moon and the sun cannot be neglected. In the case of a geostationary orbit, the sun and moon perturbation forces, respectively, amount to about 1/37 and 1/6000 of the gravitational force of the earth. The combined effect would change the orbital inclination of a satellite originally placed in a perfectly equatorial orbit by about 1 degree per year.

INFORMATION TRANSMISSION IN SPACE

In a link between two points in space, A and B, separated by the distance r , a transmitter located at A of

power P_t , radiated by an antenna having a gain G_t , produces at point B a power flux density

$$\phi = P_t G_t / (4\pi r^2) \quad (\text{Eq. 11})$$

A receiving antenna of effective aperture A_r at B intercepts a signal power

$$P_r = P_t G_t A_r / (4\pi r^2) \quad (\text{Eq. 12})$$

which equals the product of the rate of information transmission, R , and the energy, E , required to transmit one bit of information, whence

$$R = BP_t G_t A_r / (4\pi\beta Nr^2) \quad (\text{Eq. 13})$$

where β is the ratio of the energy required to transmit one bit of information and the noise power density (i.e., $\beta = E/N_0$ and $N_0 = N/B$ where N is the noise power in bandwidth B). This parameter, which is a function of the ratio of the channel bandwidth, B , and the information rate, R , characterizes the efficiency of modulation-demodulation schemes.

Equation 13 was derived by assuming an antenna of given gain at one end of the link and an antenna of given aperture at the other. In this case, frequency, or wavelength, does not appear explicitly in the equation, although the noise power density, N_0 , is in general frequency dependent.

By assuming antennas of given gain at both ends, and using the well-known relationship

$$A_e = (\lambda^2/4\pi)G$$

where A_e = effective aperture, the following result is obtained:

$$R = (\lambda^2/16\pi^2) [(BP_t G_t G_r) / (\beta Nr^2)] \quad (\text{Eq. 14})$$

In contrast, assuming antennas of given aperture at both ends of the link yields

$$R = BP_t A_t A_r / (\lambda^2 \beta Nr^2) \quad (\text{Eq. 15})$$

The above equations allow the information transmission rate, R , of a single space communications link to be computed. Since the noise encountered in space may be regarded as Gaussian, white, and additive, the Shannon model yields an upper bound on transmission rate. Equation 13 can be rewritten as

$$R = BP_t G_t A_r / (4\pi r^2 \beta N) = (B/\beta) \cdot (P_t/N) \quad (\text{Eq. 16})$$

The two factors on the right-hand side relate to the two cases of bandwidth- and power-limited transmission.

The lower bound of parameter $\beta = E/N_0$ has the value $\log_e 2 = 0.693$ in the limiting case of infinite bandwidth and ideal (Shannon) modulation-demodulation processes. Then the signal-to-noise ratio goes to zero, and the information rate takes the value

$$\lim_{B \rightarrow \infty} R = 1.44 P_t / N_0 \quad (\text{Eq. 17})$$

In physical systems, limited bandwidth and departure from ideal modulation-demodulation processes

require higher values of E/N_0 and P_t/N . The ratio of the actual information transmission rate and channel capacity (in the Shannon sense) becomes smaller than unity, and compromises amongst signal-to-noise ratio, bandwidth-to-information transmission rate ratio, and message error probabilities (or signal quality in analog systems) must be made.

THE CHOICE OF THE ORBIT

The orbit height influences the communications range, the transmission delay, and the duration of the connection, which is determined by the time interval between satellite rise and set for points on the surface of the earth separated by a given angle measured from the center of the earth. The actual coverage is also a function of the minimum elevation angle of earth-station antennas. With reference to Fig. 2, the following relationship can be used to compute the coverage angle:

$$R/(R+h) = \cos(\beta + \theta) / \cos \theta \quad (\text{Eq. 18})$$

where

- R = earth radius
- h = orbit height
- β = coverage angle
- θ = elevation angle

The transmission delay can reach a maximum

$$t_d = [2(R+h)/c](\sin \beta / \cos \theta) \quad (\text{Eq. 19})$$

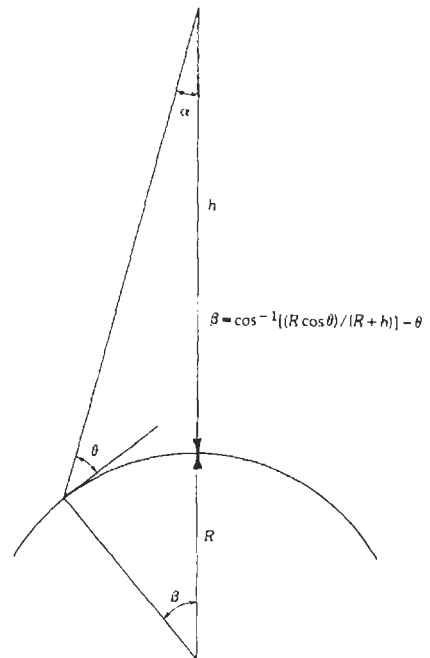


Fig. 2. Coverage and elevation angles.

Curves for the period, P , the earth coverage, and the maximum transmission delay, t_d , are shown in Fig. 3 for circular orbits.

Satellite systems can be classified in terms of the inclination angle, i , between the orbit plane and the equatorial plane:

$$i = 90^\circ, \text{ i.e., polar orbits}$$

$$0 < i < 90^\circ, \text{ i.e., inclined orbits}$$

$$i = 0, \text{ i.e., equatorial orbits}$$

If the orbit is circular and the orbital height is 35.863×10^6 meters, the period equals a sidereal day (23 hr 56 min 4 s), and if the orbital plane coincides with that of the equator of the earth, the satellite becomes geostationary. In the absence of perturbations, such satellites hover over fixed points of the equator and subtend the earth with an angle of $17^\circ 27'$. Communications can thus be established between any two points "visible" from the satellite over an area that amounts to about 40% of the entire surface of the earth. In this case, the handover problem, i.e., the passing of traffic from one satellite to another, which would be necessary with nonstationary satellites, is avoided. Other significant advantages are simplification of tracking by earth stations and near-zero Doppler effects.

The transmission delay, which ranges from 0.238 to 0.275 second for single-hop circuits, has been found acceptable for telephone communications provided that the echo at times generated in terrestrial telephone plants at the transition points between four- and two-wire circuits is kept under control. This has been done by using echo cancellers.

SOLAR ECLIPSES

A satellite is illuminated by the sun continuously except for the periods during which it is shadowed by the earth. As a spacecraft enters the shadow of the earth, the supply of electrical energy derived from solar-cell panels is cut off, and a substantial drop in the

temperature of the spacecraft body occurs; reversal of these phenomena takes place when the satellite comes out of the shadowed region. In the case of communications satellites, where continuity of service is needed during eclipses, it is clear that some alternate source of electrical energy must be provided. Storage batteries, usually of the Ni-Cd or Ni-H₂ type currently, provide this energy, being charged by a fraction of the energy obtained by sunlight solar-cell panels and discharged during eclipse periods. With regard to thermal effects, great care must be exercised in the construction of all spacecraft subsystems so that all functions are properly maintained, notwithstanding the large temperature excursions encountered.

Eclipses have known onset times and duration; i.e., they are highly predictable events calculated from ephemeris data. Although such calculations tend to be complex for satellites orbiting the earth at different altitudes and inclinations, they become fairly simple for satellites whose orbits lie in the equatorial plane of the earth, such as geostationary satellites. For this case, solar eclipses occur around the spring and autumn equinoxes. Fig. 4A illustrates what happens around the spring (vernal) equinox: Shadowing of the satellite begins at point A when it becomes positioned tangentially to the shadowed area. This occurs at the end of February. Subsequently, the duration of the eclipse increases until the spacecraft finds itself at point B on the equinox day (March 21), which yields an eclipse of maximum duration. After this, the eclipse phenomenon continues with decreased duration until point C (April 11), when the satellite position is again tangential to the shadow of the earth. The situation repeats itself from September 12 to October 14 around the autumnal equinox (September 23), at which date the eclipse duration hits a maximum again.

For a satellite having zero inclination (i.e., on the equatorial plane), the computation of the maximum duration of the eclipse is easily done, with reference to Fig. 4B under the reasonable approximation of the earth shadow as a cylinder of radius r_e equal to the radius of the earth. For a satellite orbiting at a distance R from the center of the earth, the angle α that subtends the shadowed arc of the orbit is

$$\alpha = 2 \sin^{-1} (r_e/R) \tag{Eq. 20}$$

Maximum eclipse duration T_e can then be obtained from the equation

$$T_e = P\alpha/2\pi \tag{Eq. 21}$$

where P is the orbital period. Substitution of the appropriate values for the geostationary case yields

$$\alpha = 2 \sin^{-1} (6378/42\ 164) = 0.304 \text{ rad}$$

and

$$T_e = (86\ 164.091/6.2831)0.304 = 4168 \text{ s} \\ = 69.48 \text{ min}$$

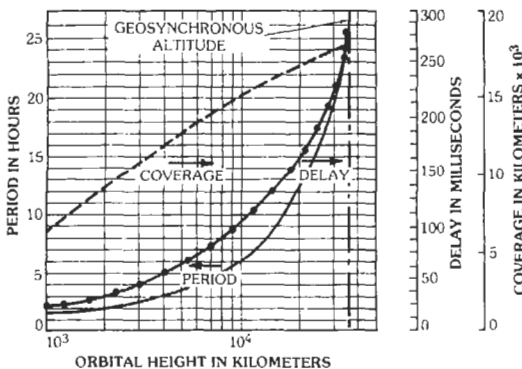
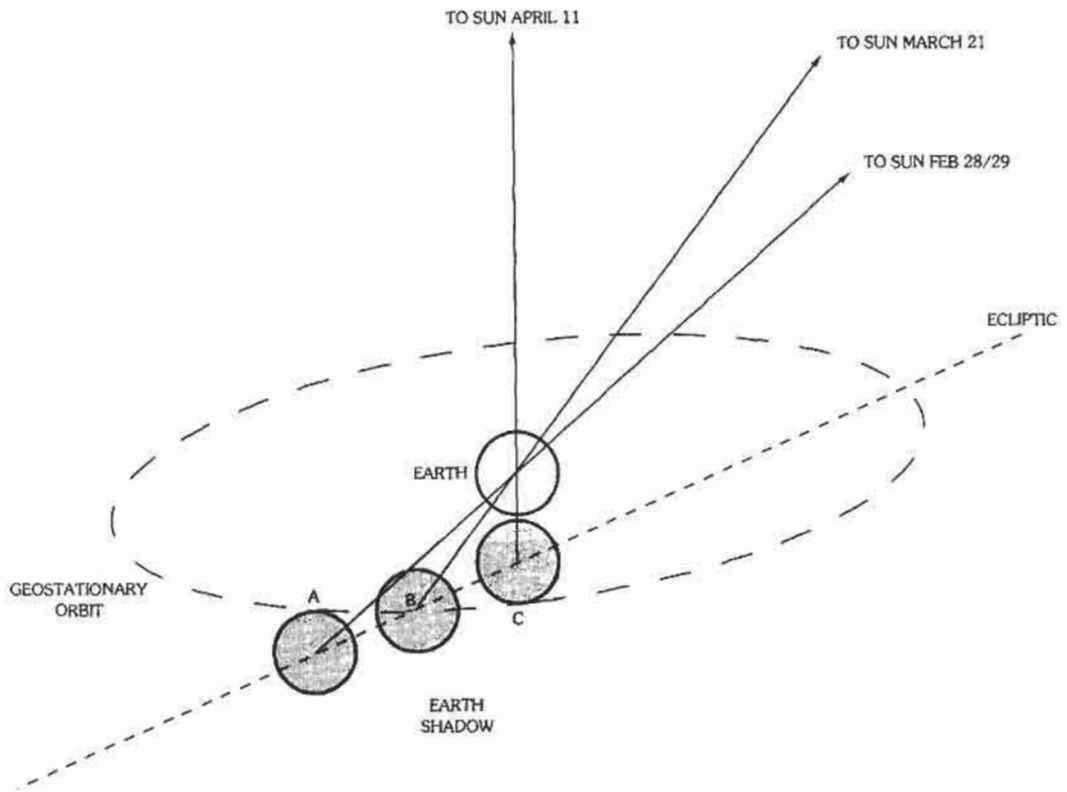
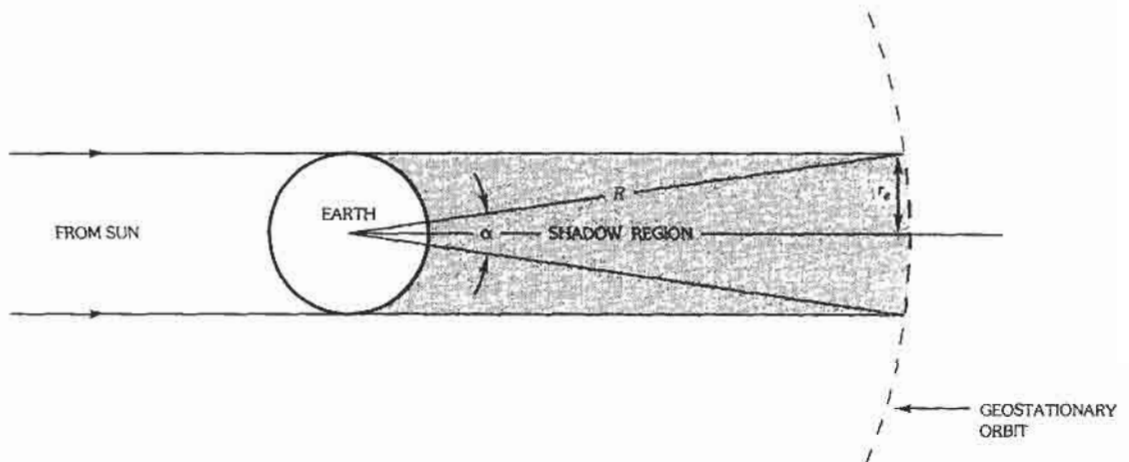


Fig. 3. Orbit period, earth coverage, and transmission delay versus orbit height.



(A) Eclipse around Spring Equinox.



(B) Geometry at equinoxes.

Fig. 4. Solar eclipse.

For a more precise computation, additional geometrical data such as the radius of the sun and the earth-sun distance are needed. The existence of umbra and

penumbra zones also must be recognized. The above simplified derivation is adequate for assessing storage-battery requirements. By keeping in mind that the

maximum eclipse duration at each of the equinoxes is about 80 minutes and adding the lesser contributions during the entire eclipse periods, it is found that the total time of darkness amounts to about 3 hours/year.

NONGEOSTATIONARY ORBITS

Since the necessary rockets and know-how related to the complex in-orbit injection maneuvers were developed, the majority of commercial communications satellites were of the geostationary type until the 1990s. However, because the elevation angle at an earth station drops at higher latitudes, geostationary satellites cannot serve the near-polar regions. For this purpose, other orbits have been used. Polar orbits are also used for non-communications-type missions such as earth observations, weather, surveillance, etc.

From Eq. 19, the maximum latitude that can be served is

$$\beta = \cos^{-1} [(R \cos \theta)/(R + h)] - \theta \quad (\text{Eq. 22})$$

Hence, for a minimum elevation angle θ of 5° , $\beta = 76^\circ 21'$. Thus, inclined orbits are required to serve those regions further north. On the basis of energy considerations, it can be shown that for a given rocket, the lower the perigee the higher the apogee. The slower satellite motion around apogee results in the dual benefit of making communications possible over greater distances on earth and for longer periods of mutual visibility; the tracking problems are also eased.

As previously mentioned, an orbit inclination of 63.5° is advantageous because of the zero rotation of the line of the apsides for this critical angle. The inclined orbit discussed above, 63.5° inclination, 12-hour period, turns out to be highly advantageous. A satellite in such an orbit spends about 8 hours near geosynchronous altitude. Thus three such satellites can provide continuous service at altitude. Further, properly adjusted, the satellites return to the same station every 24 hours. Finally, because of the inclination, at altitude the satellites can be seen in polar regions. The former USSR (and the current Russia) had extensive polar regions and used the *MOLNIYA* series of satellites in the above-mentioned inclined orbit to provide point-to-point telephony (fixed satellite services), data relay, mobile communications, TV distribution, and DBS services very successfully. Variants of this orbit will be used in the United States (for DARS) and in Japan (for multimedia).

Advances in technology permitting high levels of sophistication on board smaller satellites have led to a variety of proposals for constellations of satellites in low earth orbit (LEO) and in medium earth orbit (MEO) to provide mobile or data relay, remote sensing, radio location services, or combinations thereof. Some, such as data relay satellites (*TDRSS*), radiolocation (*LOCSTAR*), mobile communications (*Iridium*, *Globalstar*, and others) and messaging/paging (*ORB-*

COMM) are in service (in all or in part) and others, such as *ICO*, are actively under construction at the date of publication of this book. The basic reason that these systems have emerged in the 1990s is that the cost of the technology required and the reliability of the components (largely digital), combined with their reduced weight and power requirements (of VLSICS), have made such systems practical to build, launch, and maintain in orbit.

Because of the van Allen belts with their much higher radiation levels, which significantly and deleteriously affect the satellite electronics, the choice of orbits is limited. LEOs lie below the inner van Allen belt; MEOs lie between them and GEOs above them.

The LEO systems have certain advantages. They have much shorter earth-satellite paths, which leads to lower eirp requirements and much shorter delays. Further, especially of advantage in mobile/urban communications, elevation angles to the useful part of the LEO are high, thus making for less "multipath." However, their lower orbits require sizable constellations of satellites for continuous service (e.g., 66 for *Iridium*).

MEO systems are, as might be expected, intermediate in their delay and their eirp requirements, and their constellations are smaller. The trade-offs made in system design determine the altitude and inclination of the orbits and the rf frequency/elevation angle/availability considerations combine to define the constellations. (GEO systems are well known and amply described herein.)

ELEVATION AND AZIMUTH ANGLES

For the case of a geostationary satellite and an earth station having latitude ϕ and longitude λ relative to the subsatellite point on the equator, the elevation angle (i.e., the angle above the horizon) and the azimuth angle (i.e., the angle measured clockwise from the direction of true North) can be computed by using standard spherical and plane trigonometry relationships. With reference to Fig. 5 and the spherical triangle $S'EB$ at B , we have

$$\cos \beta = \cos \phi \cos \lambda \quad (\text{Eq. 23})$$

Then considering plane triangle OES , we obtain the elevation angle, θ .

$$\theta = \cos^{-1} \{ [(R + h) \sin \beta] / d \} \quad (\text{Eq. 24})$$

where

$$d = \sqrt{R^2 + (R + h)^2 - 2R(R + h) \cos \beta}$$

The azimuth angle is obtained from the spherical triangle by first using the relationship

$$\gamma = \cos^{-1} (\tan \phi \cot \beta) \quad (\text{Eq. 25})$$

which provides the angle at e ; the true azimuth is then obtained from Table 1. The earth-station quadrant is

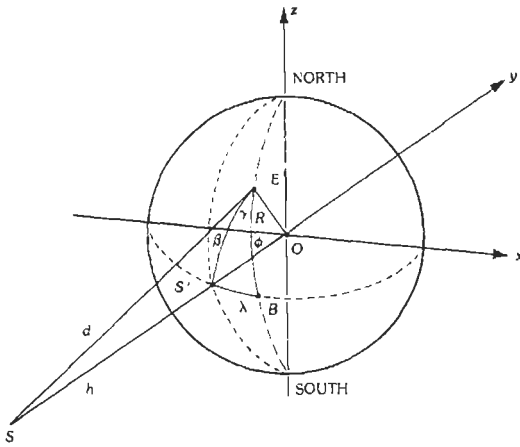


Fig. 5. Geometry of elevation and azimuth angles.

TABLE 1. TRUE AZIMUTH RELATIONSHIPS

True Azimuth	Earth-Station Quadrant
$180 - \gamma$	NW
$180 + \gamma$	NE
γ	SW
$360 - \gamma$	SE

identified with respect to the meridian passing through the subsatellite point and the equator. Fig. 6 provides the look angles to satellites in geostationary orbit.

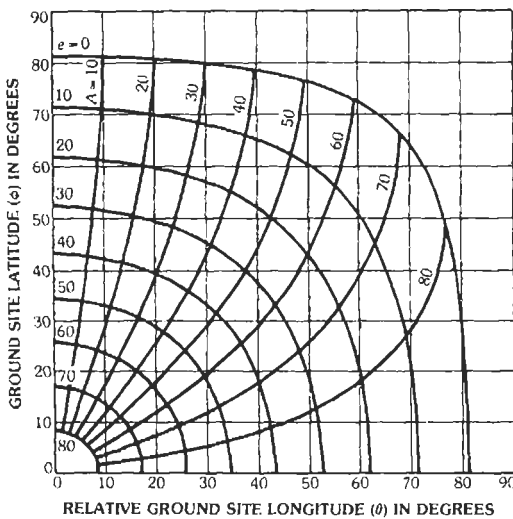


Fig. 6. Look angles to geostationary satellites. (From F. L. Smith III, "A Nomogram for Look Angles to Geostationary Satellites," *IEEE Transactions Aerospace & Electronics Systems*, May 1972, p. 394.)

THE CHOICE OF FREQUENCY

From a purely technical point of view, the choice of the operating frequency is primarily dictated by the external and internal noise sources that contribute to the noise system temperature and by the existence of spectral windows for which signal attenuation is small. External noise originating outside the solar system is designated as "cosmic," whereas sources within the solar system (terrestrial sources included) contribute to the "solar noise." At optical frequencies, the sun acts as a blackbody with temperature of about 6000 K, but below 30 GHz deviations from Planck's law result in much higher equivalent blackbody temperatures. These deviations depend on the solar activity during the 11-year sunspot cycles, with increases up to six orders of magnitude. A strong frequency dependency is observed.

The ionosphere, the troposphere, and in particular the atmosphere, contribute as "terrestrial" sources. Ionospheric effects include Faraday rotation (of linear polarized signals) and gigahertz ionospheric scintillation. These effects fall off sharply with frequency and may be neglected above 10 GHz. They are also diurnally, seasonally, and solar-cycle dependent. Tropospheric effects include tropospheric scintillation, a phenomenon that increases rapidly for decreasing elevation angle below 10° and increases with increasing frequency. It may be neglected for frequencies below 3 GHz except at the lowest elevation angles (< 3°). Hydrometeors contribute to depolarization, attenuation, and increased noise temperature. Rain- and ice-induced depolarization arises because of the nonspherical nature of the hydrometeors, resulting in differential phase shift and attenuation for vertically and horizontally polarized components of the signal. This phenomenon is important down to at least 3 GHz. Increase in attenuation through the atmosphere and in effective sky noise temperature arises from absorption and scattering by the hydrometeors, especially rain. The attenuation reduces signal strength. The increased noise temperature adds to the system noise and is particularly troublesome when a very low noise receiver is used. Increase in attenuation and noise temperature is particularly marked above 10 GHz.

As shown in Fig. 7, the "clear sky" noise temperature observable at the surface of the earth has a broad minimum from 0.8 to 8 GHz and increases rapidly at the lower frequencies on account of terrestrial contributions and at frequencies higher than 10 GHz on account of atmospheric phenomena. Fig. 7 applies to clear-sky conditions; rain and other hydrometeors introduce additional noise as well as signal attenuation. There is a water-vapor absorption-induced peak at about 23 GHz, followed by a decrease in noise temperature to about 33 GHz. There follows a sharp rise resulting from the broad oxygen absorption-induced peak at 59–60 GHz.

The effects of varying the antenna beam inclination angle are also shown. For low values of the elevation

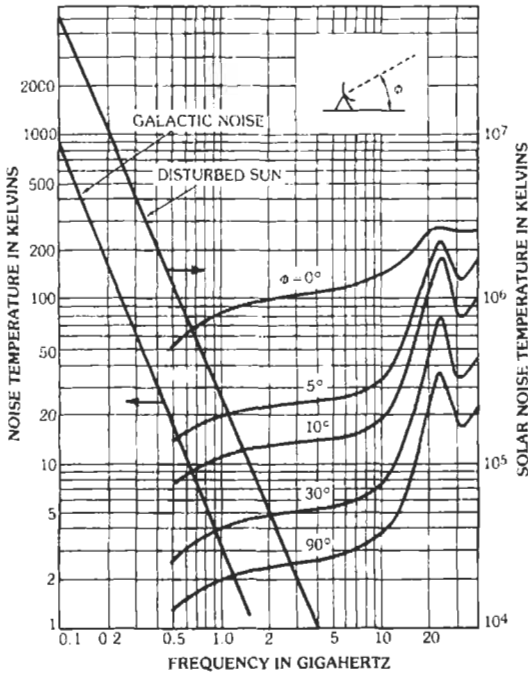


Fig. 7. Noise temperature versus frequency (clear sky).

angle, larger amounts of noise power are collected by the main beam of the antenna in the atmosphere, and, in addition, the antenna side lobes pick up noise radiated by the surface of the earth and that part of the atmospheric noise reflected by it.

A more detailed discussion of the various propagation effects is presented in the section of this chapter entitled "Propagation."

In conclusion, an optimum region for communications systems involving earth-to-space paths exists between 0.8 and 8 GHz. As it is relatively easy to build low-noise, high-gain electronic amplifiers for these frequencies, space communications systems have been preferentially designed to operate in this portion of the spectrum. In the end, it is the allocation of frequency bands to satellite communications at World Administrative Radio Conferences (WARCs) under the aegis of the ITU that determines what frequencies are available. Beyond this, to prevent interference between systems, individual system use must be coordinated by the appropriate entities—the ITU internationally and the local government agencies, such as the FCC in the United States, as well.

The expansion of satellite communications systems has resulted in overcrowding of these frequencies and has forced the opening of higher frequency bands. Allocations in the 10–15 GHz bands have been opened up for the fixed satellite service (FSS) since the late 1970s and are now used extensively. Use of allocations for the broadcasting satellite service (BSS)

in these frequency bands was begun in the 1980s. Efforts in the late 1980s in Japan and Europe and in the 1990s in the United States have opened up the very broad frequency band allocations at 20 and 30 GHz. Despite the severe rain-impairment environment, one can anticipate ever-increasing use of the 20 and 30 GHz bands for FSS and ISS. Bands above this are being used outside the United States. There is always, in this era, a shortage of bandwidth, especially for the mobile services. Allocations to the fixed satellite service, intersatellite service (ISS), broadcasting satellite service, and mobile satellite services (MSS) can be found in Chapter 1.

LINK BUDGETS

Fig. 8 shows a model of an earth-space-earth communications system that comprises an up-link from an earth station to a spacecraft and a down-link from the spacecraft to earth. In the case of space systems, the information rates in the up- and down-links are different because the up-link carries only command and control information, while the down-link carries the information gathered by the spacecraft instruments. However, in satellite communications systems, because the spacecraft functions as a relay station in space, the information flow can be the same in the up-link and the down-link. Messages (voice, video, data, etc.) generated by individual users and converted into electrical signals in the terrestrial network are routed to satellite earth stations, where, after suitable signal processing, modulated RF carriers are beamed to the satellite by a high-gain antenna. Conversely, the signals received from the satellite are processed, fed into the terrestrial network, and eventually distributed to the users.

Although direct computation of the information transmission rate can be done in terms of Eq. 16 for the up- or down-link, or for their combination, it is convenient to derive relationships for the signal-to-noise power ratio in the two links taken separately. In this manner, specific contributions to the link budgets

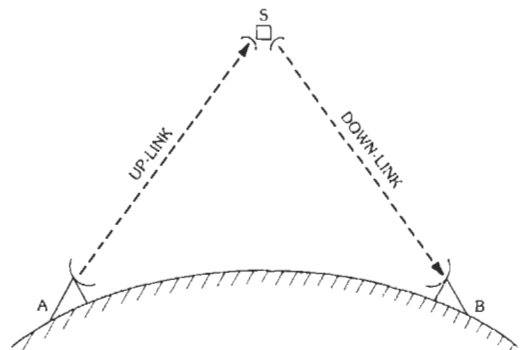


Fig. 8. Earth-space-earth links.

can be identified. In a successive step, the overall signal-to-noise power ratio is obtained by taking into account any additional noise contributions, such as interference or intermodulation. Finally, the number of channels (telephone, TV, or other) is computed for specific combinations of multiplexing, modulation, and multiple access techniques.

For the up-link, the power flux density at the satellite is given by

$$\phi = P_t G_t / 4\pi r_u^2 \quad (\text{Eq. 26})$$

Then the received signal power is

$$\begin{aligned} P_u &= \phi A_{su} = \phi G_{su} \lambda_u^2 / 4\pi \\ &= P_t G_t G_{su} \left(\frac{\lambda_u}{4\pi r_u} \right)^2 \end{aligned} \quad (\text{Eq. 27})$$

where

- A_{su} = effective area of the satellite receive antenna
- G_{su} = gain of the satellite receive antenna
- λ_u = up-link wavelength

The signal is received in the presence of thermal noise contributed by the receiver (internal noise) and the background (external noise). Both contributions can be assumed Gaussian and white; hence the total noise power density is $N_0 = kT_s$, where $k = 1.38 \times 10^{-23}$, J/K is Boltzmann's constant, and T_s is the system temperature. Thus the noise power in the RF transmission bandwidth, B , is $N = N_0 B$. Substitution in Eq. 27 yields the up-link signal-to-noise power ratio:

$$P_u / N_u = \left(\frac{P_t G_t G_{su}}{kT_s B} \right) \left(\frac{\lambda_u}{4\pi r_u} \right)^2 \quad (\text{Eq. 28})$$

and, in decibels,

$$\begin{aligned} P_u / N_u &= 10 \log P_t G_t - 20 \log (\lambda_u / 4\pi r_u) \\ &\quad - 10 \log B + 10 \log (G / T)_s - 10 \log k \end{aligned} \quad (\text{Eq. 29})$$

The first term represents the earth station eirp; the second, the "free space loss"; the third, the RF channel bandwidth; the fourth, the satellite GIT ratio (G_{su} / T_s); and the fifth, Boltzmann's constant.

A similar procedure yields the down-link budget:

$$P_d / N_d = \left(\frac{P_s G_{sd} G}{kT_c B} \right) \left(\frac{\lambda_d}{4\pi r_d} \right)^2 \quad (\text{Eq. 30})$$

and, in decibels, noting that $\frac{G_r}{T_c}$ is $(G / T)_c$,

$$\begin{aligned} P_d / N_d &= 10 \log (P_s G_{sd}) - 20 \log (\lambda_d / 4\pi r_d) \\ &\quad - 10 \log B + 10 \log (G / T)_c - 10 \log k \end{aligned} \quad (\text{Eq. 31})$$

The numerical value of the Boltzmann's-constant term in Eqs. 29 and 31 is

$$10 \log (1.38 \times 10^{-23}) = -228.6 \text{ dBW/Hz} \cdot \text{K}$$

and the "free-space loss" expressed in decibels is

$$L_{fs} = 92.44 + 20 \log r + 20 \log f \quad (\text{Eq. 32})$$

A plot of free-space loss versus frequency for the geostationary orbit appears in Fig. 9. If signal-power-to-system-noise-temperature ratios (dBW/K) are used in lieu of signal-to-noise power ratios (dB), the Boltzmann's-constant term disappears in Eqs. 28 through 31. The RF channel bandwidth also disappears when signal-power-to-noise-power-density ratios (dB · Hz) are used. The choice of the notation is immaterial, but consistency is required in the use of units. Whenever angle modulation is employed, the symbol C (for carrier power) can be used in lieu of P (for signal power), as these two quantities have equal values.

SPACECRAFT ARCHITECTURE

The communications function is supported by the following subsystems:

- A. Main structure
- B. Primary propulsion subsystem, i.e., apogee/perigee motor(s)
- C. Auxiliary propulsion subsystem for station-keeping and orientation
- D. Thermal control subsystem
- E. Power subsystem
- F. Telemetry, command, and control subsystem

The communications subsystem consists of antennas, receivers, and transmitters. Its mass plus that of the power subsystem defines the useful payload mass, M_u . The ratio of this mass and the total mass, M , of the spacecraft in orbit determines the spacecraft utilization factor:

$$u = M_u / M \quad (\text{Eq. 33})$$

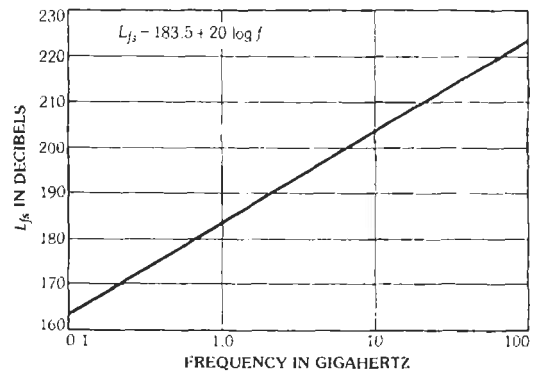


Fig. 9. Free-space loss versus frequency (GEO).

Typical values of u range from 0.35 to 0.60. The difference,

$$M_b = M - M_u = M(1 - u) \quad (\text{Eq. 34})$$

represents the mass of the bus (items A, C, D, and F of the above-mentioned support subsystems). Finally, the spacecraft mass at launch, M_l , is

$$M_l = M + M_b \quad (\text{Eq. 35})$$

where M_b represents the mass of the primary propulsion subsystem.

The three satellite parameters—transmit power, antenna gain, and bandwidth—that contribute to communications capacity are proportional to the total spacecraft mass. Designers' efforts aim at obtaining values as high as possible of the utilization factor (defined by Eq. 33), and for a given communications capacity, the total spacecraft mass should be kept as small as possible to keep launch costs down. Ever-increasing communications capacity can be achieved not just by constructing bigger and more powerful satellites but also by the introduction of advanced technologies in all spacecraft subsystems as well as in the transmission system design.

TRANSPONDERS

Transponders are microwave repeaters carried by communications satellites. Three possible configurations are shown in Fig. 10. The first case (Fig. 10A), in which no signal processing takes place other than heterodyning from the up- to the down-link frequencies and amplification, represents a class of transponders designated as "transparent." It can handle any signal whose format can fit in the transponder bandwidth. The second case (Fig. 10B) involves switching at rf, and the third case (Fig. 10C) involves switching at base band with demultiplexing, demodulation and remodulation, and possibly other signal processing such as routing, buffering, or storage.

The receive and transmit antenna configurations (including frequency reuse by orthogonal polarization and/or separate beams) depend on the earth-station topology and traffic flow. The transponder front end, which includes filters, low-noise amplifiers, local oscillators, mixers, and preamplifiers, is designed for linear operation over the entire bandwidth of the up-link to minimize volume, mass, and power consumption of the down-converter portion of the transponder. Field-effect transistors are commonly used as active circuit elements. Distortionless transmission is ensured by flat amplitude and group-delay (linear phase) responses over the passband. Equalizers (mostly for compensation of group delay distortion) are also used. However, limiters are used when anti-jam protection is needed.

Various contributions enter into the system noise temperature. If the receive antenna input port is taken

as the reference point for the configuration shown in Fig. 11, the system noise temperature is

$$T_s = T_{\text{ANT}} + [(L - 1) + L(F_{pr} - 1) + L(F_r - 1)/G_{pr}]T_{\text{ref}} \quad (\text{Eq. 36})$$

where

- T_{ANT} = antenna noise temperature
- L = transmission line (waveguide) loss
- F_{pr} = preamplifier noise factor
- G_{pr} = preamplifier power gain
- F_r = receiver (down-converter) noise figure
- T_{ref} = 290 K

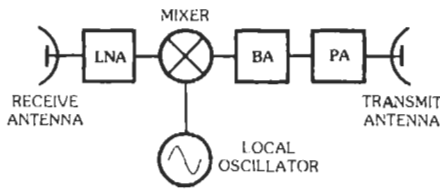
For a satellite antenna looking at the earth, the contribution of the term T_{ANT} is near 290 K. Low-loss waveguides and high-gain, low-noise preamplifiers are clearly desirable. After the signal is converted to the down-link frequency and preamplified, its level needs to be raised to the required rf power output. Redundant receivers are used to ensure survivability over the planned spacecraft lifetime.

The linearity requirements, which are fairly easy to achieve in the receiver, are difficult to meet in the power amplifiers. Transmitters having output power from 5 to 20 watts are typically used in communications satellites and in deep-space probes; hundreds of watts will be used in satellites for direct broadcasting. The high efficiency in the energy-conversion process from dc to rf that is desirable for maximum utilization of spacecraft mass and prime power conflicts with the linearity requirements.

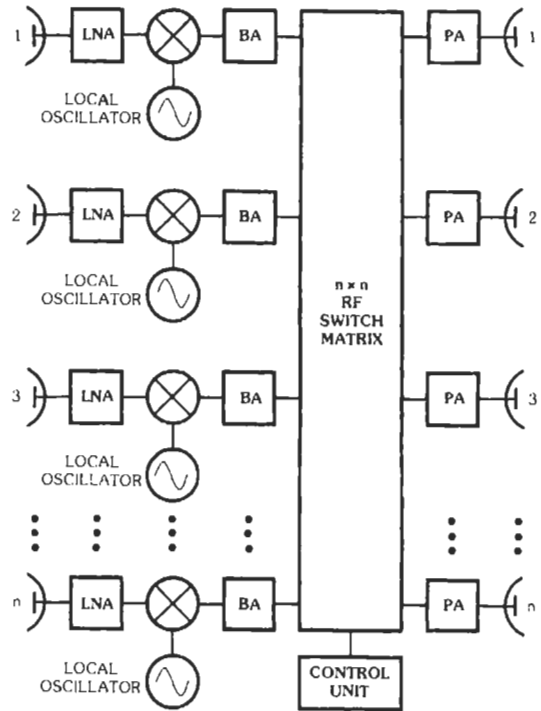
Traveling-wave tube amplifiers (TWTAs) dominated as satellite high-power amplifiers (HPAs) into the 1980s because of their wide-band, high-gain, and high-efficiency characteristics accompanied by their light weight, long life, and high reliability. They remain the satellite HPA of choice when higher power is needed, and at higher frequencies (above 15 GHz). The typical TWTA characteristics shown in Fig. 12 indicate maximum conversion efficiency at saturation with departure from linearity. Reduction of the input signal level yields linear, or quasilinear, operation at lower conversion efficiency and power output.

Traveling-wave tubes also produce AM-to-PM modulation conversion effects because the phase of the output signal is affected by the input signal amplitude. Intermodulation noise caused by amplitude nonlinearity is maximum when the TWT is driven at saturation. Reducing the input signal amplitude yields an almost linear operation, but the interaction of the electron beam and the wave advancing along the helix structure of the tube is such that AM-to-PM conversion effects are worst in the quasilinear region of the tube characteristics.

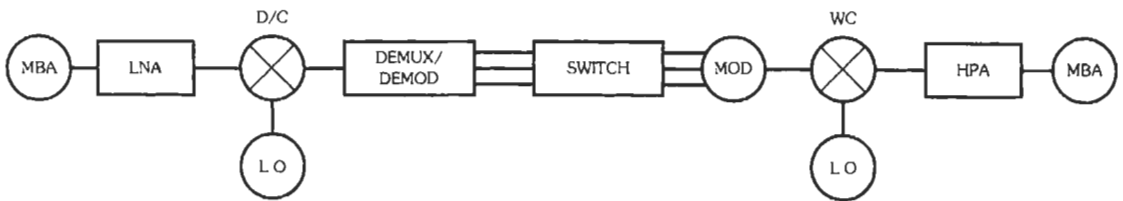
When several modulated carriers are present at the input, the resultant signal envelope fluctuates and produces phase variations in the output signal, resulting in intermodulation products. Third-order intermodulation products of the form $(f_1 + f_2 - f_3)$ and $(2f_1 - f_2)$ are



(A) Single conversion.



(B) With switching at RF.



(C) Baseband processing.

Fig. 10. Transponders.

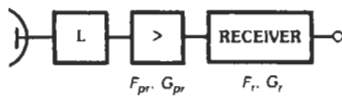


Fig. 11. Receiver front end.

objectionable because they fall inside the transponder passband; the amplitude of the $(f_1 + f_2 - f_3)$ product is about 6 dB higher than that of the $(2f_1 - f_2)$ product. For a small number of carriers, the third-order intermodulation products are mainly dictated by the coefficient of the third-power term of the power series expansion representing the nonlinear element. A 1-dB

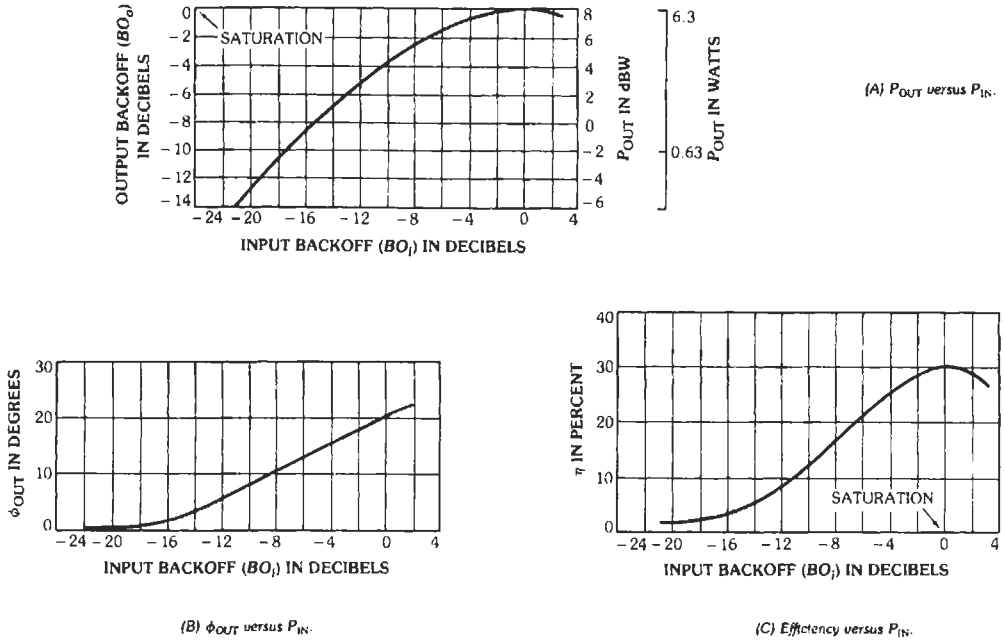


Fig. 12. TWTA characteristics.

change in the input produces a 3-dB change in the third-order intermodulation product, and hence a 2-dB change in the output carrier-to-noise ratio. When the number of carriers is increased, the intermodulation products are compressed toward the limit case of $n = \infty$, and for large n , the fifth-order product ($3f_1 - 2f_2$) becomes negligible compared with third-order products. In practice, an input level is chosen that leads to a compromise between the amplitude-to-amplitude and amplitude-to-phase distortions.

Advances in solid-state devices have made it possible to replace TWTAs with solid-state power amplifiers (SSPAs), usually using FETs at frequencies up to 15 GHz, for all but the higher-power applications. However, the limits of frequencies and power for SSPAs are continually moving upward as new structures emerge from research work in solid-state devices and production techniques. New devices include high-electron-mobility transistors (HEMTs), and new techniques include molecular beam epitaxial (MBE) deposition and electron-beam writing of circuit/transistor structures. Although the characteristics of SSPAs differ from those of TWTs, the two above-mentioned sources of signal distortion and intermodulation noise (AM to AM and AM to PM) can still be identified, and, in general, intermodulation noise is lower in solid-state amplifiers. Complementary nonlinear elements cascade-connected with nonlinear power amplifiers can provide considerable overall linearity improvement.

The down-link bandwidth is divided into sub-bands, each one of them accepting a limited number of mod-

ulated carriers to keep intermodulation products under control. Input demultiplexers separate the down-converted signals and distribute them among several transmitters, each one using a fraction of the entire down-link bandwidth. Output multiplexers are used at the interface between the transmit power amplifiers and the transmit antenna-feed waveguides.

Transponders are usually counted by the number of separate transmit channels in the down-link and the corresponding portions of the up-link wide-band receivers.

For given amounts of power and mass available for the communications subsystem, the number of transponders and their bandwidth depends on trade-offs of power, bandwidth, and intermodulation noise. With a fraction of the bandwidth lost to guard bands, the total usable bandwidth decreases when the number of transponders increases. A similar situation arises in regard to the available power per transponder, on account of the direct arithmetical proportion and also because of the increased complexity of the distribution network.

Transponders of equal bandwidth are advantageous for fixed allocations of traffic but do not provide flexibility when a network is reconfigured within a given coverage area or when a satellite is moved to serve an area characterized by different network topology and traffic requirements.

Until the mid 1970s, the prevalent use of FDM/FM/FDMA transmission led to a standardization of transponder bandwidth around 36 MHz and to the number of transponders per spacecraft of around 12 for

single-polarization satellites and 20–24 for dual-polarization satellites. The introduction of the 11/14-GHz frequency bands and of high-speed digital transmission has accelerated the use of wider-band transponders (80 and 200 MHz) in addition to those having 36–40 MHz nominal bandwidth.

The increasing number of transponders (reaching 46 in *INTELSAT VI*) needed to provide ever greater amounts of communications capacity has required advances in the design of flight-qualified microwave filters and multiplexers. Light weight has been achieved by using thin-wall Invar or graphite-fiber-reinforced plastic cavities. The requirement for flat amplitude and group time delay necessary to maintain low signal distortion and minimal crosstalk contrasts with the high skirt selectivity needed for efficient spectrum utilization. The design has evolved from conventional single-mode cascade rectangular-cavities Chebychev filters (Fig. 13) to longitudinal cross-coupled dual-mode circular-cavities filters as shown in Fig. 14, improved skirt selectivity being obtained at the expense of a finite but sufficiently large out-of-band attenuation. True elliptical response can be achieved in the canonical dual-mode filter shown in Fig. 15, in which the same physical cavity provides the input and output ports. This arrangement is also mechanically advantageous because of mass and volume reduction. Further reduction in mass and volume

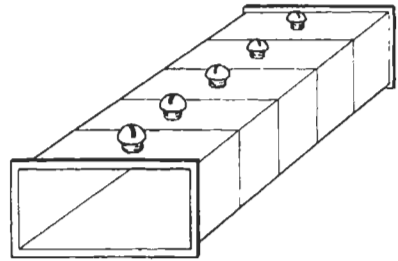


Fig. 13. Single-mode cascaded-cavities filter.

can be achieved by use of dielectric pucks, which drastically reduce cavity size. Another approach uses still higher-order-mode filters, such as those of fourth order. Finally, recent advances in high-*T_c* superconductors have permitted development of microwave filters employing such materials. These offer still greater promise, but practical realization for flight use is probably a decade off.

Technology is evolving at such a rapid pace that it is useful to look ahead a few years. One is never certain about the technological future. But, based on what is currently emerging from the R&D facilities, it is possible to indicate what looks likely to be available for commercial use in space in the next 5–10 years.

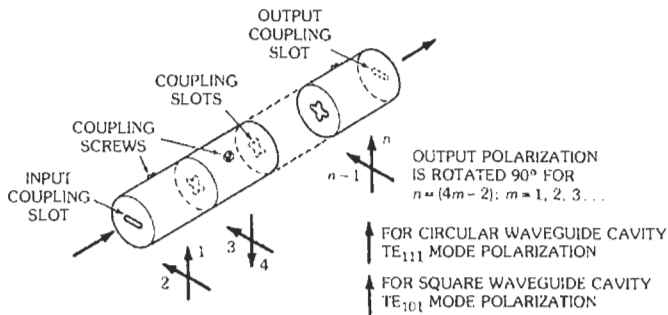


Fig. 14. Coupled dual-mode circular-cavities filter.

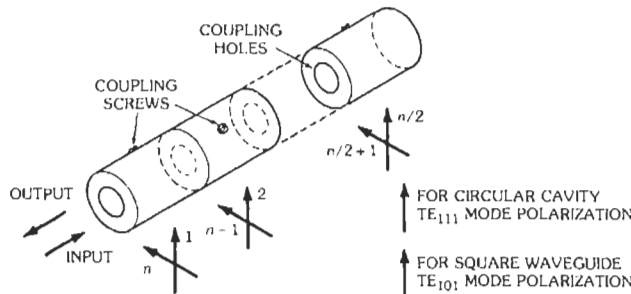


Fig. 15. Canonical dual-mode filter.

In recent years, the evolution of the transponder from a simple microwave repeater to a node, which provides a variety of communications processing functions, has been progressing rapidly. The initial phases of development of onboard processing (OBP) included dc- and remodulation and satellite-switched TDMA (detailed in its own section below), followed by deployment and service. Other forms of OBP rerouting have been developed, including OBP ATM (Asynchronous Transfer Mode) switches and packet switches. OBP error correction and de- and recoding are also in the OBP mix. An OBP ATM switch is scheduled for flight and on-orbit testing. Deployment can be anticipated. What can be anticipated for the first 10 years of the 21st century is the use of a wide variety of transponder types ranging from microwave repeaters used for single links, to multicast in their simpler forms, all the way to full OBP processor-equipped payloads. These latter will be able to receive the signals from a multibeam multifrequency antenna farm, take them down to baseband, perform error correction, reroute them as needed, encode and modulate the newly formed bit streams, and transmit them through the multibeam multifrequency antenna farms to their destinations. Every element needed for these latter payloads either has been developed or is under development. Obviously system and service needs will define the transponder types deployed in any system. But the capability will be there.

One further type of transponder needs mention at this time. Laser and optical technology has reached the point where commercial deployment of laser communications links in space and satellite applications awaits only the final demonstration. Equipment has been deployed in space for such a demonstration, for example, on the Japanese *ETS-VI*, which unfortunately failed to reach the planned orbit. Currently, the *SPOT-4* satellite carries an optical package in space, awaiting the launch of the *ARTEMIS* satellite by ESA, which will carry the already developed *SILEX* optical system with which it will communicate. Further, the Japanese *OICETS*, already ready, will be launched to also work with the *SILEX* system. One can foresee that the success of these efforts will presage the commercial use of optical transponders and links in space. The more difficult links, optically, are earth-space, because of atmospheric impairments to the signal. However, if current ongoing efforts in Japan, Europe, and the United States to develop high-altitude platform systems (HAPS) that hover above the stratosphere at altitudes above 18 km are successful, links between satellites and such platforms will be eminently feasible.

OVERALL TRANSMISSION SYSTEM CONSIDERATIONS

With transparent transponders and analog transmission, the total signal- (carrier-) to-noise power ratio can be computed from the expression

$$(C/N_{total})^{-1} = (C/N_{up})^{-1} + (C/N_{down})^{-1} + (C/N_{im})^{-1} + (C/N_{in})^{-1} \quad (\text{Eq. 37})$$

where, in addition to the up- and down-link noise, the third term represents intermodulation noise, and the fourth term represents interfering signals. In the presence of intermodulation noise, C/N_{total} can be maximized by reducing the transponder input drive. Backing off the TWTA reduces C/N_{up} and also C/N_{down} (via the power input/output relationship of the transponder), but as C/N_{im} increases rapidly when the input drive is reduced, an optimum value of C/N_{total} is obtained at a specific backoff level, as shown in Fig. 16. Interference noise is kept down by proper antenna design, transponder selectivity, and planned coordination of satellite systems sharing given frequency allocations.

Satellite communications systems involve the three interacting processes of multiplexing, modulation, and multiple access. Signals appearing at the interface of the terrestrial and the satellite transmission system include single voice channels, groups and supergroups of voice channels, analog TV channels, and digital bit streams. Multiplexing combines these diverse signals into a composite baseband signal with individual inputs distinguishable in the frequency or time domain. Modulation can be either analog (i.e., continuous) or discrete (digital), and in the latter case, sampling and quantization occur prior to the modulation

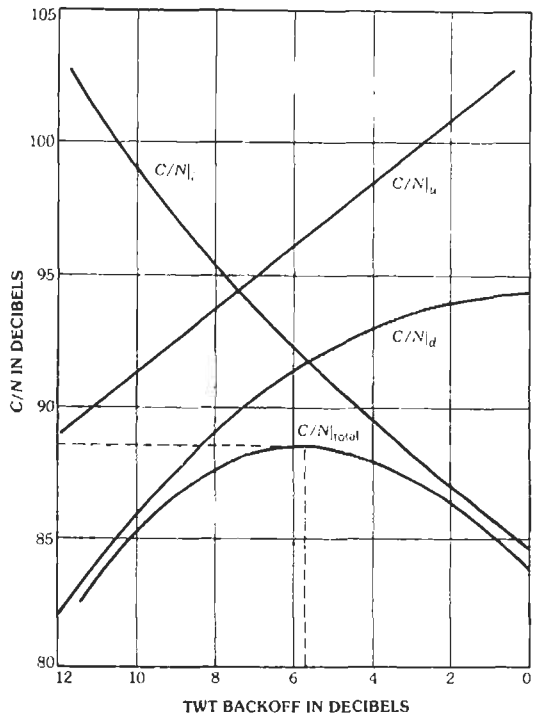


Fig. 16. C/N ratios.

process itself. Amplitude modulation has been little used; analog phase or frequency modulation or digital BPSK and QPSK modulations have been widely used.

Multiple access is a unique feature of satellite communications: When m earth stations are "visible" from a satellite, transponders can be shared by several earth-station pairs. The number of possible connections is then

$$n = m(m-1)/2 \quad (\text{Eq. 38})$$

This "n-port" network feature is very attractive in comparison with the inflexibility of the "2-port" configuration encountered in most earth-based communications systems.

Individual accesses are kept separate in a physical domain such as frequency, time, or space (i.e., by separate antenna beams) or by encoding. The accesses can be preassigned in a static mode or dynamically assigned on demand. Fig. 17 illustrates three principal modes of multiple access, and Table 2 outlines their major features.

Because modulation and multiple access are interactive processes, the choice of a specific method of multiple access depends on the best utilization of the available power and bandwidth. The number of accesses should be made as large as possible while signal impairments caused by imperfections of the transmission channel are kept within specified margins. The

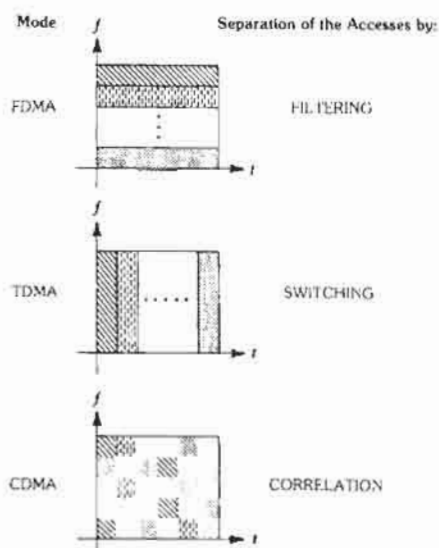


Fig. 17. Multiple-access schemes.

overall system design is also influenced by the characteristics of the terrestrial interface. Analog voice and video transmission systems predominant in the mid-1960s in terrestrial systems led to the adoption of

TABLE 2. FEATURES OF MULTIPLE-ACCESS SCHEMES

Type	Characteristics	Advantages	Disadvantages
FDMA	Constant envelope signals Angle modulation (analog or digital) Separation by filtering	Easy interfacing with terrestrial analog facilities No need for synchronization Simple earth-station equipment	Intermodulation noise in nonlinear amplifiers reduces communications capacity Backoff required Up-link power coordination is required Transmission plan is difficult to reconfigure
TDMA	Bursts from and to different stations do not overlap Only one burst present at a given time Separation by time gating	Fairly easy interfacing with terrestrial digital facilities Power amplifiers operate near saturation at higher efficiency Highest communications capacity No need for power control and coordination Transmission plan is flexible	Need for coordinated synchronization More complex earth-station equipment A/D and D/A conversion required
CDMA	Constant envelope signals Each access takes the whole transponder bandwidth Separation by correlation	Spread spectrum provides protection against jamming	Synchronization is required Communications capacity is somewhat limited Signal quality worsens when number of accesses increases Power coordination is needed

FDM/FM/FDMA techniques in most commercial satellite systems, notwithstanding the loss of communications capacity resulting from the nonlinear channel characteristics when several carriers occupy a transponder. The expansion of digital systems, accelerated by the advances in solid-state devices that have made possible rugged, reliable, and cost-effective signal-processing equipment, has led to a wider use of time-division multiple access (TDMA). As noted above, digital technology has advanced to the point where digital transmission is now dominant. Other multiple-access techniques such as code-division multiple access (CDMA) have proved to be competitive with TDMA.

Digital transmission offers the following major advantages:

- A. Guarantee of error control
- B. Reduced sensitivity to channel nonlinearity
- C. Efficient trade-off of power and bandwidth
- D. Flexibility with regard to multiplexing diverse signals
- E. Easy combination of the functions of transmission switching and routing
- F. Capability of signal regeneration
- G. Implementation with rugged hardware

In satellite systems, points A, B, C, F, and G are especially meaningful. The up- and down-link can be separately designed to provide specified amounts of bit error rate (BER). For a given type of modulation, BER is primarily a function of the carrier-to-noise power ratio, C/N (or the energy-per-bit-to-noise-power-density ratio, E/N_0), in the channel. Assuming independent noise performance of the links, the BER of their combination is the sum of the individual BERs. This condition is satisfied with repeaters capable of providing signal regeneration in the case of "transparent repeaters." Otherwise, the up-link BER must be much smaller than the down-link BER.

Digital methods are especially attractive in demand assigned multiple access (DAMA) and TDMA systems.

ANALOG TRANSMISSION SYSTEMS

In satellite systems, multiplexing, modulation, multiple-access methods, and the channel characteristics interact in the determination of communications capacity. With analog transmission, the quality of the signal delivered to the user is a function of the post-detection signal-to-noise ratio, S/N , through a general functional relationship of the form

$$S/N = f(C/N_0, CH, SIG, MUX, MOD, MA) \tag{Eq. 39}$$

where

C/N_0 = carrier-to-noise power density ratio in the RF channel

CH = RF channel characteristics (amplitude and phase, nonlinearity)

SIG = original signal descriptors

MUX = type of multiplexing

MOD = type of modulation

MA = type of multiple access

System performance can be assessed with reference to the Shannon bound (rate distortion theory and communications capacity) and, in practice, by the number of telex, telephone, or video channels that can be handled.

FDM/FM (Single-Carrier Case)

When a transponder handles a single carrier modulated by a baseband signal resulting from frequency-division multiplexing of numerous voice channels and with no interference, Eq. 37 simplifies to

$$(C/N|_{tot})^{-1} = (C/N|_{up})^{-1} + (C/N|_{down})^{-1} \tag{Eq. 40}$$

The signal-to-noise power ratio in the worst (top) voice channel at the output of the FM demodulator can be computed by using the appropriate FM equations.

First, the weighted signal-to-noise ratio for a 1-mW test tone is computed by taking into account ITU-R (formerly CCIR) recommendation on satellite systems channel noise performance, which requires the average noise power not to exceed 10 000 pW psophometrically weighted (pWp). Allowing 1000 pWp for the earth-station equipment noise, 1000 pWp for interference noise, and 500 pWp for the earth-station out-of-band noise, the balance of 7500 pWp is the noise allowed for the combination of the up- and down-links. Hence,

$$S/N = 10^{-3}/(7500 \times 10^{-12}) = 1.33 \times 10^5 \tag{Eq. 41}$$

or

$$S/N_{dB} = 51.25 \text{ dB} \tag{Eq. 42}$$

The relationship between the signal-to-noise ratios at the output and the input of the FM detector operating beyond threshold is*

$$S/N = (C/N)B/b(f_r/f_m)^2 pW \tag{Eq. 43}$$

where

S/N = weighted signal-to-noise power ratio at 1-mW test-tone level

C/N = carrier-to-noise power ratio over the RF channel bandwidth

B = RF channel bandwidth

b = audio channel bandwidth (3.1 kHz)

f_r = rms test-tone frequency deviation

* For a limiter-discriminator FM detector, the threshold is often taken as $(C/N) = 10$ to 13 dB. FMFB and PLL detectors can provide up to 3 to 5 dB of threshold extension.

- f_m = maximum baseband frequency $\approx 4.2 \times n$ kHz, with n the number of voice channels
- p = psophometric factor (1.78 or 2.5 dB)
- W = preemphasis factor (2.5 or 4.0 dB)

Carson's rule defines the RF channel bandwidth as

$$B = 2(f_p + f_m) \quad (\text{Eq. 44})$$

where

- f_p is the peak frequency deviation,
- f_m is the maximum frequency component of the baseband modulating signal,

Because this signal is Gaussianly distributed, the peak frequency deviation is undefined; however, a peak factor p , which is the ratio of the peak to the rms frequency deviations, can be assumed, with values from 3.16 (10 dB) to 8.5 (18.6 dB), the lower value corresponding to large n .

If we consider the load factor

$$L = \begin{cases} -15 + 10 \log n, & n \geq 240 \\ -1 + 4 \log n, & 12 < n < 240 \end{cases} \quad (\text{Eq. 45})$$

which yields the average power of the baseband signal in dBm0, the quantity

$$g = \text{antilog}(L/20) \quad (\text{Eq. 46})$$

is the ratio of the rms multichannel frequency deviation to the test-tone deviation. By introducing both factors in Eq. 44, the value of the rms test-tone deviation is determined:

$$f_r = (B/2 - f_m)/\rho g \quad (\text{Eq. 47})$$

As the value of B is fixed for a given transponder design, the number of voice channels per transponder can be computed by an iterative procedure. Starting with an estimate for n , computed values of f_m and f_r are introduced into Eq. 43 together with the RF channel C/N power ratio. If the resulting S/N turns out to be greater (or smaller) than the prescribed 51.25 dB, the estimate for n was too low (or too high) and needs to be changed until the correct value of S/N is found.

In practice, each voice channel requires from 40 to 50 kHz of RF bandwidth; hence the channel density falls in the range from 20 to 25 voice channels per megahertz. The use of companded FDM/FM permits greater capacity; typically a doubling of the capacity can be achieved for a given transponder bandwidth.

FDM/FM/FDMA (Multicarrier Case)

As transponder bandwidth and power are shared among the modulated carriers, the 7500 pWp must include the intermodulation noise arising from the transponder nonlinearity. Although the intermodulation noise is reduced by backing off the traveling-wave-tube power amplifier, the ultimate result is a reduction of the number of voice channels that can be handled by a transponder in proportion to the number of accesses. The smaller value of the loading factor contributes also to reducing communications capacity. Fig. 18 illustrates the performance of a typical commercial satellite using FDM/FM/FDMA in a global beam.

Within the constraints of the available bandwidth and power, and by taking into consideration interference and frequency-coordination problems, the actual value of the maximum transponder capacity depends on trade-offs among the following causes of signal impairment:

- A. Spacecraft TWTA impairments, which include in-band intermodulation products arising from both amplitude and phase nonlinearity and intelligible crosstalk due to AM-PM conversion.
- B. Frequency-modulation transmission impairments not attributable to the TWTA characteristics. These impairments include adjacent-channel interference due to spectral overlap, which produces convolution and impulse noise at baseband; dual paths between transponders on a given spacecraft; interference produced by intermodulation in adjacent transponders; co-channel interference in frequency reuse systems; and earth station RF out-of-band emission.

At the earth stations, the baseband signals frequency modulate preassigned multdestination carriers, which are transmitted via satellite to various receiving stations. After demodulation of the RF signals and demultiplexing of the baseband, the individual voice channels are recovered. This transmission method is rather inflexible with respect to changing traffic requirements, but it has been widely used for high-density trunks.

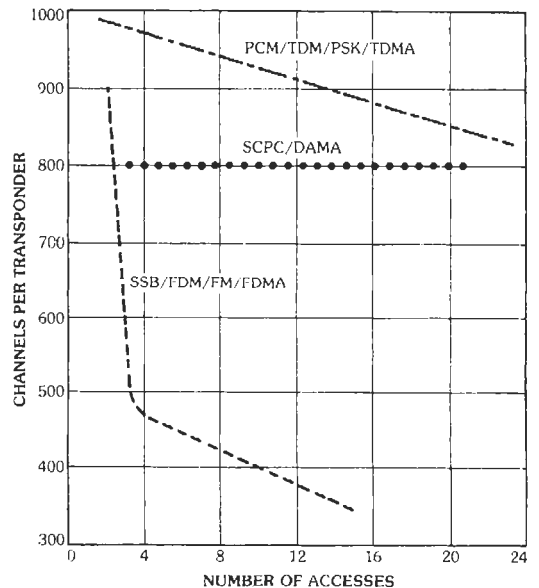


Fig. 18. Multiple-access systems performance.

DIGITAL TRANSMISSION SYSTEMS

Since

$$C/N = (E/N_0)(R/B) \tag{Eq. 48}$$

where

- C = carrier power
- N = noise power
- E = energy per bit
- N₀ = noise power density
- R = transmission rate
- B = channel bandwidth

the transmission rate in digital systems is limited either by power or by bandwidth constraints. In the power-limited case, and for a given modulation type, the bit error probability, *p_e*, is a function of *E/N*, and, hence, of *C/N₀R*; thus, the allowable rate is

$$R_p = f(C/N_0, p_e) \tag{Eq. 49}$$

where the subscript *p* indicates a power-limited situation.

In the band-limited case

$$R_b = F(B, \text{MOD}) \tag{Eq. 50}$$

where the subscript *b* indicates a band-limited situation and MOD indicates the modulation type.

Thus, in general

$$R_{\max} = \min(R_p, R_b) = f(\text{MOD}, C/N_0, B, p_e) \tag{Eq. 51}$$

Expressing bit error probability as a function of *E/N₀* determines modulation performance. In general,

$$p_e = f(E/N_0) = (1/2) \operatorname{erfc} [(E/2N_0)(1-\rho)]^{1/2} \tag{Eq. 52}$$

where

$$\operatorname{erfc}(x) = (2/\sqrt{\pi}) \int_x^\infty e^{-t^2} dt$$

and

$$-1 \leq \rho = \left(\frac{1}{T}\right) \int_{-\infty}^{+\infty} s_1(t)s_2(t) dt \leq 1$$

where the normalized correlation coefficient of the waveforms *s₁(t)* and *s₂(t)* are used to represent zeros and ones.

The most frequently used digital modulation systems are binary phase shift keying (BPSK) and quaternary phase shift keying (QPSK), for both of which

$$p_e = (1/2) \operatorname{erfc} [E/N_0]^{1/2} \tag{Eq. 53}$$

since in QPSK the symbol energy is twice the bit energy. Various forms of modulation characterized by constant envelope have been proposed to achieve more

efficient spectrum utilization and to avoid the amplitude changes occurring in BPSK and QPSK, which are objectionable in band-limited nonlinear systems. Plots of the bit error probability, *p_e*, as a function of *E/N₀* and various types of modulation appear in Fig. 19. The *R/B* ratio, which is a function of *E/N₀*, provides a measure of bandwidth utilization efficiency and, since *R* cannot exceed Shannon's channel capacity

$$C = B \log_2 [1 + C/(N_0B)]$$

a lower bound for *E/N₀* is

$$\frac{E}{N_0} \geq [\exp(0.69R/B) - 1] / (R/B)$$

In practice, since *R/C* << 1, *E/N₀* must be well above 0 dB. In power-limited situations, when the desired transmission rate cannot be achieved with a prescribed value of bit error probability, channel coding (forward error correction (FEC)) can be helpful. Error control in the channel requires the addition of redundant bits to the information bits and, hence, an increase of the overall rate and of the required bandwidth. Under certain conditions, the receiver can exploit the redundancy and significantly reduce the error probability of the recovered message.

Various kinds of codes have been proposed and used in satellite systems; the two major distinct approaches are block and convolutional codes. Although the

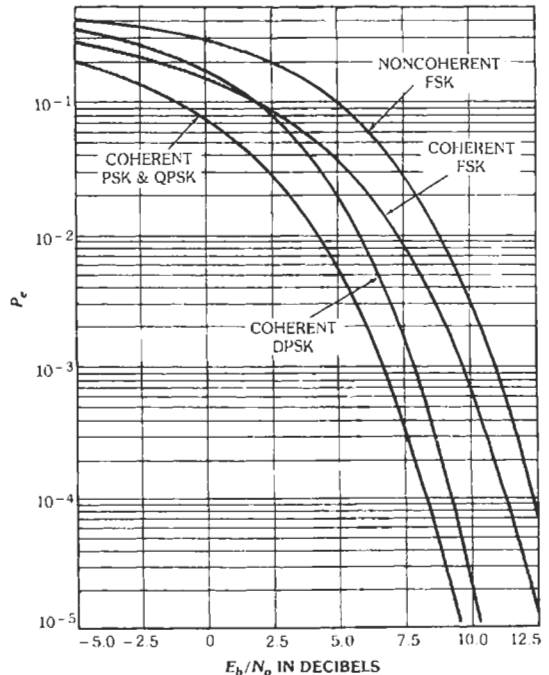


Fig. 19. Curves of BER versus *E_b/N₀*. (From H. L. Van Trees, Ed., *Satellite Communications*. New York: IEEE Press, 1979, Fig. 12, p. 76.)

former may be attractive in certain cases, the latter have been found preferable because of ease of implementation and availability of efficient schemes such as Viterbi and sequential decoding.

In the power-limited case, the link budget is calculated from Eq. 48:

$$R_p = (C/N_0) \cdot (N_0/E) \quad (\text{Eq. 54})$$

which is often written in logarithmic form as

$$\begin{aligned} R_p |_{\text{dB}} &= C/N_0 |_{\text{dB-Hz}} - E/N_0 |_{\text{dB}} \\ &= \text{eirp} + G/T - \text{free space loss} \\ &\quad + 228.6 - E/N_0 |_{\text{dB}} \end{aligned} \quad (\text{Eq. 55})$$

For the bandwidth-limited case:

$$R_p = B \cdot \log_2 m R_s / B \quad (\text{Eq. 56})$$

where

m = number of bits/symbol
 R_s = channel symbol rate (symbols/s or bauds)

Using logarithmic units

$$R_p |_{\text{dB}} = B |_{\text{dB}} + 10 \log_{10}(\log_2 m) - B/R_s |_{\text{dB}} \quad (\text{Eq. 57})$$

The lower bound of Eqs. 56 and 57 is the actual communications capacity.

The two operations of modulation and coding can be combined in a process known as "coded trellis signaling." By mating FEC coding to the modulation signal space, the minimum Euclidean distance between "words" can be maximized with the result of enhanced bit error rate performance versus E/N_0 . In this unified approach, redundancy for FEC coding is achieved without sacrificing bandwidth by first doubling the dimension of the modulation signal space to an M -ary alphabet ($M = 2^n$). Then a convolutional code of rate $(n-1)/n$ is used to obtain the coding gain that provides better utilization of the available power.

At present, QPSK is the most commonly used modulation method. The advent of reliable low-power LSICs and VLSICs and, more recently, programmable-gate arrays has made practical implementation of soft decision decoding (Viterbi decoding) with a significant (≈ 3 dB) improvement in coding gain possible. At high rates, recent developments include coded octal PSK (COPSK) and other higher-order modulation methods. Coding rates vary from rate 1/2 to rate 15/16; use of inner and outer codes at different rates permits quite efficient and high-coding-gain systems with very low error rates ($\approx 10^{-10}$ or better).

DEMAND-ASSIGNED MULTIPLE ACCESS (DAMA)

In DAMA systems, satellite circuits are assigned to earth-station pairs upon request. Each voice channel is transmitted on a separate carrier, which is activated by the speech input. Thus, transponder power and band-

width are more efficiently used than in preassigned systems. Modulation can be analog or digital, the latter being preferred, with either PCM or AM voice encoding combined with BPSK or QPSK. Carriers are taken from a pool of frequencies that can be shared by the entire network of earth stations (fully variable DAMA) or limitedly assigned to certain destinations or origins (semivariable DAMA). The choice of the approach impacts on the complexity of the earth-station equipment. These techniques are particularly applicable to thin-route traffic situations.

An early example of the use of DAMA techniques was the SPADE system of *INTELSAT*. The terminal shown in Fig. 20 comprises a terrestrial interface unit, several transmit/receive units, a signaling and switching unit, and an IF subsystem. Through the terrestrial interface unit, calls originated in the local transmit center activate the signaling and switching unit, which selects a pair of frequencies and alerts the destination. Modulation/demodulation and A/D and D/A functions were performed by QPSK modems and PCM codecs. Single-channel per carrier transmission (SCPC) was used.* The common signaling system operated in a TDMA broadcast mode at 128 kbits with BPSK modulation at an error rate of 10^{-7} , which is three orders of magnitude better than that of the 64-kbits QPSK communications channel. A variety of signaling and switching systems as used by different countries can be accommodated. Fig. 21 shows the signaling flow, and Fig. 22 shows the communications flow.

As shown in Fig. 18, the SPADE system provides 800 voice channels with an *INTELSAT IV* type transponder connected to a global beam antenna independently of the number of accesses. A voice activation factor of 0.4 implies that when all 800 channels are in use to provide up to 400 two-way conversations, only 320 channels are simultaneously active 90 percent of the time; thus, a power savings of $800/320 = 2.5$, or 4 dB, is achieved.

Although DAMA systems such as SPADE (providing FSS) have limited utility in the current environment, DAMA is a basic strategy. As such, DAMA is used as a set-up protocol, in mobile services, for example.

TIME-DIVISION MULTIPLE ACCESS (TDMA)

In time-division multiple access (TDMA), transponder power and bandwidth are shared by several earth stations. Each station transmits RF bursts at the same carrier frequency at different times; the carriers are modulated by signals coming from different sources. Because all stations are synchronized, only one burst occupies the transponder at any given time, and bursts from different stations never overlap. Power amplifiers can be operated at saturation, i.e., at maximum effi-

* SCPC/FM/FDMA transmission is also used in the *INTELSAT* preassigned network for small earth terminals.

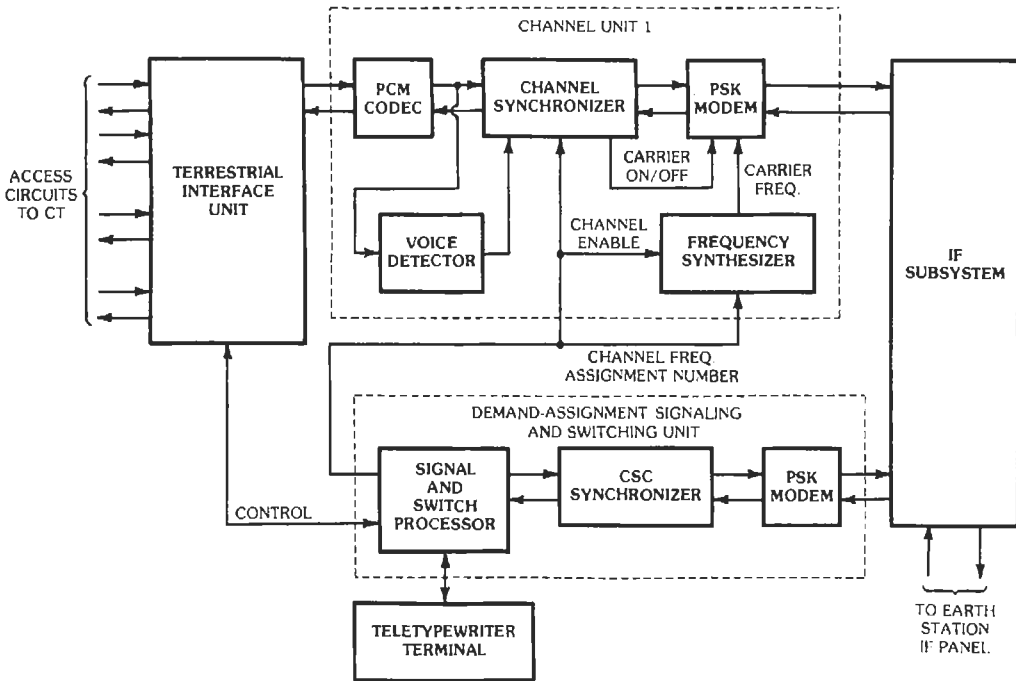


Fig. 20. SPADE terminal block diagram. (From COMSAT Technical Review, Vol. 2, No. 1, p. 226, Spring 1972.)

ciency. Consequently, since intermodulation noise is absent, communications capacity is higher than that achievable with analog FDM/FM/FDMA. The bursts are amplified by the satellite transponder and beamed down to earth to various earth stations. The down-link bursts are received and detected, and the demultiplexed signals are delivered to their destinations. Other advantages of TDMA include compatibility with terrestrial digital systems and flexibility in accommodating transmission systems growth. The requirements of up-link power control, which are very stringent in FDM/FM/FDMA systems, are greatly reduced in TDMA systems. However, when bandwidth is only slightly greater than the signal symbol rate, the combined effect of transmit-side filtering and nonlinearities produces intersymbol interference and signal distortion that cannot be eliminated by linear filtering at the receive side. In such a case, backoff of the power amplifier is still required to keep signal degradation within limits.

Earth stations have parallel steady inputs, some already in digital form and others in analog form that need to be converted to digital format. Since each station transmits periodically in bursts at a rate much higher than the bit rate of any of the input signals, memory circuits are needed for data buffering. The burst periodicity determines the TDMA frame duration; thus, for voice inputs ($f_{max} = 4 \text{ kHz}$), the minimum value of the TDMA frame duration corresponds to a

Nyquist interval of $125 \mu\text{s}$. Multiples of this minimum duration can be used, leading to increased communications capacity through trade-offs of TDMA frame duration and amounts of memory capacity. An upper bound to frame duration is determined by the transmission delay from earth station to earth station (about 270 ms for geostationary satellites at 10° elevation).

Table 3 illustrates trends of operational and planned TDMA systems. Fig. 23 shows the TDMA frame structure adopted in field trials conducted by *INTELSAT* during 1978–1979. The transmission rate of 60 Mbits/s (935 voice channels) achieved with PCM/QPSK via a 40-MHz transponder yields a channel density of 23 channels/MHz.

A reference burst consisting of 30 symbols for carrier and clock recovery, 10 symbols for unique word, and 3 symbols for station identification, transmitted by a primary reference station (two standbys were provided) permits all other stations to transmit their bursts in proper sequence, i.e., without mutual overlap. Transmission rate and channel density can be doubled with the use of digital speech interpolation (DSI).

The number of voice channels in a TDMA system is

$$n = (1/r)(R - NP/T) \quad (\text{Eq. 58})$$

where

r = voice channel bit rate

R = satellite channel bit rate (power or band limited)

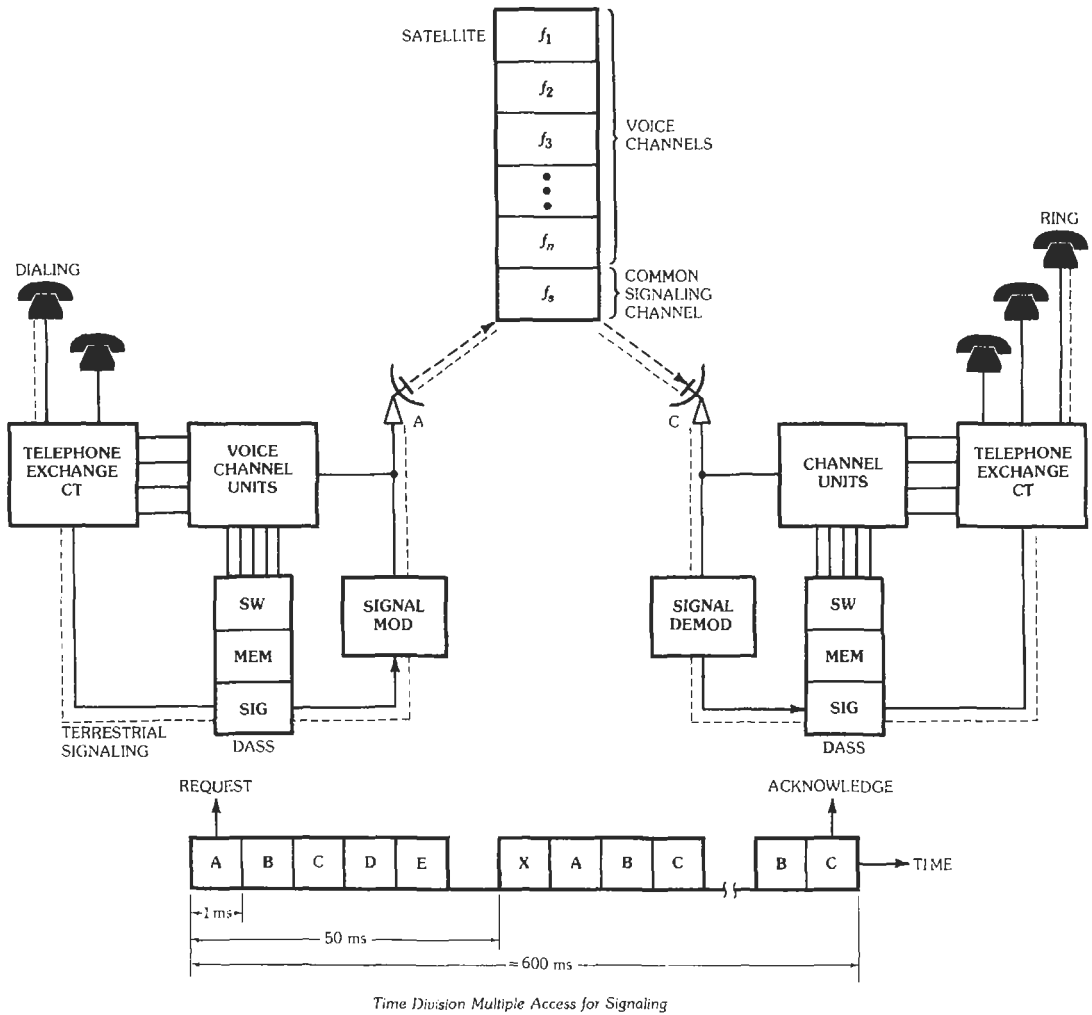


Fig. 21. SPADE signaling flow. (From COMSAT Technical Review, Vol. 2, No. 1, p. 229, Spring 1972.)

N = number of bursts in a frame
 P = number of digits in the preamble
 T = frame period

Synchronization and acquisition in TDMA systems can be implemented via closed- or open-loop techniques. In the closed-loop case, each earth station can monitor its own signals returned by the satellite via global or wide-coverage-area beam antennas. However, when narrow-beam antennas (which provide greater eirp and communications capacity) are used, open-loop techniques must be employed.

Notwithstanding the loss of communications capacity encountered with FDM/FM/FDMA when several accesses share a transponder, its simplicity and its "natural" interfacing with analog terrestrial networks

made it almost universally the primary transmission method of the first decade of commercial satellite systems. With the growth of available spacecraft power and a view of bandwidth limitations, more complex transmissions techniques such as TDMA have been used.

Fig. 18 illustrates the multiple-access performance of the three transmission systems—FDM/FM/FDMA, SCPC/DAMA, and TDM/PCM/PSK/TDMA—when used with a global-coverage transponder. Clearly, bandwidth is more efficiently used by the latter two methods, which require, however, more complex equipment at the earth stations. TDMA, whose potential was identified many years ago, is now more than competitive with FDMA and has already supplanted it operationally in many systems.

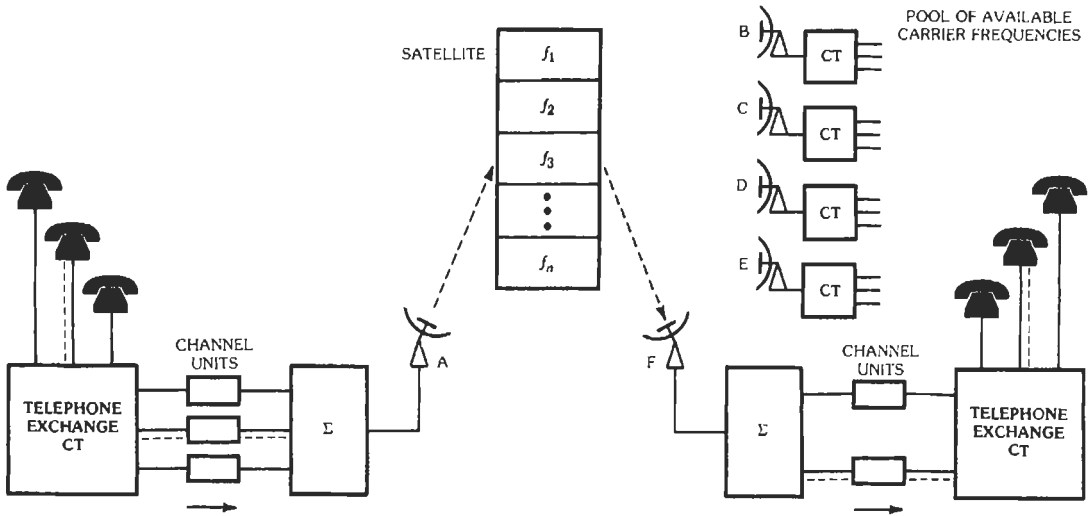


Fig. 22. SPADE communications flow. (From COMSAT Technical Review, Vol. 2, No. 1, p. 229, Spring 1972.)

TABLE 3. TDMA SYSTEMS CHARACTERISTICS

	TELESAT Canada	INTELSAT V	Advanced WESTAR	TELECOM I	INTELSAT VI
Operational Date	1975	1983	1983	1984	1989
Frequency	6/4 GHz	6/4 GHz and 14/11 GHz	14/12 GHz	14/12 GHz	6/4 GHz
Transmission Mode	TDMA	TDMA	SS-TDMA	TDMA	SS-TDMA
Transponder Hopping	No	Yes	No	Down-link switching	Yes
Modulation	QPSK	QPSK	QPSK	BPSK	QPSK
Bit Rate (Mbits/s)	61	121	250	25	120
Frame Period	250 μ s	2 ms	125-750 μ s	20 ms	2 ms
Channel Assignment	Fixed	Fixed	Fixed	Demand assignment	Fixed
Acquisition	Open-loop	Open-loop	?	Open-loop	Open-loop
Synchronization	Closed-loop	Feedback	?	Feedback	Feedback

The trend toward digital communications in terrestrial networks has been a major factor favoring the use of TDMA in satellite systems.

SATELLITE-SWITCHED TIME-DIVISION MULTIPLE ACCESS (SS-TDMA)

Narrowing satellite antenna beamwidth increases communications capacity but decreases connectivity among coverage areas. By introducing onboard switch-

ing, connectivity can be restored. Fig. 24 illustrates the basic concept of satellite-switched time-division multiple access. Spacecraft so equipped operate not only as repeaters but also as switchboards. With time-division multiplexing and digital modulation, the separation of the accesses is more effectively implemented in the time domain. It is also possible to use SS-FDMA techniques, but their more complex implementation would require heavier spacecraft.

Switching can be done at RF or at baseband. In the first case, up- and down-links are interconnected via a microwave switching matrix (MSM) whose output signals

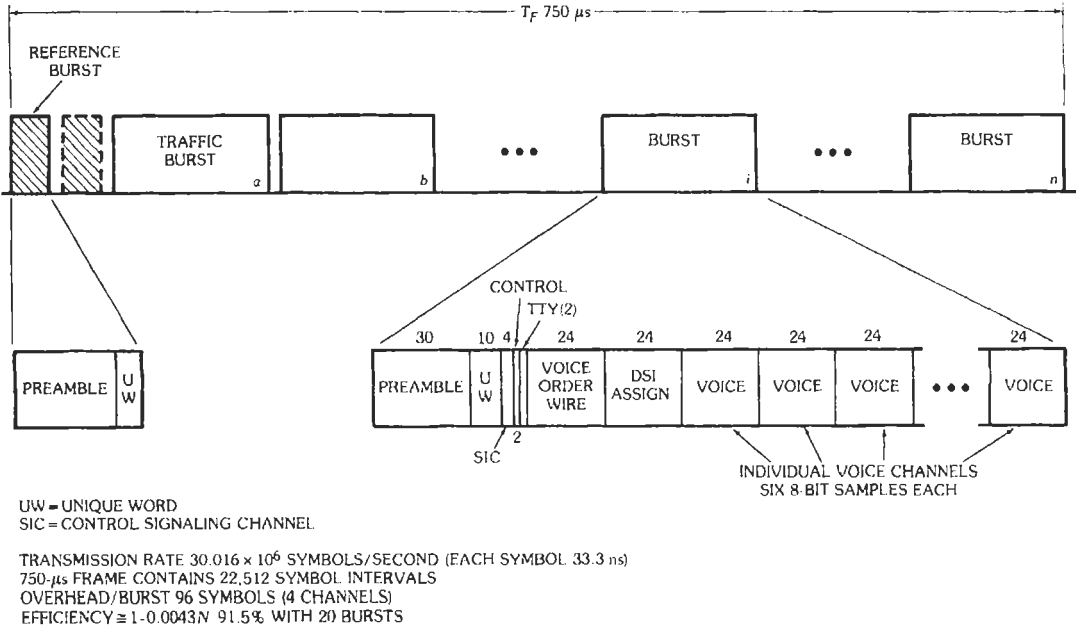


Fig. 23. TDMA frame structure.

are amplified and eventually fed to the antennas. The MSM is steered by a distribution control unit (DCU), and an acquisition and synchronization unit provides the time references necessary for TDMA operation. Configurations have been designed to provide bandwidth up to 500 MHz at 4 GHz.

Fig. 25 shows a block diagram of an SS-TDMA satellite carrying a 6×6 MSM. Data controlling the programmable cyclic switching states are transmitted via the telemetry, tracking, and command (TTC) channel for storage in the DCU memory circuits. To achieve the required reliability on board the *INTELSAT VI*, a 10×6 MSM is used in the SS-TDMA onboard package. SS-TDMA is in operational use in the *INTELSAT VI* system. The onboard package was configured to permit "transparent" use by the *INTELSAT VIVA* system TDMA and thus allow an easy transition. This highlights a major concern when introducing new access methods or any other significant change. The introduction should not interrupt service and should be transparent to the end user.

The capabilities of digital satellite transmission can be further enhanced by providing onboard signal processing. Switching is then performed at baseband after demodulation and demultiplexing of the up-link signals. The onboard switch is followed by remodulation of the down-link signals. Enhancement of overall link performance is obtained as other forms of signal processing can be introduced, such as onboard demand assignment, handling of bit streams at different rates, destination-directed packet transmission, etc.

OTHER SERVICES

Two further *INTELSAT* digital services warrant discussion. One is the International Business Service (IBS), which is the basis for private networks using the *INTELSAT* system. It carries digital traffic at bit rates from about 64 kbits/s to about 8.448 Mbits/s using QPSK and coding rates from rate 1/2 and higher. The other is the Intermediate Data Rate (IDR) service, intended originally to bridge the gap between SPADE and SCPC at 64 kbits/s (effectively 56 kbits/s) and *INTELSAT* TDMA at 120 Mbits/s, working into the public switched telephone networks (PSTN). IDR carries digital information at hierarchical rates from 64 kbits/s to 44.736 Mbits/s, using QPSK and a coding rate 3/4 (this implies transmission rates of 85.33 kbits/s to 59.776 Mbits/s). Recent advances in voice processing and digital interpolation allow up to a factor of 6 in digital circuit multiplication (DCM). Thus a 45-Mbits/s IDR channel could carry up to about 2800 voice circuits, under optimum conditions of DCM. Even higher DCM factors can be anticipated in the future.

Mobile services operate under quite different parameters than those for FSS. Such systems and their trade-offs will not be fully described here. Suffice it to say that allocated bandwidths are limited to a few tens of megahertz, at L-band (1.5 and 1.6 MHz) and S-band between satellite and mobile terminals and C-band and K_u -band (feeder link allocations in FSS band) from satellite to coastal or base earth stations. Consequently, bandwidth is a most precious resource. Only 5 kHz is

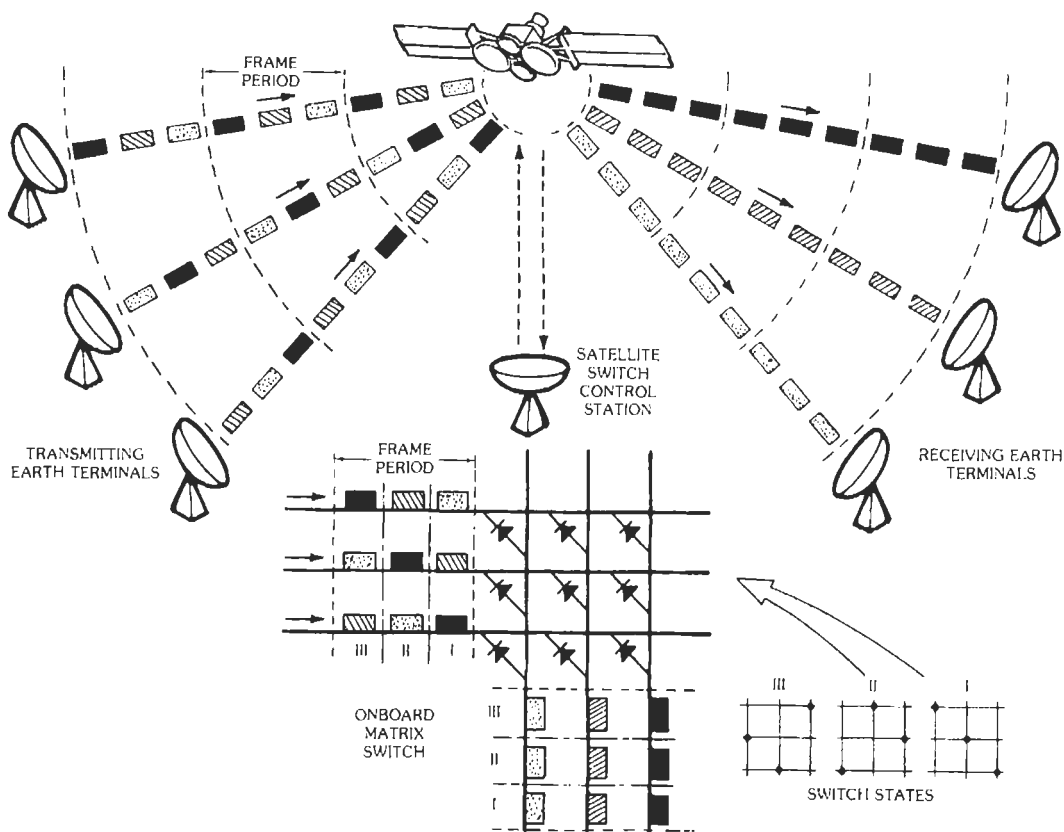


Fig. 24. Basic concept of SS-TDMA. (Courtesy S. J. Campanella.)

used per voice channel, in general, although wider channels are available for special purposes (such as higher-rate data transmission). Because mobile terminals cannot be high gain or high eirp, in general, there is a heavier burden on the satellite for the L-band links. Finally, the economics of the system are entirely different than for FSS, as even with frequency reuse and lower quality narrow-band voice channels, capacity could be at most of the order of a few thousand channels. In this regard, it should be noted that although *INTELSAT V* had a capacity above 10 000 circuits, *MARISAT* capacity was less than 50, and they were contemporaneous.

As satellite eirp grew with the introduction of zone/country coverage beams, it became evident that by using transmission in and out of a central high-gain high-eirp hub, it would be possible to send modest data rates (up to 64 kbits/s) or a single 64-kbits/s voice channel to or through that hub from earth stations with quite small antennas (down to 1 meter diameter), thus allowing very-small-aperture terminals, or VSATs. The basic configuration (shown in Fig. 26) is a star wherein the satellite (hub) is the node of the star and one or more VSAT networks are the

star points. The hub-star-point links are usually higher-speed links using some form of destination-addressed packet transmission. The hub-star-point links transmission and access methods depend on the activity anticipated. For low activity, random-access systems (with retransmission for failed access) such as ALOHA work reasonably well. For higher levels of activity, more complicated access techniques and protocols have to be used, usually entailing demand-assigned channels after a failed access is reported. Up to a few hundred VSAT terminals can be serviced in a single network, and a hub can support several such networks. VSAT antennas of 1 to 3 meters using K_u - and C-bands, transmitters at the VSAT of 1 to 5 watts (SSPAs), and PSK modulation (sometimes using FEC, sometimes not) are sufficient to support 56- to 128-kbits/s service into the hub. The hub is usually a substantial earth station with an antenna of the order of 10 to 20 meters in diameter and a TWTA with an output of the order of hundreds of watts. It usually supports outbound data rates of 512 kbits/s or less per network. The hub usually has direct access (which may even be on the premises of the HQ facility of the entity establishing the VSAT network). Much of the

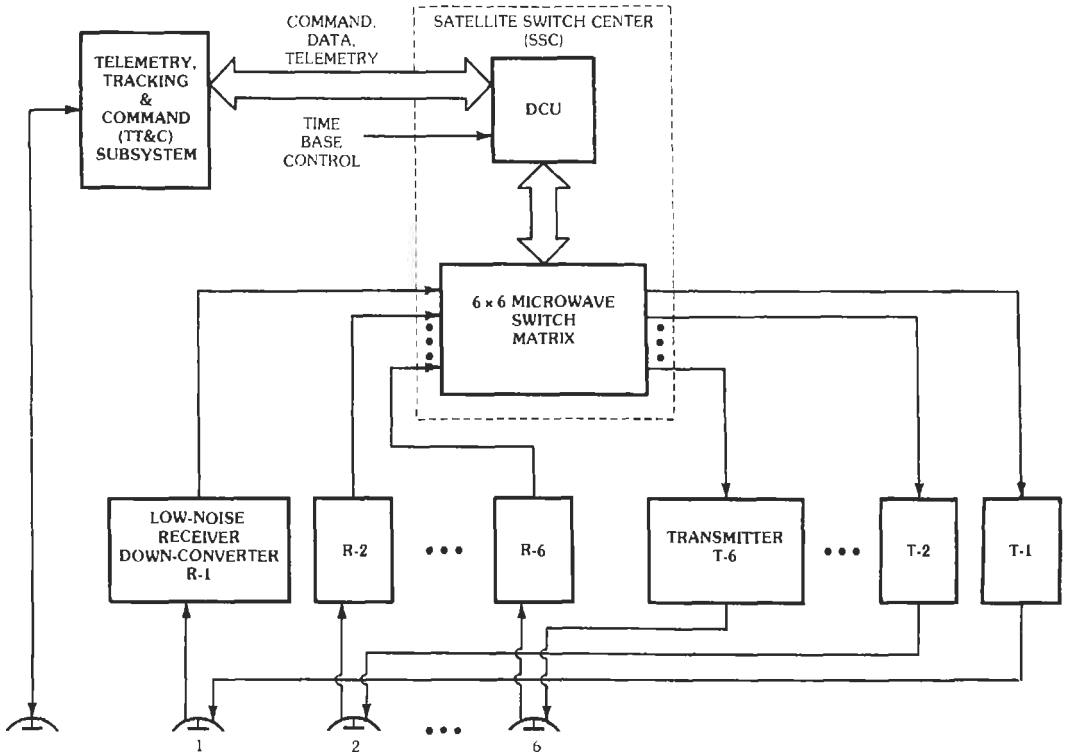


Fig. 25. Simplified block diagram of an SS-TDMA satellite.

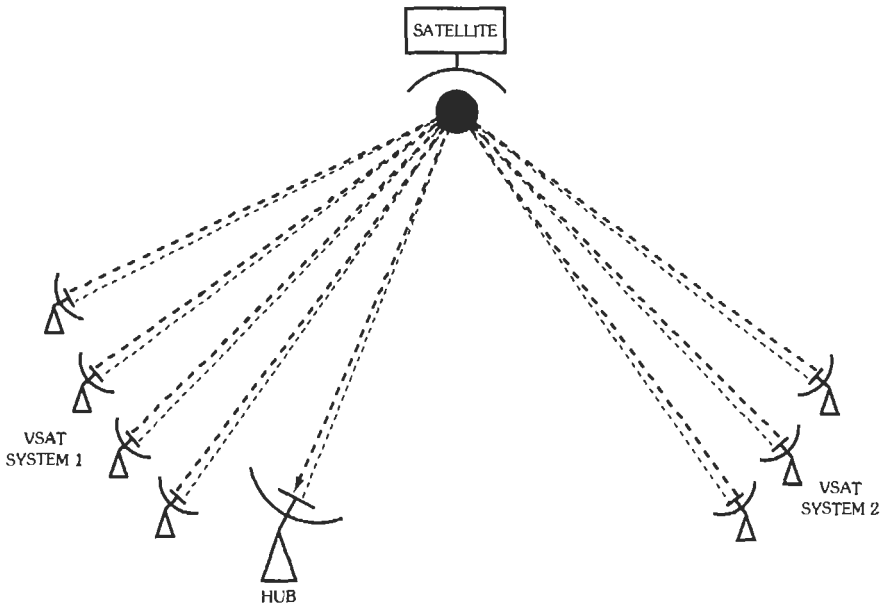


Fig. 26. VSAT network configuration.

transmitted data is broadcast to the star points or comes into the hub only. However, using double-hop, the input to the hub can be retransmitted to another star point.

As noted above, recent advances in satellite eirp have permitted the introduction of even smaller user's premises antennas and terminals, termed VSATs in FSS. In the Mobile Satellite Service (MSS), hand-held terminals have begun to appear. It is expected that Internet services will access users through these and perhaps DBS-like terminals.

In a sense, DBS-TV or DBSA can be considered an extreme form of VSAT, where the VSATs on the user's premises are truly small (down to 0.3-meter antennas) and receive-only, and the hub station transmits the program material to the satellite at feeder link frequency allocations in the FSS to be transmitted at K_u -band at very high eirps (using onboard TWAs of the order of 200 watts). Of course, there are (potentially) millions of star points; there are FM as well as digital signals; and the information rates from hub to star point are quite high (in the tens of megabits per second equivalent). DBS-TV systems have long been established in Japan, Europe, and the former USSR. Others have been implemented in Canada, the United States, South and Central America, and elsewhere.

NETWORKS

One must view the vast communications capacity of satellites not as a collection of links but as networks or parts of networks, interfacing with terrestrial networks. Thus, much attention must be placed on the interfaces, the signaling systems (such as signaling system 7), the many protocols (such as X.25 and the seven-layer ISO structure shown in Fig. 27), and services such as ISDN and B-ISDN.

SPACECRAFT ANTENNAS

Differences in the design of antennas used in space and on earth were quite marked in early communications satellite systems. Large apertures and very low noise system temperatures were then necessary at the earth stations because satellites were very much power-limited. Satellite antennas were small and had little directivity; this was especially the case with satellites in low- or medium-altitude orbits. Dipoles aligned with the spacecraft spin axis produced toroidal radiation patterns, and consequently the energy radiated outside the 18° angle subtending the earth from geosynchronous altitude was wasted. Around 1968, mechanically de-spun antennas consisting of a conical

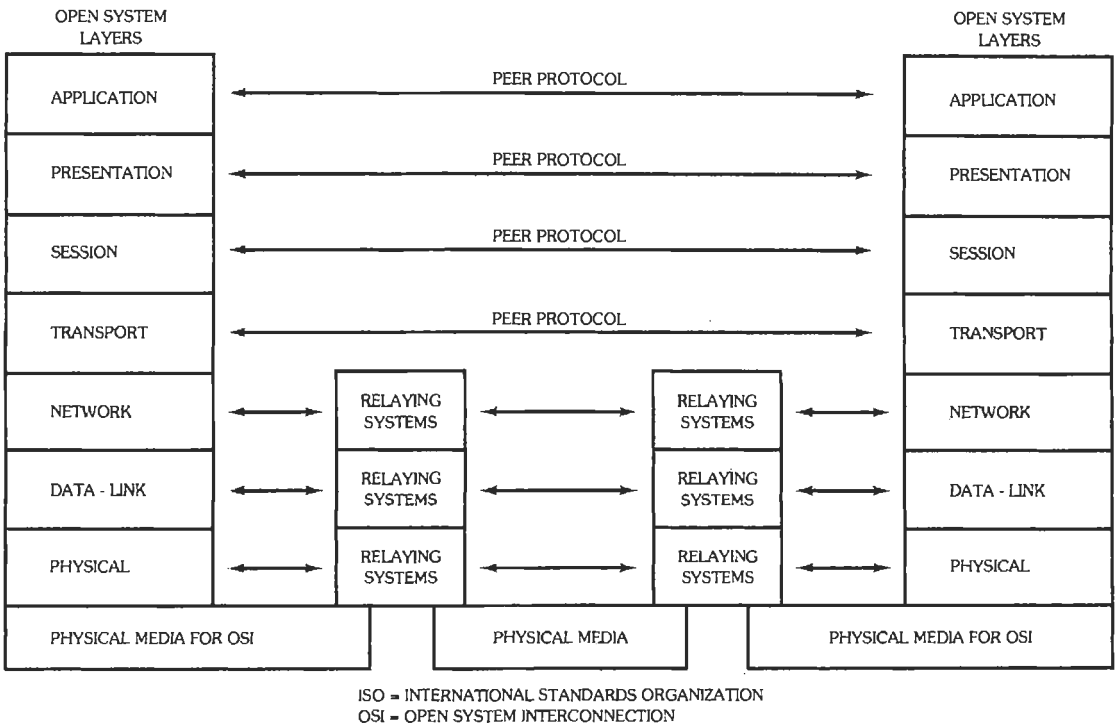


Fig. 27. ISO Reference Model.

horn and a reflector provided about 19 dB of gain (subtending an angle slightly greater than 18°) at 4 and 6 GHz (*INTELSAT III*). Progress has continued since then along two lines. First, large parabolic antennas have provided spot beam coverage down to a few degrees. Second, through frequency reuse the allocated frequency bands are utilized many times over by means of orthogonal polarizations (vertical/horizontal linear or clockwise/counterclockwise circular) and/or by means of spatially separated beams. Finally, zone beams have been synthesized to follow the contours of specific geographical areas such as national or regional boundaries, continents, etc. By assembling numerous suitably arranged feed horns, each excited with proper amplitude and phase, the radiated energy impinges upon a reflector and illuminates the desired areas on earth. It should be noted that earth-coverage horns are still used for a few C-band transponders.

In the design of contour-shaped beams and frequency-reuse systems, it is difficult to maintain isolation between dually polarized and/or spatially separate beams, especially when the number of beams is large. Until recently, parabolic reflectors with offset feed assemblies in the focal region have been adequate. The major role played by the antenna subsystem is exemplified by the design shown in Fig. 28 (*INTELSAT V*) with the coverage areas in the Atlantic Ocean Region shown in Fig. 29. In addition to beacon and TTC&M (TTC and Monitoring) antennas, circularly polarized conical horns provide global coverage at 6 GHz (receive) and 4 GHz (transmit). The 18° transmit horn is steerable up to $\pm 2^\circ$; the receive horn, which has a wider beamwidth (22°), is fixed. Two parabolic reflectors of 2.44- and 1.54-m diameter provide hemispherical and zone area coverage for reception and transmission at 6 and 4 GHz. The offset hemi/zone antenna feeds are clusters of square horns with excitation in amplitude and phase producing the required shaped beams. Hemispherical and zone coverages are simultaneously obtained with opposite-sense circular polarizations. Feed characteristics, reflector size, and focal length are selected to yield high gain, sharp beam edges, low side lobes, and the high polarization purity required to achieve isolation between beams up to 27 dB.

Nominal 1-m diameter mechanically steerable parabolic receive and transmit antennas with linear orthogonal polarizations are used at 11/14 GHz to provide circular 1.6° beams in the west spot area and elliptical 1.8° by 3.2° beams in the east spot area. Spatial discrimination at 6/4 and 14/11 GHz and polarization discrimination at 6/4 GHz provide a fourfold frequency reuse of portions of the allocated spectrum and a total maximum usable bandwidth of 2137 MHz. Thus, it is clear that communications capacity can be effectively increased by adding sophistication to the antenna subsystem.

The *INTELSAT VI* antenna farm is a logical extension of that of *INTELSAT V/VA*. Earth coverage is still achieved with horns. Two hemispheric and four zone

beams each for transmit and receive are generated using somewhat larger reflectors (3.2-m and 2.0-m diameter transmit and receive, respectively). The feed array provides six-fold frequency reuse at C-band. The K_u -band spot beams are essentially the same as in *INTELSAT V/VA*. There is, as in *INTELSAT V/VA*, a triple layer of beam-forming networks (BFNs) behind the feed arrays, switchable to select each ocean region.

As more beams (i.e., greater amounts of frequency reuse) and better isolation among beams are required, two approaches appear possible. Since feed blockage limits the use of single-offset reflectors, dual-offset reflectors with Gregorian or Cassegrain feeds could be used to allow better control of the illumination and permit the use of larger reflectors without incurring excessive focal lengths. Reflectors fed by phased arrays or phased arrays by themselves can also be used. Ultimately, size and mass limitations are encountered. The launch vehicle payload volume geometry limits the maximum diameter of fixed antennas. This limitation can be circumvented by using deployable apertures. Deployable reflectors have been used since the 1970s (e.g., *ATS-6*). Currently, deployable reflectors for the MSS at L-band are in use with diameters the order of 15 meters. Great care in structural, thermal, and mechanical design is required to maintain the specified surface tolerances after deployment, throughout eclipses and other thermal exposure, over the life in orbit. Such antennas provide many high-gain, high-erp beams. Such beams are the basis for the use of handsets that offer cellular service via satellite.

PROPAGATION

Electromagnetic waves traveling along an earth-space path encounter four distinct regions: the troposphere, inner free space, the ionosphere, and outer free space.

Free space is lossless and has unity refractive index, and the troposphere and the ionosphere have refractive indices respectively greater and smaller than unity; hence, refraction and absorption phenomena arise. In addition, the ionosphere, a magnetoionic medium, induces Faraday rotation. All these phenomena affect space and satellite communications systems. Propagation delays resulting from ray bending influence the performance of satellite navigation systems.

Faraday Rotation

A linearly polarized wave can be regarded as the sum of two counterrotating circularly polarized waves. Because the phase velocities of the two waves differ in a magnetoionic medium, the polarization plane rotates. Faraday rotation effects are negligible above 10 GHz; below 10 GHz they can be circumvented by using circular polarization. Faraday rotation can, however, degrade cross-polarization discrimination in frequency-reuse systems. As with most ionospheric

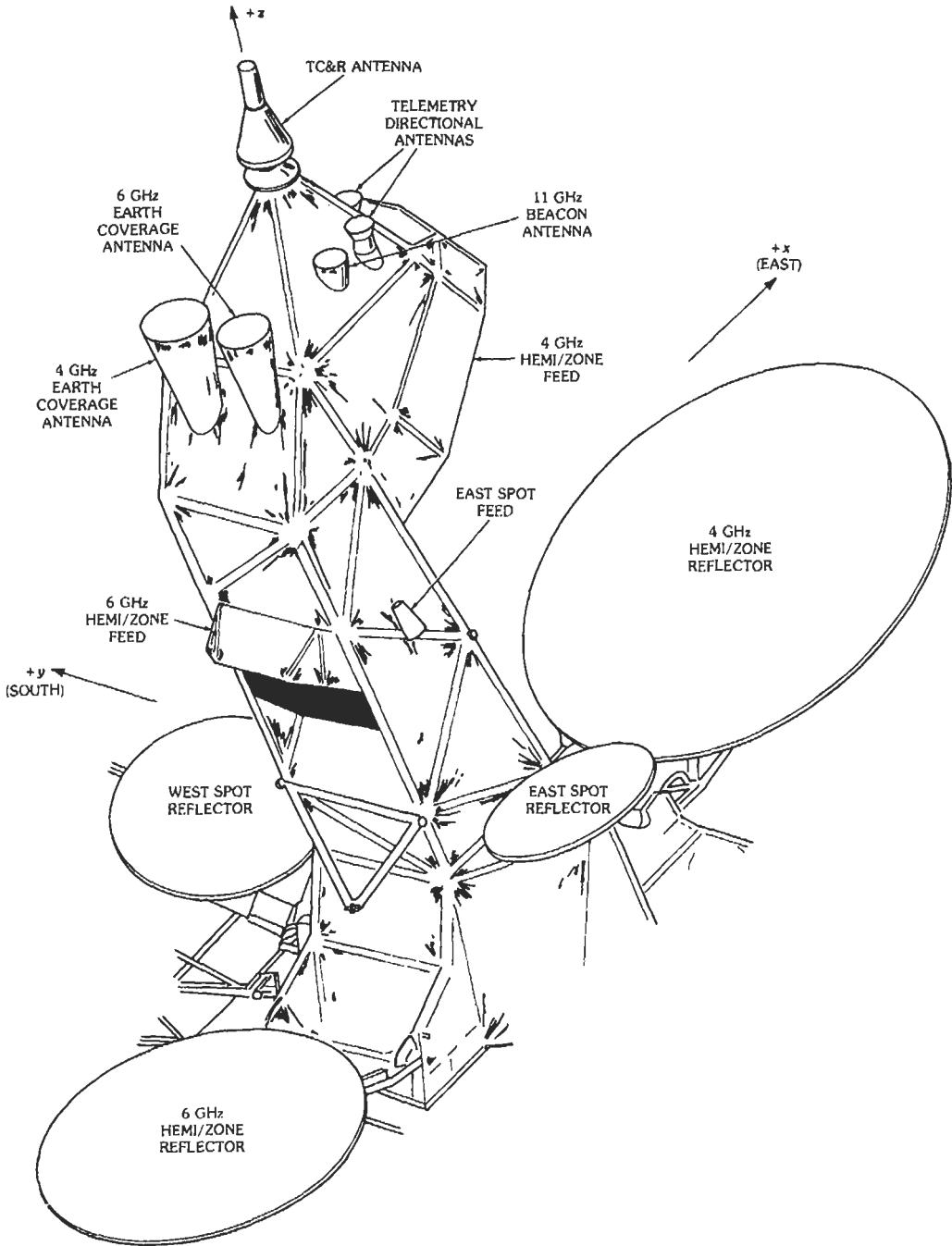


Fig. 28. *INTELSAT V* antenna subsystem. (Courtesy Ford Aerospace and Communications Corp.)

effects, there is a strong diurnal and seasonal dependence as well as increase with increasing solar activity; i.e., Faraday rotation follows the solar cycle, peaking approximately every 11 years.

Ionospheric Scintillations

Strong short-term (1- to 15-s) variations of the amplitude, phase, polarization angle, and angle of

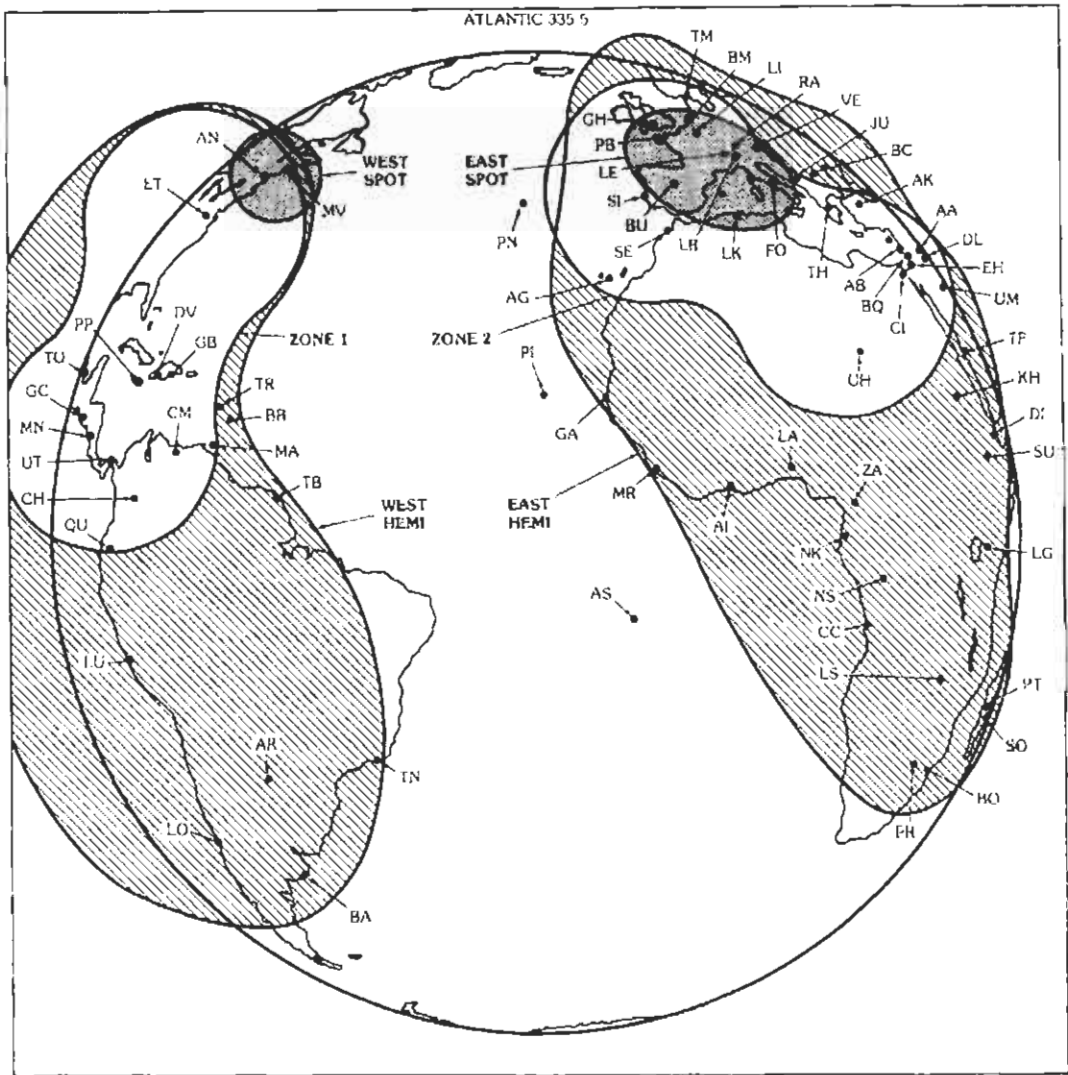


Fig. 29. INTELSAT V coverage of the Atlantic Ocean Region. (From COMSAT Technical Review, Vol. 7, No. 1, p. 314, Spring 1977.)

arrival of electromagnetic waves reaching the earth from space observed up to microwave frequencies are known under the general name of ionospheric scintillations. These are attributable to fluctuations of the electron density in the sporadic E layer and to the spreading of the F layer. Earth-station location, season of the year, local time, and amount of solar activity influence the magnitude and occurrence of ionospheric scintillations; geomagnetic latitude, solar activity, and local time are the most influential factors.

Within 120° from the geomagnetic equator, heavy scintillation can occur for a few hours just after sunset in the ionosphere on the earth-satellite path, with magnitude proportional to the number of sunspots (i.e., following the solar cycle) and maxima occurring around the equinoxes.

The scintillation index, S_4 , defined as

$$S_4 = (1/\bar{R}^2) \{ (\overline{R^2} - \bar{R}^2)^2 \}^{1/2} \quad (\text{Eq. 59})$$

has an $f^{-\alpha}$ frequency dependence with

$$\alpha = \begin{cases} 1/2 \text{ to } -1 & \text{at vhf} \\ -1 & \text{from 1.5 to 4.0 GHz} \\ -2 & \text{above 4 GHz} \end{cases}$$

Scintillations become less intense at intermediate latitudes but increase at higher latitudes beyond the boundary of the auroral region.

Additional ionospheric effects are angle-of-arrival variations and absorption. Both phenomena follow an

f^2 law; hence, their impact on communications systems operating at gigahertz frequencies is negligible.

Tropospheric Effects

Ray bending, scintillation, attenuation, and increased sky noise temperature are major effects. Tropospheric bending is opposite to ionospheric bending; it is frequency independent, whereas ionospheric bending follows an f^2 law. Because the refractive index, n , is slightly greater than one, the quantity N (refractivity) is used

$$N = (n - 1) \times 10^6$$

It is empirically given as

$$N = (77.6/T) [p + 4810\epsilon/T] \quad (\text{Eq. 60})$$

where

- T = air temperature in kelvins
- p = atmospheric pressure in millibars
- ϵ = partial water vapor pressure in millibars

Since N decreases with height, the apparent elevation of a space object is greater than the geometrical elevation angle, and the radio range is greater than the geometric distance. Eq. 60 does not hold in extremely humid climates; ducting will occur if the vertical gradient of N is high. All the above phenomena are not significant in the design of satellite communications systems except for operation at very low angles of elevation. In tropical climates, mean values of ray bending around 0.5° – 0.6° have been observed at elevation angles between 1° and 2° . Ducting can produce interference to earth-based microwave radio relays by high-power earth stations of communications satellite systems radiating at low angles.

Scintillations are induced by atmospheric turbulence with time dependence about ten times greater than that of ionospheric scintillations. The phenomenon is strongly dependent on elevation angle, increasing with decreasing elevation angle, and may be neglected above 15° . Between 5° and 10° , peak-to-peak scintillations in excess of 6 dB have been encountered at 6 GHz.

Under clear-sky conditions, water vapor and oxygen give origin to molecular resonance absorption bands whose width is affected by the atmospheric pressure (i.e., change with altitude). Clear-sky atmospheric loss (attenuation) is of the form

$$L_a = \int_0^l \{ \gamma_{O_2}(r, f) + \gamma_{H_2O}(r, f) \} dr \quad (\text{Eq. 61})$$

where

l = total propagation distance

$\gamma_{O_2}, \gamma_{H_2O}$ = absorption coefficients for oxygen and water vapor in dB/km

Plots of the attenuation per unit length as a function of frequency are given in Fig. 30.

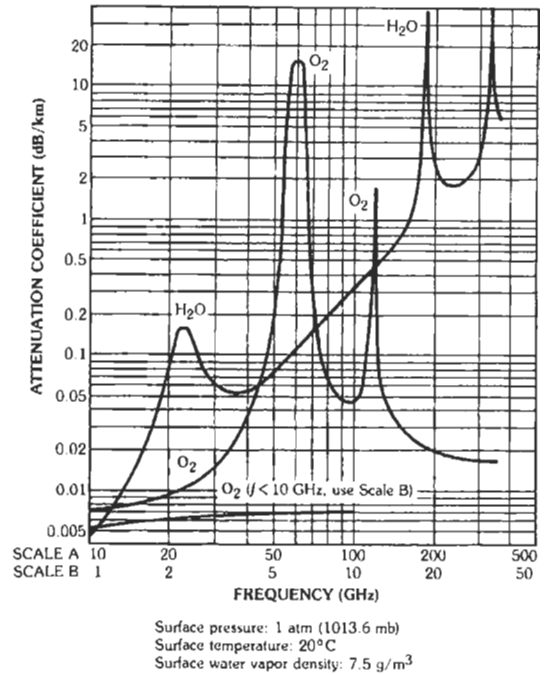


Fig. 30. Clear-sky attenuation per unit length versus frequency. (From K. Miya, Ed., *Satellite Communications Technology*. Tokyo: KDD Engineering and Consulting, Inc., 1982.)

The simplified expression

$$L_a = \gamma'_{O_2} L_{O_2} + \gamma'_{H_2O} L_{H_2O} \text{ in dB} \quad (\text{Eq. 62})$$

may be used for horizontal paths, assuming that the absorption coefficients, γ' , at the surface and the effective path length, L , for oxygen and water vapor are known. It is generally assumed that

$$L_{O_2} \cong 4 \text{ km}$$

$$L_{H_2O} \cong 2 \text{ km}$$

For a vertical path, total clear-sky loss versus frequency is shown in Fig. 31. For slant paths, a factor about equal to the cosecant of the elevation angle (valid down to about 5°) must be introduced.

Clear-Sky Noise Temperature

In the absence of hydrometeors, the absorption by O_2 and H_2O molecules contributes to clear-sky noise temperature T_{cs} :

$$T_{cs} = T_{gal}/L_{cs} + [(L_{cs} - 1)/L_{cs}] T_m \quad (\text{Eq. 63})$$

where

$$T_{gal} \cong 26/f_{GHz} \text{ kelvins}$$

$$L_{cs} = 10^{L_a/10}, \text{ clear sky loss factor}$$

$$T_m \cong 1.12 T_{surface} - 50, \text{ mean raindrop temperature along the path}$$

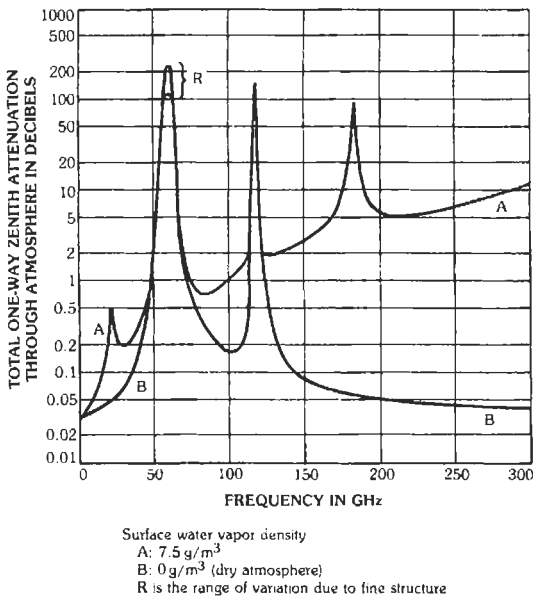


Fig. 31. Total clear-sky loss versus frequency. (From K. Miya, Ed., *Satellite Communications Technology*. Tokyo: KDD Engineering and Consulting, Inc., 1982.)

The dependence of clear-sky noise temperature on frequency and elevation angle is shown in Fig. 7.

Above 10 GHz, the term T_{gal}/L_{cs} , which represents the galactic noise contribution after passage through the clear sky, can be neglected. Then

$$T_{cs} \cong [(L_{cs} - 1)/L_{cs}]T_m = T_m(1 - 10^{-L_{cs}/10}) \quad (\text{Eq. 64})$$

Hydrometeors

On account of the basic interaction of electromagnetic waves with water in liquid form, raindrops cause absorption, scattering, and depolarization phenomena. The first two result in signal attenuation (not to be confused with either water-vapor attenuation or attenuation by rain clouds) and an increase in sky noise temperature. These effects are quite noticeable above 10 GHz. Depolarization is of little or no harm to communications systems using only one polarization, but it worsens the performance of dual-polarization systems and may also result in interference between systems. Rain may also limit signaling bandwidth.

Rain Attenuation

The amount of attenuation depends on the type of rain (stratiform, convective, cyclonic) and intensity (rain rate). As raindrop size distribution is a function of rain type and intensity, various models have been proposed to best fit specific situations. Rainfall data are available for most parts of the world; different climate types have been identified, with their regions on the

earth's surface identified and their associated rain-rate exceedance statistics defined. Cumulative rain statistics provide information about probability of exceedance, i.e., the total time that a specific rain rate will be exceeded over a sufficiently long observation period. Cumulative statistics do not provide information about the frequency of occurrence and the duration of the periods of exceedance.

The attenuation per unit length (specific attenuation), α_r (dB/km), is tied to the rain rate, R (mm/hr), by the empirically derived relationship

$$\alpha_r = a(f)R^{b(f)} \quad (\text{Eq. 65})$$

where $a(f)$ and $b(f)$ are frequency-dependent coefficients. The approximate analytic expressions given in Table 4 are adequate for engineering use. Up to 50 GHz, $b(f)$ is near unity; hence specific attenuation is essentially proportional to rain rate. However, since $a(f)$ is heavily frequency dependent, attenuation per unit length increases rapidly with frequency. Fig. 32 shows the frequency dependence of α_r for various rain rates. If the concept of equivalent path length, $L_{eq}(R)$, is introduced, the total rain attenuation loss in decibels is simply

$$A_r = \alpha_r \times L_{eq}(R) \quad (\text{Eq. 66})$$

Equivalent path length is primarily determined by the height of the melting layer, which depends on latitude and season, and by the cosecant of the elevation angle (for latitudes within $\pm 30^\circ$, the melting layer is at 4.8 km). A correction coefficient is required to take into account the distribution of rain over a long path; however, except for very light or very heavy rain, the value of the correction is near unity. Equivalent path length can be expressed in terms of rain rate and elevation angle θ as follows:

$$L_{eq}(R, \theta) = [7.413 \times 10^{-3}R^{0.766} + (0.232 - 1.803 \times 10^{-4}R) \sin \theta]^{-1} \quad (\text{Eq. 67})$$

Curves of equivalent path lengths versus elevation angle and for different rain rates are shown in Fig. 33.

Rain rate can be modeled by the Rice-Holmberg distribution

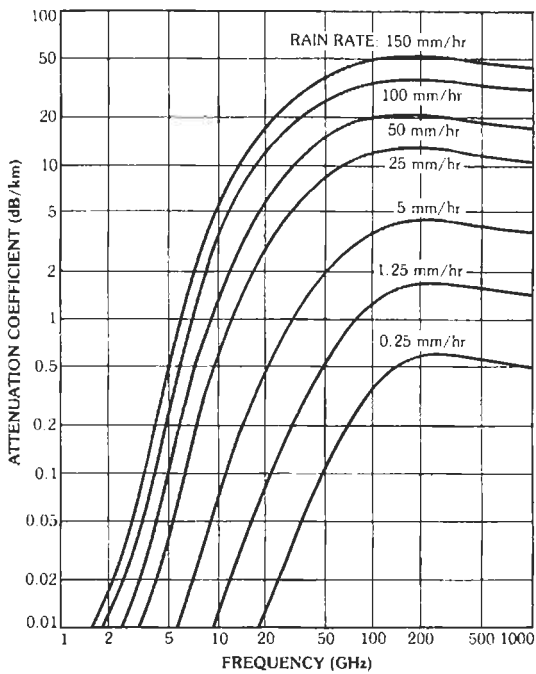
$$P(R) = ae^{-0.03R} + be^{-0.258R} + ce^{-1.63R} \quad (\text{Eq. 68})$$

with

$$\begin{aligned} a &= M\beta/2922 \\ b &= M(1 - \beta)/438.3 \\ c &= 1.86\beta \end{aligned}$$

TABLE 4. PROPAGATION COEFFICIENTS

Frequency f (GHz)	$a(f)$	$b(f)$
8.5-25	$4.21 \times 10^{-5}(f)^{2.42}$	$1.41(f)^{-0.0779}$
25-54	$4.21 \times 10^{-5}(f)^{2.42}$	$2.63(f)^{-0.272}$
54-100	$4.09 \times 10^{-2}(f)^{0.699}$	$2.63(f)^{-0.272}$



Raindrop size distribution: Laws and Parsons, 1943
 Terminal velocity of raindrops: Gunn and Kinzer, 1949
 Dielectric constant of water at 20°C: Ray, 1972

Fig. 32. Attenuation per unit length versus frequency and rain rate. (From K. Miya, Ed., *Satellite Communications Technology*. Tokyo: KDD Engineering and Consulting, Inc., 1982.)

where

- M = total mean yearly rainfall in millimeters
- β = ratio of thunderstorm rain accumulation to total accumulation

In Eq. 68, each term represents the predominant component for specific ranges of R . Comparative data of measured surface rain rate and theoretical model distributions are shown in Fig. 34. Cumulative attenuation statistics are shown in Figs. 35, 36, and 37. Formulae and constants used in them are continuously being refined. Although those given herein are adequate for many purposes, see the latest ITU-R reports for more current information.

The comparison of rain attenuation at different frequencies leads to a scaling relationship of the form

$$A_r(f_2) = A_r(f_1) (f_2/f_1)^\alpha \quad (\text{Eq. 69})$$

with $1.7 < \alpha < 2.0$.

Water-cloud attenuation is about the same as that produced by rainfall at rates less than 5 mm/hr. Ice-cloud attenuation is at least two orders of magnitude lower and hence negligible.

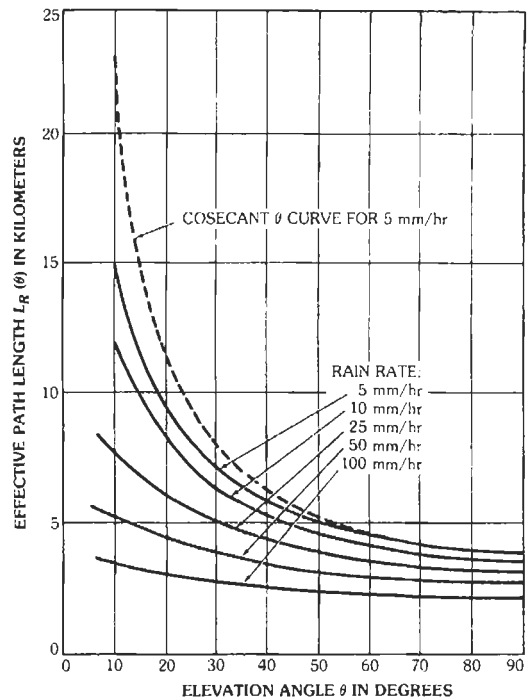


Fig. 33. Equivalent path length versus rain rate and elevation angle. (From K. Miya, Ed., *Satellite Communications Technology*. Tokyo: KDD Engineering and Consulting, Inc., 1982.)

Sky Noise Temperature with Rain

Neglecting the galactic contribution, the sky noise temperature in the presence of rain, T_s , is

$$T_s = T_{cs}/L_r + [(L_r - 1)/L_r]T_m \quad (\text{Eq. 70})$$

If the first term, which represents the contribution of the clear sky through the lossy rainy medium, is negligible compared to the second term, one can write

$$T_s \cong [(L_r - 1)/L_r]T_m = T_m(1 - 10^{-A_r/10}) \quad (\text{Eq. 71})$$

with all symbols as previously defined. Then the decibel difference (excess loss) between the attenuation caused by rain and clear sky attenuation is

$$(A_r - A_a) \Big|_{\text{dB}} = 10 \log [(T_m - T_{cs})/(T_m - T_s)] \quad (\text{Eq. 72})$$

Conversely, the excess noise temperature is

$$(T_s - T_{cs}) \Big|_{\text{kelvin}} = T_m(L_r - L_a)/L_r L_a \quad (\text{Eq. 73})$$

Corrections should be introduced when antenna sidelobe noise pickup and/or receiver front-end noise

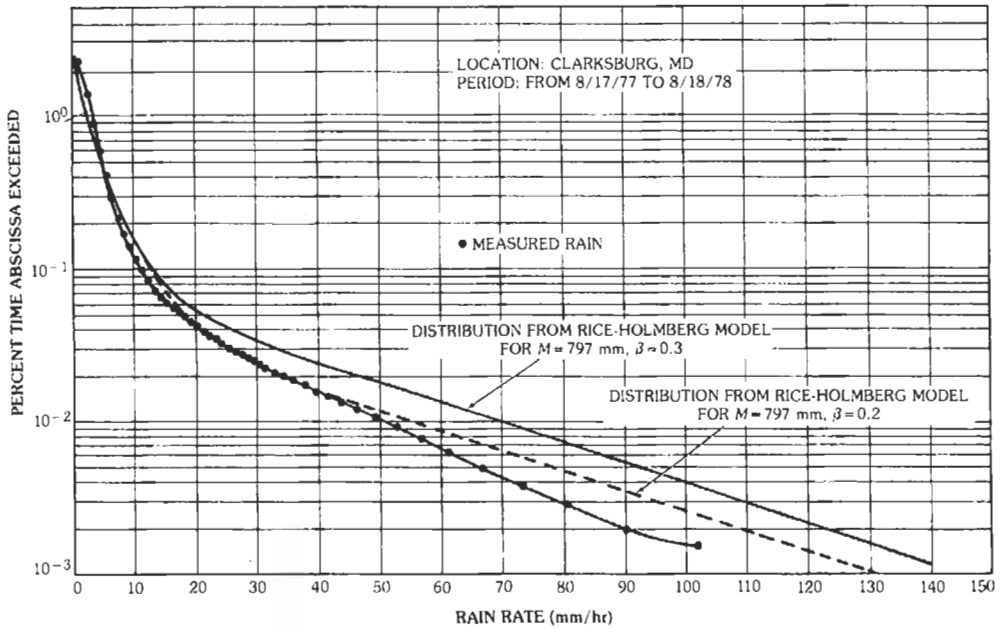


Fig. 34. Comparison of measured surface rain rate and theoretical model distributions. (From *COMSAT Technical Review*, Vol. 12, No. 1, p. 7, Spring 1982.)

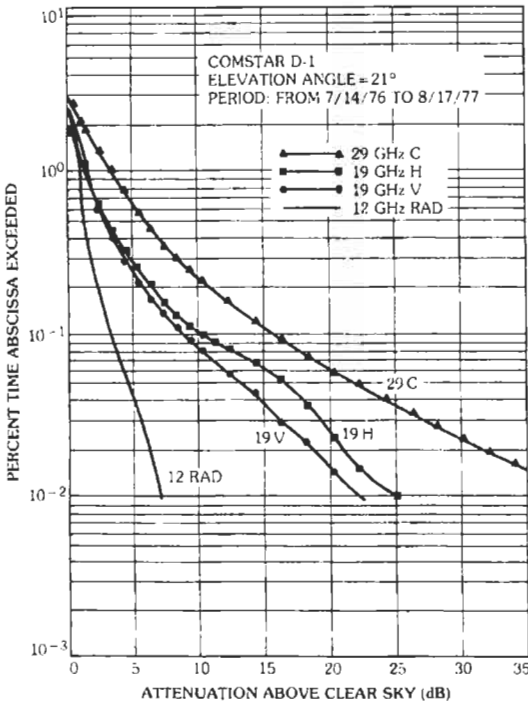


Fig. 35. Cumulative attenuation statistics (*COMSTAR D-1*). (From *COMSAT Technical Review*, Vol. 12, No. 1, p. 10, Spring 1982.)

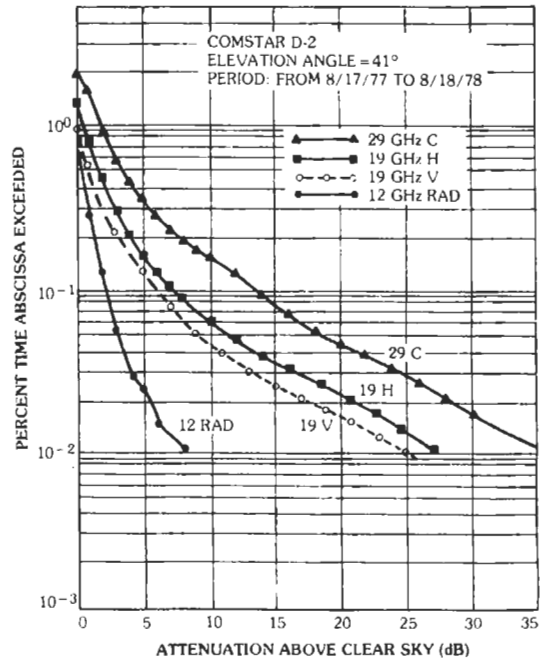


Fig. 36. Cumulative attenuation statistics (*COMSTAR D-2*). (From *COMSAT Technical Review*, Vol. 12, No. 1, p. 12, Spring 1982.)

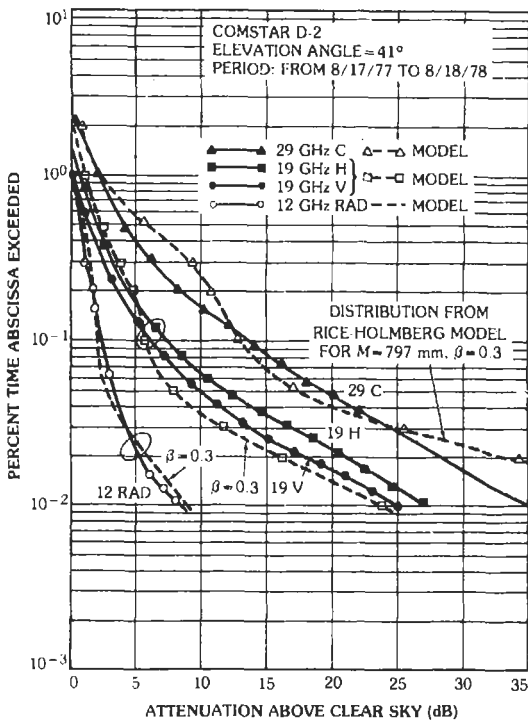


Fig. 37. Comparison of measured attenuation with theoretical distributions. (From COMSAT Technical Review, Vol. 12, No. 1, p. 13, Spring 1982.)

become nonnegligible. The following points should be emphasized:

- A. The cumulative distributions of rain rate and attenuation depend greatly on the type of climate.
- B. Rates below 10 mm/hr are usually associated with stratiform rain, which occurs over wide areas; attenuation is proportional to the cosecant of the elevation angle.
- C. Rates above 20 mm/hr are usually associated with convective rain, which occurs over smaller areas. Attenuation is proportional to less than the cosecant of the elevation angle, but the vertical extent of the rain can be quite large (up to 10 km).
- D. Very intense rates are encountered in cyclones, hurricanes, or typhoons, which may affect wide areas at times. Their probability of occurrence, although low over the year, is often seasonally high.
- E. Rain attenuation that is not a serious obstacle below 10 GHz becomes important at the higher frequencies. Up to 11–15 GHz, continuity of service can be maintained at quite high levels of availability by allowing adequate power mar-

gins (3–12 dB) in the link design. Higher power margins are uneconomical in most cases. Significant additional margin can be obtained by using coding and forward error correction.

Especially at the higher frequencies (20/30 GHz) at which rain-induced attenuation can be quite severe, unless a lower value of service continuity is acceptable, means other than power margins and coding and forward error correction may be needed.

Diversity

One proposed alternate solution is site diversity based on having two earth stations jointly operating with a given satellite. By placing the two stations at a distance for which the rain statistics along the two paths become sufficiently uncorrelated, and selecting at any time the best path, diversity gain over single-site operation can be obtained.

With reference to Fig. 38, once the attenuation statistics for a single and a double path have been acquired, diversity gain G_{div} can be defined as

$$G_{div} = A_{r1} - A_{r1,2} | p \tag{Eq. 74}$$

where the single and double subscripts refer to the same value of single and joint event probabilities. Otherwise, for a given value of the abscissa A_r , the ratio

$$I_{div} = p_1/p_{1,2} | A_r \tag{Eq. 75}$$

defines diversity advantage, i.e., the factor by which exceedance time is improved.

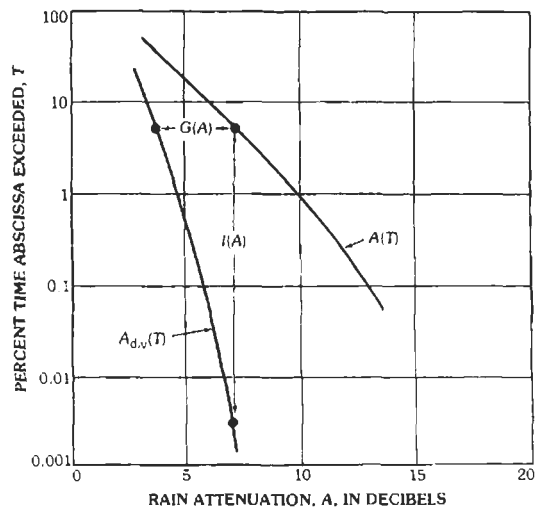


Fig. 38. Diversity and advantage gain definitions. (From L. J. Ippolito, R. D. Kout, and R. G. Wallace, Propagation Effects Handbook for Satellite Systems Design. NASA Reference Publication 1082, Dec. 1981.)

Experiments have shown that both diversity gain and advantage increase with increasing spacing of the two sites, but at a decreasing rate, leveling off at the distance for which the path statistics become independent.

Depolarization

Although very small raindrops such as those encountered in light drizzle are spherical, the larger drops associated with heavier rain depart substantially from spherical shape under the combined effects of gravity, hydrostatic forces, and aerodynamic forces. Nonspherical drops and their canting angle cause differential attenuation and phase shift for different wave components and hence depolarization.

For linear- or circular-polarized waves, the power ratio of the copolarized and cross-polarized components

$$XPD = 10 \log_{10} (P_{\text{copol}}/P_{\text{xpol}}) \quad (\text{Eq. 76})$$

gives a measure of polarization purity or cross-polarization discrimination. Once differential attenuation A and differential phase shift B are determined, cross-polarization discrimination can be computed. For circular polarization, it is

$$XPD = 20 \log_{10} \left| \frac{1 + e^{A+jB}}{1 - e^{A-jB}} \right| \quad (\text{Eq. 77})$$

and for linear polarization

$$XPD = 20 \log_{10} \left| \frac{1 + \tan^2 \chi e^{A+jB}}{1 - e^{A+jB} \tan \chi} \right| \quad (\text{Eq. 78})$$

where χ is the angle between the incident wave polarization plane and the major axis plane of the raindrops. With linear polarization, when $\chi = 0^\circ$ or 90° , XPD becomes infinite; for $\chi = \pm 45^\circ$ XPD takes a minimum value. With circular polarization, XPD is always at a minimum. At 4 and 6 GHz, rain depolarization arises from differential phase effects and can be a significant impairment in dual-polarized systems.

Proper operation of frequency reuse systems with dual polarization requires XPD of 27 dB or more. As a residual coupling usually originating in the antennas and feeds always exists, the presence of rain worsens the situation and further reduces the value of XPD . Depolarization compensation networks can be used at the earth stations to restore orthogonality.

Ice depolarization is also encountered, due to the nonspherical shape of ice particles, generally encountered just above the melting layer. This can be significant at frequencies above 15 GHz and can be identified by the absence of appreciable concurrent attenuation.

EARTH STATIONS

At the earth stations that form the ground segment of space and satellite communications systems, a variety of equipment is needed depending on: (A) function

of the station, (B) type of service, (C) frequency bands used, (D) transmitter, (E) receiver, and (F) antenna characteristics.

Three categories can be distinguished: (1) transmit-and-receive, (2) receive-only, and (3) transmit-only stations. The first is encountered in two-way communications systems. Receive-only stations are presently used in CATV systems in the United States and in great numbers in DBS systems around the world. Transmit-only earth stations are found in data collection systems and DBS feeder links. Types of service include fixed, mobile, broadcasting, and others. Spectrum use is regulated by international allocations with actual bandwidth occupancy as needed for different services, type of traffic, and modulation.

In satellite communications systems, earth-station transmitter power ranges from a few watts generated by SSPAs to about ten kilowatts generated by klystrons or TWTs. Transmitters up to 400 kW are used in the deep-space network. Receivers cover a wide range of sensitivity and bandwidth with noise temperature from a few tens of kelvins achieved with cryogenically cooled amplifiers to hundreds of kelvins in uncooled amplifiers. A highly simplified block diagram of a typical transmit-and-receive earth-station layout is shown in Fig. 39. The station has six major subsystems:

- A. The power subsystem
- B. The terrestrial interface
- C. The transmit chain
- D. The receive chain
- E. The antenna subsystem
- F. The control subsystem

Signal paths and functions are as indicated. The value of the intermediate frequency (IF) in the receive and transmit chains is usually the same, with 70 MHz being commonly used. The up-link signals may be transmitted by means of waveguides to high-power amplifiers located just behind the feed, or in IF form via coaxial cables to the up-converters and high-power amplifiers located there. It is also possible to locate all the equipment at the base of the antenna and to extend the input and output feeds via beam waveguides. High-power-amplifier (HPA) outputs can be combined through bandpass filters and circulators or by means of hybrids. In the receive chain, the weak signals from the satellite are accepted by the same feed that carries the transmitter output. These two signals, which differ in power by several orders of magnitude, are kept separate by means of orthogonal polarizations and diplexers used to enhance the separation in the frequency domain. Orthomode transducers and polarizers are employed, respectively, to couple orthogonally polarized signals into a single waveguide and to convert linear into circular polarization and vice versa. After preamplification in a low-noise amplifier (LNA), the received signals are down-converted to IF and demodulated to baseband.

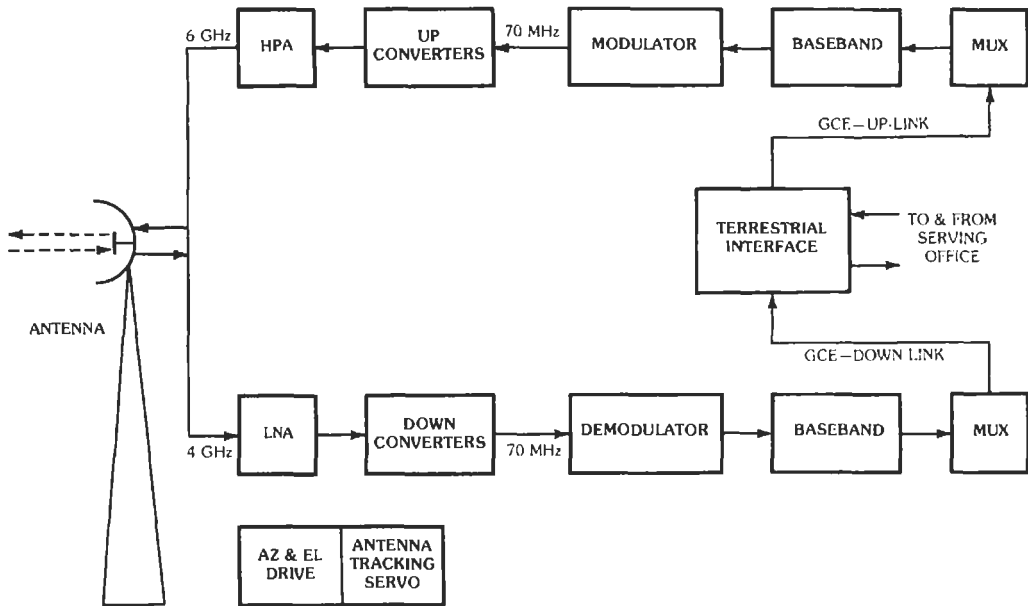


Fig. 39. Earth-station block diagram.

The antenna is a very important part of an earth station, since it affects eirp, receive figure of merit (G/T), and interference. Antenna size varies considerably, from a diameter of 0.5 meter or less used for DBS-TV receive only in private homes to much larger transmit and receive antennas (up to 32-meter diameter parabolas) used in high-capacity systems, and even larger antennas (up to 64-meter diameter) used in the deep-space network. In addition to parabolic reflectors, which are widely used, other forms of antennas such as horns, torus-shaped reflectors, Yagis, helices, and phased arrays have been used.

In terms of the general relationship,

$$G = (4\pi/\lambda^2)A_{\text{eff}} = (4\pi/\lambda^2) \cdot \eta \cdot A_{\text{geom}} \quad (\text{Eq. 79})$$

where

$$\eta = A_{\text{eff}}/A_{\text{geom}} \quad (\text{Eq. 80})$$

is the antenna efficiency. In the case of parabolic reflectors,

$$G = \pi^2(D/\lambda)^2 \cdot \eta \quad (\text{Eq. 81})$$

In practice, values of D/λ ranging from about 20 to 700 are common in communications satellite systems, and values up to 2000 are encountered in deep space systems. With efficiency in the range from 0.5 to 0.8, gains between 30 and 75 dB (dB above isotropic) have been achieved with corresponding half-power beamwidths from a few degrees to a few hundredths of a degree. Very high gain values imply electrically as

well as mechanically large antennas and consequently large and costly structures because cost is approximately proportional to aperture area.

Various forms of mounts have been used with limited or full steerability as required by the specific system characteristics. Antenna-mount types can be 2-axes (x, y or azimuth, elevation), 3-axes (x, y, z' or azimuth, elevation, and cross-elevation), or 4-axes for maritime applications.

Reflector antennas can be classified according to the number of reflecting surfaces and/or the type of feed positioning. The overall antenna efficiency factor, η , in Eq. 80 can be broken into four parts:

$$\eta = \eta_1 \times \eta_2 \times \eta_3 \times \eta_4 \quad (\text{Eq. 82})$$

where

- η_1 = illumination factor
- η_2 = spillover factor
- η_3 = blockage factor
- η_4 = surface tolerances factor

A decrease in main-lobe gain often implies an increase in an increase of side-lobe energy and consequently leads to a higher noise temperature and a decrease in interference-rejection capability. The term η_4 is due to departure of the reflector surface from the ideal. Assuming randomly distributed profile errors over the surface, efficiency is reduced by a factor

$$\eta' = e^{-(4\pi\epsilon/\lambda)^2} \quad (\text{Eq. 83})$$

where ϵ is the rms surface tolerance. As the gain loss in decibels is

$$\eta'_{\text{dB}} = 685.8 (\epsilon/\lambda)^2 \quad (\text{Eq. 84})$$

if the maximum permissible gain loss is set at 1 dB, the rms surface tolerance, ϵ , must be held to $\lambda/25$ or less. The ratio D/ϵ represents a measure of the quality of a parabolic antenna. For the above-mentioned values of the D/ϵ ratio, the corresponding range of the D/ϵ ratio is between 500 and 25 000. In electrically larger antennas intended for radio astronomy, the higher D/ϵ values have been achieved with protective means against the environment. In satellite communications systems, aside from early designs of large horn-type antennas with radome protection against the weather, open-air antennas are generally used.* The only weather protection device is electric heaters for melting snow and ice buildups.

The other three losses—blockage, spillover, and illumination—apply most significantly to folded optics systems such as the widely used Cassegrain configuration, which is characterized by a hyperbolic subreflector and parabolic main reflector. Since aperture illumination and far-field radiation are related by Fourier transforms, uniform illumination leads to a $\sin x/x$ type of radiation pattern, unity illumination factor, but relatively high-level side lobes. With tapering of the illumination, the illumination factor falls below unity, but better side-lobe control can be achieved.

Although folded optics systems such as the Cassegrain and Gregorian configurations provide perfect focusing, they do not provide it with optimum aperture illumination efficiency. Aperture efficiency can be enhanced by appropriately shaping subreflector and reflector to remove energy from the central area, blocked by the subreflector and feed, and redistributing it to achieve more uniform illumination across the remaining aperture. Such shaping can raise the gain of a reflector antenna by more than 1 dB. It is most effective for antennas with $D/\lambda > 100$.

The crowding of the geostationary orbit resulting from the expansion of satellite systems requires very strict control of antenna side lobes in order to satisfy the opposite requirements of minimum mutual interference and decreased in-orbit spacing (from 5° to 4° to 3° and even to 2°).

The side-lobe envelope of most electrically large Cassegrain antennas can be approximated by an expression of the form

$$G_{\text{dBi}} = A - B \log \theta \quad (\text{Eq. 85})$$

where A and B are constants and θ is the off-boresight angle. In 1965, the CCIR (now ITU-R) adopted the rule

$$G_{\text{dBi}} = \begin{cases} 32 - 25 \log \theta & 1^\circ \leq \theta \leq 48^\circ \\ 10 & \theta > 48^\circ \end{cases} \quad (\text{Eq. 86})$$

The antennas of the *INTELSAT* system have been standardized to follow the CCIR rule. Until 1977, side lobes beyond 1° from boresight were not to exceed the -29 dB level, but after 1977 a new rule was introduced requiring that no more than 10 percent of the side-lobe peaks exceed the envelope as defined above. A more recent recommendation would change the value of the constant A from 32 to 29 in the above expression.

Increasingly, to meet the more stringent side-lobe requirements, earth-station antennas are being built in offset configurations. By placing the feed structure below the reflector optical line of sight, the reflector aperture remains unblocked, and much lower side lobes can be realized.

The parameter that characterizes system performance is the G/T ratio, usually expressed in dB/K. All contributions to the system noise temperature must be properly taken into account once the measurement port is chosen. Antenna-feed design is of great importance with regard to both gain and side-lobe considerations (as well as polarization).

The amounts of spillover and illumination taper for both the subreflector and main reflector depend on the design of the feed system. Much progress has occurred in the transition from pyramidal to conical horns of various types such as single mode (TE_{11}), multimode ($\text{TE}_{11} + \text{TM}_{11}$), and hybrid mode (EH_{11}). In the last mentioned case, corrugated horns with $\lambda/4$ grooves have made it possible to increase bandwidth, improve symmetry of the radiation pattern, reduce the side lobes, and achieve better off-axis polarization characteristics, a feature of special importance in dual-polarization systems. In this case, the isolation between the two orthogonally polarized co-frequency channels should be as high as possible. Experience in systems using linear or circular orthogonal polarizations has confirmed that 30-dB isolation is a representative design goal.

Sophisticated orthomode junctions and an arrangement of cascaded polarizers are required to satisfy the above-mentioned requirements.

The actual isolation depends on the polarization purity of the signal source and of the antenna system. Since rain along the signal path depolarizes the signal, automatic means of depolarization correction have been successfully introduced at 6/4 GHz.

When needed, monopulse tracking systems are employed to correct continuously the pointing of the antenna in the direction of the satellite. To maintain continuity of service, an auxiliary power source is provided as well as redundant communications equipment.

* Shipboard antennas of maritime satellite communications systems are an exception.

Much of the foregoing description of earth stations applies to those used as gateways for the PSTN. In the 1980s, a marked change occurred in the application of satellite communications. Many classes of earth stations are unmanned, located on the user's premises, mobile and with low G/T , receive-only, electrically small, or combinations of these qualities. As previously mentioned, VSATs are usually at the user's premises, unmanned, and relatively small. USATs are similarly deployed. The numbers of user's premises earth stations is rising rapidly, as is the amount of traffic they carry. DBS-TV home stations are on the user's premises, receive-only, and electrically fairly small. Mobile terminals for aircraft, land mobile, and smaller ships are on the user's premises (vehicle), mobile and with low G/T (< -10 dB/K) (and thus often suffer from multipath), and quite small electrically. The emerging handset terminals are even smaller physically and electrically. Terminals used for IBS or its domestic equivalents are usually on the user's premises and often unmanned.

What has made much of this possible is the availability of highly reliable low-cost LSICs, VLSICs, and DSPs, along with advances in modulation, coding, and protocols for baseband modules, and solid-state MICs and more recently MMICs of moderate cost and high reliability at microwave frequencies. Combined with the higher eirps available in satellites in service in the 1980s and 1990s, services using earth terminals such as those described above have greatly expanded and the expansion will increase at least at the same rate in the first decade of the 21st century.

REFERENCES

1. Gordon, G. D., in *Communications Satellite Handbook*, W. L. Morgan, Ed. New York: John Wiley & Sons, 1989.
2. *Satellite Communications Systems: Systems, Techniques and Technology*, 3rd ed., Wiley Series in Communication and Distributed Systems by Gerard Maral, Michel Bousquet, and Michele Bousquet. New York: John Wiley & Sons, 1998.
3. *The Satellite Communication Applications Handbook* (Artech House Telecommunications Library) by Bruce R. Elbert. Norwood, MA: Artech House, 1997.
4. *Space Communications*, Special Issue "Global Satellite Communications Technology and Systems," Vol. 16, Nos. 2 and 3, 2000.
5. *Satellite Communications* (CRL Wave Summit Course series), T. Iida, Ed. Ohmsha, Ltd. and IOS Press, Netherlands, 2000.
6. *Low Earth Orbital Satellites for Personal Communication Networks* (Artech House Mobile Communications Library) by Abbas Jamalipour. Norwood, MA: Artech House, 1998.
7. Ippolito, L. J., *Propagation Effects Handbook for Satellite Systems Design*. NASA Reference Publication 1082(04), Feb. 1989.
8. *Proceedings of the IEEE*, Special Issue on "Satellite Communications," Vol. 78, No. 7, July, 1990.
9. Estabrook, P., and Motamedi, M. "Use of Non-geostationary Orbits for K_a -band Personal Array Satellite System." *13th AIAA COMSAT System Conference*, Los Angeles, 1990, pp. 14–24.
10. Bousquet, M., and Maral, G. "Orbital Aspects and Useful Relations from Earth-Satellite Geometry in the Frame of Future Mobile Systems." *13th AIAA COMSAT System Conference*, Los Angeles, 1990, pp. 783–789.
11. Maral, G., de Ridder, J. J., and Evans, B. G. "Low-Earth Orbit Satellite Systems for Communications." *International Journal of Satellite Communications*, Vol. 9, No. 4, July–August 1991, pp. 209–226.

28

Discrete-Time Signal Processing

*W. Kenneth Jenkins,
Douglas L. Jones, and Bill J. Hunsinger*

Introduction	28-2
Fundamentals for Discrete-Time Systems	28-2
Basic Definitions	
Finite Convolution and Difference Equations	
The Z-Transform	
The Discrete-Time Fourier Transform	
Sampling and Reconstruction	
Discrete Fourier Transform	28-7
Definitions and Properties	
FFT Algorithms	
The FFT in Spectral Analysis	
Digital-Filter Design	28-12
IIR Filters	
FIR Filters	
Digital-Filter Implementation	28-16
Network Structures	
Finite Wordlength Effects	
Adaptive Digital Signal Processing	28-22
The LMS Algorithm for FIR Adaptive Filters	
The LMS Algorithm for IIR Adaptive Filters	
Generic Applications of Adaptive Methods	
Discrete-Time Analog Technologies	28-29
Switched Capacitor Circuits	
Surface-Acoustic-Wave Filters	
SAW Filters for IS-95 Compatible CDMA Digital Receivers	
Minimum Requirements for SAW Performance	

INTRODUCTION

During the last twenty-five years, there has been rapid advancement in the theory and application of digital signal processing (DSP) in various engineering disciplines. Interest has grown in digital signal processing because, not only has the general-purpose computer become more readily available, but digital integrated circuits have become more highly integrated and cheaper, a trend that will continue into the foreseeable future. Very large scale integration (VLSI) techniques have produced high-density read-only memories (ROM) and microprocessors that provide enormous flexibilities in the design of digital hardware systems.

Digital filters offer distinct advantages over analog (continuous-time) filters in many applications, although they are not good substitutes for all analog filters. The major advantages are good numerical accuracy, programmability, stability in the presence of changing environmental conditions, suitability for multiplexing, and convenience for processing data that is directly available in binary form. Some of the disadvantages of these filters are the relatively high per-unit costs for high-frequency applications, frequency limitations imposed by the speed of the digital hardware, and the necessity for a significant amount of clocking and control circuitry to sequence the binary operations properly.

The *design* of a digital filter involves determining either a set of time-domain difference equations or a z -domain digital transfer function that satisfies given specifications. A digital filter can be obtained by first designing an analog prototype and then transforming it into a discrete-time system by a sampled-data transformation. Another approach is to use a computer optimization to place the z -domain poles and zeros so the discrete-time system will meet specifications directly. The first approach takes advantage of well known analog design techniques, while the second provides greater flexibility because it does not depend on an analog design step. Digital filter *implementation* involves choosing a network topology and hardware modules for the final network. At this stage, the designer must analyze the effects of quantization error, because error performance and network topology are closely related. Digital *hardware design* consists of designing the individual circuit elements (adders, multipliers, shift registers, etc.). If the system is to be integrated, it also includes the IC layouts.

In many applications, implementation is ultimately accomplished in software on a general-purpose computer. In these cases, the emphasis is on the design and implementation stages, since hardware design and system architecture are dictated by the general computer system. However, with the rapid advances that are now being made in the automated design and manufacture of VLSI monolithic circuits, engineers enjoy the freedom to specify custom designed digital functions and have them quickly fabricated in low-cost silicon devices. This new capability will result in digital signal processing becoming less dependent on the general

computer. New techniques for improving data rates, reducing circuit complexity, and improving reliability will become increasingly important as more custom-designed VLSI digital systems come into common usage.

Several sampled-data analog technologies exist which implement discrete-time filters. While lacking the flexibility of all-digital implementations, these technologies, which include switched-capacitor filters, surface-acoustic-wave filters, and acoustic-charge-transport devices, often represent the most cost-effective method of filtering in certain performance regimes.

FUNDAMENTALS FOR DISCRETE-TIME SYSTEMS

Basic Definitions

A *continuous-time* (CT) signal is a function, $s(t)$, that is defined for all time t contained in some interval on the real line. For historical reasons, CT signals are often called *analog signals*. If the domain of definition for $s(t)$ is restricted to a set of discrete points $t_n = nT$, where n is an integer and T is the sampling period, the signal $s(t_n)$ is called a *discrete-time* (DT) signal. Often, if the sampling interval is well understood within the context of the discussion, the sampling period is normalized by $T = 1$, and a DT signal is represented simply as a sequence $s(n)$. If the values of the sequence $s(n)$ are to be represented with a finite number of bits (as required in a finite state machine), then $s(n)$ can take on only a discrete set of values. In this case, $s(n)$ is called a *digital signal*. Much of the theory that is used in DSP is actually the theory of DT signals and DT systems, in that no amplitude quantization is assumed in the mathematics. However, all signals processed in binary machines are truly digital signals. One important question that arises in virtually every application is the question of how many bits are required in the representation of the digital signals to guarantee that the performance of the digital system is acceptably close to the performance of the ideal DT system.

Linear CT systems are characterized by the familiar mathematics of differential equations, continuous convolution operators, Laplace transforms, and Fourier transforms. Similarly, linear DT systems are described by the mathematics of difference equations, discrete convolution operators, Z -transforms, and discrete Fourier transforms. It appears that for every major concept in CT systems, there is a similar concept for DT systems (e.g., differential equations and difference equations, continuous convolution and discrete convolution, etc.). However, in spite of this duality of concepts, it is impossible to apply directly the mathematics of CT systems to DT systems, or vice versa.

Many modern systems consist of both analog and digital subsystems, with appropriate analog-to-digital (A/D) and digital-to-analog (D/A) devices at the inter-

faces. For example, it is common to use a digital computer in the control loop of an analog plant. Analytical difficulties often occur at the boundaries between the analog and digital portions of the system because the mathematics used on the two sides of the interface must be different. It is often useful to assume that a sequence $s(n)$ is derived from an analog signal $s_a(t)$ by ideal sampling, i.e.

$$s(n) = s_a(t) \Big|_{t=nT} \quad (\text{Eq. 1})$$

An alternative model for the sampled signal is denoted by $s^*(t)$ and defined by

$$s^*(t) = \sum_{n=-\infty}^{+\infty} s_a(t) \delta_a(t - nT) \quad (\text{Eq. 2})$$

where $\delta_a(t)$ is an analog impulse function. Both $s(n)$ and $s^*(t)$ are used throughout the literature to represent an ideal sampled signal. Note that even though $s(n)$ and $s^*(t)$ represent the same essential information, $s(n)$ is a DT signal and $s^*(t)$ is a CT signal. Hence, they are not mathematically identical. In fact, $s(n)$ is a "DT-world" model of a sampled signal, whereas $s^*(t)$ is a "CT-world" model of the same phenomenon.

Finite Convolution and Difference Equations

Let $y(n) = F\{x(n)\}$ define the input-output relation for a discrete-time system with input $x(n)$ and output $y(n)$. The following definitions are commonly used to define the properties of linearity, shift invariance, causality, and stability:

Linearity- F is linear if and only if $F\{ax_1(n) + bx_2(n)\} = aF\{x_1(n)\} + bF\{x_2(n)\}$, where a and b are scalar constants and $x_1(n)$ and $x_2(n)$ are two arbitrary input sequences.

Shift-Invariance- F is shift-invariant if and only if $y(n - n_0) = F\{x(n - n_0)\}$ for all integer values of n_0 .

Causality - F is causal if and only if $y(n)$ depends on samples of $x(k)$ at times $k \leq n$.

Stability (BIBO)- F is stable in the bounded-input bounded-output (BIBO) sense if and only if $y(n)$ remains bounded for all $x(n)$ that are bounded.

Whenever F is linear and shift-invariant, it is possible to express the zero state response, $y(n)$, due to an arbitrary input, $x(n)$, in terms of a discrete convolution of $h(n)$ and $x(n)$, where $h(n) = F\{\delta(n)\}$, and $\delta(n) = 0$ for $n \neq 0$ and $\delta(n) = 1$ when $n = 1$. The function $\delta(n)$ is called a unit pulse, $h(n)$ is the unit-pulse response* of F , and the finite convolution is expressed by

$$y(n) = \sum_{k=-\infty}^{+\infty} h(k) x(n - k) \quad (\text{Eq. 3})$$

* The term "impulse response" is often used interchangeably with the term "unit pulse response."

If F is causal, $h(k) = 0$ for $k < 0$, and the lower limit on the summation in Eq. 3 becomes zero.

In certain types of DT systems, the unit-pulse response is zero outside of a finite interval containing N samples; i.e., $h(n) \equiv 0$ for $n < 0$ and $n \geq N$. This type of system is a finite impulse response (FIR) DT system. If $h(n)$ is supported over an infinite length interval, then F is called an infinite impulse response (IIR) DT system. These two classes constitute the two important types of digital filters, with each class having distinct advantages and disadvantages with regard to stability, quantization error performance, and computational efficiency.

An IIR system has an infinite memory because, in general, the output depends on the input all the way into the infinite past. Such a system can also be characterized by an N th order linear difference equation, as given by

$$y(n) + a_1y(n - 1) + \dots + a_{N-1}y(n - N + 1) = b_0x(n) + \dots + b_{M-1}x(n - M + 1) \quad (\text{Eq. 4})$$

A linear difference equation is a recursive relation that can be realized quite easily with digital multipliers, adders, and unit delay registers. Hence, in some of the older literature, IIRs are referred to as recursive filters. Similarly, FIRs were often referred to as nonrecursive (or transversal) filters. However, the terms recursive and nonrecursive refer more specifically to the implementation of the filter, rather than the mathematical structure. This fact is illustrated by the frequency sampling structure,† which is a recursive realization of a FIR system. The terms IIR and FIR are used in the modern literature to eliminate any ambiguity that was caused by the older terminology.

In general, IIR systems are cheaper to implement than FIR systems because the iteration of a difference equation requires fewer arithmetic operations per output sample, as compared to calculating a finite convolution. However, IIR systems have "poles," whereas FIR systems do not. This implies that IIR systems must be carefully designed to ensure stability. Also, IIR systems suffer from limit cycle oscillations and quantization error accumulation because quantization efforts are recycled through the inherent feedback in the recursion.

The Z-Transform

The Z-transform occupies the same position of importance in DT system theory as the Laplace transform does in CT system theory. They are very similar transforms that share many common properties. However, it is important to emphasize that they cannot be used interchangeably because they apply to different types of systems.

† Reference 25.

The 2-sided Z-transform of a DT signal, $s(n)$, is defined by

$$S(z) = \mathcal{Z}\{s(n)\} = \sum_{n=-\infty}^{+\infty} s(n)z^{-n} \quad (\text{Eq. 5})$$

where $S(z)$ is said to exist for all $z \in R$ such that the infinite summation converges. Symbol R represents the region of convergence (R.O.C.) of $S(z)$. If the summation is taken from $n = 0$, rather than $n = -\infty$ the result is called the 1-sided Z-transform. It is clear that for all $s(n)$ such that $s(n) \equiv 0$ for $n < 0$, the 1-sided and 2-sided Z-transforms are equivalent.

The inverse Z-transform is defined by

$$s(n) = (1/2\pi j) \oint_C S(z)z^{n-1} dz \quad (\text{Eq. 6})$$

where C is a closed contour that encircles the origin and lies entirely within the region of convergence for $S(z)$. Eq. 6 is valid for both positive and negative values of n . In many applications, the inverse Z-transform can be found from a table of Z-transform pairs, and therefore it is seldom necessary to find the inverse Z-transform from its definition. However, for cases where it is desired to apply the definition of Eq. 6, the

contour integral can be evaluated by means of the well-known residue theorem.

$$\begin{aligned} & (1/2\pi j) \oint_C S(z) z^{n-1} dz \\ &= \sum (\text{residues of } S(z)z^{n-1} \text{ at the poles inside } C) \end{aligned} \quad (\text{Eq. 7})$$

The normal technique for finding the residues of $S(z) z^{n-1}$ is by the same method used in expanding $S(z) z^{n-1}$ in a partial fraction expansion. It is important to note that the residue theorem is applicable only when $S(z) z^{n-1}$ has singularities (poles) that occur at isolated points in the z -plane. Fortunately, virtually all one-dimensional DT systems and DT signals of practical interest produce Z-transforms that are ratios of polynomials in powers of z . The singular behavior of such systems is characterized by poles and zeros, and hence the residue theorem can be applied in virtually all cases of practical interest. Table 1 summarizes a number of Z-transform pairs that are frequently encountered in practical problems. More complete tables of Z-transforms can be found in References 16 and 30. Note that most of the sequences in Table 1 are right-

TABLE 1. Z-TRANSFORMS

$x(n)$	$X(z)$
1. $\delta(n)$ (unit pulse)	1, for all z
2. $u(n)$ (unit step)	$1/(1-z^{-1}), 1 < z $
3. $n u(n)$	$1/(1-z^{-1})^2, 1 < z $
4. $n^r a^n u(n)$ ($r = \text{integer}$)	$-z(d^r/dz^r)\{1/(1-az^{-1})\}, a < z $
5. $a^n u(n)$	$1/(1-az^{-1}), a < z $
6. $n^r a^n u(-n)$ ($r = \text{integer}$)	$-z(d^r/dz^r)\{1/(1-az^{-1}), z < a $
7. $a^{ n }$	$(a-a^2)/(1-az)(1-az^{-1}), a < z < 1/a $
8. $n^{-1} u(n)$	$-1n(1-r^{-1}), 1 < z $
9. $\cos(an) u(n)$	$\frac{(1-z^{-1}\cos a)}{(1-2z^{-1}\cos a+z^{-2})}, 1 < z $
10. $\sin(an) u(n)$	$\frac{z^{-1}\sin a}{(1-2z^{-1}\cos a+z^{-2})}, 1 < z $
11. $(K_1 \cos an + K_2 \sin an)u(n)$	$\frac{K_1 + z^{-1}(K_2 \sin a - K_1 \cos a)}{(1-2z^{-1}\cos a+z^{-2})}, 1 < z $
12. $(K_1 \cos an + G \sin an) u(n)$ where, $G = (K_2 + K_1 \cos a)/\sin a$	$\frac{K_1 + K_2 z^{-1}}{(1-2z^{-1}\cos a+z^{-2})}, 1 < z $
13. $(n+1)a^n u(n)$	$1/(1-az^{-1})^2, a < z $
14. $(1/2)(n+1)(n+2)a^n u(n)$	$1/(1-az^{-1})^3, a < z $
15. $[1/(r-1)!](n+1)(n+2) \dots (n+r-1)a^n u(n)$	$1/(1-az^{-1})^r, a < z $

sided; i.e., $x(n) = 0$ for $n < 0$. The exceptions are entry 6, which is left-sided, and entry 7, which is two-sided. For a given $X(z)$, the corresponding $x(n)$ may be right-sided, left-sided, or two-sided, depending on the nature of the region of convergence, R . Fig. 1 shows the three general types of regions of convergence for these three cases.

Table 2 summarizes some important properties of the Z-transform. Although the table is specifically given for the two-sided Z-transform, most of the prop-

erties hold for the 1-sided Z-transform also. A notable exception is property 4, the shift theorem. Since the shift theorem for 1-sided sequences is often needed for solving difference equations, it is stated here:

If $X(z)$ is the 1-sided Z-transform of $x(n)$ with a region of convergence R , then for $n_0 > 0$,

$$(i) \mathfrak{Z}\{x(n+n_0)\} = z^{n_0} \left\{ X(z) - \sum_{k=0}^{n_0-1} x(k)z^{-k} \right\}$$

$$(ii) \mathfrak{Z}\{x(n-n_0)\} = z^{-n_0} \left\{ X(z) + \sum_{k=-n_0}^{-1} x(k)z^{-k} \right\}$$

with region of convergence R .*

The Discrete-Time Fourier Transform

It is important to distinguish between the concepts of the discrete-time Fourier transform (DTFT) and the discrete Fourier transform (DFT). The DTFT is a transform-pair relationship between a DT signal and its continuous-frequency transform that is used extensively in the analysis and design of DT systems. In contrast, the DFT is a transform-pair relationship between a DT signal, and its discrete-frequency transform that is used in practical digital processing.

The DTFT of a sequence $x(n)$ is defined by

$$X(e^{j\omega}) = \sum_{n=-\infty}^{+\infty} x(n)e^{-jn\omega} \quad (\text{Eq. 8})$$

for all sequences $x(n)$ for which the sum converges. The spectrum $X(e^{j\omega})$ is a continuous-frequency function which is periodic, with period 2π . Eq. 8 is equivalent to evaluating the 2-sided Z-transform $X(z) = Z\{x(n)\}$ on the unit circle $z = e^{j\omega}$, provided that the unit circle lies within R , the region of convergence of $X(z)$. A sequence $x(n)$ can be recovered from its spectrum by

$$x(n) = (1/2\pi) \int_{-\pi}^{+\pi} X(e^{j\omega}) e^{jn\omega} d\omega \quad (\text{Eq. 9})$$

which can be interpreted on the inverse-DTFT (IDTFT). An alternate interpretation is that Eq. 8 is a Fourier series expansion of the periodic function $X(e^{j\omega})$, the $x(n)$ s are the Fourier coefficients, and Eq. 9 is an expression for finding the Fourier coefficients.

The DTFT obeys the same properties as those given for the Z-transform in Table 2, as long as the unit circle is contained within the various regions of convergence. In particular, the DTFT obeys the convolutional property (entry 11), which makes it useful for frequency analyses of DT systems. If $h(n)$ is the unit pulse response of a linear shift-invariant DT system, then

* Reference 25 describes pathological cases where R for the shifted sequence may no longer include the points $z = 0$ or $z = \infty$.

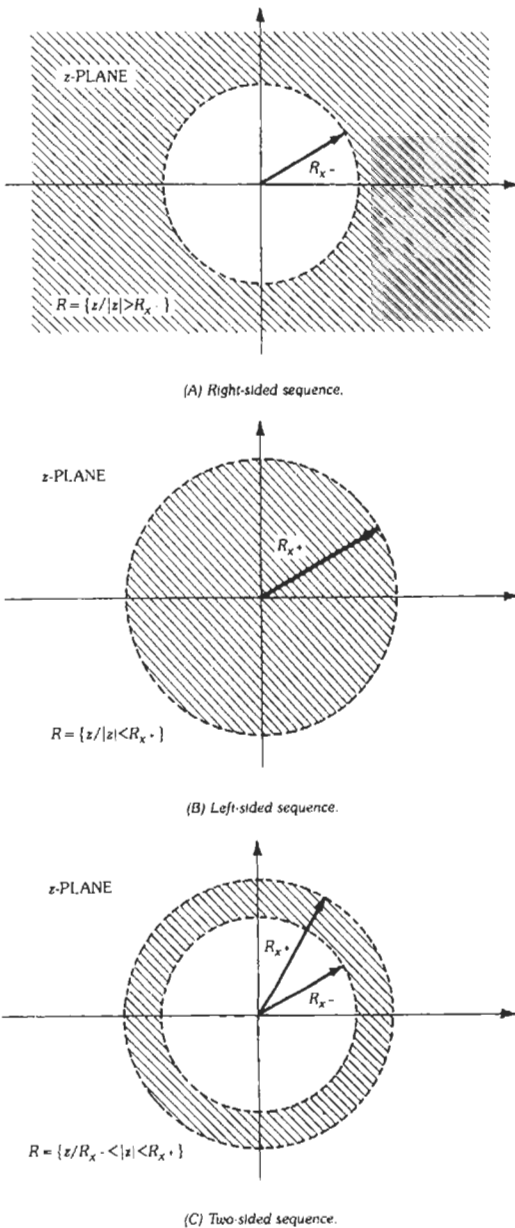


Fig. 1. Typical regions of convergence for Z-transforms.

TABLE 2. SOME IMPORTANT PROPERTIES OF THE 2-SIDED Z-TRANSFORM*

Sequence		Z-Transform
1.	$x(n)$	$X(z)$ $R_{x-} < z < R_{x+}$
2.	$y(n)$	$Y(z)$ $R_{y-} < z < R_{y+}$
3.	$ax(n) + by(n)$	$aX(z) + bY(z)$ $\max [R_{x-}, R_{y-}] < z < \min [R_{x+}, R_{y+}]$
4.	$x(n + n_0)$	$z^{n_0}X(z)$ $R_{x-} < z < R_{x+}$
5.	$a^n x(n)$	$X(a^{-1}z)$ $ a R_{x-} < z < a R_{x+}$
6.	$nx(n)$	$-z(dX(z)/dz)$ $R_{x-} < z < R_{x+}$
7.	$x^*(n)$	$X^*(z^*)$ $R_{x-} < z < R_{x+}$
8.	$x(-n)$	$X(1/z)$ $1/R_{x+} < z < 1/R_{x-}$
9.	$Re[x(n)]$	$(1/2)[X(z) + X^*(z^*)]$ $R_{x-} < z < R_{x+}$
10.	$Im[x(n)]$	$(1/2j)[X(z) - X^*(z^*)]$ $R_{x-} < z < R_{x+}$
11.	$x(n) * y(n)$	$X(z)Y(z)$ $\max [R_{x-}, R_{y-}] < z < \min [R_{x+}, R_{y+}]$
12.	$x(n)y(n)$	$(1/2\pi j) \oint X(v)Y(z/v)v^{-1}dv$ $R_{x-}R_{y-} < z < R_{x+}R_{y+}$

* From A. V. Oppenheim and R. W. Schaffer. *Digital Signal Processing*. Englewood Cliffs, N. J.: Prentice-Hall, Inc., 1975.

$H(e^{j\omega}) = \text{DTFT} \{h(n)\}$ is interpreted as the frequency response of the DT system. Therefore, if a discrete sinusoidal waveform $x(n) = A \cos(n\omega_0)$ is applied as an input to a linear shift-invariant DT system, the sinusoidal steady-state output is given by $y_{ss}(n) = A |H(e^{j\omega_0})| \cos[n\omega_0 + \angle H(e^{j\omega_0})]$.

Sampling and Reconstruction

The traditional CT Fourier transform can be applied to the ideal sampled signal, $s^*(t)$, as follows:

$$\begin{aligned} \mathcal{F}\{s^*(t)\} &= \sum_{n=-\infty}^{+\infty} s(nT) \mathcal{F}\{\delta_a(t - nT)\} \\ &= \sum_{n=-\infty}^{+\infty} s_a(nT) e^{-jnT\omega} \\ &= \text{DTFT}\{s_a(nT)\} \end{aligned} \tag{Eq. 10}$$

This verifies that the traditional Fourier transform of $s^*(t)$ is identical to the DTFT of $s(n)$, where $s(n) \equiv s_a(t)|_{t=nT}$; i.e., the sequence $s(n)$ is derived from $s_a(t)$ by ideal sampling. This proves that $s^*(t)$ and $s(n)$ are really different models of the same phenomenon, since their spectra (as computed with appropriate transforms) are identical. Suppose that $S_a(j\Omega)$ is the spec-

trum of $s_a(t)$ and $S(e^{j\omega})$ is the spectrum of $s(n)$. It can be shown that

$$S(e^{j\omega}) = (1/T) \sum_{r=-\infty}^{+\infty} S_a(j[\Omega - 2\pi r/T]) \tag{Eq. 11}$$

where $\omega = \Omega T$ is often referred to as the "normalized digital frequency." Eq. 11 shows that the DT spectrum is formed from a superposition of an infinite number of replicas of the analog signal, as illustrated in Fig. 2. As long as the sampling frequency, $\Omega_s = 2\pi/T$, is chosen so that $\Omega_s > 2\Omega_B$, where Ω_B is the highest frequency component contained in $s_a(t)$, then each period of $S(e^{j\omega})$ contains a perfect copy of $S_a(j\Omega)$, and $s_a(t)$ can be recovered exactly from $s(n)$ by ideal low-pass filtering. Sampling under these conditions is said to satisfy the Nyquist sampling criterion, since the sampling frequency exceeds the Nyquist rate, $2\Omega_B$. If the sampling rate does not satisfy the Nyquist criterion, the adjacent periods of the analog spectrum will overlap, causing a distorted spectrum (see Fig. 2). This effect, called *aliasing distortion*, is rather serious because it cannot be easily corrected once it has occurred. In general, an analog signal should be prefiltered with an analog low-pass filter prior to sampling so that aliasing distortion does not occur.

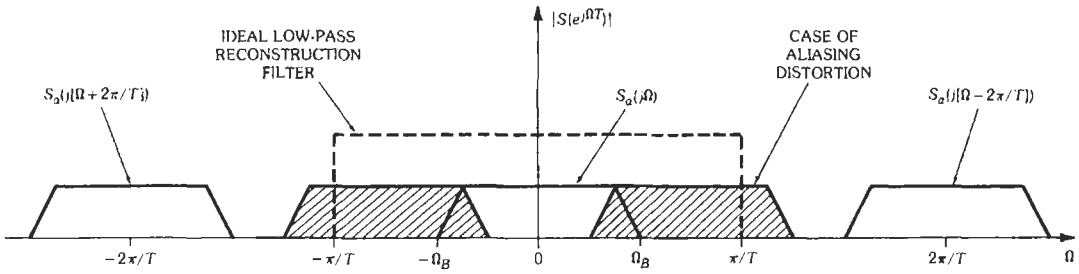


Fig. 2. Relationship between the spectrum of an analog signal and the spectrum of the ideally sampled signal.

If the Nyquist criterion has been satisfied, it is always possible to reconstruct an analog signal from its samples according to

$$s_a(t) = \sum_{k=-\infty}^{+\infty} s_a(kT) \text{sinc}\left(\left(\frac{\pi}{T}\right)(t - kT)\right) \quad (\text{Eq. 12})$$

This reconstruction formula results by filtering $s_a^*(t)$ with an ideal low-pass analog filter with a bandwidth of $\Omega_B = \pi/T$ (see Fig. 2). In general, exact reconstruction requires an infinite number of samples, although a good approximation can be obtained by using a large but finite number of terms in Eq. 12.

Most practical systems use a digital-to-analog converter for reconstruction, which results in an analog staircase approximation to the true analog signal; i.e.

$$\hat{S}_a(t) = \sum_{k=-\infty}^{+\infty} s_a(kT) [u(t - kT) - u(t - (k + 1)T)] \quad (\text{Eq. 13})$$

It can be shown that $\hat{S}_a(t)$ is obtained by filtering $s_a^*(t)$ with an analog filter whose frequency response is

$$H_a(j\Omega) = 2Te^{-j\Omega T/2} \text{sinc}(\Omega T/2) \quad (\text{Eq. 14})$$

The approximation $\hat{S}_a(t)$ is said to contain “sin x/x distortion,” which occurs because $H_a(j\Omega)$ is not an ideal low-pass filter. The $H_a(j\Omega)$ response distorts the signal by causing a droop near the band edge, as well as passing high-frequency distortion terms that “leak” through the side lobes of $H_a(j\Omega)$. Therefore, a practical D/A converter is normally followed by a postfilter

$$H_p(j\Omega) = \begin{cases} H_a^{-1}(j\Omega), & 0 \leq |\Omega| \leq \pi/T \\ 0, & \Omega \text{ otherwise} \end{cases} \quad (\text{Eq. 15})$$

which compensates for the distortion and produces the correct $s_a(t)$ analog output. Notice, however, that $H_p(j\Omega)$ can only be approximated in practice, so that the best reconstruction is necessarily an approximation. Fig. 3 shows a digital processor complete with sampling and reconstruction devices at the input and output.

DISCRETE FOURIER TRANSFORM

Definitions and Properties

Let $x(n)$ be a finite-length sequence of length N . The discrete Fourier transform (DFT) pair is defined by

$$X(k) = \sum_{n=0}^{N-1} x(n) e^{-j(2\pi/N)nk} \quad (\text{Eq. 16a})$$

$$k = 0, \dots, N - 1$$

and

$$x(n) = (1/N) \sum_{k=0}^{N-1} X(k) e^{j(2\pi/N)nk} \quad (\text{Eq. 16b})$$

$$n = 0, \dots, N - 1$$

The DFT treats $x(n)$ as though it were one period of a periodic sequence. This is an important feature that must be handled properly in signal processing to prevent the introduction of artifacts. Important properties of the DFT are summarized in Table 3. The notation $((k))_N$ denotes the least positive residue of k modulo N , and $R_N(n)$ is a rectangular window function such that $R_N(n) = 1$ for $n = 0, \dots, N - 1$, and $R_N(n) = 0$ for $n < 0$ and $n \geq N$.

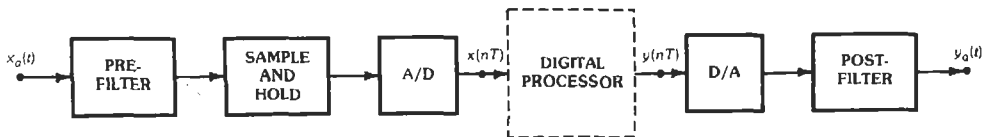


Fig. 3. Elements required for the digital processing of an analog signal.

TABLE 3. PROPERTIES OF THE DFT*

Finite-Length Sequence (length N)	DFT
1. $x(n)$	$X(k)$
2. $y(n)$	$Y(k)$
3. $ax(n) + by(n)$	$aX(k) + bY(k)$
4. $x((n+m))_N \mathfrak{R}_N(n)$	$W_N^{-km} X(k)$
5. $W_N^{kn} x(n)$	$X((k+l))_N \mathfrak{R}_N(k)$
6. $\left[\sum_{m=0}^{N-1} x((m))_N y((n-m))_N \right] \mathfrak{R}_N(n)$	$X(k)Y(k)$
7. $x(n)y(n)$	$(1/N) \left[\sum_{l=0}^{N-1} X((l))_N Y((k-l))_N \right] \mathfrak{R}_N(k)$
8. $x^*(n)$	$X^*((-k))_N \mathfrak{R}_N(k)$
9. $x^*((-n))_N \mathfrak{R}_N(n)$	$X^*(k)$
10. $\text{Re} [x(n)]$	$X_{cp}(k) = \frac{1}{2} [X((k))_N + X^*((-k))_N] \mathfrak{R}_N(k)$
11. $j \text{Im} [x(n)]$	$X_{op}(k) = \frac{1}{2} [X((k))_N - X^*((-k))_N] \mathfrak{R}_N(k)$
12. $x_{cp}(n)$	$\text{Re} [X(k)]$
13. $x_{op}(n)$	$j \text{Im} [X(k)]$
The following properties apply only when $x(n)$ is real:	
14. Any real $x(n)$	$\begin{cases} X(k) = X^*((-k))_N \mathfrak{R}_N(k) \\ \text{Re} [X(k)] = \text{Re} [X(((-k))_N) \mathfrak{R}_N(k)] \\ \text{Im} [X(k)] = \text{Im} [X(((-k))_N) \mathfrak{R}_N(k)] \\ X(k) = X(((-k))_N) \mathfrak{R}_N(k) \\ \arg [X(k)] = -\arg [X(((-k))_N) \mathfrak{R}_N(k)] \end{cases}$
15. $x_{cp}(n)$	$\text{Re} [X(k)]$
16. $x_{op}(n)$	$j \text{Im} [X(k)]$

* From A. V. Oppenheim and R. W. Schaffer, *Digital Signal Processing*. Englewood Cliffs, N. J.: Prentice-Hall, Inc., 1975.

Most of the properties given in Table 3 for the DFT are the same as those of the z -transform, although there are some important differences. Property 4, the shift property, holds for *circular* shifts of the finite-length sequence $x(n)$. Also, the multiplication of two DFTs results in the *circular convolution* of the corresponding time-domain sequences (property 6).

Suppose it is desired to implement an FIR filter

$$y(n) = \sum_{k=0}^{N-1} h(k)x(n-k)$$

by transforming $h(n)$ and $x(n)$ into $H(k)$ and $X(k)$ using a DFT, multiplying pointwise to obtain $Y(k) = H(k)X(k)$, and then using the inverse-FFT to obtain $y(n) = \text{IFFT}\{Y(k)\}$. If $x(n)$ is a finite-length sequence of length M , then the results of the circular convolution implemented by the FFT will correspond to the desired

linear convolution if and only if $N_{\text{FFT}} \geq N + M - 1$, where N_{FFT} is the block length of the FFT and $h(n)$ and $x(n)$ are each padded with zeros to create blocks of length N_{FFT} . In some applications, either the value of M is too large for the memory available, or $x(n)$ may not be finite in length, but rather it may be a continual stream of data that arrives at the filter for real-time processing. Two algorithms are available that partition $x(n)$ into smaller blocks and process the individual blocks with a smaller-length FFT: (1) overlap-save partitioning and (2) overlap-add partitioning. These are summarized below.

Overlap-Save—In this algorithm, N_{FFT} is chosen to be some convenient value with $N_{\text{FFT}} > N$. The signal, $x(n)$, is partitioned into blocks which are of length N_{FFT} and which overlap by $N - 1$ data points. Hence, the k th block is $x_k(n) = x(n + k(N_{\text{FFT}} - N + 1))$, $n = 0, \dots$

$N_{\text{FFT}} - 1$. The filter is augmented with $N_{\text{FFT}} - N$ zeros to produce

$$\hat{h}(n) = \begin{cases} h(n) & n = 0, \dots, N-1 \\ 0 & N \leq n < N_{\text{FFT}} \end{cases}$$

The FFT is then used to obtain $\hat{Y}(k) = \text{FFT}\{\hat{h}(n)\} \cdot \text{FFT}\{x_i(n)\}$, and $\hat{y}_k(n) = \text{IFFT}\{\hat{Y}(k)\}$. From the $\hat{y}_k(n)$ array, the values that correctly correspond to the linear convolution are saved; values that are “contaminated by wrap-around error” are discarded; i.e.,

$$y_k(n) = \begin{cases} \hat{y}_k(n), & N-1 \leq n \leq N_{\text{FFT}} - 1 \\ 0, & 0 \leq n < N-1 \end{cases}$$

The blocks of output data $y_k(n)$ are then assembled to produce the correct output by

$$y(n) = \sum_{k=0}^{\infty} y_k(n - k(N_{\text{FFT}} - N + 1))$$

For the overlap-save algorithm, each time a block is processed there are $N_{\text{FFT}} - N + 1$ points saved and $N - 1$ points discarded. Each block moves forward by $N_{\text{FFT}} - N + 1$ data points and overlaps the previous block by $N - 1$ points.

Overlap-Add—This algorithm is similar to the previous one except that the k th input block is defined to be

$$x_k(n) = \begin{cases} x(n + kL), & n = 0, \dots, L-1 \\ 0 & , \quad n = L, \dots, N_{\text{FFT}} - 1 \end{cases}$$

where $L = N_{\text{FFT}} - N + 1$. The filter function $\hat{h}(n)$ is augmented with zeros, as before, to create $\hat{h}(n)$, and the DFT processing is executed exactly as before. In each block $\hat{y}_k(n)$, that is obtained at the output, the first $N - 1$ points are “bad,” the last $N - 1$ points are “bad,” and the middle $N_{\text{FFT}} - 2(N - 1)$ points are “good.” However, if the last $N - 1$ points from block k are overlapped with the first $N - 1$ points of block $k + 1$ and added pointwise, correct results corresponding to linear convolution are obtained in these positions. Hence, after this addition the number of correct points produced per block is $N_{\text{FFT}} - (N - 1)$, which is the same as that for the overlap-save algorithm. The overlap-add algorithm requires essentially the same amount of computation as the overlap-save algorithm, although the addition of the overlapping blocks is extra. This feature, together with the extra delay of waiting for the next block to be finished before the previous one is complete, has resulted in more popularity for the overlap-save algorithm in practical applications.

These block-filtering algorithms make it possible to filter continual data streams in real time because the efficiency of the FFT minimizes the total computation time and can achieve reasonably high overall data rates. However, block filtering generates data in bursts; i.e., there is a delay during which no filtered data appears, and then suddenly an entire block is generated. In real-time systems, buffering must be used. The

block algorithms are particularly effective for filtering very long sequences of data that are prerecorded on tape or magnetic disk.

FFT Algorithms

“Fast Fourier transform” (FFT) is a generic name for a class of algorithms that efficiently compute the DFT (see Eqs. 16). The FFT is easily understood by examining a radix-2 FFT for the case $N = 2^3$. First, each of the indices k and n can be expressed in binary form, $k = k_2 4 + k_1 2 + k_0$ and $n = n_2 4 + n_1 2 + n_0$, where k_i and n_i are bits that take the values of either 0 or 1. If these expressions are substituted into Eq. 16a, all terms in the exponent that contain the factor $N = 8$ can be simply deleted, because $e^{j2\pi l} = 1$ for any integer l . Upon deleting such terms and regrouping, the product nk can be expressed in one of two ways:

$$nk = (4k_0)n_2 + (4k_0 + 2k_0)n_0 + (4k_2 + 2k_1 + k_0)n_0 \tag{Eq. 17a}$$

$$nk = (4n_0)k_2 + (4n_1 + 2n_0)k_1 + (4n_2 + 2n_1 + n_0)k_0 \tag{Eq. 17b}$$

Substituting Eq. 17a into Eq. 16a leads to the decimation-in-time (D-I-T) FFT, whereas substituting Eq. 17b into Eq. 16a leads to the decimation-in-frequency (D-I-F) FFT. Only the D-I-T FFT is discussed further here. The D-I-F FFT and various related forms are treated in detail in Reference 4.

The D-I-T FFT decomposes into $\log_2 N$ stages of computation, plus a stage of bit reversal. (Let $W_N = e^{j2\pi/N}$)

$$x_1(k_0, n_1, n_0) = \sum_{n_2=0}^1 x(n_2, n_1, n_0) W_8^{4k_2 n_2} \tag{stage 1}$$

$$x_2(k_0, k_1, n_0) = \sum_{n_1=0}^1 x_1(k_0, n_1, n_0) W_8^{(4k_1 + 2k_0)n_1} \tag{stage 2}$$

$$x_3(k_0, k_1, k_2) = \sum_{n_0=0}^1 x_2(k_0, k_1, n_0) W_8^{(4k_2 + 2k_1 + k_0)n_0} \tag{stage 3}$$

$$X(k_2, k_1, k_0) = x_3(k_0, k_1, k_2) \tag{bit reversal}$$

In each summation above, one of the n_i 's is summed out of the expression, while at the same time a new k_i is introduced. The notation is chosen to reflect this. For example, in stage 3, n_0 is summed out, k_2 is introduced as a new variable, and n_0 is replaced by k_2 in the result. The last operation, called bit reversal, is necessary to correctly locate the frequency samples $X(k)$ in the memory. It is easy to show that if the samples are

paired correctly, an in-place computation can be done by a sequence of butterfly computations. For example, in stage 3 the $k = 6$ and $k = 7$ samples should be paired, yielding a butterfly computation that requires one complex multiply, one complex add, and one complex subtract.

$$x_3(1,1,0) = x_2(1,1,0) + W_8^3 x_2(1,1,1)$$

$$x_3(1,1,1) = x_2(1,1,0) - W_8^3 x_2(1,1,1)$$

Therefore, samples $x_2(6)$ and $x_2(7)$ are read from memory, the butterfly is executed, and $x_3(6)$ and $x_3(7)$ are written back to memory, thereby destroying the original values of $x_2(6)$ and $x_2(7)$. In general, there are $N/2$ butterflies per stage and $\log_2 N$ stages, so the total number of butterflies is $(N/2) \log_2 N$. Since there is at most one complex multiply per butterfly, the total number of multiplies is bounded by $(N/2) \log_2 N$ (some of the multiplies involve factors of 1, and should not be counted). A direct computation of the DFT requires on the order of N^2 complex multiplies. It is obvious that for large N , the order of complexity of the FFT, $O\{(N/2) \log_2 N\}$, is much less than the order of complexity $O\{N^2\}$ of the DFT. The saving is so dramatic for large N that the FFT has made possible the solution of many DSP problems that are compute-bound and impossible to solve with the direct DFT.

Fig. 4 shows the signal-flow graph for the D-I-T FFT with $N = 8$, which is referred to as an in-place FFT with normally ordered input and bit reversed output. Minor variations that include bit reversed input and normally ordered output, and non-in-place algorithms with normally ordered inputs and outputs are

possible. Also, when N is not a power of 2, a mixed-radix algorithm can be used to reduce computation.* The mixed-radix FFT is most efficient when N is highly L composite, i.e.,

$$N = \prod_{i=1}^L P_i^{r_i}$$

where the p_i 's are small primes and the r_i 's are positive integers. It can be shown that the order of complexity of the mixed-radix FFT is

$$O\left\{N \sum_{i=1}^L r_i (p_i - 1)\right\}$$

Because of the lack of uniformity of structure among stages, this algorithm has not received much attention for hardware implementation. However, the mixed-radix FFT is often used in software applications, especially for processing data recorded in laboratory experiments where it is not convenient to restrict the block lengths to be powers of 2. The most widely used mixed-radix FFT program was published by Singleton.† Singleton's program is also contained in a collection of DSP programs published by the IEEE Signal Processing Society.

Other advanced algorithms, such as higher-radix algorithms and the prime-factor algorithm, are described in the literature. Algorithms specialized for

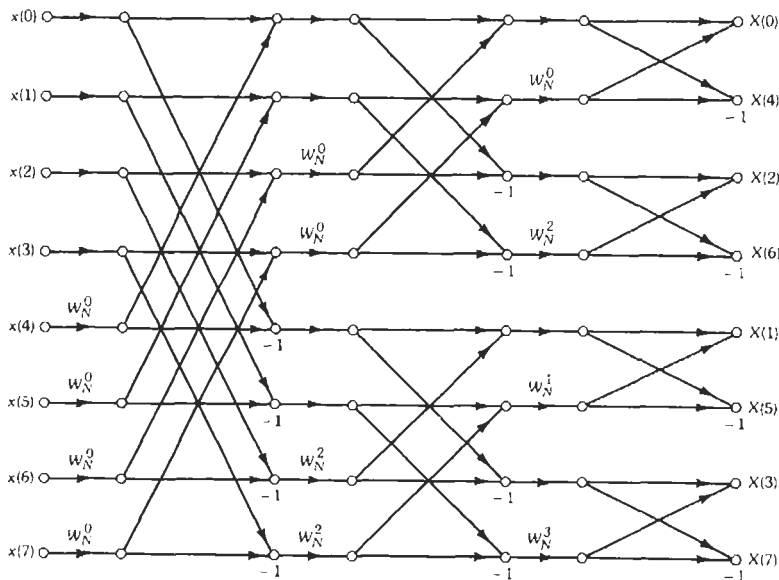


Fig. 4. Decimation-in-time FFT algorithm with normally ordered inputs and bit reversed output.

* Reference 29.

† Reference 29.

real-valued data lead to a reduction at the computational cost by a factor of two.

A radix-2 decimation-in-time FFT program, written in C, is listed in Chart 1 for reference.

The FFT in Spectral Analysis

An FFT program is often used to perform spectral analysis on signals that are sampled and recorded as part of laboratory experiments, or in certain types of data acquisition systems. There are several issues that should be addressed when spectral analysis is performed on (sampled) analog waveforms that are observed over a finite interval of time.

Windowing—The FFT treats the block of data as though it is one period of a periodic sequence. If the underlying waveform is not periodic, then harmonic distortion may occur because the periodic waveform created by the FFT may have sharp discontinuities. This effect is minimized by removing the mean of the data (it can always be reinserted) and by windowing the data so the ends of the block are smoothly tapered to zero and discontinuities do not occur when the FFT treats the windowed block as one period of a periodic sequence. A good rule of thumb is to taper 10 percent of the data on each end of the block using either a cosine taper or one of the other common windows from Table 4. An alternate interpretation of this phenomenon is that the finite length observation has already windowed the true waveform with a rectangular window that has large spectral side lobes. Hence, applying an additional window results in a more desirable window that minimizes frequency-domain distortion.

Zero-Padding—It would appear that more accuracy is produced in the spectral domain if the block

length of the FFT is increased. This can be done by (1) taking more samples within the observation interval; (2) increasing the length of the observation interval, without increasing the sampling rate; or (3) augmenting the original data set with zeros. First, it must be recognized that the finite observation interval causes a fundamental limit on the spectral resolution, even before the signals are sampled. The CT rectangular window has a sinc spectrum, which is convolved with the true spectrum of the analog signal that is being observed. Therefore, frequency resolution is limited by the width of the main lobe in the sinc spectrum, which is inversely proportional to the length of the observation interval. Sampling causes a certain degree of aliasing; although this effect can be minimized by sampling at a high enough rate. Therefore, lengthening the observation interval increases the fundamental resolution limit, while taking more samples within the same observation interval minimizes aliasing distortion and provides a better definition (more sample points) on the underlying spectrum.

Padding the data with zeros and computing a longer FFT does give more frequency domain points, but it does not improve the fundamental resolution limit, nor does it alter the effects of aliasing. The resolution limits are established by the observation interval and the sampling rate. No amount of zero padding can improve these basic limits. However, zero padding is a useful tool for providing more spectral definition; i.e., it enables one to get a better look at the (distorted) spectrum that exists after the observation and sampling effects have occurred.

Leakage and the Picket-Fence Effect—An FFT with block length N can accurately resolve only frequencies $\omega_k = (2\pi/N)k$, $k = 0, \dots, N - 1$ that are harmonics of the fundamental, $\omega_1 = 2\pi/N$. An analog waveform that is

TABLE 4. WINDOW FUNCTIONS

Name	Function	Peak Side-lobe Amplitude (dB)	Main-Lobe Width	Minimum Stop-Band Attenuation (dB)
Rectangular	$\omega(n) = 1, 0 \leq n \leq N - 1$	-13	$4\pi/N$	-21
Bartlett	$\omega(n) = \begin{cases} 2n/N, & 0 \leq n \leq (N-1)/2 \\ 2 - 2n/N, & (N-1)/2 \leq n \leq N-1 \end{cases}$	-25	$8\pi/N$	-25
Hanning	$\omega(n) = (1/2)[1 - \cos(2\pi n/N)], 0 \leq n \leq N - 1$	-31	$8\pi/N$	-44
Hamming	$\omega(n) = 0.54 - 0.46 \cos(2\pi n/N), 0 \leq n \leq N - 1$	-43	$8\pi/N$	-53
Blackman	$\omega(n) = 0.42 - 0.5 \cos(2\pi n/N) + 0.08 \cos(4\pi n/N), 0 \leq n \leq N - 1$	-57	$12\pi/N$	-74

CHART 1. AN IN-PLACE D-I-T FFT PROGRAM IN C

```

/*****
/*  fft: in-place radix-2 DIT DFT of a complex input          */
/*  input:                                                    */
/*      n:      length of FFT: must be a power of two      */
/*      m:      n= 2**m                                     */
/*  input/output:                                           */
/*      x:      float array of length n with real part of data */
/*      y:      float array of length n with imag part of data */
*****/
fft (n,m,x,y)
int    n,m;
float  x[ ],y[ ];
{
    int    i,j,k,n1,n2;
    float  c,s,e,a,t1,t2;

    j = 0;                          /* BIT-REVERSE*/
    n2 = n/2;
    for (i=1; i < n - 1; i++)       /* bit-reverse counter*/
    {
        n1 = n2;
        while ( j >= n1 )
        {
            j = j - n1;
            n1 = n1/2;
        }
        j = j + n1;

        if (i < j)                   /* swap data*/
        {
            t1 = x[i]; x[i] = x[j]; x[j] = t1;
            t1 = y[i]; y[i] = y[j]; y[j] = t1;
        }
    }

    n1 = 0; n2 = 1;                 /* FFT */

    for (i=0; i < m; i++)           /* stage loop*/
    {
        n1 = n2; n2 = n2 + n2;
        e = -6.283185307179586/n2;
        a = 0.0;

        for (j=0; j < n1; j++)     /* flight loop*/
        {
            c = cos(a); s = sin(a);
            a = a + e;

            for (k=j; k < n; k=k+n2) /* butterfly loop*/
            {
                t1 = c*x[k+n1] - s*y[k+n1];
                t2 = s*x[k+n1] + c*y[k+n1];
                x[k+n1] = x[k] - t1;
                y[k+n1] = y[k] - t2;
                x[k] = x[k] + t1;
                y[k] = y[k] + t2;
            }
        }
    }
    return;
}

```

sampled and subjected to spectral analysis may have frequency components between the harmonics. A component at frequency $\omega_{k+1/2} = (2\pi/N)(k + 1/2)$ will appear to be scattered throughout the spectrum. The effect is illustrated in Fig. 5 for a sinusoid that is observed through a rectangular window and then sampled at N points. The "picket-fence effect" means that not all frequencies can be seen by the FFT. Harmonic components are seen accurately. Other components "slip through the picket fence" while their energy is "leaked" in the harmonics.* These effects produce artifacts in the spectral domain that must be carefully monitored to assure that an accurate spectrum is obtained from FFT processing.

DIGITAL-FILTER DESIGN

A common approach to designing IIR digital† filters is first to design an analog prototype, which is then transformed into a digital filter by one of several analog-to-digital mappings. This approach allows the designer to take advantage of many well-known techniques for analog-filter design. Computer-based design techniques are also used for IIRs, as evidenced by the availability of several IIR design programs. Usually, FIRs are designed either by windowing techniques or by computer optimization algorithms. This section summarizes some of the more common design techniques.

* Reference 3.

† It is customary to refer to these as "digital" filters, although at the design stage they really are discrete-time filters.

IIR Filters

In the following discussion, $h_a(t)$ and $H_a(j\Omega)$ denote the impulse response and frequency response of an analog prototype, respectively. Similarly, $h(n)$ and $H(e^{j\omega})$ denote the unit pulse response and the frequency response of the digital filter to be designed. The objective is to obtain $H(e^{j\omega})$ (or $h(n)$) from $H_a(j\Omega)$ (or $h_a(t)$) so that the desirable features of the prototype are preserved.

Impulse Invariance—A digital filter is produced that has a unit pulse response which is exactly equal to the impulse response of the analog filter at the sampling instants; i.e., $h(n) \equiv h_a(t)|_{t=nT}$. The frequency responses are related by

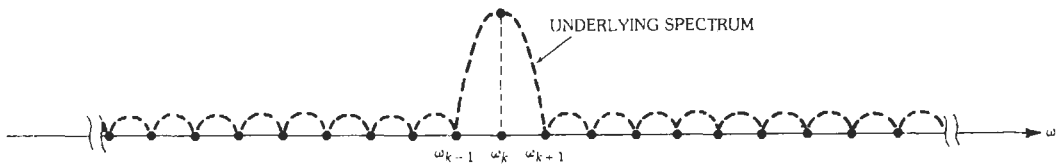
$$H(e^{j\Omega T}) = \sum_{k=-\infty}^{+\infty} H_a(j(\Omega + k2\pi/T)) \quad (\text{Eq. 18})$$

Hence, this technique preserves the frequency response when $H_a(j(\Omega))$ is band-limited, and T is chosen to satisfy the Nyquist criterion.

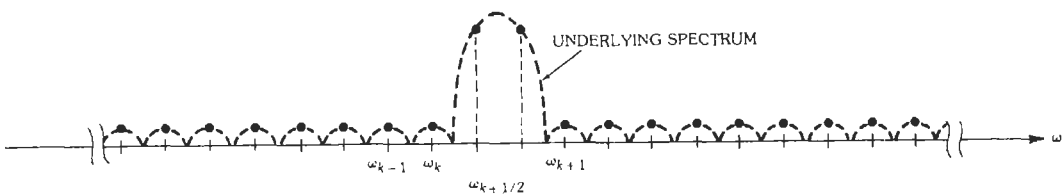
The design procedure is: Design the prototype to obtain $H_a(s)$. Then expand $H_a(s)$ into partial fractions (assuming the multiplicity of each pole is unity), i.e.

$$H_a(s) = \sum_{k=1}^N [A_k / (s - s_k)] \quad (\text{Eq. 19})$$

Each real A_k has a real s_k , which generates a first-order section in the digital filter. For every complex A_k with pole s_k , there is a conjugate term A_k^* with pole s_k^* . These can be combined to generate a second-order section in the digital filter that has only real coefficients. If



(A) FFT of a windowed sinusoid with frequency $\omega_k = 2\pi k/N$.



(B) Leakage for a nonharmonic sinusoidal component.

Fig. 5. Picket-fence effect.

there are M real poles and L complex conjugate pairs, the resulting digital filter becomes

$$H(z) = \sum_{i=1}^M \left[\frac{a_i}{1 + b_i z^{-1}} \right] + \sum_{k=1}^L \left[(c_{k0} + c_{k1} z^{-1})(1 + d_{k1} z^{-1} + d_{k2} z^{-2}) \right] \quad (\text{Eq. 20})$$

where

$$\begin{aligned} a_i &= A_i \text{ (real)} \\ b_i &= -e^{-\alpha_i T} \\ c_{k0} &= -2T \operatorname{Re} \{A_k\} \\ c_{k1} &= -2T e^{\alpha_k T} \operatorname{Re} \{A_k e^{-j\beta_k T}\} \\ d_{k1} &= -2 e^{\alpha_k T} \cos(\beta_k T) \\ d_{k2} &= e^{2\alpha_k T} \end{aligned} \quad (\text{Eq. 21})$$

and $s_k = \alpha_k + j\beta_k$ is the k th pole of $H_a(s)$. Note that $H(z)$ is completely defined in terms of a_k, s_k , and T .

Features:

- (1) Good only for band-limited $H_a(j\Omega)$.
- (2) Results initially in a parallel form.
- (3) $\omega = \Omega T \rightarrow$ linear relationship between the analog and digital frequency variables.
- (4) The impulse invariant technique is not an algebraic mapping; i.e., $H(z)$ cannot be obtained from $H_a(s)$ by a substitution of variable.

Example: The magnitude of the frequency response is shown in Fig. 6 for a fourth-order IIR digital filter designed by the impulse invariant technique with a sampling frequency $f_s = 16$ kHz. The analog prototype is a fourth-order Butterworth low-pass function with a -3 dB cutoff frequency of $f_c = 2000$ Hz.

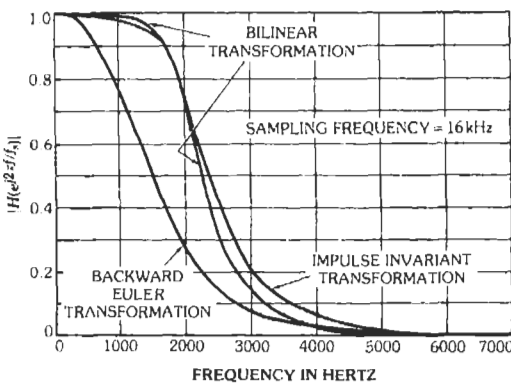


Fig. 6. Frequency response of fourth-order Butterworth analog prototype with $f_c = 2000$ Hz transformed to digital filters by various methods.

Bilinear z-Transformation—A digital filter is produced by substituting $s = (2/T)[(1 - z^{-1})/(1 + z^{-1})]$ into $H_a(s)$ to produce $H(z)$. The relation between the analog frequency variable, Ω , and the digital frequency variable, ω , is given by

$$\omega = 2 \arctan [\Omega T/2] \quad (\text{Eq. 22})$$

which implies that the entire analog frequency axis $-\infty < \Omega < \infty$ is mapped into the interval $-\pi < \omega < \pi$ so there is no aliasing distortion. Fig. 7 shows the important features of how the s -plane is mapped into the z -plane: (A) the $j\Omega$ axis is mapped into the unit circle in the z -plane; (B) the left half s -plane is mapped inside the unit circle; and (C) the right half z -plane is mapped outside the unit circle. Since the bilinear z -transform warps the frequency axis, it is necessary to prewarp the analog prototype so that the critical frequencies occur at the right places. Suppose the final digital filter has N critical frequencies (corner frequencies, band edges, zeros, etc.), denoted $\omega_i, i = 1, \dots, N$. The prewarped analog frequencies, $\Omega_i, i = 1, \dots, N$ are generated by

$$\Omega_i = (2/T) \tan [\omega_i/2], i = 1, \dots, N \quad (\text{Eq. 23})$$

and are then used to design the analog prototype.

Features:

- (1) Eliminates aliasing distortion, but causes non-linear warping of the frequency axis.
- (2) Requires prewarping of the analog prototype.
- (3) Transformation is algebraic, preserves the prototype order, and can be applied to $H_a(s)$ in any form (parallel, cascade, direct, ladder, etc.).
- (4) Particularly useful for wideband filters with piecewise constant magnitude specifications.

Example: Fig. 6 shows the magnitude of the frequency response for a fourth-order Butterworth IIR filter designed with the bilinear z -transform for a sampling frequency $f_s = 16$ kHz. Prewarping has caused the -3 dB cutoff of the digital filter to occur at 0.707. The bilinear z -transform design does not exhibit aliasing distortion, and hence its transition regions fall off faster than for the impulse invariant design.

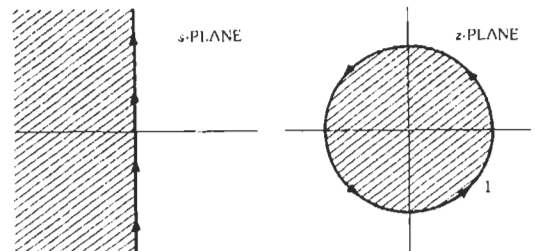


Fig. 7. Mapping of the s -plane into the z -plane by the bilinear z -transformation.

Euler Transformations—Although the Euler transformations are, in general, inferior to the above two design techniques for digital filters, they are mentioned here because they occur frequently for the discrete-time solution of differential equations. The backward and forward Euler transformations are algebraic transformations characterized in the frequency domain by $s = (1/T)(1 - z^{-1})$ and $s = (1/T)(z - 1)$, respectively. The backward Euler transformation results from a backward difference approximation to the first derivative, whereas the forward Euler transformation results from a forward difference approximation.* Both of the Euler transformations are known to distort the frequency response unless the sampling rate is very high with respect to the upper band edge of the analog prototype. The backward Euler transformation always maps stable analog prototypes into stable digital filters. In contrast, the forward Euler transformation has the potential to map a stable prototype into an unstable digital filter unless the sampling rate is chosen sufficiently high.

In Fig. 6, a backward Euler design for the fourth-order Butterworth IIR example is shown for the sampling rate $f_s = 16$ kHz. It is apparent that the backward Euler design results in considerable distortion in comparison with the other two designs.

Nonlinear Optimization Methods—Digital filter design techniques exist which treat the design as a nonlinear optimization problem. In particular, Deczky's method is often used.† Nonlinear design methods are more complicated than transformation of an analog prototype filter via the methods described above. However, direct nonlinear optimization allows much more freedom in the frequency responses that are designed and the accompanying optimization criteria.

FIR Filters

Two general approaches are available for the design of FIR digital filters: (1) designs by windowing techniques and (2) optimal designs by computer optimization techniques. In general, window designs can be carried out with the aid of a hand calculator and a table of well-known window functions. In contrast, computer optimization requires a sophisticated computer program and a considerable amount of computer time. Filters obtained by method 2 are, in general, higher quality filters. Also, more elaborate filter specifications can be approximated because the computer does most of the work.

Window Designs—Let $h(n)$ be the unit pulse response that corresponds to some desired frequency response, $H(e^{j\omega})$. If $H(e^{j\omega})$ has sharp discontinuities, such as the low-pass example shown in Fig. 8, then

$h(n)$ will represent an IIR function. The objective is to time-limit $h(n)$ in such a way as to not disturb $H(e^{j\omega})$ any more than necessary. If $h(n)$ is simply truncated, a ripple‡ occurs around discontinuities, resulting in a distorted filter, as illustrated in Fig. 8.

Suppose that $w(n)$ is a window function that time-limits $h(n)$ to create an FIR approximation, $\hat{h}(n)$; i.e., $\hat{h}(n) = w(n)h(n)$. Then if $W(e^{j\omega})$ is the spectrum of $w(n)$, $\hat{h}(n)$ has a Fourier transform given by $\hat{H}(e^{j\omega}) = W(e^{j\omega}) \otimes H(e^{j\omega})$. From this it can be seen that the ripples in $\hat{H}(e^{j\omega})$ result from the side lobes of $W(e^{j\omega})$. Ideally, $W(e^{j\omega})$ should be similar to an impulse so that $\hat{H}(e^{j\omega}) \cong H(e^{j\omega})$.

Special case (simple FIRs): Let $h(n) = \cos n\omega_0, -\infty < n < +\infty$. Then $\hat{h}(n) = w(n) \cos n\omega_0$, and

$$\hat{H}(e^{j\omega}) = (1/2) W(e^{j(\omega+\omega_0)}) + (1/2) W(e^{j(\omega-\omega_0)})$$

as depicted in Fig. 9. For this simple class, the center of the passband is controlled by ω_0 , and both the shape of the passband and the side-lobe structure are strictly determined by the choice of the window. While this simple class of FIRs does not allow very flexible designs, it is a simple technique for obtaining relatively good low-pass, bandpass, and high-pass FIRs.

General case (design procedure): Specify an ideal frequency response, $H(e^{j\omega})$, and choose samples from the curve. Use a long inverse FFT of length N to find

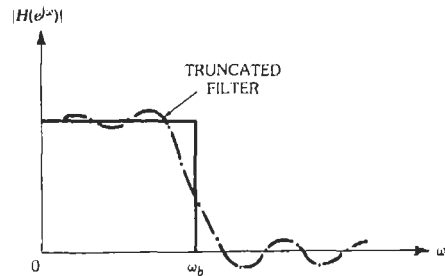


Fig. 8. Gibbs effect in a low-pass filter caused by truncating the unit pulse response.

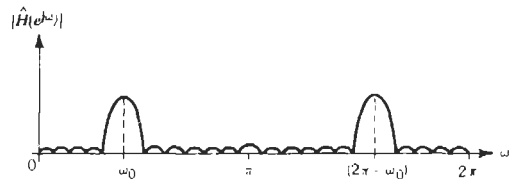


Fig. 9. Design of a simple FIR filter by windowing.

* Reference 25.

† Reference 6

‡ This ripple, known as the Gibbs phenomenon, occurs around points of discontinuity due to the truncation of the Fourier series expansion of $H(e^{j\omega})$.

an approximation to $h(n)$, where if N is the desired length of the FIR, then $\hat{N} \gg N$. Then use a carefully selected window to truncate $h(n)$ according to $\hat{h}(n) = \omega(n)h(n)$. Finally, use an FFT of length \hat{N} to find $\hat{H}(e^{j\omega})$. If $\hat{h}(e^{j\omega})$ is a satisfactory approximation to $H(e^{j\omega})$, the design is finished. If not, choose a new $H(e^{j\omega})$ or a new $\omega(n)$ and repeat. Throughout the procedure, it is important to choose $\hat{N} = kN$, with k an integer in the range [4, . . . , 10]. The k should be made as large as possible within the limits of the computer resources. This design technique is a trial-and-error procedure. The quality of the result will depend to some degree on the skill and experience of the designer.

Table 4 lists a few well-known window functions that can be used in this procedure. An extensive list of windows and their figures of merit has been published in Reference 12.

Computer Design by Equiripple Approximation—There are quite a few publications and numerous computer programs that deal with optimized designs for FIRs using equiripple approximations in the frequency domain. One popular program that has been widely distributed is the one published by McClellan, et. al.* This program is capable of designing bandpass (including low-pass and high-pass), differentiator, and Hilbert-transform FIR filters under the constraint of linear phase. Fig. 10 illustrates an equiripple approximation for an ideal low-pass response. In the program, the order, N , of the equiripple approximation (also the resulting order of the FIR design) is fixed by the designer. Also, the designer specifies the transition band, which is treated as a “don’t care” interval during the computer optimization; i.e., the optimization is unconstrained in this interval. This particular program uses the Remez Exchange Algorithm to find the best equiripple approximation to the specified response, returning to the user the FIR coef-

ficients and the stopband and passband ripple factors. If the result is unsatisfactory, the user can change the order, N , the transition width or the relative weighting factors in the various bands to improve the result. In general, this design program is easy to use, executes efficiently, and produces excellent results. In certain circumstances, the design for multiple-band bandpass filters results in undesirable resonance peaks in the “don’t-care” bands, a problem that has been addressed in the literature. Fig. 11 shows a 64th-order dual bandpass FIR filter that was designed with the program of McClellan, et. al.†

DIGITAL-FILTER IMPLEMENTATION

The fundamental building blocks for implementing digital filters are multipliers, adders, unit delay registers, and semiconductor memories. This section addresses some of the issues concerning the selection of a network topology, arithmetic number codes, and the analysis of quantization effects in digital signal processors. Much of the discussion is directed toward IIR digital filters, because the feedback in IIR structures causes the performance to be rather sensitive to network topology and the effects of quantization error.

Network Structures

The generic building block for IIR filters is the second-order section (Fig. 12). The canonical-form second-order section requires two delay registers, five multipliers, and four adders and is capable of realizing a pair of complex poles and complex zeros with real-valued multiplier coefficients. It is well known that the sensitivity of the zeros of a polynomial with respect to incremental changes in the coefficients grows large rapidly as the polynomial order increases. Therefore, with respect to transfer function sensitivity and round-off error accumulation, it is always preferable to decompose a high-order transfer function into lower-order subsystems. This results in a parallel or cascade realization.

In the parallel realization, $H(z)$ is expanded as a partial fraction expansion of the form

$$H(z) = \sum_{k=0}^{M-1} c_k z^{-k} + \sum_{j=1}^L H_j(z) \quad (\text{Eq. 24})$$

where the first M terms are called the feed-forward terms, and each $H_j(z)$ has a first-order numerator and a second-order denominator of the form

$$H_j(z) = (\gamma_{0j} + \gamma_{1j} z^{-1}) / (1 + b_{1j} z^{-1} + b_{2j} z^{-2})$$

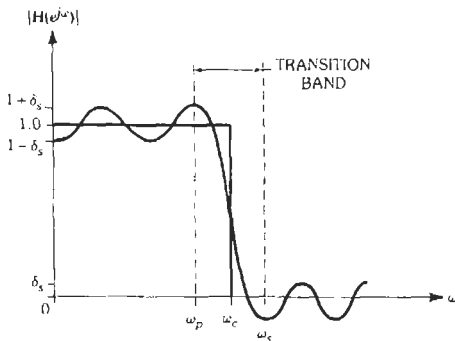


Fig. 10. Equiripple approximation to an ideal low-pass response.

* Reference 22.

Note that this program has been integrated into many commercially available DSP software design tools.

† Reference 22.

‡ Some references call this form a “series” realization, although “cascade” is now more commonly used.

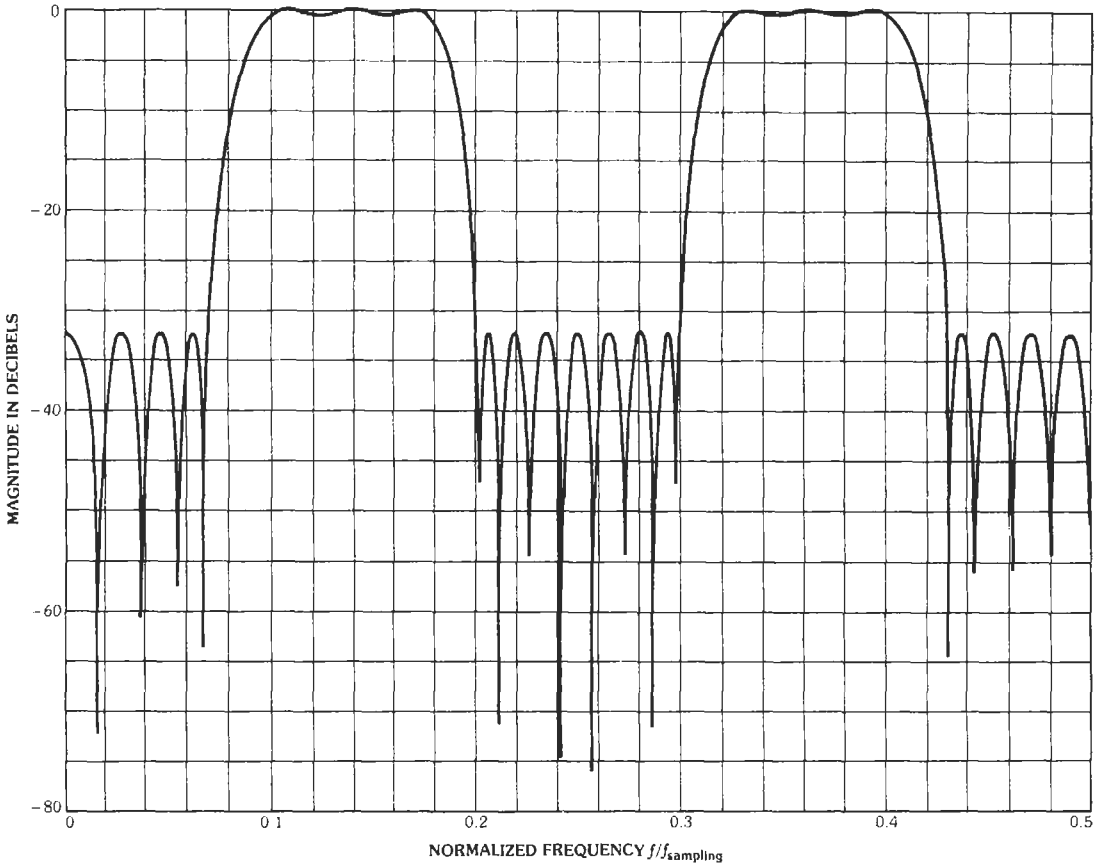


Fig. 11. Magnitude response of a 64th-order dual passband FIR digital filter.

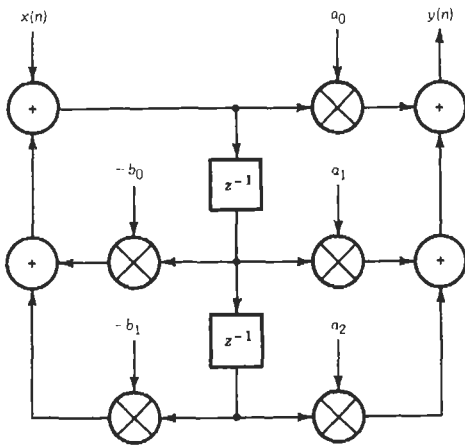


Fig. 12. Second-order section.

The network structure takes the form shown in Fig. 13. The c_k 's result by dividing the denominator of $H(z)$ into the numerator until the remainder is in proper form to expand into partial fractions.

In the cascade realization, $H(z)$ is factored into numerator and denominator factors that are at most second order, i.e.

$$H(z) = H_1(z) \cdot H_2(z) \dots H_L(z) \quad (\text{Eq. 25})$$

where each factor has the form

$$H_j(z) = (a_{0j} + a_{1j}z^{-1} + a_{2j}z^{-2}) / (1 + b_{1j}z^{-1} + b_{2j}z^{-2})$$

The cascade structure is shown in Fig. 14.

Both the parallel and cascade structures are convenient for multiplexing the hardware of one second-order section. Also, in general, the two forms require the same number of multiplications because the feed-forward multiplies make up for the lower-order numerators of the parallel second-order sections. In certain types of filter functions (elliptic), it is found that $a_{0j} = a_{2j}$, so that one multiply can be eliminated. This feature, together with the fact that extra adders are required at the output of the parallel structure, has resulted in more popularity of the cascade structure.

The order in which the factors in Eq. 25 are implemented, as well as the pairing of the numerator and denominator factors, greatly affects quantization error accumulation and represents an additional design

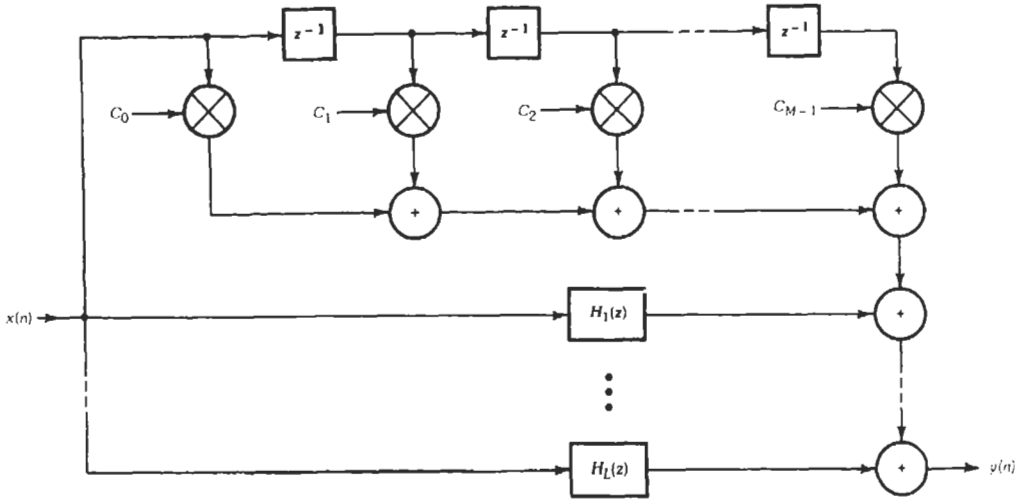


Fig. 13. Parallel realization.

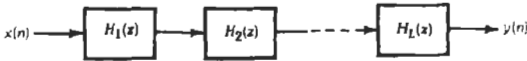


Fig. 14. Cascade realization.

parameter for cascade filters. The problem has been called pole-zero pairing, and at the present time there does not appear to be a known analytical technique to determine the correct pairing for minimum-error performance.

Two other realizations for the second-order section have been proposed for special situations. The coupled-form second-order section, shown in Fig. 15, realizes a complex pair of poles, $P_k = re^{j\omega}$ with the transfer function

$$H(z) = r \sin \theta z^{-1} / [1 - (2r \cos \theta)z^{-1} + r^2 z^{-2}] \quad (\text{Eq. 26})$$

The coupled form is a low-sensitivity second-order section that is preferable for very narrow-band low-pass filters, which are well known to be quite sensitive with respect to coefficient quantization. A second structure that has gained popularity in many applications is a multiplierless ROM-accumulator realization that was popularized by Peled and Liu.* Assume that the second-order section is characterized by the difference equation

$$y(n) = a_0 x(n) + a_1 x(n-1) + a_2 x(n-2) - b_1 y(n-1) - b_2 y(n-2) \quad (\text{Eq. 27})$$

and that each data sample is encoded as a $(b + 1)$ -bit two's-complement binary word, e.g.

$$x(n) = -x^0(n) + \sum_{j=1}^b x^j(n) 2^{-j} \quad (\text{Eq. 28})$$

* Reference 27.

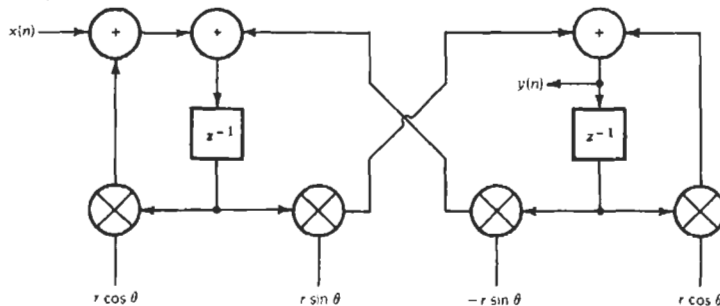


Fig. 15. A coupled-form second-order section.

If Eq. 28 is used in Eq. 27 for both the $x(n)$'s and the $y(n)$'s and the order of summations is interchanged, the output of the section can be expressed as

$$y(n) = \sum_{j=1}^b 2^{-j} \phi(A_j(n)) - \phi(A_0(n)) \quad (\text{Eq. 29})$$

where $A_j(n) = x^j(n), x^j(n-1), x^j(n-2), y^j(n-1), y^j(n-2)$ is a five-bit binary address that is used to address the stored function

$$\phi(A_j(n)) = a_0 x^j(n) + a_1 x^j(n-1) + a_2 x^j(n-2) - b_1 y^j(n-1) - b_2 y^j(n-2)$$

The structure for the resulting filter is shown in Fig. 16. The output, $y(n)$, is computed by a sequence of memory fetches, shifts, and adds. Therefore, multiplication has been entirely eliminated. This architecture is very appealing for high-speed real-time operation, as well as for filters that may be integrated as part of a

VLSI system implementation. For a second-order section, the ROM must have 32 words of memory, and the generation of each $y(n)$ requires $b + 1$ memory fetches and b add-shift cycles. The major disadvantage of this structure is that the filter coefficients are fixed in the ROM according to the function $\phi(\cdot)$. Although several ROMs can be used to select among a predetermined set of filter functions, the ROM-accumulator structure is not well suited for adaptive filters, where the a_i 's and b_i 's are changed in a continual time-varying fashion. This is due to the fact that a change of only a few coefficients would require a complete recalculation of the stored $\phi(\cdot)$ function, thereby nullifying the efficiencies of using a stored function in the first place.

The most common network structure for FIR filters is shown in Fig. 17. A length- N filter requires N multipliers, $N - 1$ delay registers, and $N - 1$ two-input adders. Normally, one hardware multiplier and one adder are time-multiplexed to minimize hardware costs, resulting in a low-cost/low-speed realization.

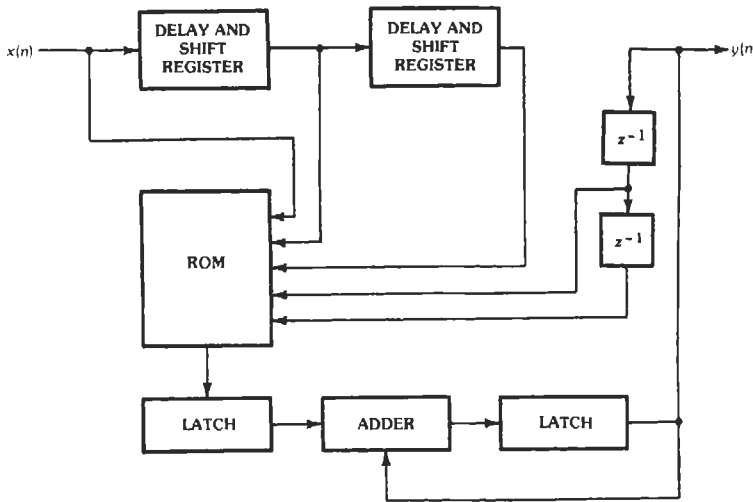


Fig. 16. ROM-accumulator realization of a second-order section.

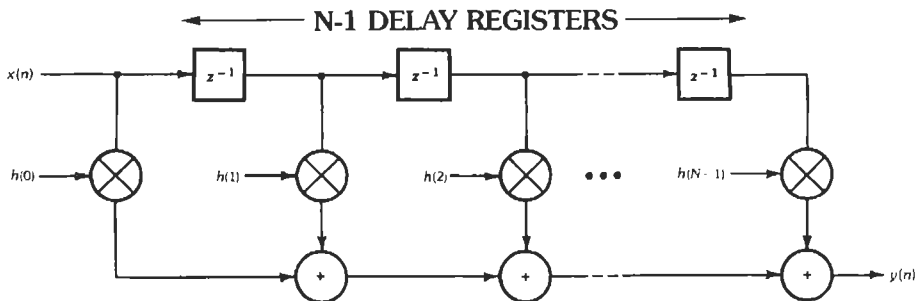


Fig. 17. Common FIR filter structure.

When higher speed is required, more computational elements are used in a pipelined configuration, resulting in a high-speed/high-cost realization. Since there is no feedback, limit cycles cannot occur, and quantization error in the output is simply the additive result of all the quantization error in the multipliers.

Finite Wordlength Effects

Three basic sources of quantization error occur in digital filters: (1) errors due to quantizing an input signal in the A/D converter; (2) errors introduced by rounding during arithmetic operations; and (3) errors in the filter response due to a finite number of bits used in the filter coefficients. The following describes quantization in fixed-point systems.

Input Error—Quantization that occurs in the A/D converter does not corrupt the filter itself, but rather acts as an additive noise component at the input of the ideal system. This noise cannot cause filter instabilities, although it does create inaccuracy in the time-domain response. A common model for A/D quantization noise is shown in Fig. 18, where $e(n)$ is usually assumed to be white noise that is uncorrelated with the signal, $x_a(nT)$, and which is uniformly distributed over the smallest quantization interval. If the probability density function has the form shown in Fig. 19, then the mean and variance of the noise at the filter input are $m_e = 0$ and $\sigma_e^2 = [2^{-2b}/12]$. Then the mean and variance of the noise in the filter output, $\epsilon(n)$, are given by

$$m_\epsilon = 0$$

$$\sigma_\epsilon^2 = [2^{-2b}/12] \sum_{n=-\infty}^{+\infty} |h(n)|^2$$

$$= [2^{-2b}/24\pi] \int_{-\pi}^{\pi} |H(e^{j\omega})|^2 d\omega$$

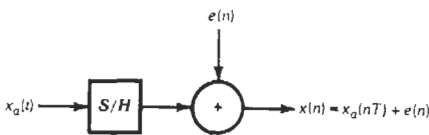


Fig. 18. Noise model for A/D quantization error.

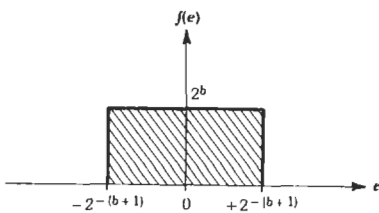


Fig. 19. PDF for rounded A/D quantization.

where $h(n)$ is the unit pulse response and $H(e^{j\omega})$ is the frequency response of the DT system. The statistics of $e(n)$ depend on the specific hardware details of the A/D converter; *i.e.*, different assumptions would be made depending on whether the quantization is implemented by rounding or truncation, and whether the digital code is one's-complement, two's-complement, or sign-magnitude.

Coefficient Error—Filter coefficients are represented by a finite number of bits, commonly in the range of 8 to 16 bits. This results in the implementation of a filter that differs slightly from the original design. The effect is particularly noticeable in IIR filters, mainly because the poles are moved when the coefficients are quantized. If the unquantized design has a pole very close to the unit circle (or on the unit circle, as in an iterative sine generator), slight coefficient changes may move a pole across the unit circle, causing instability. Pole movement also affects frequency response, so the final result must be carefully checked to ensure that coefficient quantization does not cause intolerable magnitude (phase) distortion.

Coefficient errors are fixed at the time of implementation. These errors have no inherently random properties, although one technique of predicting their effect is to treat them as random variables. This approach has been quite successful in predicting the number of bits required to keep the frequency response between specified tolerance limits.

Uncorrelated Roundoff Errors—A product resulting from the multiplication of two $b + 1$ bit two's-complement numbers requires $2b + 1$ bits for its representation. Therefore, a product must be rounded (truncated) to maintain a $b + 1$ bit data stream. Roundoff is an inherently random process when viewed over many computations. It has been rather successfully modeled as a random noise entering the system through a summation node immediately following a multiplier. If the assumption holds that the error from each multiplier is a white-noise source that is uncorrelated with other quantization errors in the structure, roundoff simply causes a jitter in the output, which can be well described by a statistical noise analysis.*

For example, a second-order section with two real poles at $z = 1/2$ and $z = 1/3$ is shown in Fig. 20 as a cascade of two first-order sections. The errors $e_1(n)$, $e_2(n)$, and $e_3(n)$ are assumed to be uncorrelated white-noise sources that represent the quantization errors due to the multipliers $1/2$, $-5/6$, and $1/3$. (Note that multiplication by 2 does not require quantization.) If it is also assumed that the number code is two's-complement with $b + 1$ bits and that the quantization is rounding, then the statistics are $m_{e_i} = 0$ and $\sigma_{e_i}^2 = [2^{-2b}/12]$, for $i = 1, 2$, and 3. If the e_i 's are uncorrelated with the data stream, the output can be expressed as $\hat{y}(n) = \hat{y}(n) + \epsilon(n)$ with

$$m_\epsilon = 0$$

* Reference 25

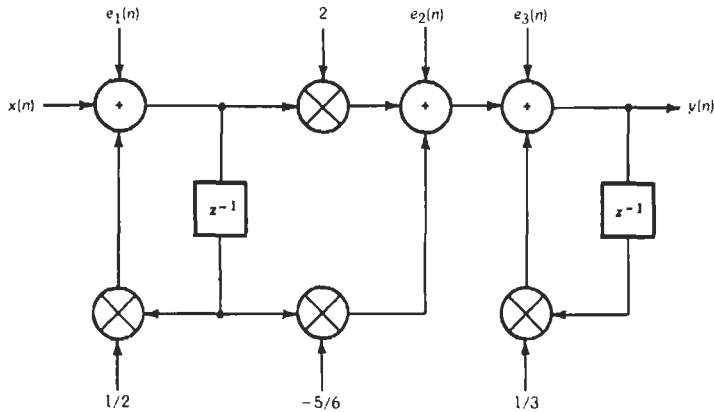


Fig. 20. Second-order section with roundoff error sources $e_1(n)$, $e_2(n)$, and $e_3(n)$.

and

$$\sigma_\epsilon^2 = \sigma_{\epsilon_1}^2 \sum_{n=0}^{\infty} |h(n)|^2 + (\sigma_{\epsilon_2}^2 + \sigma_{\epsilon_3}^2) \sum_{n=0}^{\infty} |h_2(n)|^2$$

where

$$h(n) = [(1/2)^n + (1/3)^n] u(n) \text{ is the overall system unit pulse response,}$$

$$h_2(n) = (1/3)^n u(n) \text{ is the unit pulse response of the second section.}$$

Furthermore, it is found that $\sigma_\epsilon^2 = [7.26/12]2^{-2b}$. If the order of the sections is interchanged, a similar analysis reveals that $\sigma_\epsilon^2 = [9.89/12]2^{-2b}$. This demonstrates that the error accumulation in the output of a cascade filter depends on the ordering of the sections, and that in general there is a preferred ordering (the first in this case).

This general approach to the analysis of arithmetic quantization errors can be extended to higher-order filters realized in any configuration, provided the transfer functions (unit pulse responses) from each noise source to the filter output can be determined. Noise analysis is particularly well suited for computer-aided analysis routines that obtain the required transfer function as a step during a frequency-domain analysis. Figs. 21 and 22 show the roundoff noise spectrum at the output of a fourth-order Butterworth filter realized with 18-bit two's-complement arithmetic in the parallel and cascade form. The noise spectrum for the parallel structure is somewhat smaller than that of the cascade structure, although for this filter the difference is not very great.

In filters implemented with floating-point arithmetic, a roundoff error can also occur during addition. When two floating-point numbers are added, the exponents must be made equal, usually by scaling the smaller exponent until it equals the larger. This requires a right shift of the mantissa, resulting in a loss of bits on the lower end of the word and the necessity

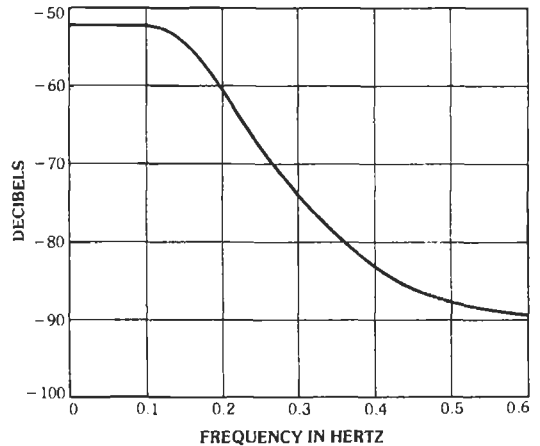


Fig. 21. Roundoff noise spectrum for fourth-order Butterworth low-pass filter in cascade form.

for rounding (truncating). Adder roundoff can also be modeled as a random process, although there is evidence that the process is not zero mean, and elementary statistical assumptions are more difficult to justify than in the case of multiplicative roundoff error.

There arise situations in which roundoff errors become correlated. This often occurs during zero input response, when all internal states are decaying toward the zero state. Roundoff forces the filter states to take on one of a finite number of distinct levels. Forcing the state into these unnatural quanta by the nonlinear quantization scheme often results in limit cycles. These small amplitude limit cycles, which have been called *deadbands*, are a manifestation of highly correlated roundoff errors. In most filters using fixed-point hardware, and in some using floating-point hardware, these deadbands result in idle channel noise, which can be unsatisfactory. Since it is a tedious job to discover all

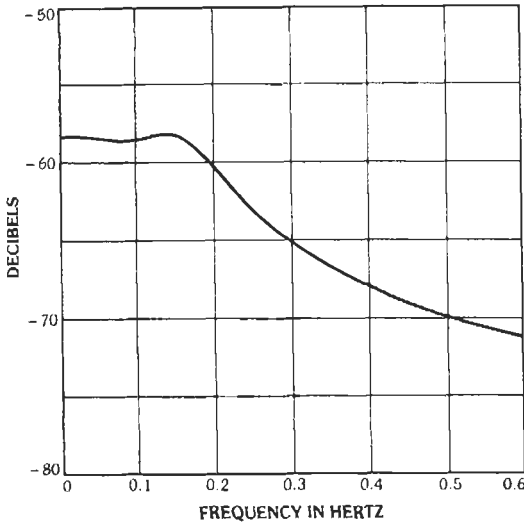


Fig. 22. Roundoff noise spectrum for fourth-order Butterworth low-pass filter in parallel form.

possible limit cycles, a great amount of attention has been focused on the computation of an absolute bound that is tight enough to be useful, and yet which bounds all possible limit cycles.

Overflow Oscillations—If a filter is not scaled properly or if a large unexpected transient should occur at the input, some of the internal states might grow too large for the dynamic range of the filter. In the two's-complement system, a number that grows too large reenters the dynamic range as a large negative number. This phenomenon occurs because the range of a two's-complement system forms a closed ring, as illustrated in Fig. 23. The transfer function of a

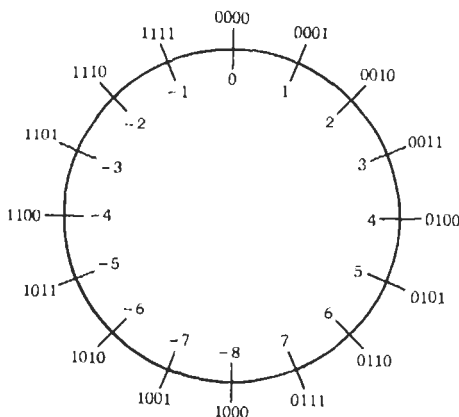


Fig. 23. Dynamic range of two's-complement number system.

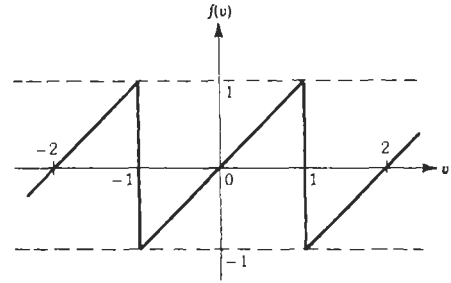


Fig. 24. The transfer function of a two's-complement adder.

two's-complement adder is shown in Fig. 24, where the adder performs $[x + y] = f(x + y)$ and $[\cdot]$ means the actual machine result. The sudden change in sign due to range overflow can result in large self-sustained limit cycles, which are called *overflow oscillations*. These oscillations are characterized as being large and disabling, contrary to the deadband effect, which is usually small and primarily a source of annoyance.

Overflow oscillations cannot be tolerated. One way to prevent them is to scale the input so that overflow at internal summation points cannot occur. A second method is to alter the adders as shown in Fig. 25 so that they saturate (or zero) when the result is outside the dynamic range. Such a scheme, called saturating (or zeroing) arithmetic, will stop overflow oscillations. Generalized forms of saturation arithmetic have been reported that allow $f(v)$ to take on a continuum of different shapes, although there appears to be no practical application of these generalizations to date.

ADAPTIVE DIGITAL SIGNAL PROCESSING

Adaptive digital signal processing is a subject that has attracted increasing attention in recent years due to demands for improved performance in high data rate digital communication systems and in wideband image/video processing systems. Adaptive system identification, adaptive noise cancellation, adaptive linear prediction, and adaptive channel equalization

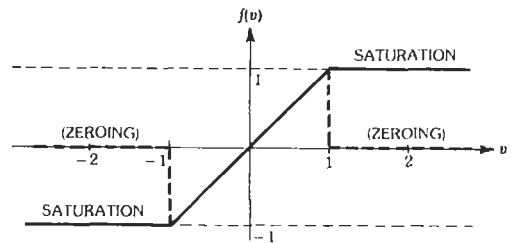


Fig. 25. Modified two's-complement adder.

are just a few of the important applications areas that have been significantly advanced with adaptive signal processing techniques. Much of the recent success in adaptive signal processing has been facilitated by improvements in VLSI digital signal processor (DSP) integrated circuit technology, which currently provides large amounts of digital signal processing power in a convenient and reliable form. Since many improvements in integrated circuit technology are still to come in the future, it is expected that adaptive techniques will assume an even more important role for high performance electronic systems of the future.

This discussion is limited to the simplest and most widely used adaptive systems, i.e., linear adaptive digital filters based on the Least Mean Square (LMS) algorithm.* Although adaptive analog filters are effective in certain applications, adaptive digital filters are far more common because of the efficiency and accuracy of DSP technology. Adaptive filters are self-designing in the sense that they continually adjust their own characteristics by means of a real-time recursive algorithm so that the filter is constantly optimized in some statistical sense while in actual operation.

The LMS Algorithm for FIR Adaptive Filters

An adaptive FIR filter consists of a digital tapped delay line with variable multiplier coefficients that are adjusted by an adaptive algorithm. The adaptive algorithm attempts to minimize a cost function that is designed to provide an instantaneous on-line estimate of how closely the adaptive filter achieves a prescribed optimum condition. The cost function most frequently used is an approximation to the expected value of the squared error, $E\{e^2(n)\}$, where $e(n) = d(n) - y(n)$ is the difference between a training signal $d(n)$ and the filter output $y(n)$, and $E\{\cdot\}$ denotes the expected value. The training signal $d(n)$ is obtained by different means in different applications, some of which are described later in the context of system identification, adaptive noise cancellation, adaptive linear prediction, and adaptive equalization.

The input vector and the coefficient weight vector of the adaptive filter at the n th iteration are defined as $\mathbf{X}(n) = [x(n), x(n-1), \dots, x(n-N+1)]^t$ and $\mathbf{W}(n) = [w_0(n), w_1(n), \dots, w_{N-1}(n)]^t$, respectively, where the superscript t denotes vector transpose. The n th output is then given by

$$y(n) = \mathbf{W}^t(n)\mathbf{X}(n) = \sum_{k=0}^{N-1} w_k(n)x(n-k) \quad (\text{Eq. 30})$$

In the following discussion, the training signal $d(n)$ and the input signal $x(n)$ are assumed to be stationary and ergodic. An adaptive filter uses an iterative method by which the tap weights $\mathbf{W}(n)$ are made to converge

to the optimal solution \mathbf{W}^* that minimizes the cost function. It is well known that \mathbf{W}^* , referred to in the literature as the Wiener solution, is given by

$$\mathbf{W}^* = \mathbf{R}_{xx}^{-1}\mathbf{R}_{xd} \quad (\text{Eq. 31})$$

where $\mathbf{R}_{xx} = E\{\mathbf{X}(n)\mathbf{X}^t(n)\}$ is the autocorrelation matrix of the input and $\mathbf{R}_{xd} = E\{\mathbf{X}(n)d(n)\}$ is the cross correlation vector between the input and the desired response. The most common iterative approach is to update each tap weight according to a steepest descent strategy; i.e., the tap weight vector is incremented in proportion to the gradient $\nabla \mathbf{w}$ according to

$$\mathbf{W}(n+1) = \mathbf{W}(n) - \mu \nabla \mathbf{w} \quad (\text{Eq. 32})$$

where μ is the step size, $\nabla \mathbf{w} = [\nabla_{w_0}(n), \dots, \nabla_{w_{N-1}}(n)]^t$, and $\nabla_{w_k}(n) = \delta E\{e^2(n)\}/\delta w_k$ is the partial derivative of the cost function with respect to $w_k(n)$, for $k = 0, \dots, N-1$. However, the precise value of the cost function is not known, nor is the gradient known explicitly, so it is necessary to make some simplifying assumptions that allow gradient estimates to be computed on-line. Different approaches to estimating the cost function and/or the gradient lead to different adaptive algorithms, such as the Least Mean Squares (LMS), the Recursive Least Squares (RLS), and the quasi-Newton (q-N) algorithms.†

The LMS algorithm‡ makes the simplifying assumption that the expected value of the squared error is approximated by the squared error itself:

$$E\{e^2(n)\} \approx e^2(n)$$

In deriving the algorithm, the error squared is differentiated with respect to \mathbf{W} to give a term that approximates the true gradient. Differentiating the squared error with respect to each filter coefficient w_k yields

$$\nabla_{w_k}(n) = \delta [e^2(n)]/\delta w_k = -2e(n)x(n-k) \quad (\text{Eq. 33})$$

for each $k = 0, \dots, N-1$. Thus in vector notation the complete LMS update relation becomes

$$\mathbf{W}(n+1) = \mathbf{W}(n) + 2\mu e(n)\mathbf{X}(n) \quad (\text{Eq. 34})$$

The value of μ , usually determined experimentally, greatly affects both the convergence rate of the adaptive process and the minimum mean squared error after convergence. To ensure stability and guarantee convergence of both the coefficients and the mean squared error estimate, μ must satisfy the condition

$$0 < \mu < 1/NE\{x^2(n)\}$$

where $E\{x^2(n)\} = (1/N)\text{tr}[R_{xx}]$ is the average input signal power, which can be determined directly from the input signal, or can be calculated from R_{xx} if the input autocorrelation matrix is known. When μ is properly chosen, the weight vector converges to an estimate of

* Reference 35.

† Reference 18.

‡ Reference 35.

the Wiener solution. In practice, because of gradient noise, the weight vector converges to a solution that constantly fluctuates about \mathbf{W}^* . It has been shown that the estimate of \mathbf{W}^* achieved by the LMS algorithm is an unbiased estimate of the true Wiener solution.

The LMS Algorithm for IIR Adaptive Filters

The MSE approximation that led to the conventional LMS algorithm for FIR filters has also been applied to the general class of IIR filter.* Recall that a direct form IIR digital filter is characterized by a difference equation,

$$y(n) = \sum_{k=0}^{N_b} b_k x(n-k) + \sum_{k=1}^{N_a} a_k y(n-k) \quad (\text{Eq. 35})$$

where the b_k s are the coefficients that define the zeros of the filter and the a_k s define the poles. The LMS adaptive algorithm for IIR filters is derived in a similar manner as in the FIR case, although the recursive relation of Eq. 35 is used instead of the convolutional sum to characterize the input-output relationship of the filter. The IIR derivation is more complicated because the recursive terms on the right side of Eq. 35 depend on past values of the filter coefficients.

If the derivatives of the squared error function are calculated using the chain rule, so that all first order dependencies are taken into account, the result is

$$\begin{aligned} \nabla_{\epsilon[e^2]} &= \left[\frac{\partial (e^2(n))}{\partial \mathbf{a}}, \frac{\partial (e^2(n))}{\partial \mathbf{b}} \right]^T \\ &= \left[2e(n) \frac{\partial e(n)}{\partial \mathbf{a}}, 2e(n) \frac{\partial e(n)}{\partial \mathbf{b}} \right]^T \\ &= \left[-2e(n) \frac{\partial y(n)}{\partial \mathbf{a}}, -2e(n) \frac{\partial y(n)}{\partial \mathbf{b}} \right]^T \end{aligned}$$

where

$$\frac{\partial y(n)}{\partial b_k} = x(n-k) + \sum_{j=1}^N a_j(n) \frac{\partial y(n-j)}{\partial b_k} \quad k = 0, \dots, N_b \quad (\text{Eq. 36a})$$

and

$$\frac{\partial y(n)}{\partial a_k} = y(n-k) + \sum_{j=1}^N a_j(n) \frac{\partial y(n-j)}{\partial a_k} \quad k = 0, \dots, N_a \quad (\text{Eq. 36b})$$

This procedure does not generate a closed form expression for the gradient as it did in the FIR case, but it does result in a mechanism by which the gradient can be generated recursively using Eq. 36. Note that Eq. 36a consists of N_b "all-pole" filters, where the input to the k -th filter is simply the input signal shifted

by k time steps; similarly, Eq. 36b represents a similar bank of filters, but where the input to the k -th filter is the shifted output signal $y(n-k)$. Strictly speaking, Eqs. 36a and b require $N_a + N_b$ distinct filters, one for each of the gradient terms. However, in practice it is usually assumed that the filter coefficients are slowly varying in comparison to the signal frequency content, so that within a shift of N_a or N_b time samples, the filters can be treated as though they are time-invariant. With this assumption all of the gradient components represented by Eq. 36a can be generated by a single all-pole gradient filter whose input is $x(n)$.

Similarly, the gradient terms in Eq. 36b can be produced with a second filter that is identical to the first, but which uses $y(n)$ as the input signal. Then time-shifted outputs from these filters will serve as accurate approximations to the true gradient components represented by Eqs. 36a and b. This assumption of "stationarity" in the input and output gradient filters is justified by imposing slow parameter variation by choosing a small step size. In the adaptive filtering literature, the filter represented by Eq. 36a is typically called the input gradient filter, whereas the one represented by Eq. 36b is the output gradient filter. Similar concepts are found in the fields of automatic control and analog circuit theory, where such filters are typically called input and output sensitivity filters

It is known that the use of the output error in the formulation of the cost function prevents bias in the solution due to noise in the desired signal. However, the effect of this recursion is to make the problem nonlinear in terms of the coefficient parameters. The current filter parameters now depend upon previous filter coefficients, which are time-varying. This can lead to MSE surfaces that are not quadratic in nature. There are many examples in the literature for which the MSE surface demonstrates one or more local minima, in addition to the global minimum. In these cases the LMS algorithm, being a gradient search technique, may converge to a local minimum, resulting in poor performance when used in practical applications. Therefore, in general, practicing engineers are often reluctant to use the LMS IIR adaptive filter in practical applications due to its uncertain performance.

Generic Applications of Adaptive Methods

There are four basic configurations in which adaptive filters are typically used to solve practical engineering problems:†

System Identification Configuration—An adaptive filter is said to be used in the system identification configuration when both the adaptive filter and an unknown system are excited by the same input signal $x(n)$, the system outputs are compared to form the

* Reference 35.

† Reference 13.

error signal $e(n) = d(n) - y(n)$, and the parameters of the adaptive filter are iteratively adjusted to minimize some specified function of the error $e(n)$. In the system identification configuration, shown in Fig. 26, the desired signal is produced as the output of an unknown plant whose input is accessible for excitation. When the minimum of the cost function is achieved and the adaptive filter parameters have converged to stable values, the adaptive filter provides a model of the unknown system, in the sense that the adaptive process has formed the best approximation it can in the MSE sense using the structure imposed by the adaptive system. The converged coefficients provide good estimates of the model parameters.

In order for the adaptive system to form a good model of the unknown system at all frequencies, it is important that the input signal have sufficiently rich spectral content. For example, if the adaptive filter is an FIR filter structure with N adjustable coefficients, the input signal must contain at least N distinct frequency components in order to uniquely determine the set of coefficients that minimizes the MSE. A white noise input signal is ideal because it excites all frequencies with equal power. A broadband colored noise input will also provide a good excitation signal, in the sense of driving the adaptive filter to the minimum MSE solution. Although, in general, the convergence rate of the learning process will be slower than for white noise inputs because the frequencies that are excited with small power levels will converge slowly. Many adaptive algorithms attempt to normalize (or whiten) the input power spectrum in order to improve the convergence rate of the learning process.*

The system identification configuration is a fundamental adaptive filtering concept that underlies many applications of adaptive filters. The major attraction of the system identification configuration is that the training signal is automatically generated as the output of

the unknown system. The disadvantage is that the input of the unknown system must be accessible to be excited by an externally applied input noise signal. In some applications obtaining a model of the unknown system is the desired result, and the accuracy of the adaptive coefficients is a primary concern. In other applications it is not necessary that the unknown system be identified explicitly, but rather that the adaptive filter is required to model the unknown system only to generate accurate estimates of its output signal. This leads to the next configuration called adaptive noise cancellation, which is really a variation on the fundamental theme of system identification.

Adaptive Noise Cancelling Configuration—A block diagram for an adaptive noise canceller is shown in Fig. 27, where it is seen that the unknown system in this configuration is not shown explicitly, nor is it desired to identify the implicit unknown system in a direct way. The primary signal is assumed to be the sum of an information bearing signal $s(n)$ and an additive noise component $N_0(n)$, which is uncorrelated with $s(n)$. The primary signal is used to train the adaptive noise canceller, so that $d(n) = s(n) + N_0(n)$ and the error signal becomes $e(n) = d(n) - y(n) = s(n) + N_0(n) - y(n)$. The reference signal, which is used as the input to the adaptive filter, should be a reference noise $N_1(n)$ that is uncorrelated to $s(n)$, but which is correlated in an unknown way with $N_0(n)$. The adaptive filter forms an estimate of $N_0(n)$ and subtracts this estimate from the primary input signal, thereby forming in a good estimate of the information signal at the output of the noise canceller. Note that

$$E[|e(n)|^2] = E[|s(n)|^2] + E[|N_0(n) - y(n)|^2] \quad (\text{Eq. 37})$$

so that minimizing $E[|e(n)|^2]$ will also minimize $E[|N_0(n) - y(n)|^2]$ because the first term in Eq. 37 is dependent only on the information signal $s(n)$ and its mean squared value cannot be affected by the adaptive filter as long as $s(n)$ and $N_1(n)$ are uncorrelated. After

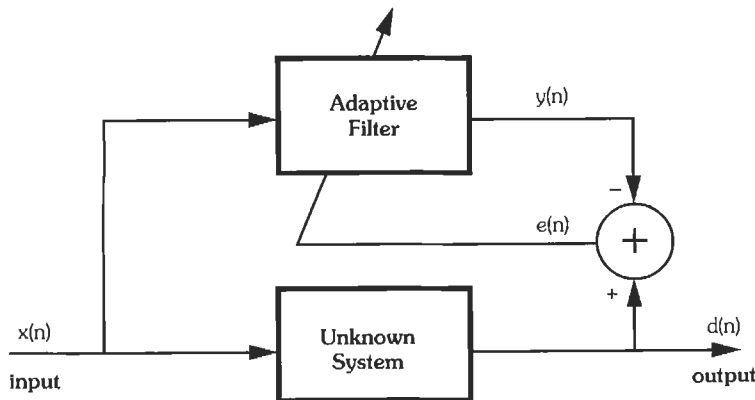


Fig. 26. System identification configuration.

* Reference 13.

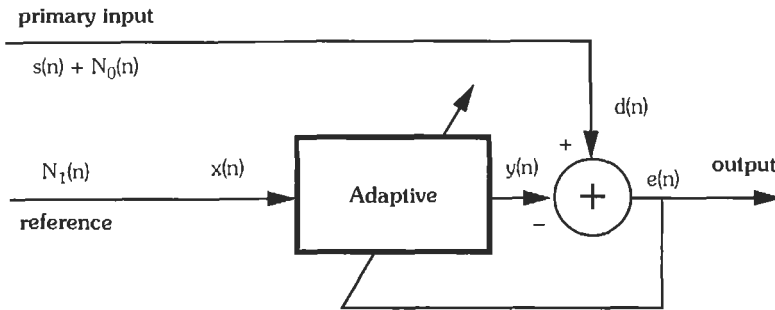


Fig. 27. Noise cancelling configuration.

the adaptive filter converges $y(n)$ becomes the best estimate of $N_0(n)$ according to the MSE criterion.

Since the unknown system in the adaptive noise-cancelling configuration is implicit, there is no need for access to its input in this configuration. However, it is necessary to find a suitable reference signal that does not contain any significant amount of the information signal $s(n)$. If the reference contains even small levels of $s(n)$, then some part of the primary signal $s(n)$ will be canceled and the overall signal-to-noise ratio will degrade.

Adaptive Linear Prediction Configuration—

Adaptive linear prediction is a very important and well developed subject that spans many different areas of engineering. A block diagram of this configuration is shown in Fig. 28. In this configuration the input vector is delayed, usually by one time sample, and the delayed input vector $x(n-1) = [x(n-1), x(n-2), \dots, x(n-N)]^T$ is then used to predict $x(n)$, the current value of the input. The prediction error is given by $e(n) = d(n) - y(n) = x(n) - y(n)$. Sometimes the entire system of Fig. 28 from the input $x(n)$ to output 1 is considered to be a single complete system, in which case it is referred to as a prediction error filter. Whenever the mean squared prediction error is minimized $e(n)$ will become uncorrelated with $x(n)$, while $y(n)$ remains highly correlated with $x(n)$. Therefore, since $d(n) = y(n) + e(n)$, the prediction filter decomposes the

input signal into two components, one that is uncorrelated to the input and one that is highly correlated to the input. In this sense the linear predictor is a type of correlation filter.

Note that two distinct outputs, output 1 and output 2, are labeled in Fig. 28 to give access to both the correlated and uncorrelated components. Output 1 is used in applications such as adaptive linear predictive coding (LPC) for speech analysis and synthesis, and in adaptive differential pulse code modulation (ADPCM) for speech (and image) waveform compression. Since the prediction error is a difference between the actual value of $x(n)$ and its predicted value $y(n)$, the dynamic range needed for accurately encoding $e(n)$ is usually much smaller than $x(n)$ itself. This is the fundamental mechanism by which a linear prediction filter is able to compress waveforms. Alternately, output 2 produces a filtered version of $x(n)$ with the uncorrelated noise component removed. When used in this mode, the adaptive linear predictor becomes a line enhancer, which is capable of removing broadband noise from a narrow band information signal, a function frequently needed in communication systems.

The linear predictor described above is called a forward prediction filter because it uses the N "past" samples contained in $x(n-1)$ to predict the "future" sample $x(n)$. In this case the prediction error is called the forward prediction error, denoted by $e_f(n)$, and the overall filter from the input to output 1 is called a for-

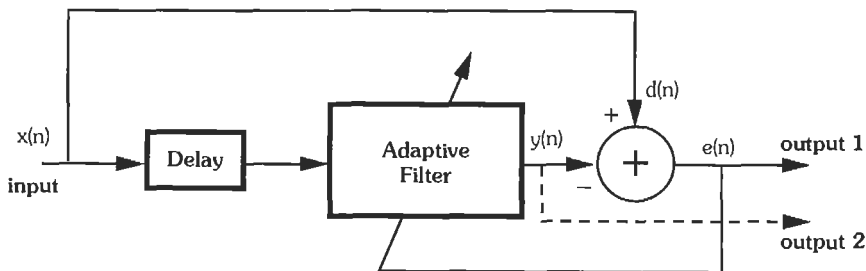


Fig. 28. Linear prediction configuration.

ward prediction error filter. The linear prediction problem can also be formulated as backward linear prediction, in which case the filter is used to estimate the sample $x(n-N)$ from the N "future" samples contained in $x(n) = [x(n), x(n-1), \dots, x(n-N+1)]^T$. In this case, the prediction filter is called the *backward prediction filter*, the prediction error is called the *backward prediction error*, denoted $e_b(n)$, and the overall filter is called the *backward prediction error filter*. A combination of forward and backward prediction is used in the conventional adaptive lattice filters, where the uncorrelated properties of the prediction error leads to excellent learning characteristics.*

Inverse System Configuration—The fourth adaptive filtering configuration is the inverse system configuration shown in Fig. 29. In this configuration the adaptive filter is placed in series with an unknown system and the output $y(n)$ is driven by the adaptive algorithm to form the best MSE approximation to a delayed copy of the input signal. When the adaptive filter reaches convergence, the series combination of the unknown and adaptive systems forms an overall frequency response that approximates a pure delay, i.e., the overall system approximates a flat magnitude response and a linear phase characteristic across the usable bandwidth of the excited spectrum. In this case the adaptive filter estimates $H^{-1}(j\omega)$, where $H(j\omega)$ is the frequency response of the unknown system. The inverse system configuration is the basis for adaptive equalization, in which non-ideal communication channels are equalized in order to reduce dispersion and to eliminate intersymbol interference in high-speed digital communications (see Reference 10). The adaptive equalizer forms an essential component of most state-of-the-art modems today, where the equalization function is required to maintain acceptable bit error rates when binary information is transmitted across narrow-band (4 kHz) telephone channels. Equalizers have also

been used to equalize the dispersive channel that a computer faces when transferring high-speed digital data to a magnetic recording medium (disk or tape). It has been shown that properly designed equalizers will permit symbols to be more densely written on the magnetic recording medium due to the reduction in intersymbol interference. This methodology has attracted the attention of many disk manufacturers due to its ability to effectively increase the capacity of the disk.

The training of an adaptive equalizer in the inverse system configuration raises a number of problems that are unique to this configuration. Note that by the nature of the network configuration, the input to the adaptive filter has already been filtered by the unknown system. Hence in most situations the input to the equalizer cannot be a white noise signal, and depending on the severity of the channel imperfections, the equalizer may experience trouble converging quickly. In a communication system, the transmitter and the receiver are typically located at separate physical locations, so it may not be a simple matter to provide a training signal that is an exact delayed copy of the transmitted waveform. For this reason, channel equalizers are often trained during prescribed "hand shaking" intervals, during which time a pseudorandom binary sequence with known spectral characteristics is transmitted. Once the equalizer has converged to equalize the present characteristics of the unknown channel, the parameters of the equalizer are frozen and held at their converged values during the data transfers that follow.

Due to the difficulty in obtaining a suitable training reference, there has been a great deal of interest in combining certain blind equalization schemes with decision feedback equalizers. In these cases the blind equalization technique is used to bring the equalizer into the neighborhood of proper convergence, at which point the scheme is switched over to a decision feedback algorithm which works very well as long as the

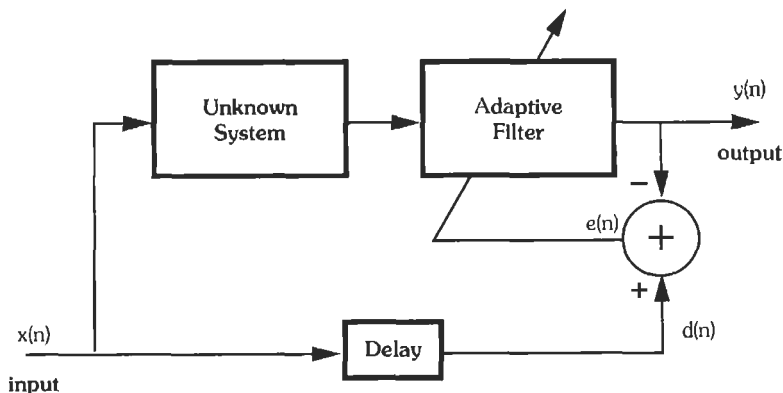


Fig. 29. Inverse system configuration.

* Reference 13.

equalizer remains in the neighborhood of its optimum solution.

Fig. 30 shows an adaptive filter used for channel equalization. The idea behind this application is that the degrading effects of a linear dispersive communications channel can be removed by using the adaptive filter to approximate the inverse of the channel characteristics. In this application, a delayed replica of the transmitted signal must be available as a training

sequence, and the approximate delay Δ in the channel must be known. When the adaptive filter is adjusted properly, the frequency response of the combined channel and adaptive equalizer will have an approximately flat spectrum with a linear phase characteristic. The linear phase models the combined delay in the channel and the equalizer. Fig. 31 shows a computer simulation of an adaptive equalizer designed to operate with the LMS algorithm. Note that the impulse

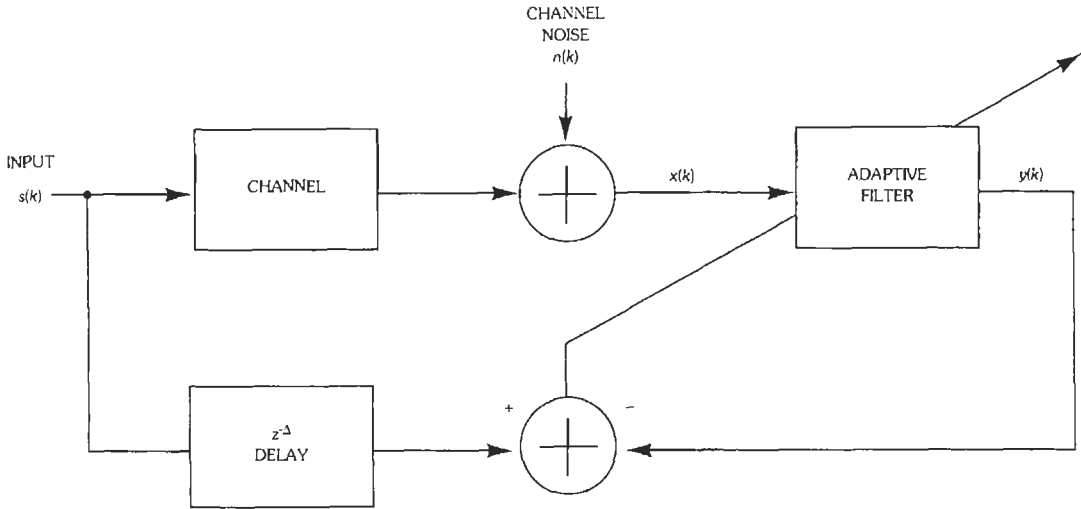


Fig. 30. An adaptive filter used for channel equalization.

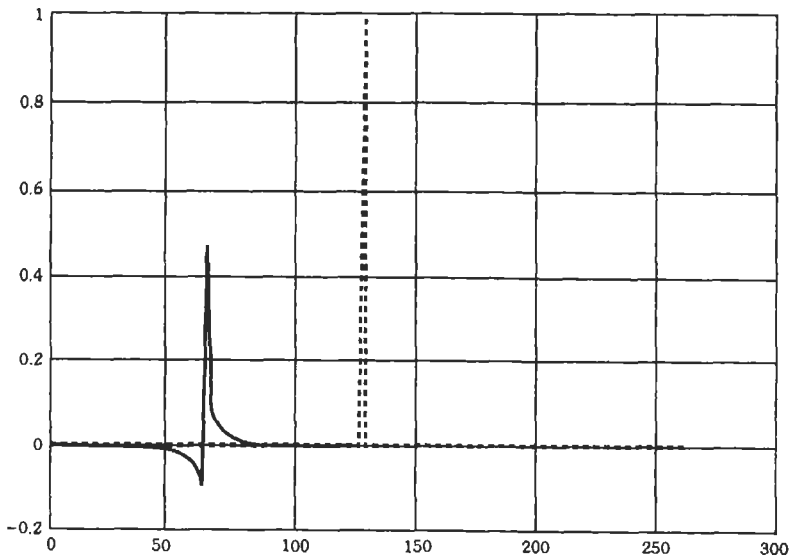


Fig. 31. Computer simulation of adaptive equalization. Solid line: impulse response of the unequalized channel. Dashed line: impulse response of the equalized channel.

response of the channel model (solid line) in Fig. 31 shows a broadening of the pulse and loss of height at the peak due to channel dispersion. The simulated impulse response of the equalized channel (dashed line) shows a much sharper impulse with a higher peak. Adaptive channel equalizers are now found in high-data-rate state-of-the-art communications equipment, such as high speed modems used for computer communications and in certain communication channels that link computer CPUs to magnetic disk storage devices.

DISCRETE-TIME ANALOG TECHNOLOGIES

The most obvious implementation of discrete-time filters is with digital technology. However, filters operating at very large bandwidths can require very high computational rates, and digital implementations become very expensive in terms of cost, size, and power consumption. Several discrete-time analog technologies have been developed which compute discrete-time filters via analog means, thus avoiding explicit analog-to-digital conversion. Although they lack some of the flexibility of a digital implementation, these technologies have found application in a number of specialized cost-sensitive or high-bandwidth applications. Filter coefficients for these technologies are obtained by using the filter design methods described earlier. The basic filtering operation can also accomplish other tasks such as radar pulse detection and compression, communication channel equalization, signal generation, and pattern matching.

Two discrete-time analog technologies are discussed in this section. Switched Capacitor Filters (SCF), which implement IIR filters at relatively low bandwidths (up to about 10 MHz), are often a cost-effective alternative to an all-digital implementation. SCFs inherently sample the analog signal and thus exhibit spectral aliasing. Surface Acoustic Wave (SAW) devices implement FIR filters with nonadjustable coefficients at very high bandwidths, and are widely used

in very high frequency systems such as radars and televisions. SAW devices process the continuous signal without sampling and thus produce no aliasing errors, but the filter is discrete and the frequency response exhibits spectral periodicity.

The following subsections provide an overview of these two discrete-time analog technologies and their operating characteristics; devices are most likely to be purchased from a manufacturer specializing in these technologies, since a number of proprietary techniques are involved in manufacturing high-quality SCF or SAW devices. However, it is important that the system designer know the capabilities and limits of these technologies to evaluate whether these technologies are attractive for a given application. Table 5 summarizes some of the key parameters of the various technologies at the present time.

Switched Capacitor Circuits

Switched capacitor (SC) circuits became popular with circuit designers and IC technologists when it was discovered that SC techniques could be used for the design and manufacture of high precision monolithic filters. SC circuits are designed to depend only on capacitor ratios, which can be controlled quite accurately in MOS technology to realize high precision circuits.

Widespread interest in SC circuits developed during the mid 1970s, although a patent by Baker* in 1967 already had established that a filter characteristic could be specified entirely in terms of capacitor ratios. Baker's patent disclosure described circuits containing only switches, capacitors, and op amps, although his invention did not directly combine the concept with MOS IC technology. It now seems clear that much of the new interest that sprang up in the mid 1970s did not result so much from new advancements in either

* Reference 2.

TABLE 5. KEY PARAMETERS OF DISCRETE-TIME FILTER TECHNOLOGIES

	Technologies		
	Digital	SCF	SAW
Programmable	Yes	yes	no
Highest freq.	A/D limited	10 MHz	200 MHz
Lowest freq.	Dc	dc	5 MHz
Filter type	FIR, IIR	IIR	FIR
Filter length	—	32	4096
Passband ripple	—	0.001	0.3
Stopband rejection	—	80 dB	60 dB
Dynamic range	A/D limited	80 dB	70 dB

circuit theory or IC technology alone, but rather from coupling these two areas in an effective way.

Basic Principles of SC Circuits—The earliest SC filter designs were derived from active filters by means of resistor replacement. The three SC circuits shown in Fig. 32 represent different types of discrete-time integrators that result by replacing the input resistor in an RC active integrator with a series, parallel, or “toggled” switched capacitor, respectively. The concept of resistor replacement is easily understood by considering a linear capacitor C as it is charged by a voltage V

that is held approximately constant during each interval of a clock sequence of period T . The voltage V causes the charge on the plates to change by the amount Δq during the time interval T . Since the capacitor obeys the linear law $\Delta q = CV$, or $\Delta q/T = C/T V$, it is easily shown that the term on the left represents the average current that flows in the capacitor terminals during the interval T . Thus the switched capacitor obeys the approximate relation $i = (C/T)V$. As long as the clock frequency is high enough relative to the signal frequency content so that the constant voltage assumption holds, the switched capacitor behaves

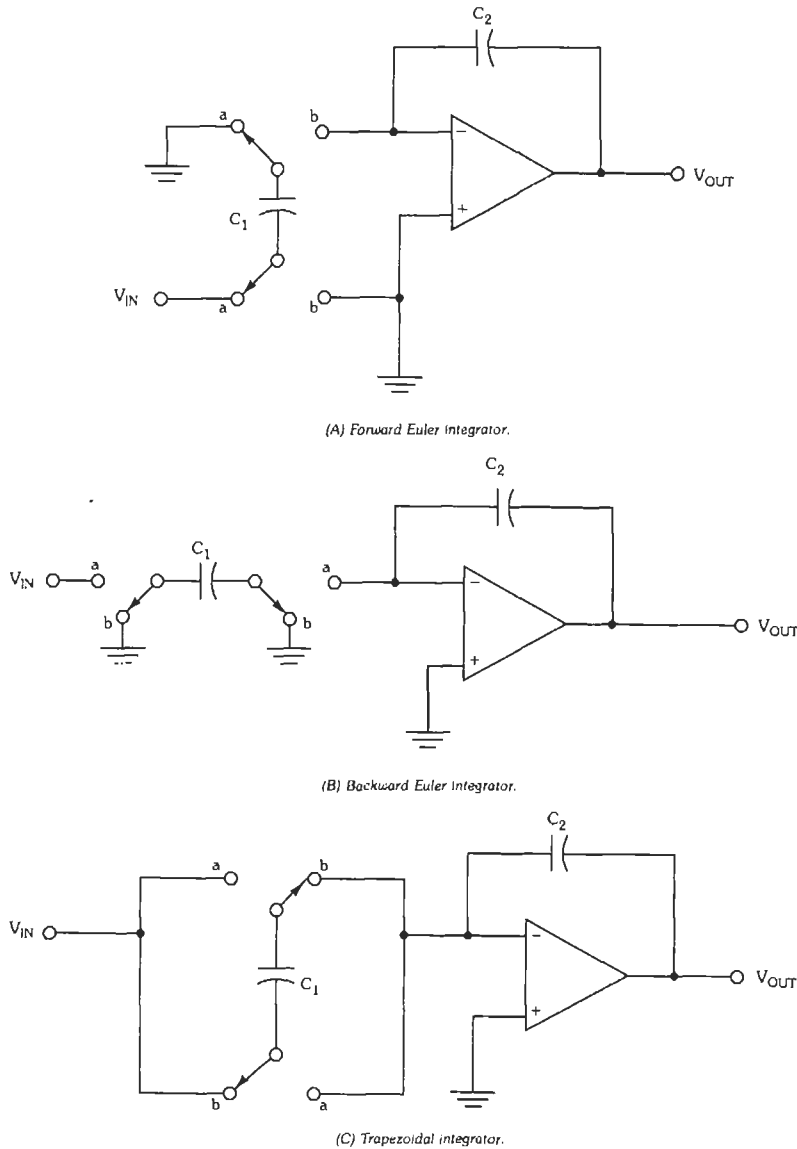


Fig. 32. SC integrators obtained by resistor replacement.

approximately like a resistor of value $R_{eq} = (T/C)$. When the frequency of the clock approaches the Nyquist limit relative to the signal bandwidth, this approximation begins to break down, and it is necessary to consider the details of the switch phasing in order to determine an accurate model for the performance of the SC circuit.

As it became better understood that resistor replacement techniques are inaccurate for switching rates near the Nyquist limit, SC circuit designers began to characterize charge transfer relationships in SC circuits with discrete-time difference equations and to apply the theory of discrete time systems to analyze performance. It was learned that if the output voltage is sampled on the "a" switch phase, then the network shown in Fig. 32a behaves as a forward Euler integrator, the network of Fig. 32b as a backward Euler integrator, and that of Fig. 32c as a trapezoidal integrator. Subsequently, Broderson et al.* showed how high quality SC filters can be derived from the "leapfrog" active filter structure, a particular form of a doubly terminated LC ladder, by using SC integrator replacements.

If a toggled switched capacitor, like the one used to replace the input resistor in the trapezoidal integrator of Fig. 32c, is used to replace every resistor in an arbitrary RC active filter, the analog structure is transformed into an SC structure according to the bilinear- z transformation discussed earlier. If each phase of the biphasic clock has period T , the effective sampling rate becomes $2/T$ because the toggle switched capacitor has identical behavior on each phase.

In MOS transistors, there normally is a rather large nonlinear parasitic capacitance C_s between the bottom plate and the substrate, as well as a considerably smaller parasitic capacitance C_p between the top plate and the substrate. C_s can be as much as 20% of the fabricated capacitor, and is rather difficult to predict or compensate. C_p is normally an order of magnitude smaller, although its presence can disturb a sensitive design significantly. It is important for SC circuits to be designed with the bottom plate of all capacitors grounded or switched between a voltage source and ground so that excess charge accumulated by the parasitic during the charging phase is harmlessly discharged to ground and does not affect the characteristics of the circuit. It is clear that a toggled switched capacitor cannot be made parasitic insensitive in general, since the top and bottom plates must perform the same function on alternate phases. Despite efforts to compensate for parasitics in the "bilinear switched capacitor," it was largely rejected by IC technologists as impractical because it led to too much uncontrolled variation in the characteristics of the final SC circuits.

Note that the forward and backward Euler integrators of Fig. 32 are insensitive with respect to both the top plate and bottom plate parasitics. For this reason

these types of integrators became more successful than the bilinear integrator, and were subsequently used in many design methods that were eventually published. Broderson et al.† described a set of practical considerations for MOS implementation of SC circuits that became standard rules of thumb for SC circuit designers: (1) switched capacitor "resistors" should not close an on amp feedback path; (2) there should be no floating nodes; (3) at least one plate of every capacitor should be connected to a voltage source or switched between voltage sources; and (4) the noninverting op amp input should be kept at a constant voltage (usually grounded) or switched between voltage sources. Although these should be guidelines rather than strict rules, they have withstood the test of time in the sense that most IC technologists tend to reject designs that violate these conditions. The later book by Ghausi and Laker‡ incorporated these rules into the design procedure, representing a change of philosophy from the resistor and integrator replacement techniques that dominated earlier designs.

Switched Capacitor Filters for Voice Quality CODECs—One successful application of switched capacitor circuit technology has been in the realization of analog pre- and post-filters required for digital voice transmission in modern telephone channels. The A/D and D/A conversion, as well as some additional signal conditioning (companding, etc.), is implemented with an IC referred to as a CODEC (coder-decoder). A single IC contains the "transmit" filter to handle transmission in the forward direction, the "receive" filters to handle the returned voice signal, as well as other circuits required for signal conditioning. The industry standard for the digital sampling rate is 8 kHz, which is designed to prevent aliasing on a 4-kHz telephone channel. Prior to SC circuit technology, the transmit and receive filtering was done with an active filter that was implemented with expensive hybrid fabrication techniques. With switched capacitor circuit techniques, the goal of integrating the transmit and receive filters together with the CODEC, while meeting the rather strict system specifications, became reachable. The philosophy is to switch the SC filters at a very high rate (usually 256 kHz) relative to the 8-kHz sampling rate of the CODEC so that sampled data effects in the SC circuits are negligible, and for all intents and purposes the SC pre- and post-filters function as analog filters. The transmit filter must have a low-pass characteristic within a 0 ± 0.01 dB tolerance in the pass-band up to a ripple cutoff frequency of 3.4 kHz. In the receive channel, a post-filter is required following the D/A converter to smooth the staircase waveform that results from the reconstruction process. The receive filter is a low-pass filter with a cutoff at 3.4 kHz, which is essentially the same as the antialiasing transmit filter.

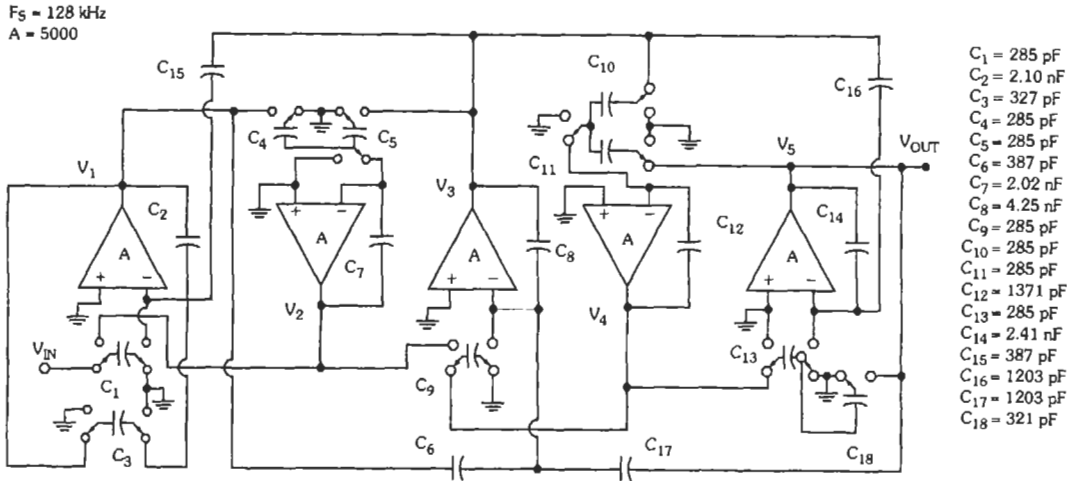
* Reference 5.

† Reference 5.

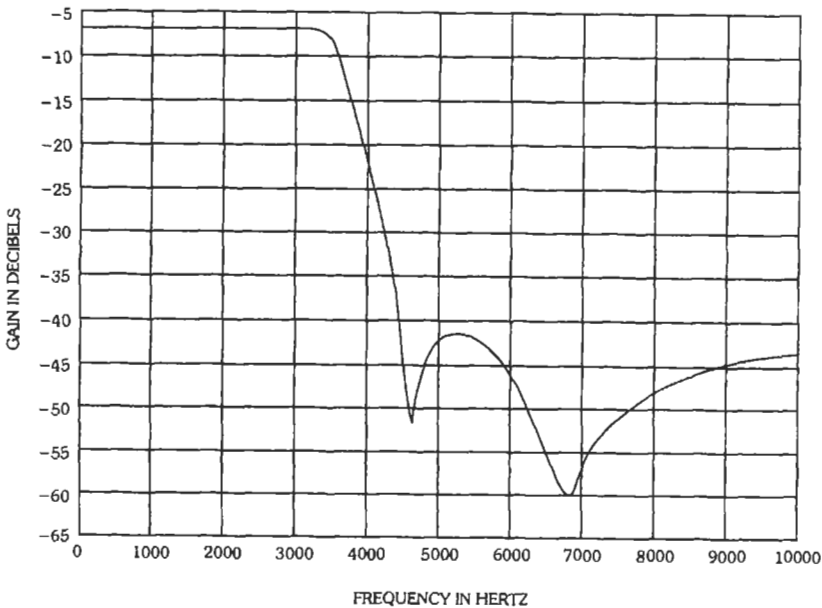
‡ Reference 11.

The SC ladder structure shown in Fig. 33 is capable of realizing the low-pass characteristic required in a CODEC by means of a fifth order elliptic low-pass design. (The example presented here is for the idealized low-pass filter alone; i.e., the design does not include the transmit notch and receive $\sin x/x$ compensation.) This circuit is a variation of a state variable ladder derived by an LD1 (lossless discrete integrator)

transformation, although certain switches have been rearranged to guarantee a parasitic insensitive circuit. It can be seen from Fig. 33 that this design approximates the desired response very accurately. Although many different SC circuits have been designed that meet the CODEC requirements, the state variable ladder is a preferred design because of its low sensitivities to variations in capacitor ratios.



(A) Circuit based on a state variable realization.



(B) Frequency response.

Fig. 33. A fifth order elliptic low-pass state variable ladder realization of a CODEC antialiasing filter. (Same filter is used as the D/A analog post-filter.)

Surface-Acoustic-Wave Filters

A surface-acoustic-wave (SAW) filter consists of an input transducer that generates a surface-acoustic-wave replica of the signal, a polished piezoelectric substrate that acts as the propagation path for the waves, and an output transducer that converts the delayed acoustic waves back to an electrical signal. The input and output transducers each have an impulse response with a finite duration, and the construction is such that the transducer electrodes have a one-to-one correspondence with individual cycles in the impulse response. The impulse response center frequency is set by electrode spacing, and its envelope is set by the electrode length. The frequency-dependent transfer function of each transducer is the Fourier transform of its impulse response. The surface-acoustic-wave device connected between a voltage source and a low-impedance load operates as a filter with an overall transfer function that is the product of the input and output transducer frequency responses and a delay proportional to the separation between the centers of the transducers.

and its envelope is set by the electrode length. The frequency-dependent transfer function of each transducer is the Fourier transform of its impulse response. The surface-acoustic-wave device connected between a voltage source and a low-impedance load operates as a filter with an overall transfer function that is the product of the input and output transducer frequency responses and a delay proportional to the separation between the centers of the transducers.

Consider the SAW device shown in Fig. 34, which uses a single electrode placed between a pair of grounded electrodes as the SAW generation transducer. When a pulse is applied, it creates a uniform

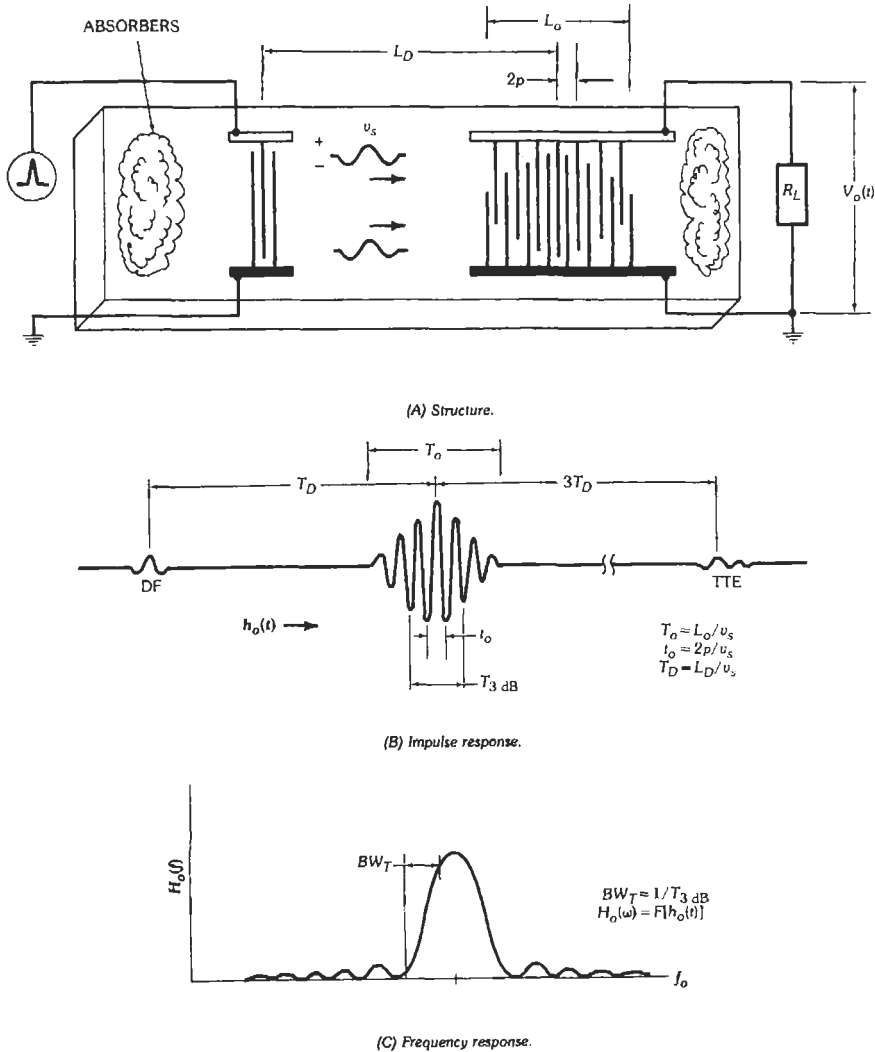


Fig. 34. Surface-acoustic-wave FIR filter.

wavefront acoustic disturbance that propagates along the surface both to the left and to the right. The wave propagating to the left is absorbed, to prevent it from bouncing back and creating an echo. Since the substrate is piezoelectric, the disturbance propagating to the right carries an electric field with it. The output transducer in this example consists of seven electrode pairs with center-to-center spacing p . The electrodes vary in length and are connected interdigitally to two bus bars, one of which is grounded and the other connected to a load. As the SAW pulse passes through the output transducer, it creates an oscillating impulse response, $h_0(t)$, with an overall duration, T_0 , equal to the output transducer length, L_0 , divided by the SAW velocity, v_s (Fig. 34B). The center frequency, f_0 , is equal to the SAW velocity divided by the electrode-to-electrode spacing, p , and the amplitude of each cycle, h_i , is proportional to the corresponding transducer electrode length, l_i . The delay, T_D , is equal to the center-to-center distance of the input and output transducers, L_D , divided by the SAW velocity. The frequency response of the output transducer, $H_0(f)$, is equal to the Fourier transform of the impulse response (Fig. 34C). The device has a transition bandwidth, BW_T that is approximately equal to the reciprocal of the impulse response duration as measured at the 3-dB points, T_{3dB} .

Response—Operation of SAW filters is limited to a specific range of frequencies, fractional bandwidths, transition bandwidths, and delays because the SAW propagation velocity is approximately 3000 meters per second (plus or minus 20 percent) for most useful materials.

Frequency: The highest center frequency achievable with the SAW device is dependent on the ability to print fine lines. If an electrode-to-electrode spacing of 0.75 micron can be achieved, then the devices will operate with a 2-GHz center frequency. At low frequencies, the wavelengths become very long (3 mm at 1 MHz). Since the device must be hundreds or thousands of wavelengths in length and be much thicker than a wavelength, SAW devices become impractical below 5 MHz.

Delay: The substrate size for practical devices is limited to approximately 6 cm, which means that the maximum usable delay is in the order of 20 microseconds. If the application calls for very long delays and cost is not a particularly strong consideration, then delays can be expanded to the 100-microsecond range by using long substrates. The problem of achieving very short delays is the problem of getting the transducers physically close to each other. The minimum possible delay that can be achieved in a filter is twice the reciprocal of the transition bandwidth.

Transition Bandwidth: The impulse response duration is limited by substrate size to periods less than 10 microseconds so that transition bandwidths are limited to values greater than 100 kHz.

Inband Ripple: The SAW devices are designed with the same mathematical techniques used to design digital FIR filters. In the case of SAW devices, the design is realized by adjusting the length and position of the electrodes, and this can be done with an accuracy of approximately one percent. In addition, internally generated echoes add ripples to the frequency response. The measured amplitude-vs-frequency response of SAW bandpass filters generally matches the theoretical design specifications within ± 0.5 dB. In special cases where the application can tolerate higher insertion losses and warrants extra development costs, the response deviations can be reduced to ± 0.2 dB.

Out-of-Band Rejection: The SAW out-of-band rejection depends on the accuracy with which the transducers are implemented and the transducer isolation built into the device. Typical state of the art SAW filters provide approximately 50–70 dB of out-of-band rejection. In certain circumstances, the out-of-band rejection can be raised to 80 dB. It is also possible to cascade SAW filters to improve rejection, if the cost, losses, and inband deviations of two filters can be tolerated.

Echo Levels: The main source of SAW-filter-generated noise is signals (called echoes) that precede and/or follow the main impulse response of the device. The two echoes in a SAW filter are the direct feed signals that are coupled directly to the output without any delay and the triple-transit echo that emerges with a delay three times that of the main signal. These echoes are highly dependent on the method by which the devices are interfaced into the supporting electronic circuitry and are typically 45 dB below the main signal.

Temperature Dependence: Both the center frequency and the delay of SAW filters are temperature dependent because the substrate dimensions and the SAW velocity vary with temperature. These two effects are combined into a single term called the SAW *temperature coefficient of delay* (TCD). The temperature coefficient of delay is approximately 90 ppm/°C (parts per million per degree Celsius) for the LiNbO₃ substrates used in large-percentage-bandwidth devices, and less than 1 ppm/°C for the quartz material used in low-percentage-bandwidth filters.

Dynamic Range: A 50–70-dB dynamic range is typical, and it can be expanded to 100 dB in certain cases. The largest possible signal is limited by burnout, and the smallest input signal that a SAW device can process is limited by the transducer losses and the thermal noise generated in the load.

Noise Figure: The noise figure of a SAW filter is equal to the noise figure of the output amplifier plus the device insertion loss. The effective noise figure is high (typically 15–35 dB), and that is the reason SAW devices are presently being used primarily in the IF sections of receivers.

SAW Filters for IS-95 Compatible CDMA Digital Receivers*

Standing Acoustic Wave (SAW) filters are commonly used as band-pass filters at the intermediate frequency in wideband superheterodyne receivers due to their sharp magnitude characteristics and excellent linear phase response. These SAW filters are one of the largest and most expensive components at the front end of the receiver, making them less desirable for portable, consumer applications. Studies have shown that minimum phase, finite impulse response SAW filters (FIRs) can obtain similar performance with fewer tap weights than their linear phase counterparts. They also lead to the shortest possible propagation, or energy, delays. In SAW filters a minimum phase design has an additional advantage that the magnitude response is less sensitive to errors in the tap weights. Therefore the potential exists for reducing overall size and cost by selecting a minimum phase design. However, the phase response of a minimum phase FIR filter is non-linear and does not meet linear phase constraints required by IS-95 CDMA specifications. This phase non-linearity can be controlled by using an adaptive all-pass phase equalizer that maintains overall phase linearity when the center frequency drifts due to environmental temperature variations.

Minimum Requirements for SAW Performance

Commercially available transversal bandpass SAW filters meet stringent stability requirements and are widely used as the IF filter in superheterodyne digital receivers. They are designed to have a very flat magnitude and linear-phase frequency response over the bandpass region. Two commercially available SAW filters which meet specifications for the IS-95-A spread spectrum standard are the Thomson Microsonics FB E872 [32] and the SAWTEK 855292 [28]. In the following discussion, the specifications of these components, as summarized in Table 6, are to represent those of a typical SAW filter.

TABLE 6. MINIMUM MAGNITUDE RESPONSE SPECIFICATIONS

-5 dB one-sided bandwidth	≥ 630 kHz
-33 dB one-sided bandwidth	≤ 900 kHz
Amplitude ripple within ± 300 kHz	≤ 0.6 dB _p
Phase linearity within ± 630 kHz	$\pm 3^\circ$ rms

* The authors would like to acknowledge contributions of T. Wilborn [36] to the section on "SAW Filters for IS-95 Compatible CDMA Digital Receivers."

Because the center frequency of a SAW filter depends on the physical dimensions of the quartz substrate, it drifts slightly over the operating temperature range of -30 to 80°C . To compensate for this drift, SAW filters are designed to have a two-sided -5 dB bandwidth greater than the minimum requirement of $2 \cdot 630$ kHz given in the table above. Thus, data signals will pass through the filter even though the IF local oscillator frequency does not exactly match the center frequency f_c of the SAW filter.

Fig. 35 illustrates the difference between the impulse responses of two SAW filters with identical magnitude responses, where one has a minimum phase characteristic, and the other, linear phase. For SAW filters, the minimum phase design has the additional advantage that the magnitude response is less sensitive to errors in the tap weights. Therefore, the potential exists for reducing the overall cost of the SAW filter, as well as its size by using a minimum phase SAW filter.

Fig. 36 shows the frequency response of a minimum phase SAW filter which meets the magnitude response specifications given in Table 6. But Fig. 37 illustrates that the rms phase is outside of those specifications, regardless of the center frequency of the filter.

Note that the degree of phase nonlinearity changes with respect to f_c . Therefore, an equalizer that adjusts with changes in f_c is necessary to maintain the maximum SNR over variations in temperature.

An IIR adaptive phase equalizer is capable of providing enough non-linear phase compensation to meet the IS-95-A linear phase constraints. Examples presented in Reference 36 showed learning characteristics of one-section IIR phase equalizers trained with both the decision-directed algorithm (DDA) and the constant modulus algorithm (CMA) adaptive phase equalizers, respectively. Both were able to minimize their corresponding cost functions, achieving local minima that brought the cascade combination of the minimum phase SAW filter and the first order IIR adaptive equalizer into compliance with the magnitude and phase specifications of Table 6. It was seen that the DDA demonstrates considerably faster learning characteristics than the CMA, providing the probability of error is low and correct decisions are made most of the time. The CMA learning characteristics showed a slower learning (convergence) rate, and a large minimum mean square error. Fig. 38 shows phase characteristic of the combined minimum phase SAW filter ($f_c = 0$) and the allpass IIR equalizing filter using six first order all-pass IIR sections. The equalized minimum phase filter has a slightly larger propagation delay than the unequalized minimum phase SAW filter, but it achieves a much smaller propagation delay than the linear phase SAW filter at a greatly reduced cost. This example demonstrates how analog and digital signal processing technologies can be combined to achieve efficient mixed mode solutions to certain practical problems.

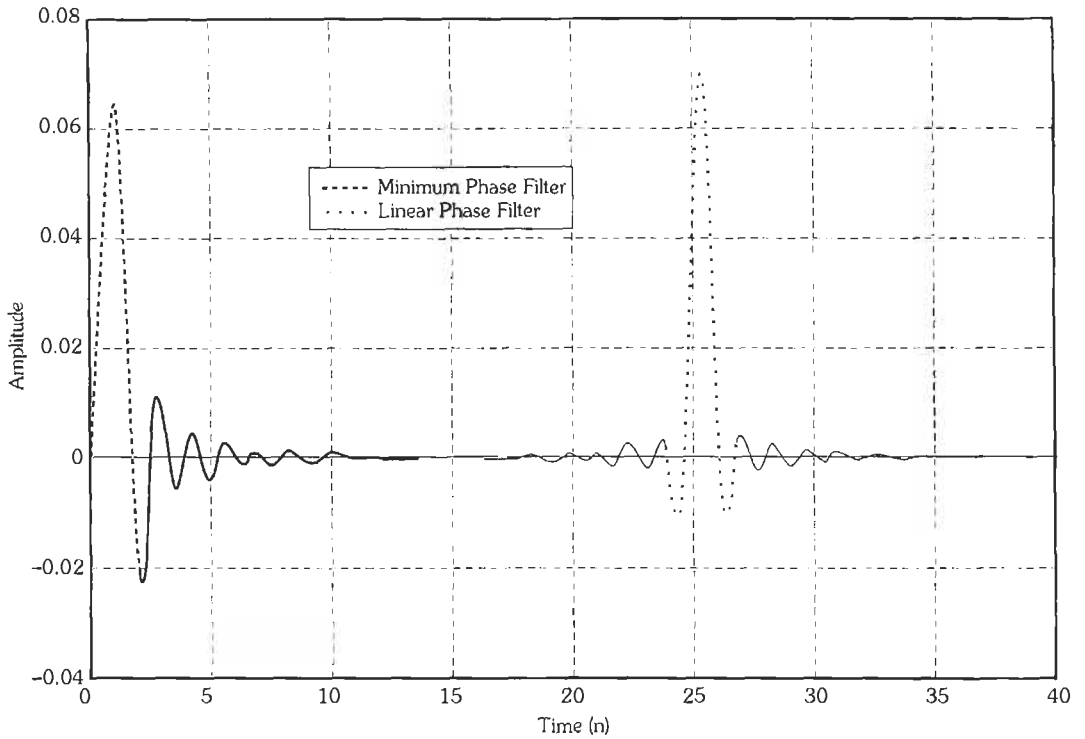


Fig. 35. Impulse responses of minimum phase and linear phase filters with identical magnitude responses.

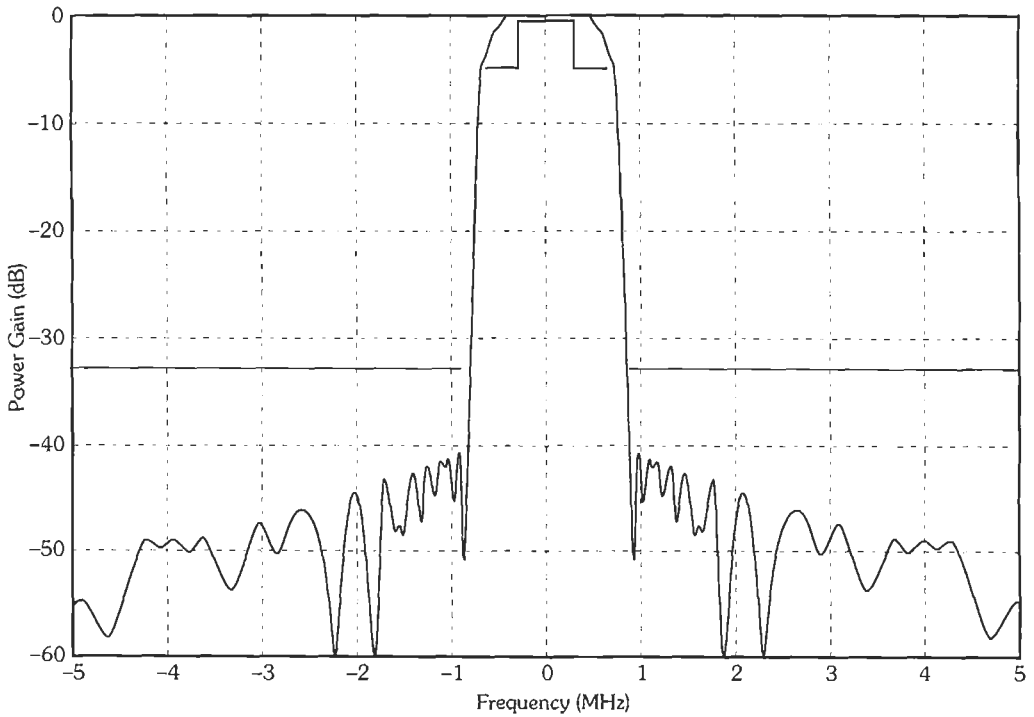


Fig. 36. Magnitude response for both minimum phase and linear phase SAW filters.

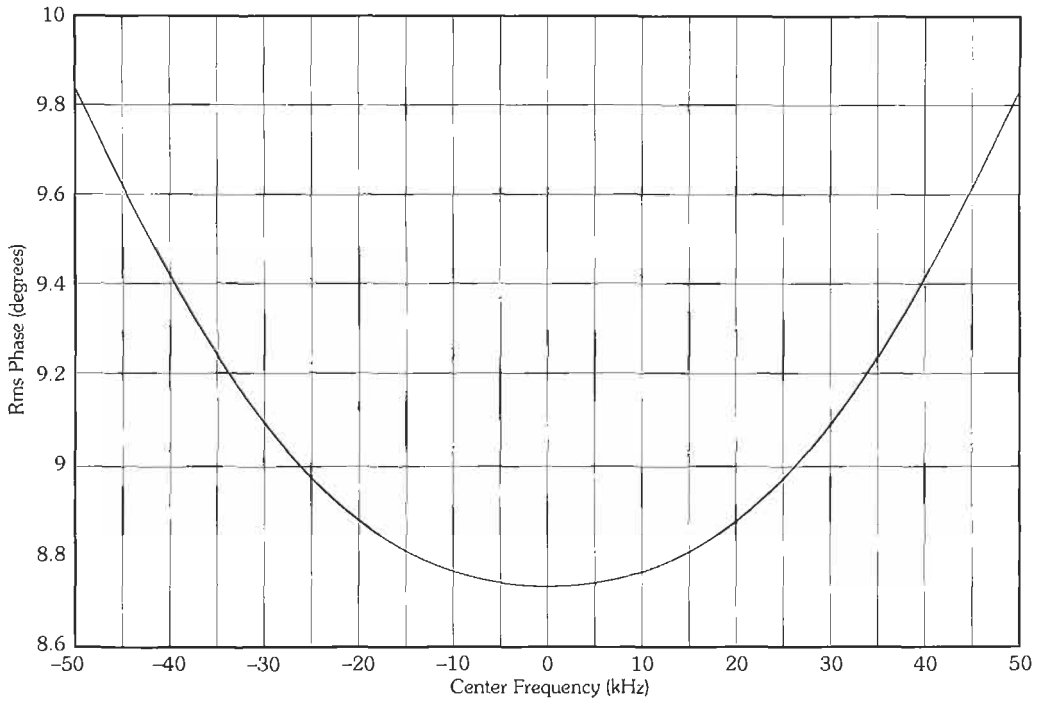


Fig. 37. RMS phase of the minimum phase filter measured over the passband region as a function of center frequency f_c .

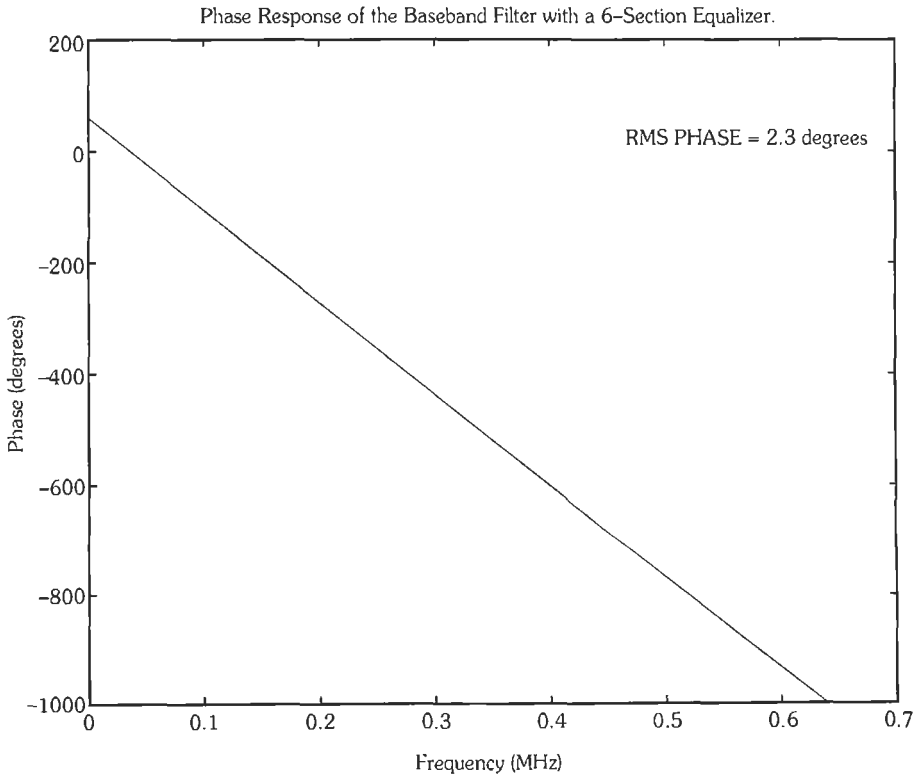


Fig. 38. Phase response of the combined minimum phase SAW filter ($f_c = 0$) and the 6-section all-pass adaptive phase equalizer.

REFERENCES

1. Allen, P. E., and Sanchez-Sinencio, E. *Switched Capacitor Circuits*. New York: Van Nostrand Reinhold Co., 1984.
2. Baker, L. "Dynamic Transfer Networks." US Patent No. 3 469 213, September 23, 1969 (filed May 16, 1967).
3. Bergiad, G. D. "A Guided Tour of the Fast Fourier Transform." *IEEE Spectrum*, Vol. 6, July 1969, pp. 41-52.
4. Brigham, E. O. *The Fast Fourier Transform*. Englewood Cliffs, NJ: Prentice-Hall, Inc., 1974.
5. Broderson, R. W., Gray, P. R., Hodges, D. A., Allstott, D., and Jacobs. G. "MOS Switched Capacitor Filters." *Proc. of the IEEE*, Vol. 67, No. 1, January 1979, pp. 61-75.
6. Deczky, A. G. "Synthesis of Recursive Digital Filters Using the Minimum P Error Criterion." *IEEE Trans. Audio Electroacoustics*, Vol. 20, Oct. 1972, pp. 257-263.
7. De Man, H., Rabacy, J., Arnout, G., and Vandewalle, J. "Practical implementation of a general computer aided design technique for switched capacitor circuits." *IEEE J. Solid-State Circuits*, Vol. SC-15, April 1980, pp. 190-200.
8. Elliott, S., Mierzewski, M., and Planting, P. "The Production of Surface Acoustic Wave Resonators." *1981 IEEE Ultrasonics Symposium*, Vol. 1 81CH1689-9, pp. 89-93.
9. Feldmann and Henaff, J. *Surface Acoustic Waves for Signal Processing*. Paris, France: Masson and CNET, 1986.
10. Gersho, A. "Adaptive equalization of highly dispersive channels for data transmission." *Bell Syst. Tech. J.*, Vol. 48, Jan. 1969, pp. 55-70.
11. Ghausi, M. S., and Laker, K. R. *Modern Filter Design-Active RC and Switched Capacitor*. Englewood Cliffs, NJ: Prentice Hall, Inc., 1981.
12. Harris, F. J. "On the Use of Windows for Harmonic Analysis with the Discrete Fourier Transform." *Proc. of the IEEE*, Vol. 66, No. 1, January 1978, pp. 51-83.
13. Haykin, S. *Adaptive Filter Theory*. Prentice-Hall, Englewood Cliffs, NJ, 1996.
14. Honig, M., and Messerschmitt, D. *Adaptive Filters: Structures, Algorithms, and Applications*. Hingham, MA: Kluwer Academic Publishers, 1984.
15. Hoskins, M. J., Morkoc, H., and Hunsinger, B. J. "Charge Transport by Surface Acoustic Waves in GaAs." *Applied Physics Letters*, Vol. 41, August 1982.
16. Jenkins, W. K. "Fourier Series, Fourier Transforms and the DFT." Chapter 4 in *Mathematics for Circuits and Filters*, ed. by Wai-Kai Chen, CRC Press, 2000, pp 83-111.
17. Jenkins, W. K., and Marshall, D. F. "Transform Domain Adaptive Filters." In *The DSP Handbook*, Vijay Madisetti, ed., CRC Press, pp 22.1-22.20, 1998.
18. Jenkins, W. K., Hull, A. W., Schnauffer, B. A., Strait, J. A., and X, Li. *Advanced Concepts in Adaptive Filtering*. Kluwer Press, 1996.
19. Jenkins, W. K. "Finite Arithmetic Concepts." Chapter 9 in *Handbook on Digital Signal Processing*, J. F. Kaiser and S. K. Mitra, eds., John Wiley and Sons, Inc., 1993.
20. Jenkins, W. K., and Leon, B. J. "The Use of Residue Number Systems in the Design of Finite Impulse Response Digital Filters." *IEEE Trans. on Circuits and Systems*, Vol. CAS-24, No. 4, April 1977, pp. 191-201.
21. Kansy, R. J., Hoskins, M. J., Fleisch, D. A., and Hunsinger, B. J. "Acoustic Charge Transport Signal Processors." *Microwave Journal*, November 1988.
22. McClellan, J. H., Parks, T. W., and Rabiner, R. "A Computer Program for Designing Optimum FIR Linear Phase Digital Filters." *IEEE Trans. Audio and Electroacoust.*, Vol. AU-21, No. 6, December 1973, pp. 506-526.
23. McClellan, J. H., Schafer, R. W., and Yoder, M. A., *DSP First*. Upper Saddle River, NJ: Prentice Hall, 1998.
24. Mitra, S. K. *Digital Signal Processing*. McGraw-Hill, New York, 1998.
25. Oppenheim, A. V., and Schafer, R. W. *Discrete-Time Signal Processing*. Englewood Cliffs, NJ: Prentice-Hall, 1989.
26. Oppenheim, A. V., S. Willsky, and Young, I. T. *Signals and Systems*. Englewood Cliffs, NJ: Prentice-Hall, 1982.
27. Peled, A., and Liu, B. "A New Hardware Realization of Digital Filters." *IEEE Trans. Acoust., Spch., and Sig. Proc.*, Vol. ASSP-22, No. 6, December 1974, pp. 456-462.
28. Sawtek Inc., "85.38 MHz Filter for CDMA Subscriber Applications." Part Number 855292, August 1997.
29. Singleton, R. C. "An Algorithm for Computing the Mixed Radix Fast Fourier Transform." *IEEE Trans. Audio and Electroacoust.*, Vol. AU-17, June 1969, pp. 93-103.
30. Soderstrand, M. A., Jenkins, W. K., Jullien, G. A., and Taylor, F. J. (Eds.), *Modern Applications of Residue Number System Arithmetic to Digital Signal Processing*. New York: IEEE Press, 1986.
31. Stribling, Sidney Nira. "Surface Acoustic Wave Filter for a Television Transmitter." MS Thesis, University of Illinois, Urbana, 1975.

32. Thomson Microsonics, (FB E872), SAW Bandpass Filter, Application: CDMA Mobile, Preliminary Specifications, December 1998.
33. TIA/EIA/IS-95-A, Mobile Station-Base Station Compatibility Standard for Dual-Mode Wideband Spread Spectrum Cellular System. Arlington, VA: Telecommunications Industry Association, 1995.
34. Trick, T. N., and Jenkins, W. K. "Uncorrelated Roundoff Noise in Digital Filters." In *Digital Signal Processing*, J. K. Agarwal, ed. North Hollywood, CA: Point Lobos Press, 1979.
35. Widrow, B., and Stearns, S. *Adaptive Signal Processing*. Englewood Cliffs, NJ: Prentice Hall, Inc., 1985.
36. Wilborn, T. "Adaptive Allpass Phase Equalizer for Digital Receivers: A Case Study," MS Thesis, Department of Electrical and Computer Engineering, University of Illinois, Urbana-Champaign, 1999.

29

Transmission Lines

*Revised by
Tatsuo Itoh*

Rule of Subscripts and Sign Conventions	29-3
Symbols	29-3
Fundamental Quantities and Line Parameters	29-4
Voltage and Current	29-5
Impedance and Admittance	29-6
Lines Open- or Short-Circuited at the Far End	29-7
Voltage Reflection Coefficient and Standing-Wave Ratio	29-8
Power and Efficiency	29-10
Transformation of Impedance on Lines With High SWR	29-12
Mismatch and Transducer Loss	29-13
One End Mismatched	
Generator and Load Mismatched	
Notes on Equation (3)	
Examples	
Attenuation and Resistance of Transmission Lines at Ultrahigh Frequencies	29-15
Resonant Lines	29-15
Quarter-Wave Matching Sections	29-19

Impedance Matching With Shorted Stub	29-19
Impedance Matching With Open Stub	29-19
Length of Transmission Line	29-21
Characteristic Impedance of Lines	29-21
Microstrip Lines	29-25
Quasi-TEM Characteristics	
Attenuation	
Frequency-Dependent Characteristics	
Power-Handling Capacity	
Strip Transmission Lines	29-27
Coplanar Transmission Lines	29-28
Attenuation and Power Rating of Lines and Cables	29-29
Attenuation	
Power Rating	
Army-Navy List of Preferred Radio-Frequency Cables	29-30

The equations and charts of this chapter are for transmission lines operating in the TEM mode.* At the beginning of several of the sections (e.g., "Fundamental Quantities and Line Parameters," "Voltage and Current," "Impedance and Admittance," "Voltage Reflection Coefficient and Standing-Wave Ratio") there are accurate equations, according to conventional transmission-line theory. These are applicable from the lowest power and communication frequencies, including direct current, up to the frequency where a higher mode begins to appear on the line.

Following the accurate equations are others that are specially adapted for use in radio-frequency problems. In cases of small attenuation, the terms $\alpha^2 x^2$ and higher powers in the expansion of $exp \alpha x$, etc., are neglected. Thus, when $\alpha x = (\alpha/\beta)\theta = 0.1$ neper (or about 1 decibel), the error in the approximate equations is of the order of 1 percent.

Much of the information is useful also in connection with special lines that function in a quasi-TEM mode (e.g., microstrip).

It should be observed that Z_0 and Y_0 are complex quantities and the imaginary part cannot be neglected in the accurate equations, unless preliminary examination of the problems indicates the contrary. Even when attenuation is small, $Z_0 = 1/Y_0$ must often be taken at its complex value, especially when the standing-wave ratio is high. In the first few pages of equations, the symbol R_0 is used frequently. However, in later charts and special applications, the conventional symbol Z_0 is used where the context indicates that the quadrature component need not be considered for the moment.

RULE OF SUBSCRIPTS AND SIGN CONVENTIONS

The equations for voltage, impedance, etc., are generally for the quantities at the input terminals of the line in terms of those at the output terminals (Fig. 1). In case it is desired to find the quantities at the output in terms of those at the input, it is simply necessary to interchange the subscripts 1 and 2 in the equations and to place a minus sign before x or θ . The minus sign may then be cleared through the hyperbolic or circular functions; thus

$$\sinh(-\gamma x) = -\sinh \gamma x$$

and so on.

SYMBOLS

Voltage and current symbols usually represent the alternating-current complex sinusoid, with magnitude equal to the root-mean-square value of the quantity.

* The information through the section "Length of Transmission Line" (page 29-21) is valid for single-mode waveguides in general, except for equations where the symbols R , L , G , or C per unit length are involved.

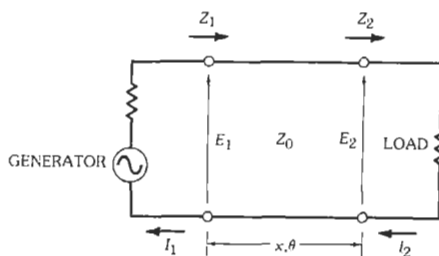


Fig. 1. Transmission line with generator and load.

Certain quantities, namely C , c , f , L , T , v , and ω are shown with an optional set of units in parentheses. Either the standard units or the optional units may be used, provided the same set is used throughout. (For the physical significance of C , G , L , and R , refer to Fig. 2.)

- $A = 10 \log_{10}(1/\eta) =$ dissipation loss in a length of line in decibels
- $A_0 = 8.686 \alpha x =$ normal or matched-line attenuation of a length of line in decibels
- $B_0 =$ susceptive component of Y_0 in mhos
- $C =$ capacitance of line in farads/unit length (microfarads/unit length)
- $c =$ velocity of light in vacuum in units of length/second (units of length/microsecond). See table of physical constants in Chapter 3.

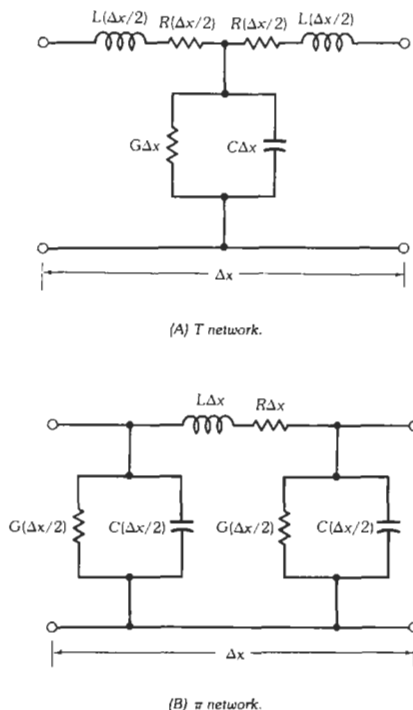


Fig. 2. Equivalent networks of a short section, Δx , of a transmission line.

E = voltage (root-mean-square complex sinusoid) in volts
 ${}_fE$ = voltage of forward wave, traveling toward load
 ${}_rE$ = voltage of reflected wave
 $|E_{\text{nat}}|$ = root-mean-square voltage when standing-wave ratio = 1.0
 $|E_{\text{max}}|$ = root-mean-square voltage at crest of standing wave
 $|E_{\text{min}}|$ = root-mean-square voltage at trough of standing wave
 e = instantaneous voltage
 $F_p = G/\omega C$ = power factor of dielectric
 f = frequency in hertz (megahertz)
 G = conductance of line in mhos/unit length
 G_0 = conductive component of Y_0 in mhos
 $g_a = Y_a/Y_0$ = normalized admittance at voltage standing-wave maximum
 $g_b = Y_b/Y_0$ = normalized admittance at voltage standing-wave minimum
 I = current (root-mean-square complex sinusoid) in amperes
 ${}_fI$ = current of forward wave, traveling toward load
 ${}_rI$ = current of reflected wave
 i = instantaneous current
 L = inductance of line in henrys/unit length (microhenrys/unit length)
 P = power in watts
 R = resistance of line in ohms/unit length
 R_0 = resistive component of Z_0 in ohms
 $r_a = Z_a/Z_0$ = normalized impedance at voltage standing-wave maximum
 $r_b = Z_b/Z_0$ = normalized impedance at voltage standing-wave minimum
 $S = |E_{\text{max}}/E_{\text{min}}|$ = voltage standing-wave ratio
 T = delay of line in seconds/unit length (microseconds/unit length)
 v = phase velocity of propagation in units of length/second (units of length/microsecond)
 X_0 = reactive component of Z_0 in ohms
 x = distance between points 1 and 2 in units of length (also used for normalized reactance = X/Z_0)
 $Y_1 = G_1 + jB_1 = 1/Z_1$ = admittance in mhos looking toward load from point 1
 $Y_0 = G_0 + jB_0 = 1/Z_0$ = characteristic admittance of line in mhos
 $Z_1 = R_1 + jX_1$ = impedance in ohms looking toward load from point 1
 $Z_0 = R_0 + jX_0$ = characteristic impedance of line in ohms
 Z_{oc} = input impedance of a line open-circuited at the far end
 Z_{sc} = input impedance of a line short-circuited at the far end
 α = attenuation constant = nepers/unit length = $0.1151 \times$ decibels/unit length
 β = phase constant in radians/unit length
 $\gamma = \alpha + j\beta$ = propagation constant

ϵ = base of natural logarithms = 2.718; or dielectric constant of medium (relative to air), according to context
 $\eta = P_2/P_1$ = efficiency (fractional)
 $\theta = \beta x$ = electrical length or angle of line in radians
 $\theta^\circ = 57.3\theta$ = electrical angle of line in degrees
 λ = wavelength in units of length
 λ_0 = wavelength in free space
 $\rho = |\rho| \angle 2\psi$ = voltage reflection coefficient
 $\text{pdB} = -20 \log_{10}(1/\rho)$ = voltage reflection coefficient in decibels
 ϕ = time phase angle of complex voltage at voltage standing-wave maximum
 ψ = half the angle of the reflection coefficient = electrical angle to nearest voltage standing-wave maximum on the generator side
 $\omega = 2\pi f$ = angular velocity in radians/second (radians/microsecond).

FUNDAMENTAL QUANTITIES AND LINE PARAMETERS

$$dE/dx = -(R + j\omega L)I$$

$$d^2E/dx^2 = \gamma^2 E$$

$$dI/dx = -(G + j\omega C)E$$

$$d^2I/dx^2 = \gamma^2 I$$

$$\gamma = \alpha + j\beta = [(R + j\omega L)(G + j\omega C)]^{1/2}$$

$$= j\omega(LC)^{1/2} \times [(1 - jR/\omega L)(1 - jG/\omega C)]^{1/2}$$

$$\alpha = \left(\frac{1}{2}\right)[(R^2 + \omega^2 L^2)(G^2 + \omega^2 C^2)]^{1/2} + RG - \omega^2 LC)^{1/2}$$

$$\beta = \left(\frac{1}{2}\right)[(R^2 + \omega^2 L^2)(G^2 + \omega^2 C^2)]^{1/2} - RG + \omega^2 LC)^{1/2}$$

$$Z_0 = 1/Y_0 = [(R + j\omega L)/(G + j\omega C)]^{1/2}$$

$$= (L/C)^{1/2} [(1 - jR/\omega L)/(1 - jG/\omega C)]^{1/2}$$

$$= R_0(1 + jX_0/R_0)$$

$$Y_0 = 1/Z_0 = G_0(1 + jB_0/G_0)$$

$$\alpha = \frac{1}{2}(R/R_0 + G/G_0)$$

$$\beta B_0/G_0 = \frac{1}{2}(R/R_0 - G/G_0)$$

$$R_0 = [M/2(G^2 + \omega^2 C^2)]^{1/2}$$

$$G_0 = [M/2(R^2 + \omega^2 L^2)]^{1/2}$$

$$B_0/G_0 = -X_0/R_0 = (\omega RC - \omega LG)/M$$

where

$$M = [(R^2 + \omega^2 L^2)(G^2 + \omega^2 C^2)]^{1/2} + RG + \omega^2 LC$$

$$1/T = v = f\lambda = \omega/\beta$$

$$\beta = \omega/v = \omega T = 2\pi/\lambda$$

$$\gamma x = \alpha x + j\beta x = (\alpha/\beta)\theta + j\theta$$

$$\theta = \beta x = 2\pi x/\lambda = 2\pi fTx$$

$$\theta^\circ = 57.3\theta = 360x/\lambda = 360fTx$$

(A) Special case—distortionless line: When $R/L = G/C$, the quantities Z_0 and α are independent of frequency.

$$X_0 = 0$$

$$\alpha = R/R_0$$

$$Z_0 = R_0 + j_0 = (L/C)^{1/2}$$

$$\beta = \omega(LC)^{1/2}$$

(B) For small attenuation: $R/\omega L$ and $G/\omega C$ are small.

$$\gamma = j\omega(LC)^{1/2}\{1 - j[(R/2\omega L) + (G/2\omega C)]\}$$

$$= j\beta(1 - j\alpha/\beta)$$

$$\beta = \omega(LC)^{1/2} = \omega L/R_0 = \omega CR_0$$

$$T = 1/v = (LC)^{1/2} = R_0 C$$

$$\alpha/\beta = (R/2\omega L) + (G/2\omega C) = (R/2\omega L) + \frac{1}{2}F_p$$

$$= (Rv/2\omega R_0) + \frac{1}{2}F_p$$

= attenuation in nepers/radian

= (decibels per 100 feet)(wavelength in line in meters)/1663

$$\alpha = \frac{1}{2}R(C/L)^{1/2} + \frac{1}{2}G(L/C)^{1/2}$$

$$= (R/2R_0) + \pi(F_p/\lambda)$$

$$= (R/2R_0) + \frac{1}{2}(F_p\beta)$$

where R and G vary with frequency, while L and C are nearly independent of frequency.

$$Z_0 = 1/Y_0$$

$$= (L/C)^{1/2}\{1 - j[(R/2\omega L) - (G/2\omega C)]\}$$

$$= R_0(1 + jX_0/R_0)$$

$$= 1/[G_0(1 + jB_0/G_0)]$$

$$= (1/G_0)(1 - jB_0/G_0)$$

$$R_0 = 1/G_0 = (L/C)^{1/2}$$

$$B_0/G_0 = -(X_0/R_0) = (R/2\omega L) - \frac{1}{2}F_p = (\alpha/\beta) - F_p$$

$$X_0 = -[R/2\omega(LC)^{1/2}] + (G/2\omega C)(L/C)^{1/2}$$

$$= -(R\lambda/4\pi) + (\frac{1}{2}F_p)R_0$$

(C) With certain exceptions, the following few equations are for ordinary lines (e.g., not spiral delay lines) with the field totally immersed in a uniform dielectric of dielectric constant ϵ (relative to air). The exceptions are all the quantities not including the symbol ϵ , these being good also for special types such as spiral delay lines, microstrip, etc.

$$L = 1.016R_0(\epsilon^{1/2}) \times 10^{-3} \text{ microhenrys/foot}$$

$$= \frac{1}{3}R_0(\epsilon^{1/2}) \times 10^{-4} \text{ microhenrys/centimeter}$$

$$C = 1.016[(\epsilon^{1/2})/R_0] \times 10^{-3} \text{ microfarads/foot}$$

$$= [(\epsilon^{1/2})/3R_0] \times 10^{-4} \text{ microfarads/centimeter}$$

$$v/c = 1016/R_0 C' = \epsilon^{-1/2}$$

= velocity factor (with capacitance C' in picofarads/foot)

$$\lambda = \lambda_0 v/c = c/f(\epsilon^{1/2}) = \lambda_0/(\epsilon^{1/2})$$

$$T = 1/v = R_0 C' \times 10^{-6} = 1.016 \times 10^{-3}/(v/c)$$

= 1.016 $\times 10^{-3} \epsilon^{1/2}$ microseconds/foot (with capacitance c' in picofarads/foot)

The line length is

$$x/\lambda = xf(\epsilon^{1/2})/984 \text{ wavelengths}$$

$$\theta = 2\pi x/\lambda = xf(\epsilon^{1/2})/156.5 \text{ radians}$$

where xf is the product of feet times megahertz.

VOLTAGE AND CURRENT

$$E_1 = jE_1 + rE_1 = jE_2\epsilon^{\gamma x} + rE_2\epsilon^{-\gamma x}$$

$$= E_2\{[(Z_2 + Z_0)/2Z_2]\epsilon^{\gamma x} + [(Z_2 - Z_0)/2Z_2]\epsilon^{-\gamma x}\}$$

$$= \frac{1}{2}(E_2 + I_2 Z_0)\epsilon^{\gamma x} + \frac{1}{2}(E_2 - I_2 Z_0)\epsilon^{-\gamma x}$$

$$= E_2[\cosh \gamma x + (Z_0/Z_2) \sinh \gamma x]$$

$$= E_2 \cosh \gamma x + I_2 Z_0 \sinh \gamma x$$

$$= [E_2/(1 + \rho_2)](\epsilon^{\gamma x} + \rho_2 \epsilon^{-\gamma x})$$

$$I_1 = I_1 + I_1 = I_1 \epsilon^{\gamma x} + I_1 \epsilon^{-\gamma x}$$

$$= Y_0(jE_2 \epsilon^{\gamma x} - I_1 \epsilon^{-\gamma x})$$

$$= I_2 \{[(Z_0 - Z_2)/2Z_0] \epsilon^{\gamma x} + [(Z_0 + Z_2)/2Z_0] \epsilon^{-\gamma x}\}$$

$$= \frac{1}{2}(I_2 + E_2 Y_0) \epsilon^{\gamma x} + \frac{1}{2}(I_2 - E_2 Y_0) \epsilon^{-\gamma x}$$

$$= I_2 [\cosh \gamma x + (Z_2/Z_0) \sinh \gamma x]$$

$$= I_2 \cosh \gamma x + E_2 Y_0 \sinh \gamma x$$

$$= [I_2/(1 - \rho_2)] (\epsilon^{\gamma x} - \rho_2 \epsilon^{-\gamma x})$$

$$E_1 = AE_2 + BI_2$$

$$I_1 = CE_2 + DI_2$$

where the general circuit parameters are $A = \cosh \gamma x$, $B = Z_0 \sinh \gamma x$, $C = Y_0 \sinh \gamma x$, and $D = \cosh \gamma x$.

Refer to section entitled "Matrix Algebra" in Chapter 47.

(A) When point 2 is at a voltage maximum or minimum, x' is measured from voltage maximum and x'' from voltage minimum (similarly for currents).

$$E_1 = E_{\max} (\cosh \gamma x' + S^{-1} \sinh \gamma x')$$

$$= E_{\min} (\cosh \gamma x'' + S \sinh \gamma x'')$$

$$I_1 = I_{\max} (\cosh \gamma x' + S^{-1} \sinh \gamma x')$$

$$= I_{\min} (\cosh \gamma x'' + S \sinh \gamma x'')$$

When attenuation is neglected

$$E_1 = E_{\max} (\cos \theta' + j S^{-1} \sin \theta')$$

$$= E_{\min} (\cos \theta'' + j S \sin \theta'')$$

(B) Letting Z_l = impedance of load, l = distance from load to point 2, and x_l = distance from load to point 1

$$E_1 = E_2 \frac{\cosh \gamma x_l + (Z_0/Z_l) \sinh \gamma x_l}{\cosh \gamma l + (Z_0/Z_l) \sinh \gamma l}$$

$$I_1 = I_2 \frac{\cosh \gamma x_l + (Z_l/Z_0) \sinh \gamma x_l}{\cosh \gamma l + (Z_l/Z_0) \sinh \gamma l}$$

(C)
$$e_1 = \sqrt{\frac{2}{Z_0}} |I_2| \epsilon^{\alpha x} \sin \omega t + 2\pi(x/\lambda) - \psi_2 + \phi + \sqrt{\frac{2}{Z_0}} |E_2| \epsilon^{-\alpha x} \sin[\omega t - 2\pi(x/\lambda) + \psi_2 + \phi]$$

$$i_1 = \sqrt{2} |I_2| \epsilon^{\alpha x} \times \sin[\omega t + 2\pi(x/\lambda) - \psi_2 + \phi + \tan^{-1}(B_0/G_0)] + \sqrt{2} |I_2| \epsilon^{-\alpha x} \times \sin[\omega t - 2\pi(x/\lambda) + \psi_2 + \phi + \tan^{-1}(B_0/G_0)]$$

(D) For small attenuation

$$E_1 = E_2 \{ [1 + (Z_0/Z_2)\alpha x] \cos \theta + j[(Z_0/Z_2) + \alpha x] \sin \theta \}$$

$$I_1 = I_2 \{ [1 + (Z_2/Z_0)\alpha x] \cos \theta + j[(Z_2/Z_0) + \alpha x] \sin \theta \}$$

(E) When attenuation is neglected

$$E_1 = E_2 \cos \theta + jI_2 Z_0 \sin \theta$$

$$= E_2 [\cos \theta + j(Y_2/Y_0) \sin \theta]$$

$$= I_2 \epsilon^{j\theta} + I_2 \epsilon^{-j\theta}$$

$$I_1 = I_2 \cos \theta + jE_2 Y_0 \sin \theta$$

$$= I_2 [\cos \theta + j(Z_2/Z_0) \sin \theta]$$

$$= Y_0 (I_2 \epsilon^{j\theta} - I_2 \epsilon^{-j\theta})$$

General circuit parameters are

$$A = \cos \theta$$

$$B = jZ_0 \sin \theta$$

$$C = jY_0 \sin \theta$$

$$D = \cos \theta$$

IMPEDANCE AND ADMITTANCE

$$\frac{Z_1}{Z_0} = \frac{Z_2 \cosh \gamma x + Z_0 \sinh \gamma x}{Z_0 \cosh \gamma x + Z_2 \sinh \gamma x}$$

$$\frac{Y_1}{Y_0} = \frac{Y_2 \cosh \gamma x + Y_0 \sinh \gamma x}{Y_0 \cosh \gamma x + Y_2 \sinh \gamma x}$$

(A) The input impedance of a line at a position of maximum or minimum voltage has the same phase angle as the characteristic impedance.

$$Z_1/Z_0 = Z_b/Z_0 = Y_0/Y_b = r_b + j0 = S^{-1}$$

at a voltage minimum (current maximum).

$$Y_1/Y_0 = Y_a/Y_0 = Z_0/Z_a = g_a + j0 = S^{-1}$$

at a voltage maximum (current minimum).

(B) When attenuation is small

$$\frac{Z_1}{Z_0} = \frac{[(Z_2/Z_0) + \alpha x] + j[1 + (Z_2/Z_0)\alpha x] \tan \theta}{[1 + (Z_2/Z_0)\alpha x] + j[(Z_2/Z_0) + \alpha x] \tan \theta}$$

For admittances, replace $Z_0, Z_1,$ and Z_2 by $Y_0, Y_1,$ and $Y_2,$ respectively.

(C) When attenuation is neglected

$$\frac{Z_1}{Z_0} = \frac{Z_2/Z_0 + j \tan \theta}{1 + j(Z_2/Z_0) \tan \theta} = \frac{1 - j(Z_2/Z_0) \cot \theta}{Z_2/Z_0 - j \cot \theta}$$

and similarly for admittances.

(D) When attenuation $\alpha x = \theta \alpha / \beta$ is small and the standing-wave ratio is large (say > 10) (Note: The complex value of Z_0 or Y_0 must be used in computing the resistive component of Z_1 or Y_1 .): For θ measured from a voltage minimum

$$\begin{aligned} Z_1/Z_0 &= [r_b + (\alpha/\beta)\theta](1 + \tan^2 \theta) + j \tan \theta \\ &= [r_b + (\alpha/\beta)\theta](\cos^2 \theta)^{-1} + j \tan \theta \end{aligned}$$

(See Note 1)

$$\begin{aligned} Z_0/Z_1 &= Y_1/Y_0 \\ &= [r_b + (\alpha/\beta)\theta](1 + \cot^2 \theta) - j \cot \theta \\ &= [r_b + (\alpha/\beta)\theta](\sin^2 \theta)^{-1} - j \cot \theta \end{aligned}$$

(See Note 2)

For θ measured from a voltage maximum

$$\begin{aligned} Z_0/Z_1 &= Y_1/Y_0 = [g_a + (\alpha/\beta)\theta](1 + \tan^2 \theta) + j \tan \theta \\ & \text{(See Note 1)} \\ Z_1/Z_0 &= [g_a + (\alpha/\beta)\theta](1 + \cot^2 \theta) - j \cot \theta \\ & \text{(See Note 2)} \end{aligned}$$

Note 1: Not valid when $\theta \approx \pi/2, 3\pi/2,$ etc., due to approximation in denominator $1 + (r_b + \theta \alpha / \beta)^2 \cot^2 \theta = 1$ (or with g_a in place of r_b).

Note 2: Not valid when $\theta \approx 0, \pi, 2\pi,$ etc., due to approximation in denominator $1 + (r_b + \theta \alpha / \beta)^2 \cot^2 \theta = 1$ (or with g_a in place of r_b). For open- or short-circuited line, valid at $\theta = 0$.

(E) When x is an integral multiple of $\lambda/2$ or $\lambda/4$: For $x = n\lambda/2,$ or $\theta = n\pi$

$$\frac{Z_1}{Z_0} = \frac{(Z_2/Z_0) + \tanh n\pi(\alpha/\beta)}{1 + (Z_2/Z_0) \tanh n\pi(\alpha/\beta)}$$

For $x = n\lambda/2 + \lambda/4,$ or $\theta = (n + \frac{1}{2})\pi$

$$\frac{Z_1}{Z_0} = \frac{1 + (Z_2/Z_0) \tanh(n + \frac{1}{2})\pi(\alpha/\beta)}{(Z_2/Z_0) + \tanh(n + \frac{1}{2})\pi(\alpha/\beta)}$$

(F) For small attenuation, with any standing-wave ratio: For $x = n\lambda/2,$ or $\theta = n\pi,$ where n is an integer

$$\begin{aligned} \frac{Z_1}{Z_0} &= \frac{(Z_2/Z_0) + n\pi(\alpha/\beta)}{1 + (Z_2/Z_0)n\pi(\alpha/\beta)} \\ g_{a1} &= \frac{g_{a2} + \alpha n\lambda/2}{1 + g_{a2}\alpha n\lambda/2} = S_1^{-1} \end{aligned}$$

For $x = (n + \frac{1}{2})\lambda/2,$ or $\theta = (n + \frac{1}{2})\pi,$ where n is an integer or zero

$$\begin{aligned} \frac{Z_1}{Z_0} &= \frac{1 + (Z_2/Z_0)(n + \frac{1}{2})\alpha(\lambda/2)}{(Z_2/Z_0) + (n + \frac{1}{2})\alpha(\lambda/2)} \\ g_{b1} &= \frac{1 + g_{a2}(n + \frac{1}{2})(\alpha/\beta)\pi}{g_{a2} + (n + \frac{1}{2})(\alpha/\beta)\pi} = S_1 \end{aligned}$$

Subscript a refers to the voltage-maximum point and b to the voltage minimum. In the above equations, subscripts a and b may be interchanged, and/or r may be substituted in place of $g,$ except for the relationships to standing-wave ratio.

LINES OPEN- OR SHORT-CIRCUITED AT THE FAR END

Point 2 is the open- or short-circuited end of the line, from which x and θ are measured.

(A) Voltages and Currents: Use the equations of the "Voltage and Current" section, with the following conditions.

Open circuited line:

$$\begin{aligned} \rho_2 &= 1.00 \angle 0^\circ = 1.00 \\ rE_2 &= jE_2 = E_2/2 \\ rI_2 &= -jI_2 \\ I_2 &= 0 \\ Z_2 &= \infty \end{aligned}$$

Short-circuited line:

$$\begin{aligned} \rho_2 &= 1.00 \angle 180^\circ = -1.00 \\ rE_2 &= -jE_2 \\ E_2 &= 0 \\ rI_2 &= jI_2 = I_2/2 \\ Z_2 &= 0 \end{aligned}$$

(B) Impedances and admittances:

$$\begin{aligned} Z_{oc} &= Z_0 \coth \gamma x \\ Z_{sc} &= Z_0 \tanh \gamma x \\ Y_{oc} &= Y_0 \tanh \gamma x \\ Y_{sc} &= Y_0 \coth \gamma x \end{aligned}$$

(C) For small attenuation: Use the equations for large swr in (D) of the preceding section, with the following conditions.

Open-circuited line:

$$g_a = 0$$

Short-circuited line:

$$r_b = 0$$

(D) When attenuation is neglected:

$$\begin{aligned} Z_{oc} &= -jR_0 \cot \theta \\ Z_{sc} &= jR_0 \tan \theta \\ Y_{oc} &= jG_0 \tan \theta \\ Y_{sc} &= -jG_0 \cot \theta \end{aligned}$$

(E) Relationships between Z_{oc} and Z_{sc} :

$$\begin{aligned} (Z_{oc}Z_{sc})^{1/2} &= Z_0 \\ \pm(Z_{sc}/Z_{oc})^{1/2} &= \tanh \gamma x \\ \pm(Z_{oc}/Z_{sc})^{1/2} &= \coth \gamma x \end{aligned}$$

(F) When attenuation is small (except for $\theta \approx n\pi/2$, $n = 1, 2, 3, \dots$)

$$\begin{aligned} \pm(Z_{sc}Z_{oc})^{1/2} &= \pm(Y_{oc}/Y_{sc})^{1/2} \\ &= \pm j[-(C_{oc}/C_{sc})]^{1/2} \\ &\quad \times [1 - j\frac{1}{2}(G_{oc}/\omega C_{oc} - G_{sc}/\omega C_{sc})] \end{aligned}$$

where $Y_{oc} = G_{oc} + j\omega C_{oc}$ and $Y_{sc} = G_{sc} + j\omega C_{sc}$. The + sign is to be used before the radical when C_{oc} is positive, and the - sign when C_{oc} is negative.

(G) $R/|X|$ component of input impedance of low-attenuation nonresonant line:

Short-circuited line (except when $\theta \approx \pi/2, 3\pi/2$, etc.)

$$\begin{aligned} R_1/|X_1| &= G_1/|B_1| \\ &= |(\alpha/\beta)\theta(\tan\theta + \cot\theta) + (B_0/G_0)| \\ &= |(\alpha/\beta)(2\theta/\sin 2\theta) + (B_0/G_0)| \end{aligned}$$

Open-circuited line (except when $\theta \approx \pi, 2\pi$, etc.)

$$\begin{aligned} R_1/|X_1| &= G_1/|B_1| \\ &= |(\alpha/\beta)\theta(\tan\theta + \cot\theta) - (B_0/G_0)| \\ &= |(\alpha/\beta)(2\theta/\sin 2\theta) - (B_0/G_0)| \end{aligned}$$

VOLTAGE REFLECTION COEFFICIENT AND STANDING-WAVE RATIO

$$\begin{aligned} \rho &= \frac{rE}{fE} = -\frac{rI}{fI} = (Z - Z_0)/(Z + Z_0) \\ &= (Y_0 - Y)/(Y_0 + Y) = |\rho| \angle 2\psi \end{aligned}$$

where ψ is the electrical angle to the nearest voltage maximum on the generator side of the point where ρ is measured (Figs. 3 and 4).

$$\rho_1 = \rho_2 e^{-2\alpha x} \angle -2\theta$$

$$|\rho_1| = |\rho_2| / 10^{A_0/10}$$

Voltage reflection coefficient in decibels

$$\rho_{dB} = -20 \log_{10} |1/\rho|$$

The minus sign is frequently omitted.

$$|\rho_{dB} \text{ at input}| = |\rho_{dB} \text{ at load}| + 2A_0$$

These two relationships and standing-wave ratio ver-

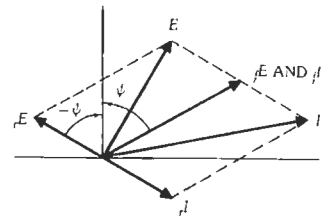


Fig. 3. Voltages and currents at time $t = 0$ at point ψ electrical degrees toward the load from a voltage standing-wave maximum.

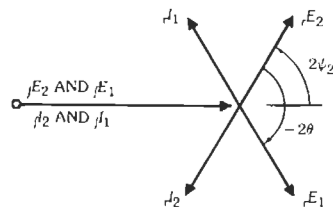


Fig. 4. Abbreviated diagram of a line with zero attenuation.

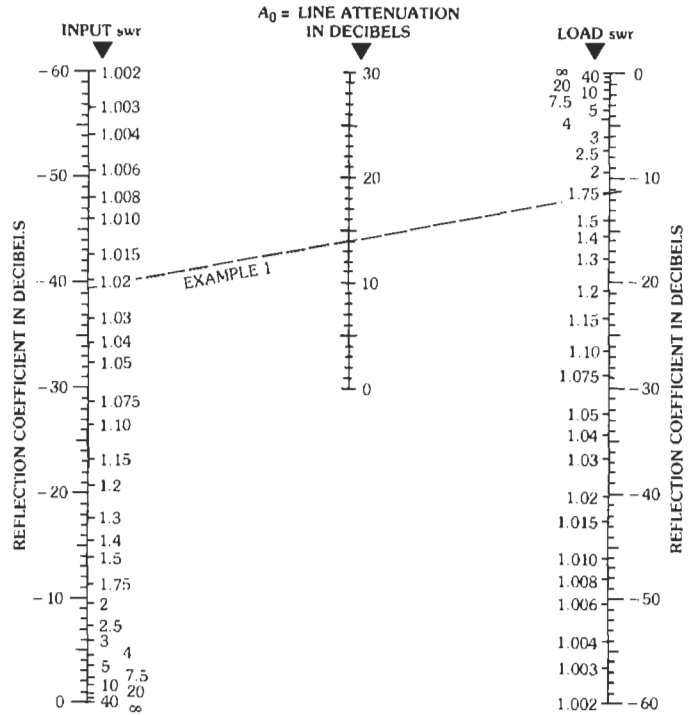


Fig. 5. Line attenuation and voltage reflection coefficient for low swr.

sus reflection coefficient in decibels are shown in Figs. 5 and 6.

$$Z = E/I = ({}_rE + {}_lE)/({}_rI + {}_lI)$$

$$= Z_0[(1 + \rho)/(1 - \rho)]$$

$$Z/Z_0 = (1 + \rho)/(1 - \rho)$$

$$= \frac{1 + jS \cot \psi}{S + j \cot \psi}$$

$$S = |E_{\max}/E_{\min}| = |I_{\max}/I_{\min}|$$

$$= \frac{|{}_rE| + |{}_lE|}{|{}_rE - |{}_lE|} = \frac{|{}_rI| + |{}_lI|}{|{}_rI| - |{}_lI|}$$

$$= \frac{1 + |\rho|}{1 - |\rho|} = r_a = g_a^{-1} = g_b = r_b^{-1}$$

$$|\rho| = (S - 1)/(S + 1)$$

$$1/S_1 = \tanh[\alpha x + \tanh^{-1}(1/S_2)]$$

$$= \tanh[0.1151A_0 + \tanh^{-1}(1/S_2)]$$

(A) For high standing-wave ratio: When the ratio S_1 is greater than 6/1, then with 1 percent accuracy or better

$$1/S_1 = 1/S_2 + \alpha x = 1/S_2 + 0.115A_0$$

$$|\rho_{dB}| = 17.4/S$$

Subject to the conditions below, the standing-wave ratio is given by one or the other of

$$S \approx (1 + x^2)/r$$

$$S \approx (1 + b^2)/g$$

where

$$r + jx = Z/Z_0 = (1/R_0)[R - (B_0/G_0)X + jX]$$

$$g + jb = Y/Y_0 = (1/G_0)[G + (B_0/G_0)B + jB]$$

Conditions, for 1-percent accuracy

$$r < 0.1 |x + 1/x| \quad \text{when } |x| > 0.3$$

$$g < 0.1 |b + 1/b| \quad \text{when } |b| > 0.3$$

The boundary of the 1-percent-error region can be plotted on the Smith chart (Fig. 7) by use of the equation (for impedances)

$$|\cot \psi| = 0.1S^2/(S^2 - 1)^{1/2}$$

The same boundary line on the chart holds when reading admittances. The area outside the solid heart-shaped curve is where the swr equation is accurate to within 1 percent. The area outside the dashed curve is where the reciprocal of $r + jx$ lies in the permitted region.

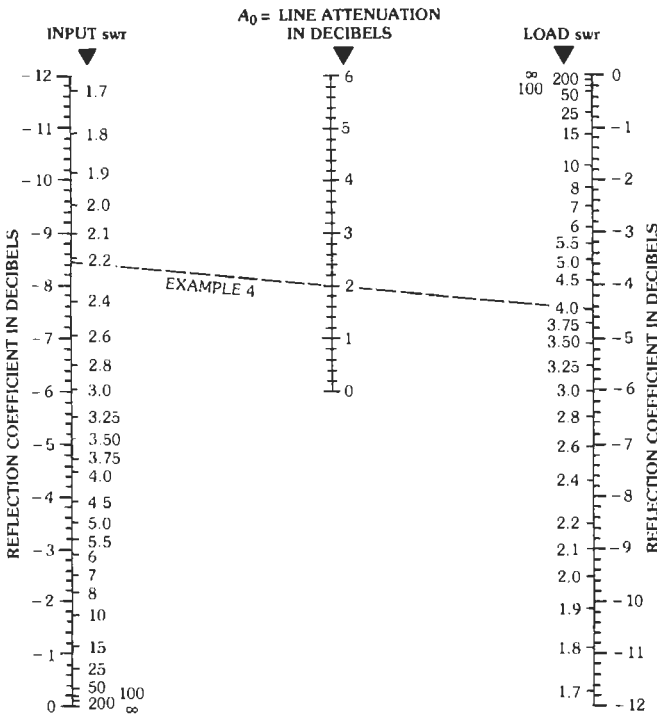


Fig. 6. Line attenuation and voltage reflection coefficient for high swr.

POWER AND EFFICIENCY

The net power flowing toward the load is

$$P = |jE|^2 G_0 [1 - |\rho|^2 + 2|\rho| (B_0/G_0) \sin 2\psi]$$

(A) When the angle B_0/G_0 of the characteristic

admittance is negligibly small, the net power flowing toward the load is given by

$$P = G_0 (|jE|^2 - |rE|^2) = |jE|^2 G_0 (1 - |\rho|^2)$$

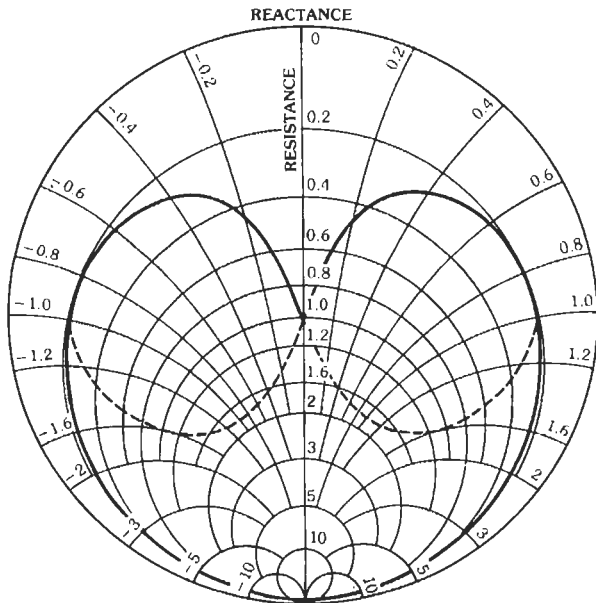


Fig. 7. Permitted region for use of equation $S \approx (1 + x^2)/r$. (From W. W. Macalpine, "Computation of Impedance and Efficiency of Transmission Line with High Standing-Wave Ratio," Trans. of the AIEE, Vol. 72, Part 1, p. 336, Fig. 2; July 1953.)

$$= |E_{\max} E_{\min}| / R_0$$

$$P_1 = |E_2|^2 G_0 (\epsilon^{2(\alpha/\beta)\theta} - |\rho_2|^2 \epsilon^{-2(\alpha/\beta)\theta})$$

(B) Efficiency, when B_0/G_0 is negligibly small:

$$\begin{aligned} \eta &= P_2/P_1 = \frac{1 - |\rho_2|^2}{\epsilon^{2(\alpha/\beta)\theta} - |\rho_2|^2 \epsilon^{-2(\alpha/\beta)\theta}} \\ &= \eta_{\max} \frac{1 - |\rho_2|^2}{1 - |\rho_2|^2 \eta_{\max}^2} = \frac{1 - |\rho_2|^2}{1 - |\rho_1|^2} \epsilon^{-2\alpha x} \\ &= \frac{1/|\rho_2| - |\rho_2|}{1/|\rho_1| - |\rho_1|} = \frac{S_1 - 1/S_1}{S_2 - 1/S_2} \end{aligned}$$

The maximum error in the above expressions is

$$\begin{aligned} &\pm 100(S_2 - 1/S_2)B_0/G_0 \text{ percent} \\ &\pm 4.34(S_2 - 1/S_2)B_0/G_0 \text{ decibels} \end{aligned}$$

When the ratio S_1 is greater than 6/1:

$$\eta \approx S_1/S_2 \approx (1 + 0.115A_0S_2)^{-1}$$

When the load matches the line, $\rho_2 = 0$ and the efficiency is accurately

$$\eta_{\max} = \exp[-2(\alpha/\beta)\theta] = \exp(-2\alpha x) = 10^{-A_0/10}$$

$$A - A_0 = 10 \log_{10}(\eta_{\max}/\eta)$$

Fig. 8 is drawn from the expressions in this paragraph.

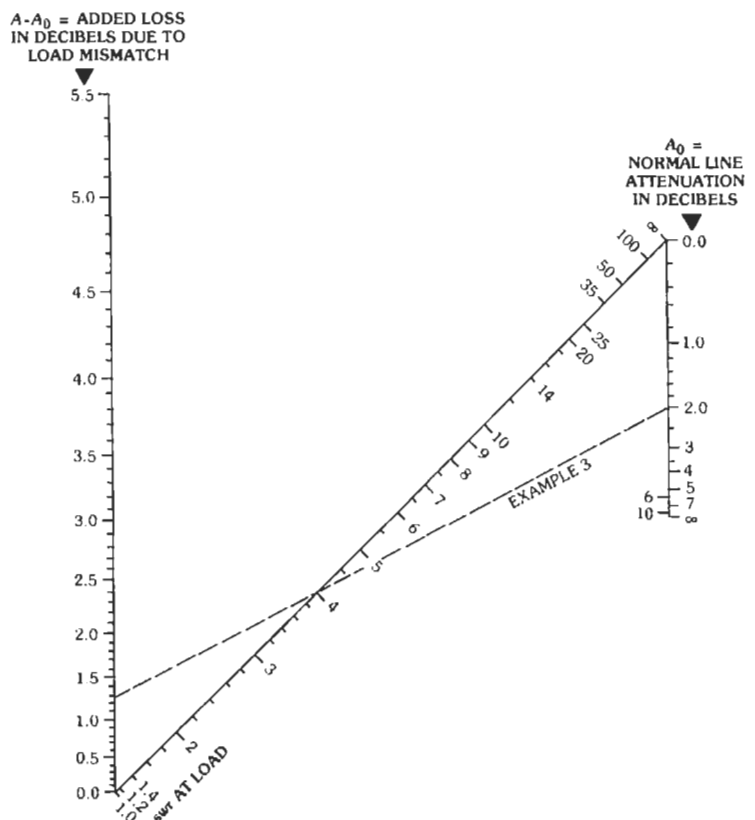
(C) Efficiency, when swr is high:

$$\begin{aligned} \eta &= \frac{P_2}{P_1} = \frac{R_2}{R_1} \left(\frac{1 + x_1^2}{1 + x_2^2} \right) = \frac{G_2}{G_1} \left(\frac{1 + b_1^2}{1 + b_2^2} \right) \\ &= \frac{R_2}{R_0^2 G_1} \left(\frac{1 + b_1^2}{1 + x_2^2} \right) = \frac{R_0^2 G_2}{R_1} \left(\frac{1 + x_1^2}{1 + b_2^2} \right) \end{aligned}$$

where R is the ohmic resistance and x is the normalized reactance, and similarly for G and b . It is important that the R s and G s be computed properly, using equations in the following section, headed "Transformation of Impedance on Lines with High SWR." Note the identity of the efficiency equations with the left-hand terms of the impedance equations. The conditions for accuracy are the same as stated for the impedance equations for high standing-wave ratio.

Example: Physical significance of the equation for efficiency at high standing-wave ratio: Subject to stated conditions, approximately, $x = \cot\psi$ and $I = I_{\max}\sin\psi$.

Fig. 8. Standing-wave loss factor. Due to load mismatch, an increase of loss in decibels as read from this figure must be added to normal line attenuation to give total dissipation loss in the line. This does not include mismatch loss due to any difference of line input impedance from the conjugate of the generator impedance [Paragraph (B) in section headed "Power and Efficiency"].



I_{max} = current standing-wave maximum, practically constant along line when standing-wave ratio > 6. Then

$$P = I^2 R = I_{max}^2 R / (1 + x^2)$$

When line length is greater than $\frac{1}{3}$ wavelength, then

$$\eta \approx [1 + 0.115 A_0 (1 + x_2^2) (R_0 / R_2)]^{-1}$$

(D) Loss in nepers = $\frac{1}{2} \log_e (P_1 / P_2) = 0.1151 \times$ (loss in decibels).

For a matched line, loss = attenuation = $(\alpha / \beta) \theta = \alpha x$ nepers.

Loss in decibels = $10 \log_{10} (P_1 / P_2) = 8.686 \times$ (loss in nepers).

When $2(\alpha / \beta) \theta$ is small

$$P_1 / P_2 = 1 + 2(\alpha / \beta) \theta \frac{1 + |\rho_2|^2}{1 - |\rho_2|^2}$$

and

decibels/wavelength

$$= 10 \log_{10} \left(1 + 4\pi(\alpha / \beta) \frac{1 + |\rho_2|^2}{1 - |\rho_2|^2} \right)$$

(E) For the same power flowing in a line with standing waves as in a matched or "flat" line:

$$P = |E_{flat}|^2 / R_0$$

$$|E_{max}| = |E_{flat}| S^{1/2}$$

$$|E_{min}| = |E_{flat}| / S^{1/2}$$

$$|_f E| = \frac{1}{2} |E_{flat}| (S^{1/2} + S^{-1/2})$$

$$|_r E| = \frac{1}{2} |E_{flat}| (S^{1/2} - S^{-1/2})$$

When the loss is small, so that S is nearly constant over the entire length, then per half wavelength

$$(\text{power loss}) / (\text{loss for flat line}) \approx \frac{1}{2} (S + 1/S)$$

(F) The power dissipation per unit length, for unity standing-wave ratio, is

$$\Delta P_d / \Delta x = 2\alpha P$$

$$\frac{(\text{dissipation in watts/foot})}{(\text{line power in kilowatts})} = 2.30 \text{ (decibels/100 feet)}$$

where the decibels/100 feet is the normal attenuation for a matched line.

When $swr > 1$, the dissipation at a current maximum is S times that for $swr = 1$, assuming the attenuation to be due to conductor loss only. The multiplying factor for local heating reaches a minimum value equal to $(S +$

$1/S)/2$ all along the line when conductor loss and dielectric loss are equal.

(G) Further considerations on power and efficiency are given in the section headed "Mismatch and Transducer Loss" (p. 29-13).

TRANSFORMATION OF IMPEDANCE ON LINES WITH HIGH SWR*

When standing-wave ratio is greater than 10 or 20, resistance cannot be read accurately on the Smith chart, although it is satisfactory for reactance.

Use the equation (Fig. 9)

$$R_1 = R_2 \frac{1 + x_1^2}{1 + x_2^2} + R_0 (1 + x_1^2) \times \left[(\alpha / \beta) \theta + (B_0 / G_0) \left(\frac{x_1}{1 + x_1^2} - \frac{x_2}{1 + x_2^2} \right) \right]$$

where,

R = ohmic resistance,

$x = X / R_0$ = normalized reactance.

When admittance is given or required, similar equations can be written with the aid of the following tabulation. The top row shows the terms in the above equation.

R_1	R_2	x_1^2	x_2^2	R_0	x_1	$-x_2$
G_1	G_2	b_1^2	b_2^2	$1/R_0$	$-b_1$	b_2
R_1	$G_2 R_0^2$	x_1^2	b_2^2	R_0	x_1	b_2
G_1	R_2 / R_0^2	b_1^2	x_2^2	$1/R_0$	$-b_1$	$-x_2$

For transforming R to G or vice versa:

$$R = R_0^2 G |x/b|$$

where x and b are read on the Smith chart in the usual manner for transforming impedances to admittances.

* W. W. Macalpine, "Computation of Impedance and Efficiency of Transmission Lines with High Standing-Wave Ratio," *Transactions of the AIEE*, Vol. 72, Part I, pp. 334-339; July 1953; also *Electrical Communication*, Vol. 30, pp. 238-246; September 1953.

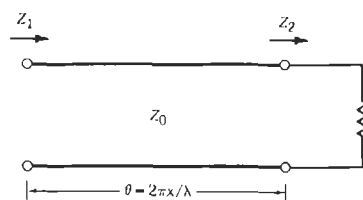


Fig. 9. Transmission-line impedances.

The conditions for roughly 1-percent accuracy of the equations are:

Standing-wave ratio is greater than 6/1 at input; $|B_0/G_0| < 0.1$; $r + jx$ or $g + jb$ (whichever is used, at each end of line) meet the requirements stipulated in paragraph (A) of the section headed "Voltage Reflection Coefficient and Standing-Wave Ratio" (p. 29-8); and the line parameters and given impedance are known to 1-percent accuracy.

When line length is greater than $\frac{1}{3}$ wavelength, then

$$R_1 \approx R_2[(1 + x_1^2)/(1 + x_2^2)] + 0.115A_0R_0(1 + x_1^2)$$

$$\frac{R_1/R_0}{1 + x_1^2} \approx \frac{R_2/R_0}{1 + x_2^2} + (\alpha/\beta)\theta$$

The equation for resistance transformation is derived from expressions for high swr in paragraph (A), just referred to.

Example: A load of $0.4 - j2000$ ohms is fed through a length of RG-218/U cable at a frequency of 2.0 megahertz. What are the input impedance and the efficiency for a 24-foot length of cable and for a 124-foot length?

For RG-218/U, the attenuation at 2.0 megahertz is 0.095 decibel/100 feet (see Fig. 29). The dielectric constant $\epsilon = 2.26$, and F_p is negligibly small. Then, by equations in (B) and (C), p. 29-5:

$$B_0/G_0 = \alpha\beta = (\text{dB}/100 \text{ ft}) (\lambda_{\text{meters}})/1663$$

$$= [0.095 \times 150/(2.26)^{1/2}]/1663 = 0.0057$$

$$x/\lambda = x\epsilon^{1/2}/984 = 24 \times 2.0 \times 1.5/984 = 0.073$$

$$\theta = 2\pi x/\lambda = 0.46 \text{ radian for 24-foot length}$$

while

$$x/\lambda = 0.38 \text{ and } \theta = 2.4 \text{ for 124-foot length}$$

$$Z_2/Z_0 \approx (0.4 - j2000)/50 = 0.008 - j40$$

For the 24-foot length, by the Smith chart

$$x_1 = X_1/Z_0 = -1.9, \text{ or } X_1 = -95 \text{ ohms}$$

The conditions for accuracy of the resistance transformation equation are satisfied. Now

$$1 + x_1^2 = 1 + (1.9)^2 = 4.6$$

$$1 + x_2^2 = 1 + (40)^2 = 1600$$

$$R_1 = 0.4(4.6/1600) + 50 \times 4.6 \times 0.0057$$

$$\times [0.46 - (1.9/4.6) + (40/1600)]$$

$$= 0.0012 + 0.105 = 0.106 \text{ ohm}$$

The efficiency equation in paragraph (C) on p. 29-11 gives

$$\eta = 0.0012/0.106 = 0.0113, \text{ or } 1.1 \text{ percent}$$

where the 0.0012 figure is taken directly from the first quantity on the right-hand side of the computation of the value of R_1 .

Similarly, for the 124-foot length, $x_1 = 1.1$, $X_1 = 55$ ohms, $1 + x_1^2 = 2.21$, $R_1 = 0.00055 + 1.83 = 1.83$ ohms

$$\eta = 0.00055/1.83 = 3.1 \times 10^{-4}, \text{ or } 0.03 \text{ percent}$$

Tabulating the results,

Length (feet)	Input Impedance (ohms)	Efficiency (%)	Loss (dB)
24	$0.106 - j95$	1.1	19.6
124	$1.8 + j55$	0.03	35

The considerably greater loss for 124 feet compared with 24 feet is because the transmission passes through a current maximum where the loss per unit length is much higher than at a current minimum.

MISMATCH AND TRANSDUCER LOSS

Figs. 5, 6, and 8, plus the equations in this section, permit the calculation of loss when impedance mismatch exists in a transmission-line system; also, the change in standing-wave ratio along a line due to attenuation can be determined.

One End Mismatched

When either generator or load impedance is mismatched to the Z_0 of the line and the other is matched (Fig. 10)

$$(\text{mismatch loss}) = P_m/P$$

$$= 1/(1 - |\rho|^2) = (S + 1)^2/4S \tag{Eq. 1}$$

where,

P = power delivered to load,

P_m = power that would be delivered were system matched,

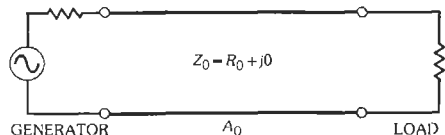


Fig. 10. Transmission line with generator and load.

S = standing-wave ratio of mismatched impedance referred to Z_0 .

Compared with an ideal transducer (ideal matching network between generator and load)

$$(\text{transducer loss}) = A_0 + 10 \log_{10}(P_m/P) \text{ decibels} \quad (\text{Eq. 2})$$

where A_0 = normal attenuation of line.

Generator and Load Mismatched

$$|X_0/R_0| \ll 1$$

When mismatches exist at both ends of the system (Fig. 11)

(mismatch loss at input)

$$\begin{aligned} &= P_m/P \\ &= [(R_g + R_1)^2 + (X_g + X_1)^2] / 4R_gR_1 \quad (\text{Eq. 3}) \end{aligned}$$

$$\begin{aligned} (\text{transducer loss}) &= (A - A_0) + A_0 \\ &+ 10 \log_{10}(P_m/P) \text{ decibels} \quad (\text{Eq. 4}) \end{aligned}$$

where $(A - A_0)$ = standing-wave loss factor obtained from Fig. 8 for S = standing-wave ratio at load.

Notes on Equation (3)

Equation (3) reduces to Eq. (1) when X_g and/or X_1 is zero.

In Eq. (3), the impedances can be either ohmic or normalized with respect to any convenient Z_0 .

When determining input impedance $R_1 + jX_1$ on a Smith chart, adjust the radius arm for S at the input, determined from that at the output by the aid of Figs. 5 and 6.

For the junction of two admittances, use Eq. (3) with G and B substituted for R and X , respectively.

Equation (3) is valid for a junction in any linear passive network. The same is true of Eq. (1) when at least one of the impedances concerned is purely resistive.

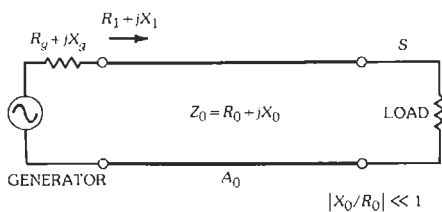


Fig. 11. Transmission line with mismatch at both ends.

tive. Determine S as if one impedance were that of a line.

Examples

Example 1: The swr at the load is 1.75, and the line has an attenuation of 14 decibels. What is the input swr ?

Using Fig. 5, set a straightedge through the 1.75 division on the "load swr " scale and the 14-decibel point on the middle scale. Read the answer on the "input swr " scale, which the straightedge intersects at 1.022.

Example 2: Readings on a reflectometer show the reflected wave to be 4.4 decibels below the incident wave. What is the swr ?

Using Fig. 6, locate the reflection coefficient 4.4 (or -4.4) decibels on either outside scale. Beside it, on the same horizontal line, read $swr = 4.0+$.

Example 3: A 50-ohm line is terminated with a load of $200 + j0$ ohms. The normal attenuation of the line is 2.00 decibels. What is the loss in the line?

Using Fig. 8, align a straightedge through the points $A_0 = 2.0$ and $swr = 4.0$. Read $A - A_0 = 1.27$ decibels on the left-hand scale. Then the transmission loss in the line is

$$A = 1.27 + 2.00 = 3.27 \text{ decibels.}$$

This is the dissipation, or heat loss, as opposed to the mismatch loss at the input (example 4).

Example 4: In the preceding example, suppose the generator impedance is $100 + j0$ ohms, and the line is 5.35 wavelengths long. What is the mismatch loss between the generator and the line?

According to example 3, the load $swr = 4.0$ and the line attenuation is 2.0 decibels. Then, using Fig. 6, the input swr is found to be 2.22. On the Smith chart, locate the point corresponding to 0.35 wavelength toward the generator from a voltage maximum, and $swr = 2.22$. Read the input normalized impedance as $0.62 + j0.53$ with respect to $Z_0 = 50$ ohms. Now the mismatch loss at the input can be determined by use of Eq. (3). However, since the generator impedance is nonreactive, Eq. (1) can be used if desired. Refer to the following paragraph and to the "Notes on Equation (3)" above.

With respect to $100 + j0$ ohms, the normalized impedance at the line input is $0.31 + j0.265$, which gives $swr = 3.5$ according to the Smith chart. Then by Eq. (1), $P_m/P = 1.45$, giving a mismatch loss of 1.62 decibels. The transducer loss is found by using the results of examples 3 and 4 in Eq. (4). This is

$$1.27 + 2.00 + 1.62 = 4.9 \text{ decibels}$$

ATTENUATION AND RESISTANCE OF TRANSMISSION LINES AT ULTRAHIGH FREQUENCIES

The normal or matched-line attenuation in decibels/100 feet is

$$A_{100} = 4.34R_t/Z_0 + 2.78f\epsilon^{1/2}F_p$$

where the total line resistance/100 feet (for perfect surface conditions of the conductors) is, for copper coaxial line

$$R_t = 0.1(1/d + 1/D)f^{1/2}$$

and for copper 2-wire open line

$$R_t = (0.2/d)f^{1/2}$$

where,

- D = diameter of inner surface of outer coaxial conductor in inches,
- d = diameter of conductors (coaxial-line center conductor) in inches,
- f = frequency in megahertz,
- ϵ = dielectric constant relative to air,
- F_p = power factor of dielectric at frequency f .

For other conductor materials, the resistance of a conductor of diameter d (and similarly for D) is

$$0.1(1/d) (f\mu_r\rho/\rho_{Cu})^{1/2} \text{ ohms/100 feet}$$

where,

- μ_r = relative permeability of material (1 for nonmagnetic materials),
- ρ = resistivity of material at any temperature,
- ρ_{Cu} = resistivity of copper at 20 °C (1.724 microhms-centimeter),
- f = frequency in megahertz.

RESONANT LINES

Symbols:

- f_0 = resonance frequency in megahertz
- G_a = conductance load in mhos at voltage standing-wave maximum, equivalent to some or all of the actual loads
- k = coefficient of coupling
- n = integral number of quarter wavelengths
- $\rho = k^2 Q_{1s} Q_{2s}$ = load transfer coefficient or matching factor
- P_c = power converted into heat in resonator
- P_m = power available from generator in watts
= $E_{oc}^2/4R_{gn}$

- P_x = power transferred when load is directly connected to generator (for single resonators); or an analogous hypothetical power (for two coupled resonators)
- Q = figure of merit of a resonator as it exists, whether loaded or unloaded
- Q_d = doubly loaded Q (all loads being included)
- Q_s = singly loaded Q (all loads included except one). For a pair of coupled resonators, Q_{1s} is the value for the first resonator when isolated from the other. (Similarly for Q_{2s})
- Q_u = unloaded Q
- R_b = resistance load in ohms at voltage standing-wave minimum, equivalent to some or all of the actual loads
- R_u = resistance similar to R_b except for unloaded resonator
- R_1 = generator resistance, referred to short-circuited end
- R_2 = load resistance
- $S_x = R_1/R_2$ or R_2/R_1 = mismatch factor between generator and load
- Z_{10} = characteristic impedance of the first of a pair of resonators
- θ_1 = electric angle from a voltage standing-wave minimum point

(A) Q of a resonator (electrical, mechanical, or any other) is

$$Q = 2\pi \frac{\text{(energy stored)}}{\text{(energy dissipated per cycle)}}$$

$$= 2\pi f \frac{\text{(energy stored)}}{\text{(power dissipation)}}$$

In a freely oscillating system, the amplitude decays exponentially.

$$I = I_0 \exp(-\pi ft/Q)$$

(B) Unloaded Q of a resonant line:

$$Q_u = \beta/2\alpha$$

the line length being n quarter-wavelengths, where n is a small integer. The losses in the line are equivalent to those in a hypothetical resistor at the short-circuited end (D) in the section headed "Impedance and Admittance"

$$R_u = n\pi Z_0/4Q_u$$

(C) Loaded Q of a resonant line (Fig. 12):

$$Q^{-1} = Q_u^{-1} + (4R_b/n\pi Z_0) + (4G_a/n\pi Y_0)$$

$$= (4/n\pi Z_0) (R_u + R_b + G_a/Y_0^2)$$

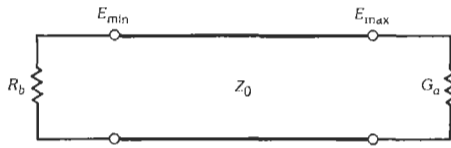


Fig. 12. Quarter-wave line with loadings at nominal short-circuit and open-circuit points.

All external loads can be referred to one end and represented by either R_b or G_a as in Fig. 13.

The total loading is the sum of all the individual loadings.

General conditions:

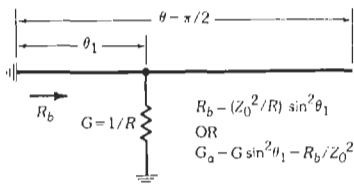
$$R_b/Z_0 = G_a/Y_0 \ll 1.0$$

or, roughly, $Q > 5$.

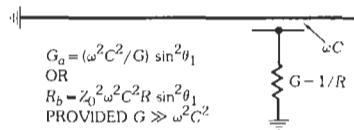
(D) Input admittance and impedance:

The converse of the equations for Fig. 13 can be used at the resonance frequency. Then R or G is the input impedance or admittance, while

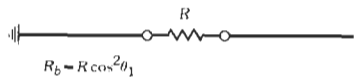
$$R_b = n\pi Z_0 / 4Q_s$$



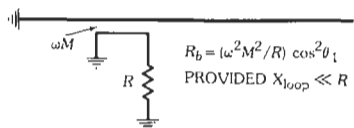
(A) Shunt or tapped load.



(B) Probe coupling.



(C) Series load.



(D) Loop coupling.

Fig. 13. Typical loaded quarter-wave sections with apparent R_b equivalent to the loading at distance θ_1 from voltage-minimum point of the line. Outer conductor not shown.

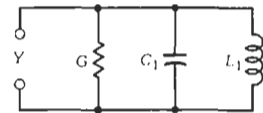
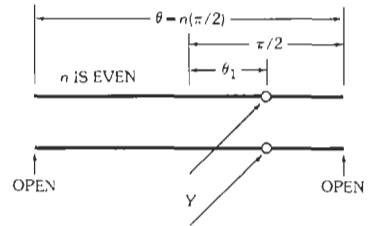
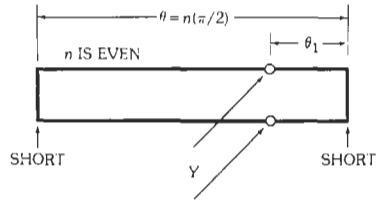
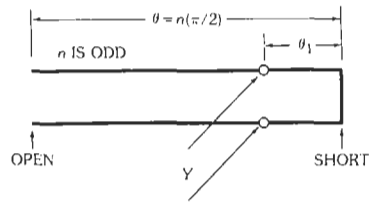


Fig. 14. Resonant transmission lines and their equivalent lumped circuit.

where $Q_s =$ singly loaded Q with the losses and all the loads considered except that at the terminals where input R or G is being measured.

In the vicinity of the resonance frequency, the input admittance when looking into a line at a tap point θ_1 in Fig. 14 is approximately

$$Y = G + jB = \frac{n\pi Y_0}{4 \sin^2 \theta_1} \left(Q_s^{-1} + j2 \frac{f - f_0}{f_0} \right)$$

provided

$$|f - f_0| / f_0 \ll 1.0$$

and

$$|[\theta(f - f_0)/f_0] \cot \theta_1| \ll 1.0$$

where $\theta = n\pi/2 =$ length of line at f_0 . The equation is

not valid when $\theta_1 \approx 0, \pi, 2\pi$, etc., except that it is good near the short-circuited end when $f - f_0 \approx 0$.

Such a resonant line is approximately equivalent to a lumped LCG parallel circuit, where

$$\omega_0^2 L_1 C_1 = (2\pi f_0)^2 L_1 C_1 = 1$$

Admittance of the equivalent circuit is

$$Y = G + j[\omega C_1 - (1/\omega L_1)]$$

$$\approx \omega_0 C_1 \{Q_s^{-1} + j2[(f - f_0)/f_0]\}$$

Then, subject to the conditions stated above

$$L_1 = (4 \sin^2 \theta_1) / n\pi \omega_0 Y_0$$

$$C_1 = n\pi Y_0 / (4\omega_0 \sin^2 \theta_1) = nY_0 / (8f_0 \sin^2 \theta_1)$$

$$G = n\pi Y_0 / (4Q_s \sin^2 \theta_1)$$

$$Q_s = \omega_0 C_1 / G = 1/\omega_0 L_1 G$$

Similarly, the input impedance at a point in series with the line (Figs. 13C and D) is

$$Z = R + jX = \frac{n\pi Z_0}{4 \cos^2 \theta_1} \left(Q_s^{-1} + j2 \frac{f - f_0}{f_0} \right)$$

provided

$$|f - f_0| / f_0 \ll 1.0$$

and

$$|\theta[(f - f_0)/f_0] \tan \theta_1| \ll 1.0$$

The equation is not valid when $\theta_1 \approx \pi/2, 3\pi/2$, etc.

The voltage standing-wave ratio at resonance, on the generator (Fig. 15), is

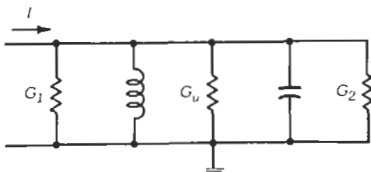
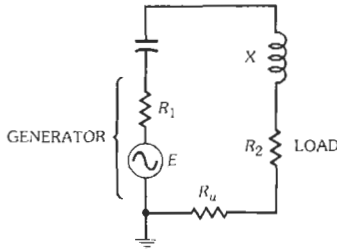


Fig. 15. Equivalent circuits of a resonant line (or a lumped tuned circuit) as seen at the short-circuited and open-circuited ends. All the power equations are good for either lumped or distributed parameters.

$$S = (R_2 + R_u) / R_1$$

$$= \frac{(R_2/R_1)Q_u + Q_d}{Q_u - Q_d}$$

When $R_1 = R_2$

$$S = \frac{1 + Q_d/Q_u}{1 - Q_d/Q_u}$$

$$\rho = Q_d/Q_u$$

(E) Insertion loss (Fig. 15): At resonance, for either a distributed or a lumped-constant device

(dissipation loss)

$$= 10 \log_{10}(P_x/P_{out})$$

$$= 20 \log_{10}[1/(1 - Q_d/Q_u)]$$

$$\approx 20 \log_{10}(1 + Q_d/Q_u)$$

$$\approx 8.7Q_d/Q_u \text{ decibels}$$

(mismatch loss)

$$= 10 \log_{10}(P_m/P_x)$$

$$= 10 \log_{10}[(1 + S_x)^2/4S_x] \text{ decibels}$$

The dissipation loss also includes a small additional mismatch loss due to the presence of the resonator. The error in the form $20 \log_{10}(1 + Q_d/Q_u)$ is about twice that of the form $8.7Q_d/Q_u$. The last expression ($8.7Q_d/Q_u$) is in error compared with the first, $20 \log_{10}[1/(1 - Q_d/Q_u)]$, by roughly $-50(Q_d/Q_u)$ percent for $(Q_d/Q_u) < 0.2$.

The selectivity is given on page 8-7, where $Q = Q_d$. That equation is accurate over a smaller range of $(f - f_0)$ for a resonant line than it is for a single tuned circuit.

At resonance*

$$P_{out}/P_{in} = R_2 / (R_u + R_2)$$

$$= \frac{Q_u - Q_d}{Q_u + (R_1/R_2)Q_d} = 1 - Q_s/Q_u$$

where Q_s is for the resonator loaded with R_2 only.

The maximum power transfer, for fixed Q_u, Q_d , and Z_0 occurs when $R_1 = R_2$. Then

* When the line is resonated by a reactive load ($\theta \neq n\pi/2$), it is frequently preferable to use the resistance form of the equation. Compute R_u by the method in the section "Transformation of Impedance on Lines With High SWR," or the section "Impedance and Admittance," where $Z_0 = R_0(1 - jB_0/G_0)$.

$$P_{out}/P_{in} = (Q_u - Q_d)/Q_u + Q_d = 1 - Q_s/Q_u$$

$$P_{out}/P_m = (1 - Q_d/Q_u)^2$$

$$P_{in}/P_m = 1 - (Q_d/Q_u)^2$$

When the generator R_1 or G_1 is negligibly small (then $Q = Q_s = Q_d$)

$$(P_{in}/P_{out})_s = Q_u/(Q_u - Q)$$

(F) Power dissipation (= P_c):

$$P_c/P_m = \frac{4(Q_d/Q_u)(1 - Q_d/Q_u)}{1 + R_2/R_1}$$

For matching input and output ($R_1 = R_2$)

$$P_c/P_m = 2(Q_d/Q_u)(1 - Q_d/Q_u) \approx 2Q_d/Q_u \quad (\text{for } Q_d \ll Q_u)$$

$$P_c/P_{out} = 2Q_d/(Q_u - Q_d)$$

$$P_c/P_{in} = 2Q_d/(Q_u + Q_d)$$

For generator matched by load plus cavity

$$P_c/P_m = 2Q_d/Q_u$$

When the generator R_1 or G_1 is negligibly small

$$(P_c/P_{out})_s = Q/(Q_u - Q)$$

$$(P_c/P_{in})_s = Q_s/Q_u$$

(G) Voltage and current: At the current-maximum point of an n -quarter-wavelength resonant line

$$I_{sc} = 4 \left[\frac{P_m Q_d (1 - Q_d/Q_u)}{(1 + R_2/R_1) n \pi Z_0} \right]^{1/2}, \text{ rms amperes}$$

$$= 4 \left[\frac{P_m Q_d}{\{1 + [(R_2 + R_u)/R_1]\} n \pi Z_0} \right]^{1/2}$$

When the generator R_1 or G_1 is negligibly small

$$I_{sc} = 2 \left[\frac{P_s Q_s}{n \pi Z_0 (R_2 + R_u)/R_s} \right]^{1/2}$$

where,

P_s = rated power of generator,
 R_s = rated load impedance as transformed into current-maximum point of cavity.

$$I = I_{sc} \cos \theta_1$$

$$E = Z_0 I_{sc} \sin \theta_1$$

The voltage and current are in quadrature time phase. When $R_1 = R_2 + R_u$ and $n = 1$

$$I_{sc} \approx (8 P_m Q_d / \pi Z_0)^{1/2}$$

In a lumped-constant tuned circuit

$$I = 2 \left[\frac{P_m Q_d (1 - Q_d/Q_u)}{(1 + R_2/R_1) X} \right]^{1/2}$$

(H) Pair of coupled resonators (Fig. 16):

With inductive coupling near the short-circuited end of a pair of quarter-wave resonant lines

$$k = (4/\pi) \omega M / (Z_{10} Z_{20})^{1/2}$$

For coupling through a lossless quarter-wavelength line, inductively coupled near the short-circuited ends of the resonators (Fig. 16D):

$$k = \frac{4 \omega^2 M_1 M_2}{\pi Z_0 (Z_{10} Z_{20})^{1/2}}$$

Probe coupling near top (Fig. 16C):

$$k = (4/\pi) \omega C_{12} (Z_{10} Z_{20})^{1/2} \sin \theta_1 \sin \theta_2$$

For lumped-constant coupled circuits, p and k are defined on pp. 8-6 and 8-2, respectively.

In either lumped or distributed resonators

(dissipation loss)

$$= 10 \log_{10} (P_s/P_{out})$$

$$= 10 \log_{10} [1/(1 - Q_{1s}/Q_{1u})(1 - Q_{2s}/Q_{2u})]$$

$$\approx 20 \log_{20} [1/(1 - Q_s/Q_u)]$$

$$\approx 20 \log_{10} (1 + Q_s/Q_u)$$

$$\approx 8.7 Q_s/Q_u \text{ decibels}$$

where

$$Q_s/Q_u = [(Q_{1s}/Q_{1u})(Q_{2s}/Q_{2u})]^{1/2}$$

provided (Q_{1s}/Q_{1u}) and (Q_{2s}/Q_{2u}) do not differ by a ratio of more than 4 to 1, and neither exceeds 0.2.

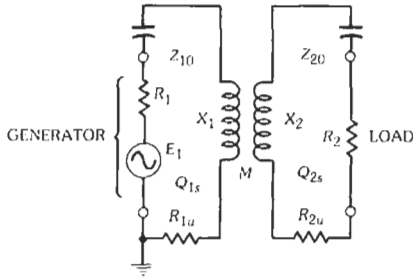
(mismatch loss at f_0)

$$= 10 \log_{10} (P_m/P_s)$$

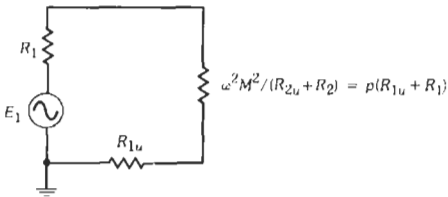
$$= 10 \log_{10} [(1 + p)^2/4p] \text{ decibels.}$$

Equations and curves for selectivity are given on pp. 8-7 and 8-8, where $Q = Q_s$.

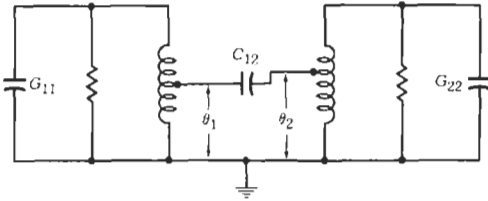
At the peaks, when $p \geq 1$, the mismatch loss is zero,



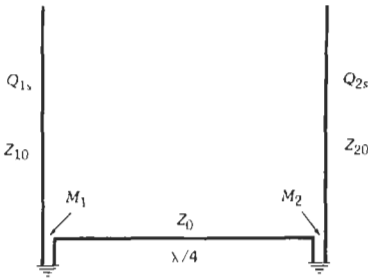
(A) Equivalent circuit with resistances as seen at the short-circuited end.



(B) Equivalent circuit of first resonator at resonance frequency.



(C) Probe-coupled or aperture-coupled resonators.



(D) Quarter-wavelength line coupling.

Fig. 16. Two coupled resonators.

except for some that is included in the dissipation loss.

Input voltage standing-wave ratio at f_0 for equal or unequal resonators:

$$S = \frac{p + Q_{1s}/Q_{1u}}{1 - Q_{1s}/Q_{1u}}$$

At the peak frequencies ($p \geq 1$) for equal or nearly equal resonators:

$$S = \frac{1 + Q_{1s}/Q_{1u}}{1 - Q_{1s}/Q_{1u}}$$

Similar equations, using subscript 2 instead of 1, apply at the output.

When the resonators are isolated, each one presents to the generator or load an swr of

$$S = (Q_u/Q_s) - 1$$

The power dissipation in either lumped or distributed (quarter-wave) devices, where the two resonators are not necessarily identical, but $Q_s \ll Q_u$, is

$$P_{1c} = I_{1sc}^2 R_{1u} [4/(1+p)^2] P_m Q_{1s}/Q_{1u}$$

$$P_{2c} = [4p/(1+p)^2] P_m Q_{2s}/Q_{2u}$$

These equations and those below for the currents assume that P_m is concentrated at f_0 .

The currents in quarter-wave resonant lines, when $Q_s \ll Q_u$.

$$I_{1sc} = [4/(1+p)] (P_m Q_{1s}/\pi Z_{10})^{1/2}$$

$$I_{2sc}/I_{1sc} = (p Z_{10} Q_{2s}/Z_{20} Q_{1s})^{1/2}$$

Similarly, for a pair of tuned circuits at resonance, when $Q_s \ll Q_u$.

$$I_1 = [2/(1+p)] (P_m Q_{1s}/X_1)^{1/2}$$

$$I_2/I_1 = (p X_1 Q_{2s}/X_2 Q_{1s})^{1/2}$$

QUARTER-WAVE MATCHING SECTIONS

Fig. 17 shows how voltage-reflection coefficient and standing-wave ratio (swr) vary with frequency f when quarter-wave matching lines are inserted between a line of characteristic impedance Z_0 and a load of resistance R . The symbol f_0 represents the frequency for which the matching sections are exactly one-quarter wave-length ($\lambda/4$) long.

IMPEDANCE MATCHING WITH SHORTED STUB

The use of a shorted stub for impedance matching is illustrated in Fig. 18.

IMPEDANCE MATCHING WITH OPEN STUB

The use of an open stub for impedance matching is illustrated in Fig. 19.

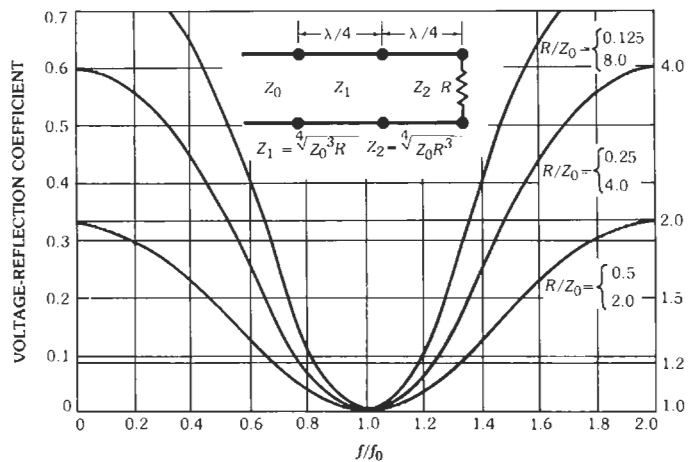
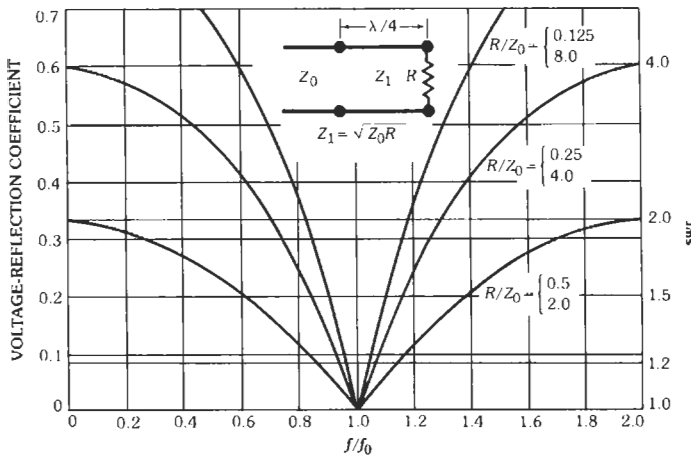
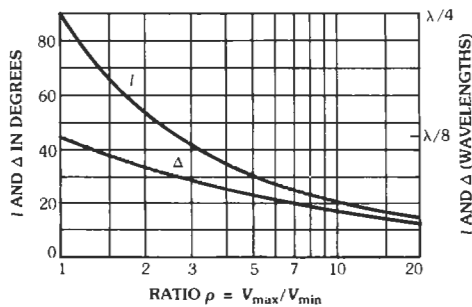


Fig. 17. Quarter-wave matching sections.



l = LENGTH OF SHORTED STUB
 Δ = LOCATION OF STUB MEASURED FROM V_{min} TOWARD LOAD

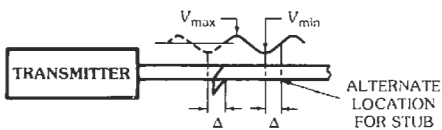
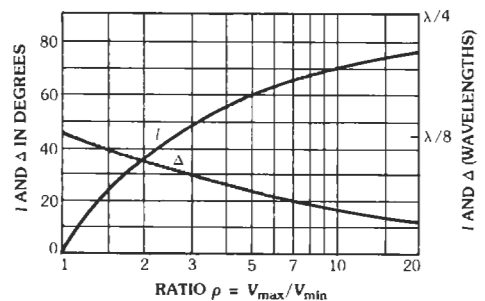


Fig. 18. Impedance matching with shorted stub.



l = LENGTH OF OPEN STUB
 Δ = LOCATION OF STUB MEASURED FROM V_{min} TOWARD TRANSMITTER

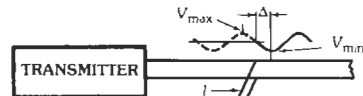


Fig. 19. Impedance matching with open stub.

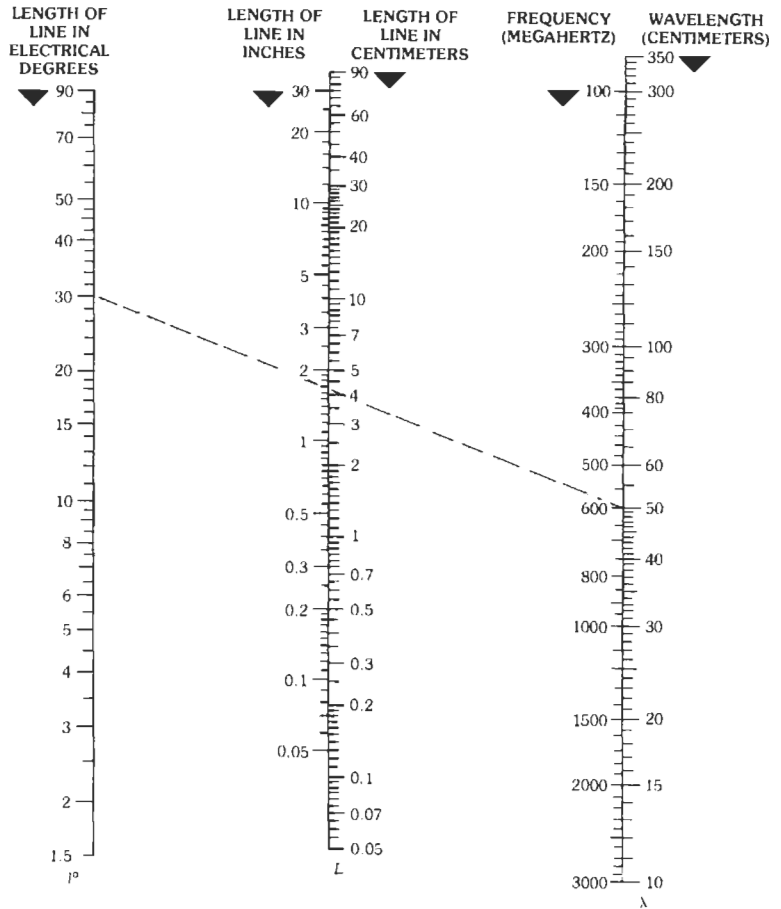


Fig. 20. Determination of length of transmission line.

LENGTH OF TRANSMISSION LINE

Fig. 20 relates the actual length of line in centimeters and inches to the length in electrical degrees and the frequency, provided that the velocity of propagation on the transmission line is equal to that in free space. The length is given on the *L*-scale intersection by a line between λ and l° , where

$$l^\circ = \frac{360L \text{ in centimeters}}{\lambda \text{ in centimeters}}$$

Example: $f = 600$ megahertz, $l^\circ = 30$, length $L = 1.64$ inches or 4.2 centimeters.

$$Z_0 = (138/\epsilon^{1/2}) \log_{10} (D/d)$$

$$= (60/\epsilon^{1/2}) \log_e (D/d)$$

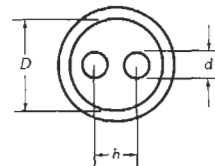
$\epsilon =$ dielectric constant
 $= 1$ in air

CHARACTERISTIC IMPEDANCE OF LINES

A. Single coaxial line (See also Fig. 21.)



B. Balanced shielded line

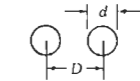
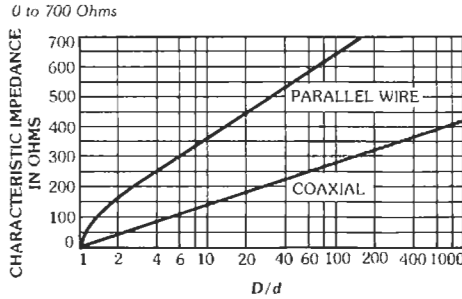
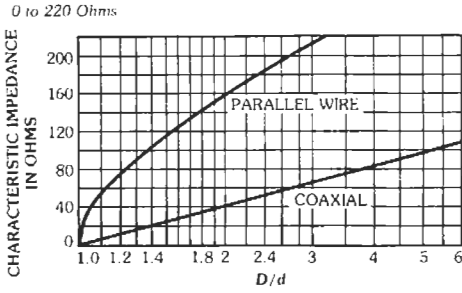


For $D \gg d, h \gg d$

$$Z_0 = (276/\epsilon^{1/2}) \log_{10} \{2v[(1 - \sigma^2)/(1 + \sigma^2)]\}$$

$$= (120/\epsilon^{1/2}) \log_e \{2v[(1 - \sigma^2)/(1 + \sigma^2)]\}$$

$$v = h/d \quad \sigma = h/D$$



Parallel Wires in Air

$$Z_0 = 120 \cosh^{-1} D/d$$

FOR $D \gg d$

$$Z_0 \approx 276 \log_{10} 2D/d$$



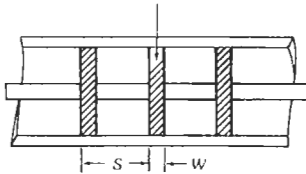
Coaxial

$$Z_0 = (138/\sqrt{\epsilon}) \log_{10} D/d$$

CURVE IS FOR $\epsilon = 1.00$

Fig. 21. Characteristic impedance of transmission lines.

C. Beads—dielectric ϵ_1

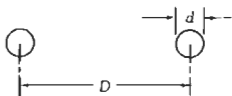


For lines A and B, if insulating beads are used at frequent intervals, call new characteristic impedance Z_0'

$$Z_0' = Z_0 \{1 + [(\epsilon_1/\epsilon) - 1](W/S)\}^{1/2}$$

$$W \ll S \ll \lambda/4$$

D. Open two-wire line in air (See also Fig. 21.)

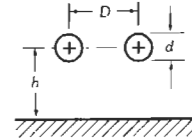


$$Z_0 = 120 \cosh^{-1}(D/d)$$

$$\approx 276 \log_{10}(2D/d)$$

$$\approx 120 \log_e(2D/d)$$

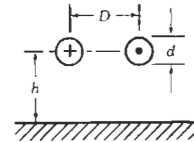
E. Wires in parallel, near ground



For $d \ll D, h$

$$Z_0 = (69/\epsilon^{1/2}) \log_{10}\{(4h/d)[1 + (2h/D)^2]^{1/2}\}$$

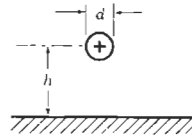
F. Balanced, near ground



For $d \ll D, h$

$$Z_0 = (276/\epsilon^{1/2}) \log_{10}\{(2D/d) [1 + (D/2h)^2]^{-1/2}\}$$

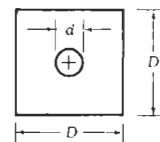
G. Single wire, near ground



For $d \ll h$

$$Z_0 = (138/\epsilon^{1/2}) \log_{10}(4h/d)$$

H. Single wire, square enclosure



$$Z_0 \approx [138 \log_{10} \rho + 6.48 - 2.34A - 0.48B - 0.12C] \epsilon^{-1/2}$$

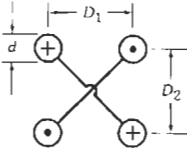
where $\rho = D/d$

$$A = (1 + 0.405\rho^{-4}) / (1 - 0.405\rho^{-4})$$

$$B = (1 + 0.163\rho^{-8}) / (1 - 0.163\rho^{-8})$$

$$C = (1 + 0.067\rho^{-12}) / (1 - 0.067\rho^{-12})$$

I. Balanced 4-wire



For $d \ll D_1, D_2$

$$Z_0 = (138/\epsilon^{1/2}) \log_{10}\{(2D_2/d)[1 + (D_2/D_1)^2]\}^{-1/2}$$

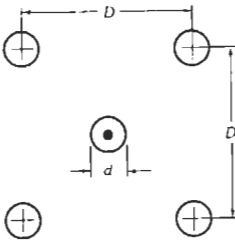
J. Parallel-strip line



$$wl < 0.1$$

$$Z_0 \approx 377(w/l)$$

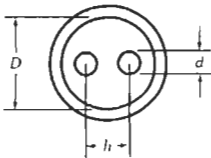
K. Five-wire line



For $d \ll D$

$$Z_0 = (173/\epsilon^{1/2}) \log_{10}(D/0.933d)$$

L. Wires in parallel—sheath return



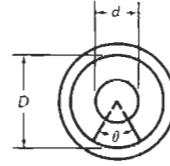
For $d \ll D, h$

$$Z_0 = (69/\epsilon^{1/2}) \log_{10}[(v/2\sigma^2)(1 - \sigma^4)]$$

$$\sigma = h/D$$

$$v = h/d$$

M. Air coaxial with dielectric supporting wedge

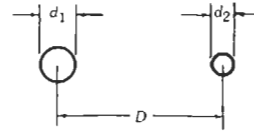


$$Z_0 \approx \frac{138 \log_{10}(D/d)}{[1 + (\epsilon - 1)(\theta/360)]^{1/2}}$$

ϵ = dielectric constant of wedge

θ = wedge angle in degrees

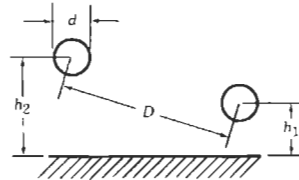
N. Balanced 2-wire—unequal diameters



$$Z_0 = (60/\epsilon^{1/2}) \cosh^{-1}N$$

$$N = \frac{1}{2}[(4D^2/d_1d_2) - (d_1/d_2) - (d_2/d_1)]$$

O. Balanced 2-wire near ground



For $d \ll D, h_1, h_2$

$$Z_0 = (276/\epsilon^{1/2}) \log_{10}\{(2D/d)[1 + (D^2/4h_1h_2)]\}^{-1/2}$$

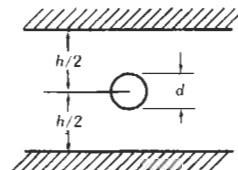
Holds also in either of the following special cases:

$$D = \pm(h_2 - h_1)$$

or

$$h_1 = h_2 \text{ (see F above)}$$

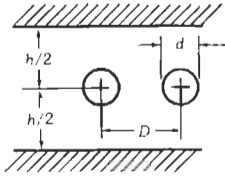
P. Single wire between grounded parallel planes—ground return



For $d/h < 0.75$

$$Z_0 = (138/\epsilon^{1/2}) \log_{10}(4h/\pi d)$$

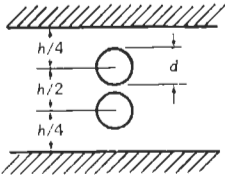
Q. Balanced line between grounded parallel planes



For $d \ll D, h$

$$Z_0 = (276/\epsilon^{1/2}) \log_{10} \left(\frac{4h \tanh(\pi D/2h)}{\pi d} \right)$$

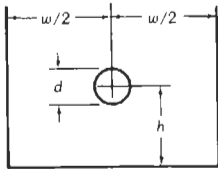
R. Balanced line between grounded parallel planes



For $d \ll h$

$$Z_0 = (276/\epsilon^{1/2}) \log_{10}(2h/\pi d)$$

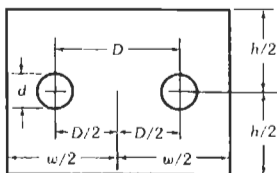
S. Single wire in trough



For $d \ll h, w$

$$Z_0 = (138/\epsilon^{1/2}) \log_{10} \left(\frac{4w \tanh(\pi h/w)}{\pi d} \right)$$

T. Balanced 2-wire line in rectangular enclosure



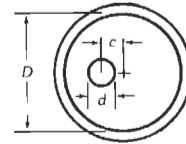
For $d \ll D, w, h$

$$Z_0 = (276/\epsilon^{1/2}) \left[\log_{10} \left(\frac{4h \tanh(\pi D/2h)}{\pi d} \right) - \sum_{m=1}^{\infty} \log_{10} \left(\frac{1 + u_m^2}{1 - v_m^2} \right) \right]$$

where

$$u_m = \frac{\sinh(\pi D/2h)}{\cosh(m\pi w/2h)} \quad v_m = \frac{\sinh(\pi D/2h)}{\sinh(m\pi w/2h)}$$

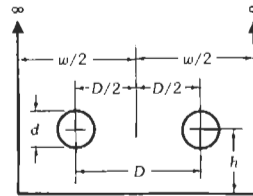
U. Eccentric line



$$Z_0 = (60/\epsilon^{1/2}) \cosh^{-1} U$$

$$U = \frac{1}{2} [(D/d) + (d/D) - (4c^2/dD)]$$

V. Balanced 2-wire line in semi-infinite enclosure



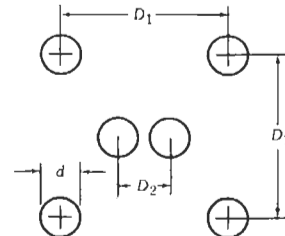
For $d \ll D, w, h$

$$Z_0 = (276/\epsilon^{1/2}) \log_{10}[2w/\pi d(A^{1/2})]$$

where

$$A = \operatorname{cosec}^2(\pi D/w) + \operatorname{cosech}^2(2\pi h/w)$$

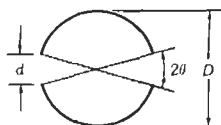
W. Outer wires grounded, inner wires balanced to ground



$$Z_0 \approx (276/\epsilon^{1/2}) \left\{ \log_{10}(2D_2/d) \right.$$

$$\left. - \left[\log_{10} \frac{1 + (1 + D_2/D_1)^2}{1 + (1 - D_2/D_1)^2} \right]^2 \left[\log_{10}(2D_1\sqrt{2}/d) \right]^{-1} \right\}$$

X. Split thin-walled cylinder



$$Z_0 \approx \frac{129}{\log_{10}[\cot^2 \frac{1}{2} \theta + (\cot^2 \frac{1}{2} \theta - 1)^{1/2}]}$$

For θ small:

$$Z_0 \approx 129 / \log_{10}(4D/d)$$

Courtesy of Electronic Engineering

Y. Slotted air line



When a slot is introduced into an air coaxial line for measuring purposes, the increase in characteristic im-

pedance in ohms, compared with a normal coaxial line, is less than

$$\Delta Z = 0.03 \theta^2$$

where θ is the angular opening of the slot in radians.

MICROSTRIP LINES

Microstrip line consists of a conductor strip placed on a dielectric substrate (relative dielectric constant ϵ_r), which is in turn backed by a conducting ground plane (Fig. 22). At lower microwave frequencies, the modal field is considered almost TEM. However, as the frequency is increased, the dispersion effect becomes more obvious, and the characteristic impedance and the phase velocity defined under the quasi-TEM analysis must be modified.

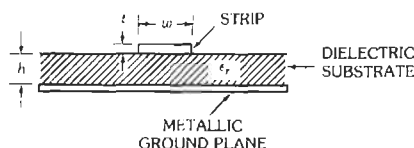


Fig. 22. Cross section of microstrip line.

Quasi-TEM Characteristics*

The characteristic impedance Z_0 and the effective dielectric constant ϵ_e are functions of structure and the dielectric constant ϵ_r .

$$Z_0(w, t, h, \epsilon_r) = Z_{01}(U_r) / [\epsilon_e(U_r, \epsilon_r)]^{1/2}$$

$$\beta = \frac{Z_{01}(U_r)}{Z_{01}(U_r)} [\epsilon_e(U_r, \epsilon_r)]^{1/2} \omega/c$$

where,

$$Z_{01}(x) = (376.73/2\pi) \ln \left\{ \frac{f(x)}{x} + [1 + (2/x)^2]^{1/2} \right\}$$

$$f(x) = 6 + (2\pi - 6) \exp \left[-(30.666/x)^{0.7528} \right]$$

$$\epsilon_e(x, \epsilon_r) = [(\epsilon_r + 1)/2] + [(\epsilon_r - 1)/2](1 + 10/x)^y$$

$$y = -a(x)b(\epsilon_r)$$

$$a(x) = 1 + (1/49) \ln \left[\frac{x^4 + (x/52)^2}{x^4 + 0.432} \right] + (1/18.7) \ln \left[1 + (x/18.1)^3 \right]$$

$$b(\epsilon_r) = 0.564 [(\epsilon_r - 0.9)/(\epsilon_r + 3)]^{0.053}$$

* These formulas are from E. Hammerstadt and O. Jensen, 1980 IEEE International Microwave Symposium Digest, pp. 407-409, June 1980. For the effects of an additional ground plane over the structure, see S. March, Microwaves, pp. 83-94, December 1981.

x is either

$$U_r = w/h + (t/\pi) \ln \left[1 + \frac{4 \exp(1)}{t \coth^2 (6.517 w/h)^{1/2}} \right] \left[1 + 1/\cosh (\epsilon_r - 1)^{1/2} \right] / 2$$

or

$$U_1 = w/h + (t/\pi) \ln \left[1 + \frac{4 \exp(1)}{t \coth^2 (6.517 w/h)^{1/2}} \right]$$

depending on which is called for in the above equations. Note that $t \rightarrow 0$, both U_r and U_1 approach w/h .

Fig. 23 shows Z_0 as a function of w/h .

Attenuation

Dielectric loss (dB/unit length)

$$\alpha_D = 4.34 \beta \tan \delta / [1 + (1/\epsilon_r)(F - 1)/(F + 1)]$$

$$F = (1 + 10h/w)^{1/2}$$

The factor $\tan \delta$ is the loss tangent of the substrate material. ** The formula for β is given in the preceding subsection, "Quasi-TEM Characteristics."

Conductor loss is†

** M. V. Schneider, "Dielectric Loss in Hybrid Integrated Circuits," *Proc. IEEE*, Vol. 57, pp. 1206-1207, 1969.

† R. A. Pucel, D. J. Massé, and C. P. Hartwig, "Losses in Microstrip," *IEEE Trans. Microwave Theory and Techniques*, Vol. MTT-16, pp. 342-350, 1968.

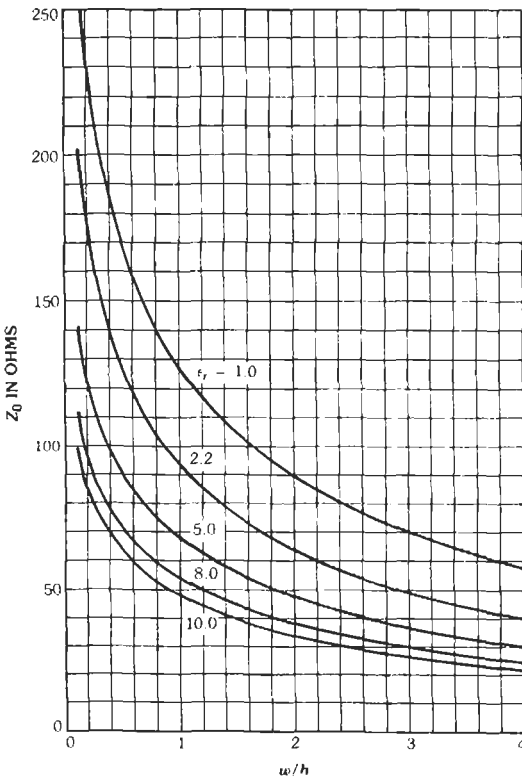


Fig. 23. Characteristic impedance of microstrip line. Calculated from quasi-tem formulas by Hammerstadt and Jensen. (Curves assume t is negligible; i.e., $U_r = U_e = w/h$.)

$$\alpha_c = (4.34R_s/\pi h Z_0)[1 - (w'/4h)^2] \{1 + h/w' + (h/\pi w') [\ln(4\pi w/t + 1) - (1 - t/w)/(1 + t/4\pi w)]\}$$

$w/h \leq 1/2\pi$

$$= (4.34R_s/\pi h Z_0) [1 - (w'/4h)^2] \{1 + h/w' + (h/\pi w')[\ln(2h/t + 1) - (1 + t/h)/(1 + t/2h)]\}$$

$1/2\pi < w/h \leq 2$

$$= \frac{8.68 R_s \left[\frac{w'}{h} + \frac{w'/\pi h}{w'/2h + 0.94} \right]}{h Z_0 \{w'/h + (2/\pi) \ln [2\pi\epsilon (w'/2h + 0.94)]\}^2} \{1 + h/w' + (h/\pi w')[\ln(2h/t + 1) - (1 + t/h)/(1 + t/2h)]\}$$

$w/h > 2$

where

$$w' = w + \Delta w$$

and for $2t/h < w/h, 1/2\pi$

$$\Delta w = (t/\pi) \ln(4\pi w/t + 1)$$

$w/h \leq 1/2\pi$

$$= (t/\pi) \ln(2h/t + 1)$$

$w/h \geq 1/2\pi$

The quantity Z_0 is the characteristic impedance discussed before, and R_s is the surface resistivity. Typically, $R_s = 2.61 \times 10^{-7} \sqrt{f}$ for copper.

Frequency-Dependent Characteristics

As the microstrip mode is not purely TEM, both Z_0 and ϵ_e are functions of frequency.[‡]

$$Z_0(f) = Z_0(0)[\epsilon_e(0)/\epsilon_e(f)]^{1/2} [\epsilon_e(f) - 1]/[\epsilon_e(0) - 1]$$

$$\epsilon_e(f) = \epsilon_r - |\epsilon_r - \epsilon_e(0)|/[1 + G(f/f_p)^2]$$

where,

$$f_p = Z_0(0)/(2\mu_0 h)$$

μ_0 = free-space permeability.

$$G = (\pi^2/12)[(\epsilon_r - 1)/\epsilon_e(0)][Z_0(0)/60]^{1/2}$$

In the above, $Z_0(0)$ and $\epsilon_e(0)$ are dc values and are obtained from the formulas in "Quasi-TEM Characteristics."

Power-Handling Capacity

For a microstrip line composed of a strip $7/32$ -inch wide on a Teflon-impregnated fiberglass base $1/16$ -inch thick:

(A) At 3000 megahertz with 300 watts cw, the temperature under the strip conductor has been measured at 50° Celsius rise above 20° Celsius ambient.

(B) Under pulse conditions, corona effects appear at the edge of the strip conductor for pulse power of roughly 10 kilowatts at 9000 megahertz.

STRIP TRANSMISSION LINES*

Strip transmission lines differ from microstrip in that a second ground plane is placed above the strip and the space between the two ground planes is filled completely with a homogeneous dielectric (Fig. 24). The characteristic impedance is shown in Fig. 25 and the attenuation in Fig. 26.

Dielectric loss in decibels/unit length

$$\alpha_d = 27.3F_p(\epsilon_r)^{1/2}/\lambda_0$$

where,

λ_0 = free-space wavelength,
 F_p = power factor or loss angle.

Conductor loss in decibels/unit length

$$\alpha_c = (y/b) (f_{GHz}\epsilon_r\mu_r\rho/\rho_{Cu})^{1/2}$$

* S. B. Cohn, "Problems in Strip Transmission Lines," *Transactions of the IRE Professional Group on Microwave Theory and Techniques*, Vol. M1T3, pp. 119-126, March 1955. Other papers on strip-type lines also appear in that issue of the journal.

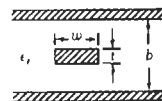


Fig. 24. Cross section of a strip transmission line.

‡ E. Hammerstadt and O. Jensen, *1980 IEEE International Microwave Symposium Digest*, pp. 407-409, June 1980.

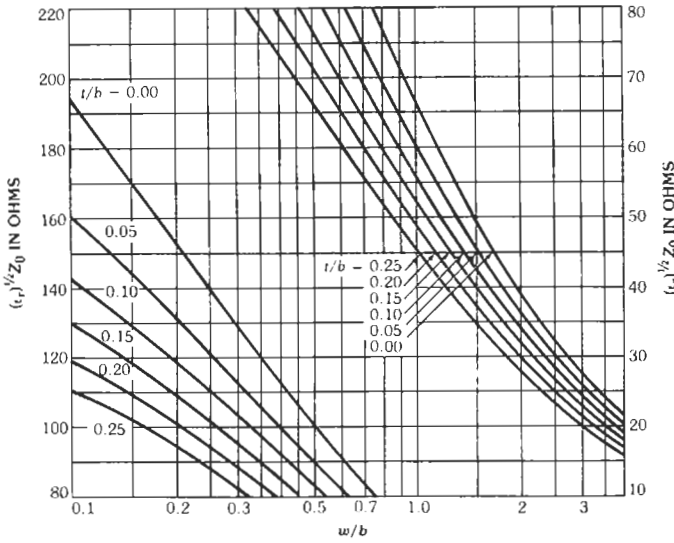


Fig. 25. Plot of strip-transmission-line Z_0 versus w/b for various values of t/b . For lower-left family of curves, refer to left-hand ordinate values; for upper-right curves, use right-hand scale. (Courtesy of Transactions of the IRE Professional Group on Microwave Theory and Techniques.)

where,

- y = ordinate from Fig. 26,
- ρ/ρ_{Cu} = resistivity relative to copper.

The unit of length in α_d is that of λ_0 , and in α_c it is that of b .

$$K_r = (1/\pi) \ln [2(1 + \sqrt{k})/(1 - \sqrt{k})]$$

$$0.707 \leq k \leq 1$$

$$K_r = \pi \ln [2(1 + \sqrt{k'})/(1 - \sqrt{k'})]$$

$$0 \leq k \leq 0.707$$

COPLANAR TRANSMISSION LINES

A quasi-TEM propagation takes place in a coplanar transmission line (Fig. 27). When $t = 0$, the phase constant and the characteristic impedance under the quasi-TEM approximation are**

$$\beta = \sqrt{\epsilon_c} \omega/c$$

$$Z_0 = (30\pi/\sqrt{\epsilon_c}) K_r$$

** An approximation for a ratio of complete elliptic integrals by W. Helberg, "From Approximation to Exact Relations for Characteristic Impedances," *IEEE Trans. Microwave Theory and Technique*, Vol. MTT-17, pp. 259-265, 1969, is used in the original expression for Z_0 by C.P. Wen, "Coplanar Waveguide: A Surface Strip Transmission Line Suitable for Non-Reciprocal Gyromagnetic Device Applications," *IEEE Trans. Microwave Theory and Techniques*, Vol. MTT-17, pp. 1087-1090, 1969.

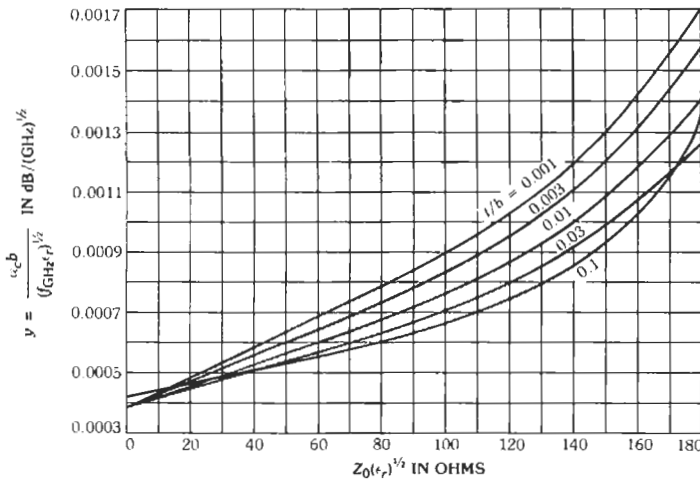


Fig. 26. Theoretical attenuation of copper-shielded strip transmission line in dielectric medium ϵ_r . (Courtesy of Transactions of the IRE Professional Group on Microwave Theory and Techniques.)

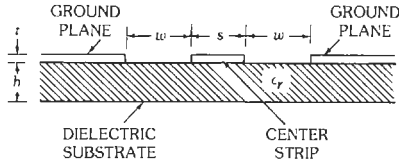


Fig. 27. Cross section of coplanar transmission line.

$$k = s/(2w + s)$$

$$k' = (1 - k^2)^{1/2}$$

where for the range of practical interest $1 \leq h/w < \infty$. The value of ϵ_e is given by†

$$\epsilon_e = [(\epsilon_r + 1)/2] \{ \tanh[0.775 \ln(h/w) + 1/75] + k(w/h)[0.04 - 0.7k + 0.01(1 - 0.1\epsilon_r)(0.25 + k)] \}$$

If t is not negligible, k and ϵ_e in the expressions for $t = 0$ must be replaced with†

$$k^e = k + (1 - k^2) (1.25t/2w\pi)[1 + \ln(4\pi s/t)]$$

$$\epsilon_e^t = \epsilon_e - 0.7(\epsilon_e - 1) (t/w)/(K_r + 0.7t/w)$$

Fig. 28 is a graph of characteristic impedance calculated from the formulas in this section for $t = 0$.

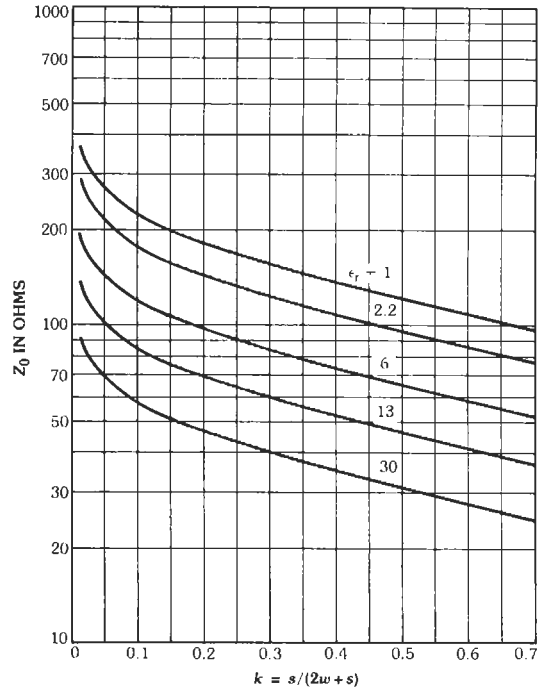


Fig. 28. Graph showing characteristic impedance of coplanar transmission line.

ATTENUATION AND POWER RATING OF LINES AND CABLES

Attenuation

Fig. 29 illustrates the attenuation of general-purpose radio-frequency lines and cables up to their practical upper frequency limit. Most of these are coaxial-type lines, but waveguide and microstrip are included for comparison.

The following notes are applicable to this figure.

(A) For the RG-type cables, only the number is given (for instance, the curve for RG-218/U is labeled 218. Refer to the table of radio-frequency cables.) The data on RG-type cables are taken mostly from "RF Transmission Lines and Fittings," MIL-HDBK-216, 4 January 1962, revised 18 May 1965, and from "Solid Dielectric Transmission Lines," Electronic Industries Association Standard RS-199, December 1957.

Some approximation is involved in order to simplify the figure. Thus, where a single curve is labeled with several type numbers, the actual attenuation of each individual type may be slightly different from that shown by the curve.

(B) The curves for rigid copper coaxial lines are labeled with the diameter of the line only, as $7/8"$ C. These have been computed for the lines listed in "Rigid Coaxial Transmission Lines, 50 Ohms," Electronic Industries Association Standard RS-225, August 1959. The computations considered the copper losses only, on the basis of a resistivity $\rho = 1.724$ microhm-centimeters; a derating of 20 percent has been applied to allow for imperfect surface, presence of fittings, etc., in long installed lengths. Relative attenuations of the different sizes are

$$A_{6\frac{1}{8}"} \approx 0.13A_{7\frac{1}{8}"}''$$

$$A_{3\frac{1}{8}"} \approx 0.26A_{7\frac{1}{8}"}''$$

$$A_{1\frac{1}{8}"} \approx 0.51A_{7\frac{1}{8}"}''$$

(C) Typical curves are shown for three sizes of 50-ohm semirigid cables such as Styroflex, Spirolinc, Heliac, Alumispline, etc. These are labeled by size in inches, as $7/8"$ S.

(D) The microstrip curve is for Teflon-impregnated fiberglass dielectric $1/16$ -inch thick and conductor strip $7/32$ -inch wide.

(E) Shown for comparison is the attenuation in the $TE_{1,0}$ mode of five sizes of brass waveguide. The resistivity of brass was taken as $\rho = 6.9$ microhm-

† K. C. Gupta, et. al., *Microstrip Lines and Slotlines*, (Dedham, Mass.: Artech House, 1969).

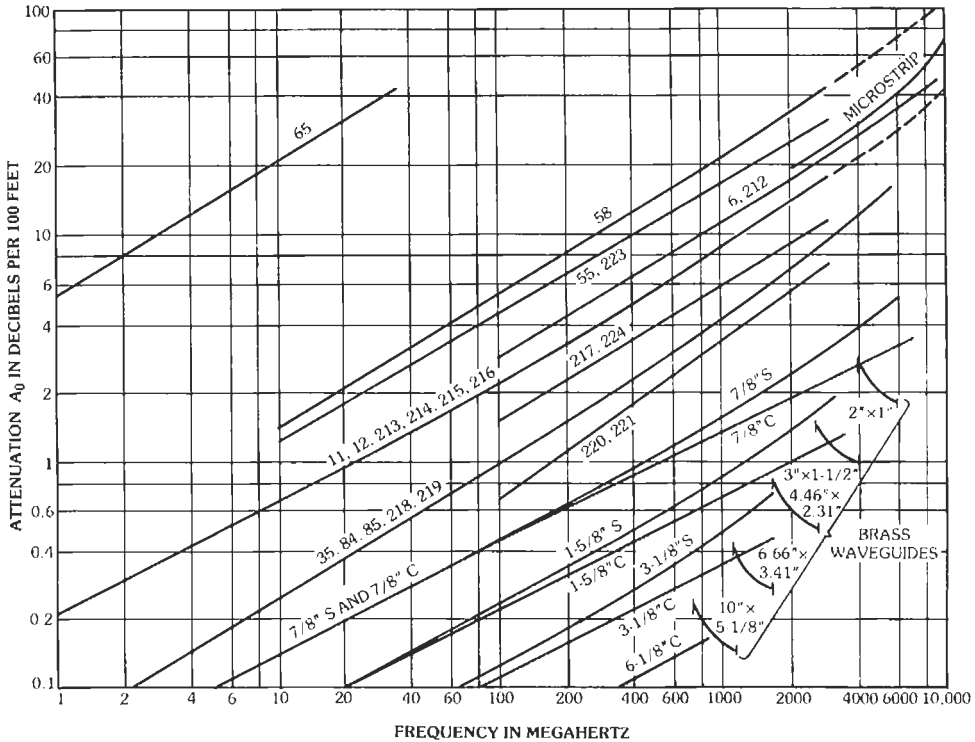


Fig. 29. Attenuation of cables.

centimeters, and no derating was applied. For copper or silver, attenuation is about half that for brass. For aluminum, attenuation is about two-thirds that for brass.

Power Rating

Fig. 30 shows the approximate power transmitting capabilities of various coaxial-type lines. The following notes are applicable.

(A) Identification of the curves for the RG-type cables is the same as in Fig. 29. The data for these cables are from the same sources. For polyethylene cables, an inner-conductor maximum temperature of 80 degrees Celsius is specified (G). For high-temperature cables (types 211, 228; 225, 227) the inner-conductor temperature is 250 degrees Celsius.

(B) The curves for 50-ohm rigid coaxial line are labeled with the diameter of the line only, as 1/8" C.

ARMY-NAVY LIST OF PREFERRED

Class of Cables		JAN Type RG-	Inner Conductor [‡]	Dielectric Material (Note 1)	Nominal Diameter of Dielectric (in.)	Shielding Braid
50 ohms	Single braid	58C/U	19/0.0071" tinned copper	A	0.116	Tinned copper
		213/U	7/0.0296" copper	A	0.285	Copper
		215/U	7/0.0296" copper	A	0.285	Copper
		218/U	0.195" copper	A	0.680	Copper

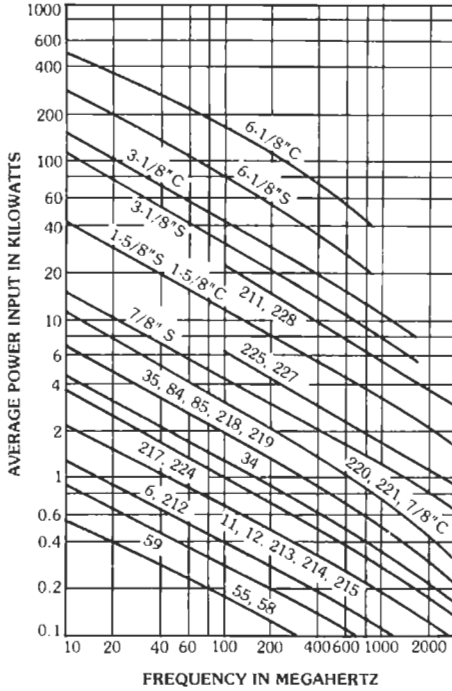


Fig. 30. Power rating of cables.

These are rough estimates based largely on miscellaneous charts published in catalogs.

(C) For Styroflex, Spiroline, Heliac, Alumipline, etc., cables, refer to (C) above.

(D) The curves are for unity voltage standing-wave ratio. Safe operating power is inversely proportional to swr expressed as a numerical ratio greater than unity. Do not exceed maximum operating voltage (see table of radio-frequency cables).

TABLE 1. DERATING FACTOR

Ambient Temperature (°C)	Maximum Allowable Center-Conductor Temperature (°C)			
	80	75	70	65
40	1.0	0.86	0.72	0.59
50	0.72	0.59	0.46	0.33
60	0.46	0.33	0.22	0.10
70	0.20	0.09	0	—
80	0	—	—	—

(E) An ambient temperature of 40 degrees Celsius is assumed.

(F) The four curves meeting the 100-watt ordinate may be extrapolated: at 3000 megahertz for 55, 58, power is 28 watts; for 59, power is 44 watts; and for 6, 212, power is 58 watts.

(G) Electronic Industries Association Standard RS-199 states that operation of a polyethylene dielectric cable at a center-conductor temperature in excess of 80 degrees Celsius is likely to cause permanent damage to the cable. Where practicable, and particularly where continuous flexing is required, it is recommended that a cable be selected which, in regular operation, will produce a center-conductor temperature not greater than 65 degrees Celsius. Rating factors for various operating temperatures are given in Table 1. Multiply points on the power-rating curve by the factors in the table to determine power rating at operating rcsct above to add one line conditions.

RADIO-FREQUENCY CABLES*

Protective Covering (Note 2)	Nominal Overall Diameter (in.)	Weight (lb/ft)	Nominal Capacitance (pF/ft)	Maximum Operating Voltage (rms)	Remarks
IIa	0.195	0.029	28.5	1900	Small-size flexible cable
IIa	0.405	0.120	29.5	5000	Medium-size flexible cable (formerly RG-8A/U)
IIa, with armor	0.475 max.	0.160	29.5	5000	Same as RG-213/U, but with armor (formerly RG-10A/U)
IIa	0.870	0.491	29.5	11,000	Large-size low-attenuation high-power transmission line (formerly RG-17A/U)

Continued on next page.

ARMY-NAVY LIST OF PREFERRED

Class of Cables	JAN Type RG-	Inner Conductor†	Dielectric Material (Note 1)	Nominal Diameter of Dielectric (in.)	Shielding Braid		
50 ohms (Continued)	219/U	0.195" copper	A	0.680	Copper		
	220/U	0.260" copper	A	0.910	Copper		
	221/U	0.260" copper	A	0.910	Copper		
Double braid	55B/U	0.032" silvered copper	A	0.116	Silvered copper		
	212/U	0.0556" silvered copper	A	0.185	Silvered copper		
	214/U	7/0.0296" silvered copper	A	0.285	Silvered copper		
	217/U	0.106" copper	A	0.370	Copper		
	223/U	0.035" silvered copper	A	0.116	Silvered copper		
	224/U	0.106" copper	A	0.370	Copper		
	75 ohms	Single braid	11A/U	7/26 AWG tinned copper	A	0.285	Copper
12A/U			7/26 AWG tinned copper	A	0.285	Copper	
34B/U			7/0.0249" copper	A	0.460	Copper	
35B/U			0.1045" copper	A	0.680	Copper	
59B/U			0.0230" copper-covered steel	A	0.146	Copper	
84A/U			0.1045" copper	A	0.680	Copper	
85A/U			0.1045" copper	A	0.680	Copper	
164/U			0.1045" copper	A	0.680	Copper	
307A/U			17/0.0058" silvered copper	A Foamed	0.029	Silvered copper	
Double braid			6A/U	21 AWG copper-covered steel	A	0.185	Inner: silvered copper. Outer: copper
			216/U	7/0.0159" tinned copper	A	0.285	Copper

RADIO-FREQUENCY CABLES—CONT.

Class of Cables	JAN Type RG-	Inner Conductor†	Dielectric Material (Note 1)	Nominal Diameter of Dielectric (in.)	Shielding Braid
IIa, with armor	0.945 max.	0.603	29.5	11,000	Same as RG-218/U, but with armor (formerly RG-18A/U)
IIa	1.120	0.745	29.5	14,000	Very-large low-attenuation high-power transmission cable (formerly RG- 19A/U)
IIa, with armor	1.195 max.	0.925	29.5	14,000	Same as RG-220/U, but with armor (formerly RG-20A/U)
IIIa	0.206	0.032	28.5	1900	Small-size flexible cable
IIa	0.332	0.093	28.5	3000	Small-size microwave cable (formerly RG-5B/U)
IIa	0.425	0.158	30.0	5000	Special medium-size flexible cable (for- merly RG-9B/U)
IIa	0.545	0.236	29.5	7000	Medium-size power transmission line (formerly RG-14A/U)
IIa	0.216	0.036	28.5	1900	Small-size flexible cable (formerly RG- 55A/U)
IIa, with armor	0.615 max.	0.282	29.5	7000	Same as RG-217/U, but with armor (formerly RG-74A/U)
IIa	0.412	—	20.5	5000	Medium-size flexible video cable
IIa, with armor	0.475	—	20.5	5000	Similar to RG-11A/U, but with armor
IIa	0.630	0.231	21.5	6500	Large-size high-power low-attenuation flexible cable
IIa, with armor	0.945 max.	0.480	21.5	10,000	Large-size high-power low-attenuation video and communication cable
IIa	0.242	—	21.0	2300	General-purpose small-size video cable
IIa, with lead sheath	1.000	1.325	21.5	10,000	Same as RG-35B/U, except lead sheath instead of armor for underground in- stallations
IIa, with lead sheath and special armor	1.565 max.	2.910	21.5	10,000	Same as RG-84A/U, with special armor for underground installations
IIa	0.870	0.490	21.5	10,000	Same as RG-35B/U, except without armor
IIIa	0.270	—	20	400	
IIa	0.332	—	20.0	2700	Small-size video and communication cable
IIa	0.425	0.121	20.5	5000	Medium-size flexible video and com- munication cable (formerly RG- 13A/U)

Continued on next page.

ARMY-NAVY LIST OF PREFERRED

Class of Cables		JAN Type RG-	Inner Conductor†	Dielectric Material (Note 1)	Nominal Diameter of Dielectric (in.)	Shielding Braid	
High tem- perature	Single braid	144/U	7/0.0179" silvered copper- covered steel	F1	0.285	Silvered copper	
		178B/U	7/0.004" silvered copper- covered steel	F1	0.034	Silvered copper	
		179B/U	Same as above	F1	0.063	Silvered copper	
		180B/U	Same as above	F1	0.102	Silvered copper	
		187A/U	7/0.004" annealed silvered copper-covered steel	F1	0.060	Silvered copper	
		195A/U	Same as RG-187A/U	F1	0.102	Silvered copper	
		196A/U	Same as RG-187A/U	F1	0.034	Silvered copper	
		211A/U	0.190" copper	F1	0.620	Copper	
		228A/U	0.190" copper	F1	0.620	Copper	
		302/U	0.025" silvered copper-covered steel	F1	0.146	Silvered copper	
		303/U	0.039" silvered copper-covered steel	F1	0.116	Silvered copper	
		304/U	0.059" silvered copper-covered steel	F1	0.185	Silvered copper	
		316/U	7/0.0067" annealed silvered copper-covered steel	F1	0.060	Silvered copper	
		Double braid	115/U	7/0.028" silvered copper	F2	0.250	Silvered copper
			142B/U	0.039" silvered copper-covered steel	F1	0.116	Silvered copper
225/U	7/0.0312" silvered copper		F1	0.285	Silvered copper		
226/U	19/0.0254" silvered copper wire		F2	0.370	Copper		
227/U	7/0.0312" silvered copper		F1	0.285	Silvered copper		
Pulse	Single braid	26A/U	19/0.0117" tinned copper	E	0.288	Tinned copper	

RADIO-FREQUENCY CABLES—CONT.

Protective Covering (Note 2)	Nominal Overall Diameter (in.)	Weight (lb/ft)	Nominal Capaci- tance (pF/ft)	Maximum Operating Voltage (rms)	Remarks
Teflon-tape moisture seal with double-braid type-V jacket	0.410	0.120	20.5	5000	Similar to RG-11A/U, except cable core is Teflon. $Z = 75$ ohms
IX	0.075 max.	—	29.0	1000	$Z = 50$ ohms
IX	0.105	—	20.0	1200	
IX	0.145	—	15.5	1500	$Z = 95$ ohms
VII	0.110	—	—	1200	Miniaturized cable. $Z = 75$ ohms
VII	0.155	—	—	1500	Miniaturized cable. $Z = 95$ ohms
VII	0.080	—	—	1000	Miniaturized cable. $Z = 50$ ohms
Same as RG-144/U	0.730	0.450	29.0	7000	Semiflexible cable operating at -55 °C to $+200$ °C (formerly RG-117A/U). $Z = 50$ ohms
Teflon-tape moisture seal with double-braid type-V jacket, with armor	0.795	0.600	29.0	7000	Same as RG-211A/U, but with armor (formerly RG-118A/U). $Z = 50$ ohms
IX	0.206	—	21.0	2300	$Z = 75$ ohms
IX	0.170	—	28.5	1900	$Z = 50$ ohms
IX	0.280	—	28.5	3000	$Z = 50$ ohms
IX	0.102	—	—	1200	Miniaturized cable. $Z = 50$ ohms
Same as RG-144/U	0.375	—	29.5	5000	Medium-size cable for use where expansion and contraction are a major problem. $Z = 50$ ohms
IX	0.195	—	28.5	1900	Small-size flexible cable. $Z = 50$ ohms
Same as RG-144/U	0.430	0.176	29.5	5000	Semiflexible cable operating at -55 °C to $+200$ °C (formerly RG-87A/U). $Z = 50$ ohms
Same as RG-144/U	0.500	0.247	29.0	7000	Medium-size cable for use where expansion and contraction are a major problem (formerly RG-94A/U). $Z = 50$ ohms
Same as RG-228A/U	0.490	0.224	29.5	5000	Same as RG-225/U, but with armor (formerly RG-116/U). $Z = 50$ ohms
IV, with armor	0.505	0.189	50.0	10,000	High-voltage cable. $Z = 48$ ohms

Continued on next page.

ARMY-NAVY LIST OF PREFERRED

Class of Cables		JAN Type RG-	Inner Conductor†	Dielectric Material (Note 1)	Nominal Diameter of Dielectric (in.)	Shielding Braid
Pulse (Continued)	Single braid	27A/U	19/0.0185" tinned copper	D	0.455	Tinned copper
	Double braid	25A/U	19/0.0117" tinned copper	E	0.288	Tinned copper
		28B/U	19/0.0185" tinned copper	D	0.455	Inner: tinned copper. Outer: galvanized steel
		64A/U	19/0.0117" tinned copper	E	0.288	Tinned copper
		156/U	7/21 AWG tinned copper	First layer A; second layer H	0.285	Inner: tinned copper
		157/U	19/24 AWG tinned copper	First layer H; second layer A;	0.455	Outer: galvanized steel.
		158/U	37/21 AWG tinned copper	third layer H	0.455	Tinned copper outer shield
		190/U	19/0.0117" tinned copper		0.380	Same as above
		191/U	30 AWG tinned copper; single braid over supporting elements; 0.485" max.	First layer H; second layer J; third layer H	1.065	Same as above
	Four braids	88/U	19/0.0117" tinned copper	E	0.288	Tinned copper
Low capacitance	Single braid	62A/U	0.0253" solid copper-covered steel	A	0.146	Copper
		63B/U	0.0253" copper-covered steel	A	0.285	Copper
		79B/U	0.0253" copper-covered steel	A	0.285	Copper
	Double braid	71B/U	0.0253" copper-covered steel	A	0.146	Tinned copper
High attenuation	Single braid	301/U	7/0.0203" Karma wire	F1	0.185	Karma wire
High delay	Single braid	65A/U	No. 32 Formex F. Helix diameter 0.128"	A	0.285	Copper
Twin conductor	Single braid	57A/U	Each conductor 7/0.0285" plain copper	A	0.472	Tinned copper
		130/U		A	0.472	Tinned copper
		131/U		A	0.472	Tinned copper

RADIO-FREQUENCY CABLES—CONT.

Protective Covering (Note 2)	Nominal Overall Diameter (in.)	Weight (lb/ft)	Nominal Capaci- tance (pF/ft)	Maximum Operating Voltage (rms)	Remarks
IV, with armor	0.670	0.304	50.0	15,000 peak	Large-size cable. $Z = 48$ ohms
IV	0.505	0.205	50.0	10,000	High-voltage cable. $Z = 48$ ohms
IV	0.750	0.370	50.0	15,000 peak	Large-size cable. $Z = 48$ ohms
IV	0.475 max.	0.205	50.0	10,000	Medium-size cable. $Z = 48$ ohms
IIa	0.540	0.211	30.0	10,000	Taped inner layers, first layer type K and second layer type A-IR, between the outer braid of the outer conductor and the tinned copper shield Triaxial pulse cables. $Z = 50$ ohms
IIa	0.725	0.317	38.0	15,000	
IIa	0.725	0.380	78.0	15,000	Same as above, except $Z = 25$ ohms
VIII over one wrap of type K	0.700	0.353	50.0	15,000	Taped inner layers, two wraps of type K and two wraps of type L between the outer braid and the tinned copper shield. Pulse cable. $Z = 50$ ohms
Same as above	1.460	1.469	85.0	15,000	Same as RG-190/U, except $Z = 25$ ohms
IIa	0.515	—	50.0	10,000	Medium-size multishielded high-voltage cable. $Z = 48$ ohms
IIa	0.242	0.382	14.5	750	$Z = 93$ ohms
IIa	0.405	0.082	10.0	1000	Medium-size low-capacitance air-spaced cable. $Z = 125$ ohms
IIa, with armor	0.475 max.	0.138	10.0	1000	Same as RG-63B/U, but with armor. $Z = 125$ ohms
IIIa	0.250 max.	—	14.5	750	Low-capacitance cable. $Z = 93$ ohms
IX	0.245	—	29.0	3000	High-attenuation cable. $Z = 50$ ohms
IIa	0.405	0.096	44.0	1000	High-impedance video cable; high-delay line. $Z = 950$ ohms. (Refer to Note 3.)
IIa	0.625	0.225	17.0	3000	$Z = 95$ ohms
I	0.625	0.220	17.0	8000	Same as RG-57A/U, except inner conductors are twisted to improve flexibility. $Z = 95$ ohms
I, with aluminum armor	0.710	0.295	17.0	8000	Same as RG-130/U, but with armor. $Z = 95$ ohms

Continued on next page.

ARMY-NAVY LIST OF PREFERRED

Class of Cables		JAN Type RG-	Inner Conductor†	Dielectric Material (Note 1)	Nominal Diameter of Dielectric (in.)	Shielding Braid
Twin con- ductor (Continued)	Double braid	22B/U	Each conductor 7/0.0152" copper	A	0.285	Tinned copper
		111A/U		A	0.285	Tinned copper
	Twin coaxial	181/U	Each conductor 7/26 AWG copper	A	0.210	Copper inner braids and common braid

* From "RF Transmission Lines and Fittings," MIL-HDBK-216, 4 January 1962, revised 18 May 1965. Requirements for listed cables are in Specification MIL-C-17.

† Diameter of strands given in inches, as, 7/0.0296" = 7 strands, each 0.0296-inch diameter.

Note 1-Dielectric materials: A = Polyethylene, D = Layer of synthetic rubber between two layers of conducting rubber, E = Layer of conducting rubber plus two layers of synthetic rubber, F1 = Solid polytetrafluoroethylene (Teflon), F2 = Semisolid or taped polytetrafluoroethylene (Teflon), H = Conducting synthetic rubber, and J = Insulating butyl rubber.

RADIO-FREQUENCY CABLES—CONT.

Protective Covering (Note 2)	Nominal Overall Diameter (in.)	Weight (lb/ft)	Nominal Capaci- tance (pF/ft)	Maximum Operating Voltage (rms)	Remarks
IIa	0.420	0.116	16.0	1000	Small-size balanced twin-conductor cable. $Z = 95$ ohms
IIa, with armor	0.490 max.	0.145	16.0	1000	Same as RG-22B/U, but with armor. $Z = 95$ ohms
IIa	0.640	—	12	3500	Filled-to-round, unbalanced transmission cable. Twin coaxial. $Z = 125$ ohms

Note 2—Jacket types: I = Polyvinyl chloride (colored black), IIa = Noncontaminating synthetic resin, IIIa = Noncontaminating synthetic resin (colored black), IV = Chloroprene, V = Fiberglass, silicone-impregnated varnish, VII = Polytetrafluoroethylene, VIII = Polychloroprene, and IX = Fluorinated ethylene propylene.

Note 3—For RG-65A/U, delay = 0.042 microsecond per foot at 5 megahertz; dc resistance = 7.0 ohms/foot.

30 Waveguides and Resonators

*Revised by
Tatsuo Itoh*

Propagation of Electromagnetic Waves in Hollow Waveguides	30-3
Rectangular Waveguides	30-3
Circular Waveguides	30-5
Square Waveguides	30-7
Attenuation in a Waveguide Beyond Cutoff	30-7
Standard Waveguides	30-7
Ridged Waveguides	30-8
Flexible Waveguides	30-8
Waveguide Losses	30-14
Waveguide Material and Modes	
Waveguide Flange Leakage	
Flange Resistance and Bolt Torque	
Flange Insertion Loss	
Losses and Noise Temperature	
Waveguide Circuit Elements	30-18
Hybrid Junctions	30-19
Resonant Cavities	30-20
Rectangular Cavity of Dimensions a , b , $2h$	
Cylindrical Cavities of Radius a and Length $2h$	

Spherical Resonators of Radius a
Additional Cavity Equations
Effect of Temperature and Humidity on Cavity Tuning
Coupling to Cavities and Loaded Q
Equations for Coupling Through a Cavity
Cavity Coupling Techniques
Simple Waveguide Cavity
Resonant Irises

Surface-Wave Transmission Line 30-24

 Symbols

 Dielectric Other Than Polyethylene

Dielectric-Rod Waveguides 30-26

Rectangular Dielectric Guides and Image Guides 30-27

Dielectric Resonators 30-28

Slot Lines 30-29

Fin-Lines 30-31

PROPAGATION OF ELECTROMAGNETIC WAVES IN HOLLOW WAVEGUIDES

For propagation of energy at microwave frequencies through a hollow metal tube under fixed conditions, the following different types of waves are available.

TE Waves: Transverse-electric waves, sometimes called *H waves*, characterized by the fact that the electric vector (*E* vector) is always perpendicular to the direction of propagation. This means that

$$E_z \equiv 0$$

where *z* is the direction of propagation.

TM Waves: Transverse-magnetic waves, also called *E waves*, characterized by the fact that the magnetic vector (*H* vector) is always perpendicular to the direction of propagation. This means that

$$H_z \equiv 0$$

where *z* is the direction of propagation.

Note—TEM Waves: Transverse-electromagnetic waves. These waves are characterized by the fact that both the electric vector (*E* vector) and the magnetic vector (*H* vector) are perpendicular to the direction of propagation. This means that

$$E_z \equiv H_z \equiv 0$$

where *z* is the direction of propagation. This is the mode commonly excited in coaxial and open-wire lines. It cannot be propagated in a waveguide.

The solutions for the field configurations in waveguides are characterized by the presence of the integers *m* and *n*, which can take on separate values from 0 or 1 to infinity. Only a limited number of these different *m, n* modes can be propagated, depending on the dimensions of the guide and the frequency of excitation. For each mode there is a definite lower limit or cutoff frequency below which the wave is incapable of being propagated. Thus, a waveguide is seen to exhibit definite properties of a high-pass filter.

The propagation constant, $\gamma_{m,n}$, determines the amplitude and phase of each component of the wave as it is propagated along the length of the guide. With *z* = (direction of propagation) and $\omega = 2\pi \times$ (frequency), the factor for each component is

$$\exp[j\omega t - \gamma_{m,n}z]$$

Thus if $\gamma_{m,n}$ is real, the phase of each component is constant, but the amplitude decreases exponentially with *z*. When $\gamma_{m,n}$ is imaginary, it is said that no propagation takes place. The frequency is considered below cutoff. Actually, propagation with high attenuation does take place for a small distance, and a short length of guide below cutoff is often used as a calibrated attenuator.

When $\gamma_{m,n}$ is imaginary, the amplitude of each component remains constant, but the phase varies with *z*. Hence, propagation takes place. The value of $\gamma_{m,n}$ is purely imaginary only in a lossless guide. In the practical case, $\gamma_{m,n}$ usually has both a real part, $\alpha_{m,n}$, which is the attenuation constant, and an imaginary part, $\beta_{m,n}$, which is the phase propagation constant. Then $\gamma_{m,n} = \alpha_{m,n} + j\beta_{m,n}$.

RECTANGULAR WAVEGUIDES

Fig. 1 shows a rectangular waveguide and a rectangular system of coordinates, disposed so that the origin falls on one of the corners of the waveguide; *z* is the direction of propagation along the guide, and the cross-sectional dimensions are *y*₀ and *x*₀.

For the case of perfect conductivity of the guide walls with a nonconducting interior dielectric (usually air), the equations for the TM_{*m,n*} or E_{*m,n*} waves in the dielectric are

$$E_x = -A \frac{\gamma_{m,n}}{\gamma_{m,n}z + \omega^2 \mu \epsilon} (m\pi/x_0) \sin[(n\pi/y_0)y] \times \cos[(m\pi/x_0)x] \exp(j\omega t - \gamma_{m,n}z)$$

$$E_y = -A \frac{\gamma_{m,n}}{\gamma_{m,n}z + \omega^2 \mu \epsilon} (n\pi/y_0) \cos[(n\pi/y_0)y] \times \sin[(m\pi/x_0)x] \exp(j\omega t - \gamma_{m,n}z)$$

$$E_z = A \sin[(n\pi/y_0)y] \sin[(m\pi/x_0)x] \times \exp(j\omega t - \gamma_{m,n}z)$$

$$H_x = -A \frac{j\omega \epsilon}{\gamma_{m,n}z + \omega^2 \mu \epsilon} (n\pi/y_0) \cos[(n\pi/y_0)y] \times \sin[(m\pi/x_0)x] \exp(j\omega t - \gamma_{m,n}z)$$

$$H_y = A \frac{j\omega \epsilon}{\gamma_{m,n}z + \omega^2 \mu \epsilon} (m\pi/x_0) \sin[(n\pi/y_0)y] \times \cos[(m\pi/x_0)x] \exp(j\omega t - \gamma_{m,n}z)$$

$$H_z \equiv 0$$

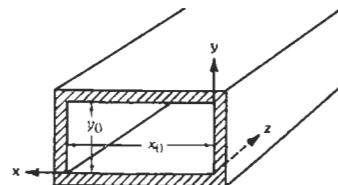


Fig. 1. Rectangular waveguide.

where ϵ is the dielectric constant and μ the permeability of the dielectric material in meter-kilogram-second (rationalized) units.

Constant A is determined solely by the exciting voltage. It has both amplitude and phase. Integers m and n may individually take values from 1 to infinity. No TM waves of the 0,0 type or 1,0 type are possible in a rectangular guide, so neither m nor n may be 0.

Equations for the TE_{*m,n*} waves or H_{*m,n*} waves in a dielectric are

$$E_x = -B \frac{j\omega\mu}{\gamma_{m,n}z + \omega^2\mu\epsilon} (n\pi/y_0) \sin[(n\pi/y_0)y] \times \cos[(m\pi/x_0)x] \exp(j\omega t - \gamma_{m,n}z)$$

$$E_y = B \frac{j\omega\mu}{\gamma_{m,n}z + \omega^2\mu\epsilon} (m\pi/x_0) \cos[(n\pi/y_0)y] \times \sin[(m\pi/x_0)x] \exp(j\omega t - \gamma_{m,n}z)$$

$$E_z = 0$$

$$H_z = B \frac{\gamma_{m,n}}{\gamma_{m,n}z + \omega^2\mu\epsilon} (m\pi/x_0) \cos[(n\pi/y_0)y] \times \sin[(m\pi/x_0)x] \exp(j\omega t - \gamma_{m,n}z)$$

$$H_y = B \frac{\gamma_{m,n}}{\gamma_{m,n}z + \omega^2\mu\epsilon} (n\pi/y_0) \sin[(n\pi/y_0)y] \times \cos[(m\pi/x_0)x] \exp(j\omega t - \gamma_{m,n}z)$$

$$H_x = B \cos[(n\pi/y_0)y] \cos[(m\pi/x_0)x] \times \exp(j\omega t - \gamma_{m,n}z)$$

Constant B depends only on the original exciting voltage and has both magnitude and phase; m and n individually may assume any integer value from 0 to infinity. The 0,0 type of wave where both m and n are 0 is not possible; all other combinations are.

As stated previously, propagation takes place only when propagation constant $\gamma_{m,n}$ is imaginary.

$$\gamma_{m,n} = [(m\pi/x_0)^2 + (n\pi/y_0)^2 - \omega^2\mu\epsilon]^{1/2}$$

This means, for any m,n mode, propagation takes place when

$$\omega^2\mu\epsilon > (m\pi/x_0)^2 + (n\pi/y_0)^2$$

or, in terms of frequency f and velocity of light c , when

$$f > \frac{c}{2\pi(\mu_1\epsilon_1)^{1/2}} [(m\pi/x_0)^2 + (n\pi/y_0)^2]^{1/2}$$

where μ_1 and ϵ_1 are the relative permeability and relative dielectric constant, respectively, of the dielectric material with respect to free space.

The wavelength in the air-filled waveguide is always greater than the wavelength in free space. The wavelength in the dielectric-filled waveguide may be less than the wavelength in free space. If λ is the wavelength in free space and the medium filling the waveguide has a relative dielectric constant ϵ ,

$$\lambda_{g(m,n)} = \frac{\lambda}{[\epsilon - (m\lambda/2x_0)^2 - (n\lambda/2y_0)^2]^{1/2}} = \frac{\lambda}{[\epsilon - (\lambda/\lambda_c)^2]^{1/2}}$$

where $(1/\lambda_c)^2 = (m/2x_0)^2 + (n/2y_0)^2$.

The phase velocity within the guide is also always greater than in an unbounded medium. The phase velocity, v , and group velocity, u , are related by

$$u = c^2/v$$

where the phase velocity is given by $v = c\lambda_g/\lambda$, and the group velocity is the velocity of propagation of the energy.

To couple energy into waveguides, it is necessary to understand the configuration of the characteristic electric and magnetic lines. Fig. 2 shows the field configuration for a TE_{1,0} wave. Fig. 3 shows the instantaneous field configuration for a higher mode, a TE_{2,1} wave.

In Fig. 4 are shown only the characteristic E lines for the TE_{1,0}, TE_{2,0}, TE_{1,1}, and TE_{2,1} waves. The arrows on the lines indicate their instantaneous relative directions. To excite a TE wave, it is necessary to insert a probe to coincide with the direction of the E lines. Thus, for a TE_{1,0} wave, a single probe projecting from the side of the guide parallel to the E lines would be sufficient to couple into it. Two ways of coupling from a coaxial line to a rectangular waveguide to excite the TE_{1,0} mode are shown in Fig. 5. With structures such as these, it is possible to make the standing-wave ratio due to the junction less than 1.15 over a 10- to 15-percent frequency band.

Fig. 6 shows the instantaneous configuration of a TM_{1,1} wave; Fig. 7 shows the instantaneous field configuration for a TM_{2,1} wave. Coupling to this type of wave may be accomplished by inserting a probe that is parallel to the E lines, or by means of a loop oriented to link the lines of flux.

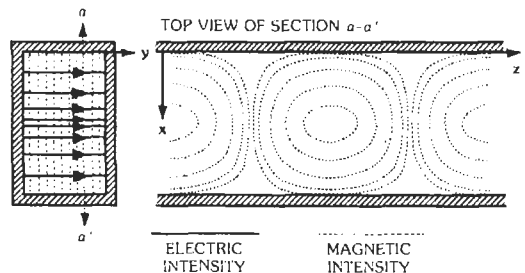


Fig. 2. Field configuration for a TE_{1,0} wave.

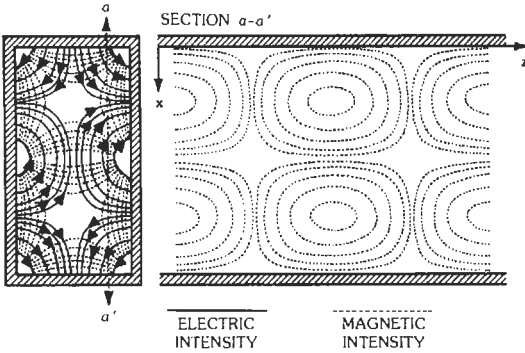


Fig. 3. Field configuration for a $TE_{2,1}$ wave.

CIRCULAR WAVEGUIDES

The usual coordinate system is ρ, θ, z , where ρ is the radial direction, θ is the angle, and z is in the longitudinal direction.

TM Waves (E Waves): $H_z \equiv 0$

$$E_\rho = H_\theta \eta (\lambda / \lambda_{g(m,n)}) \exp(j\omega t - \gamma_{m,n} z)$$

$$E_\theta = -H_\rho \eta (\lambda / \lambda_{g(m,n)}) \exp(j\omega t - \gamma_{m,n} z)$$

$$E_z = A J_n(k_{m,n} \rho) \cos n\theta \exp(j\omega t - \gamma_{m,n} z)$$

$$H_\rho = -jA (2\pi / \lambda k_{m,n}^2 \eta \rho) J_n'(k_{m,n} \rho) \sin n\theta$$

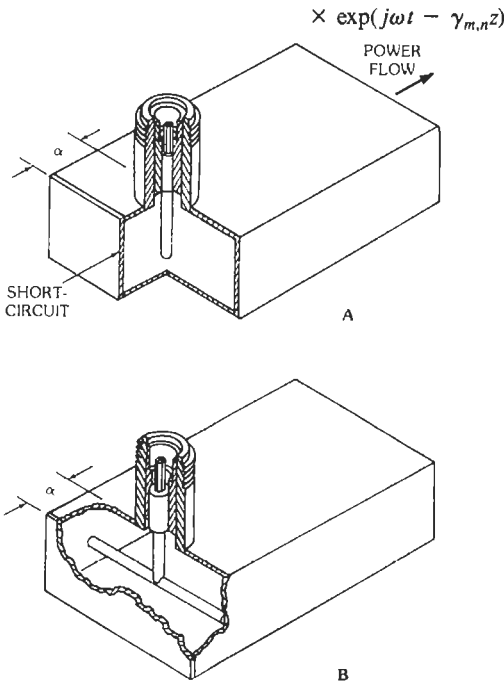


Fig. 5. Methods of coupling to $TE_{1,0}$ mode ($\alpha \approx \lambda_g/4$).

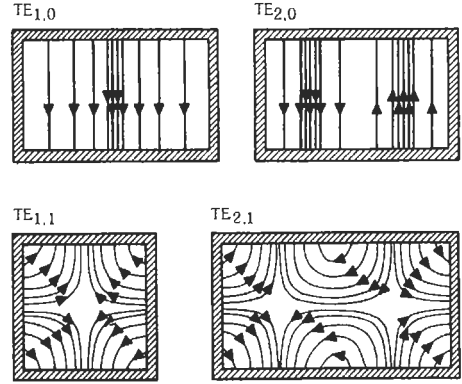


Fig. 4. Characteristic E lines for TE waves.

$$H_\theta = -jA (2\pi / \lambda k_{m,n} \eta) J_n'(k_{m,n} \rho) \cos n\theta$$

$$\times \exp(j\omega t - \gamma_{m,n} z)$$

where $\eta = (\mu/\epsilon)^{1/2}$, with μ and ϵ in absolute units.

By the boundary conditions, $E_z = 0$ when $\rho = a$, the radius of the guide. Thus, the only permissible values of k are those for which $J_n(k_{m,n} a) = 0$, because E_z must be zero at the boundary.

The numbers m, n take on all integral values from zero to infinity. The waves are seen to be characterized by the numbers m and n , where n gives the order of the Bessel functions, and m gives the order of the root of $J_n(k_{m,n} a)$. The Bessel function has an infinite number of roots, so there are an infinite number of k 's that make $J_n(k_{m,n} a) = 0$.

TE Waves (H Waves): $E_z \equiv 0$

$$E_\rho = jB (2\pi \eta / \lambda k_{m,n}^2 \rho) J_n(k_{m,n} \rho) \sin n\theta$$

$$\times \exp(j\omega t - \gamma_{m,n} z)$$

$$E_\theta = jB (2\pi \eta / \lambda k_{m,n}) J_n'(k_{m,n} \rho) \cos n\theta$$

$$\times \exp(j\omega t - \gamma_{m,n} z)$$

$$H_\rho = -E_\theta (\lambda_{g(m,n)} / \eta \lambda) \exp(j\omega t - \gamma_{m,n} z)$$

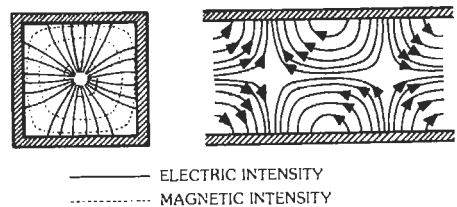


Fig. 6. Instantaneous field configuration, $TM_{1,1}$ wave.

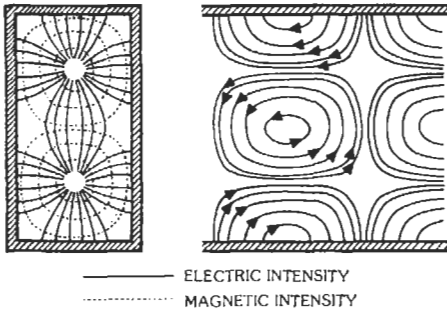


Fig. 7. Instantaneous field configuration, $TM_{2,1}$ wave.

$$H_{\theta} = E_{\rho}(\lambda_{g(m,n)}/\eta\lambda) \exp(j\omega t - \gamma_{m,n}z)$$

$$H_z = BJ_n(k_{m,n}\rho) \cos n\theta \exp(j\omega t - \gamma_{m,n}z)$$

Again, n takes on integral values from zero to infinity. The boundary condition $E_{\theta} = 0$ when $\rho = a$ still applies. To satisfy this condition, k must be such as to make $J_n'(k_{m,n}a)$ equal to zero [where the superscript indicates the derivative of $J_n(k_{m,n}a)$]. It is seen that m takes on values from 1 to infinity, since there are an infinite number of roots of $J_n'(k_{m,n}a)$.

For circular waveguides, the cutoff frequency for the m,n mode is

$$f_{c(m,n)} = ck_{m,n}/2\pi$$

where c = velocity of light, and $k_{m,n}$ is evaluated from the roots of the Bessel functions.

$$k_{m,n} = U_{m,n}/a \text{ or } U_{m,n}'/a$$

where a = radius of guide or pipe, and $U_{m,n}$ is the root of the particular Bessel function of interest (or its derivative).

The wavelength in any guide filled with a homogeneous dielectric ϵ (relative) is

$$\lambda_g = \lambda_0/[\epsilon - (\lambda_0/\lambda_c)^2]^{1/2}$$

where λ_0 is the wavelength in free space, and λ_c is the free-space cutoff wavelength for any mode under consideration.

Tables 1 and 2 are useful in determining the values of k . For TE waves, the cutoff wavelengths are given in

TABLE 1. TE WAVES, VALUES OF λ_c/a
(WHERE a = RADIUS OF GUIDE)

$\begin{matrix} m \\ n \end{matrix}$	0	1	2
1	1.640	3.414	2.057
2	0.896	1.178	0.937
3	0.618	0.736	0.631

TABLE 2. TM WAVES, VALUES OF λ_c/a

$\begin{matrix} m \\ n \end{matrix}$	0	1	2
1	2.619	1.640	1.224
2	1.139	0.896	0.747
3	0.726	0.618	0.541

Table 1, and for TM waves the cutoff wavelengths are given in Table 2, where n is the order of the Bessel function and m is the order of the root.

Fig. 8 shows λ_0/λ_g as a function of λ_0/λ_c . From this, λ_g may be determined when λ_0 and λ_c are known.

The pattern of magnetic force of TM waves in a circular waveguide is shown in Fig. 9. Only the maximum lines are indicated. To excite this type of pattern, it is necessary to insert a probe along the length of the waveguide and concentric with the H lines. For instance, in the $TM_{0,1}$ type of wave, a probe extending down the length of the waveguide at its very center would provide the proper excitation. This method of excitation is shown in Fig. 10. Corresponding methods of excitation may be used for the other types of TM waves shown in Fig. 9.

Fig. 11 shows the patterns of electric force for TE waves. Again, only the maximum lines are indicated. This type of wave may be excited by an antenna that is parallel to the electric lines of force. The $TE_{1,1}$ wave may be excited by means of an antenna extending across the waveguide. This is illustrated in Fig. 12.

Propagating E waves have a minimum attenuation at $(3)^{1/2}f_c$. The $H_{1,1}$ wave has minimum attenuation at the frequency $2.6(3)^{1/2}f_c$.

The $H_{0,1}$ wave has the interesting and useful property that attenuation decreases as the frequency increases. This has made the $H_{0,1}$ mode exceedingly useful in the transmission of microwave signals over long distances.

Table 3 presents some of the important equations for various guides.

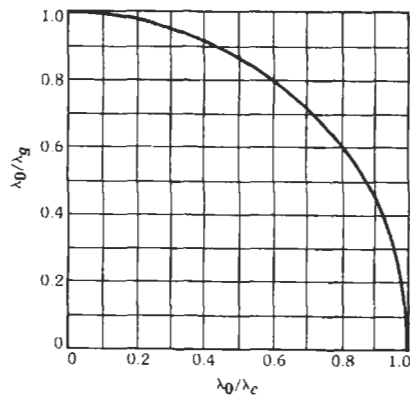


Fig. 8. Chart for determining guide wavelength.

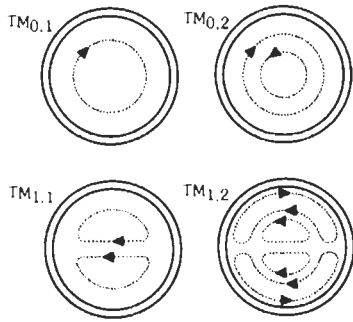


Fig. 9. Patterns of magnetic force of TM waves in circular waveguides.

SQUARE WAVEGUIDES

Waveguide having interior dimensions $x_0 = y_0$ (Fig. 1) has found increasingly important application in dual-polarized horn feeds and waveguide multiplexers. Usually these involve simultaneous propagation of the orthogonally oriented dominant modes, $TE_{1,0}$ and $TE_{0,1}$. These modes are theoretically capable of propagation without cross coupling, at the same frequency, in lossless waveguide of square cross section. In practice, wall losses, surface irregularities, and unequal transverse interior dimensions give rise to $TE_{1,0}$ and $TE_{0,1}$ mode cross-conversion. This occurs continually along the waveguide so that unless special care is taken, long lengths of dual-polarized waveguide exhibit a deteriorated mode isolation as a function of length of guide.

Most important in establishing the initial mode isolation is proximity of the operating frequency to the cutoff frequency of the $TE_{1,0}$ mode in the square waveguide so that the total operating frequency band is above $TE_{1,0}$ cutoff and well below $TE_{1,1}$ cutoff. The lowest operating frequency should be approximately 25% above the $TE_{1,0}$ cutoff frequency. Thus, a dual-polarized feed propagating a 4400-MHz signal should use a square waveguide having internal dimensions of about 1.68 inches. If the internal dimensions are arrived at by using the 1.87 internal dimension of WR 187 waveguide, the dual-mode isolation will probably not

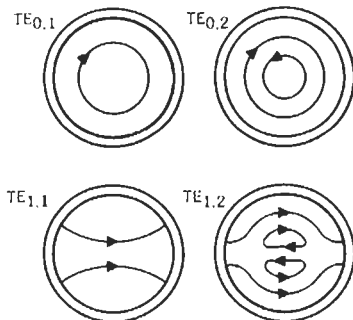


Fig. 11. Patterns of electric force of TE waves in circular waveguides.

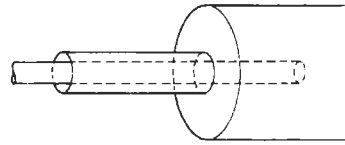


Fig. 10. Method of coupling to circular waveguide for $TM_{0,1}$ wave.

exceed 35 dB. By operating about 25% above the $TE_{1,0}$ mode cutoff frequency and well below $TE_{1,1}$ cutoff, the isolation can exceed 50 dB.*

ATTENUATION IN A WAVEGUIDE BEYOND CUTOFF

When a waveguide is used at a wavelength greater than the cutoff wavelength, there is no real propagation and the fields are attenuated exponentially. The attenuation, L , in a length, d , is given by

$$L = 54.5(d/\lambda_c)[1 - (\lambda_c/\lambda)^2]^{1/2} \text{ decibels}$$

where λ_c = cutoff wavelength and λ = operating wavelength.

Note that for $\lambda \gg \lambda_c$, attenuation is essentially independent of frequency and

$$L = 54.5d/\lambda_c \text{ decibels}$$

where λ_c is a function of geometry.

STANDARD WAVEGUIDES

Table 4 lists some properties and dimensions of standard rectangular waveguides. For other than theoretical vacuum performance, consider the relative value of ϵ for sea level, 20° C air, as approximately 1.0006. Rounded inner corners also modify performance slightly.†

* D. J. LeVine and W. Sichak, "Dual-Mode Horn Feed for Microwave Multiplexing." *Electronics*, September 1954.

† M. M. Brady, "Cutoff Wavelengths and Frequencies of Standard Rectangular Waveguides." *I.E.E. Electronics Letters*, Vol. 5, No. 17, 21 August 1969.

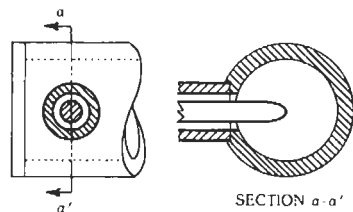
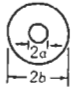
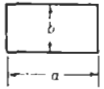
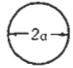


Fig. 12. Method of coupling to circular waveguide for $TE_{1,1}$ wave.

TABLE 3. CUTOFF WAVELENGTHS AND ATTENUATION FACTORS*

Type of Guide (Copper)†	Cutoff Wavelength λ_c $f_c = (c/\lambda_c)$	Attenuation Constant α (dB/100 ft)
Coaxial line‡ TEM 	0	$\frac{3.58 \times 10^{-6}(f)^{1/2}(1/b)[1 + (b/a)]}{\ln (b/a)}$ Note: Fig. 29 in Chapter 29 includes a derating factor of 20% applied to calculated α .
Rectangular pipe TE _{m,0} or H _{m,0} 	2a/m	$\frac{1.107}{a^{3/2}} \times \frac{1/2(a/b) (ff_c)^{3/2} + (ff_c)^{-1/2}}{[(ff_c)^2 - 1]^{1/2}}$
Circular pipe: 		
TM _{0,1} or E _{0,1}	2.613a	$\frac{0.485}{a^{3/2}} \times \frac{(ff_c)^{3/2}}{[(ff_c)^2 - 1]^{1/2}}$
TE _{1,1} or H _{0,1}	3.412a	$\frac{0.423}{a^{3/2}} \times \frac{(ff_c)^{-1/2} + (1/2.38) (ff_c)^{3/2}}{[(ff_c)^2 - 1]^{1/2}}$
TE _{0,1} or H _{0,1}	1.640a	$\frac{0.611}{a^{3/2}} \times \frac{(ff_c)^{-1/2}}{[(ff_c)^2 - 1]^{1/2}}$

* Dimensions are in inches and frequencies in hertz; vacuum dielectric.

† For other metals multiply α by the square root of ratio of resistivity relative to that of copper.

‡ Inner and outer conductors same material.

RIDGED WAVEGUIDES

To lower the cutoff frequency of a waveguide for use over a frequency band wider than normal, ridges may be used. By proper choice of dimensions, it is possible to obtain as much as a 4:1 ratio between cutoff frequencies for the TE_{2,0} and TE_{1,0} modes.

Tables 5 and 6 and Figs. 13 and 14 give the essential characteristics of single- and double-ridged guides. Figs. 15 and 16 show the relationship between the cutoff wavelength and the critical dimensions. Figs. 17 and 18 describe the bandwidth (ratio of cutoff wavelengths of the 0,1 and 0,2 modes). The price paid for increased

bandwidth is an increase in attenuation relative to the equivalent rectangular guide (Figs. 19 and 20).

Coaxial line can be coupled to either the single- or double-ridged guide by partly inserting the probe into the open space between the ridge and the opposite wall or ridge (Fig. 5A) or connecting it directly across the gap.

FLEXIBLE WAVEGUIDES

Flexible waveguide is used to join rigid sections or components that cannot be accurately dimensioned, positioned, or rendered immobile in space. Thus, rather than attempt to specify and hold precisely every bend,

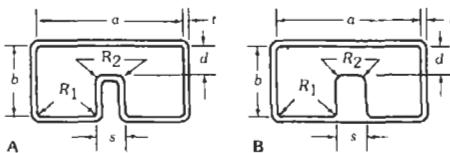


Fig. 13. Single-ridged waveguides (refer to Table 5).

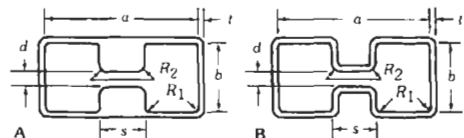


Fig. 14. Double-ridged waveguides (refer to Table 6).

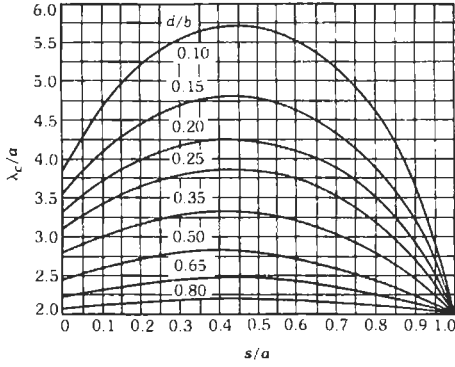


Fig. 15. Cutoff wavelength, single-ridged guide. (From S. Hopfer, "The Design of Ridged Waveguides," Transactions of the I.R.E., Vol. MTT-3, No. 5, 1955; Fig. 5.)

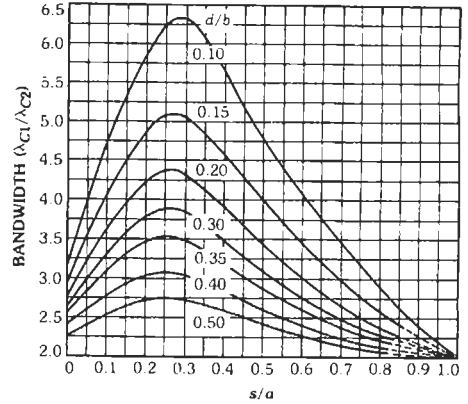


Fig. 18. Bandwidth curves, double-ridged guide. (From S. Hopfer, "The Design of Ridged Waveguides," Transactions of the I.R.E., Vol. MTT-3, No. 5, 1955; Fig. 10.)

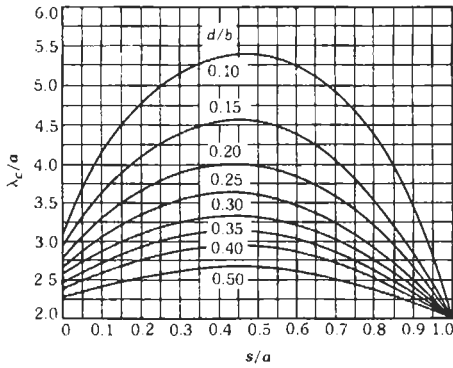


Fig. 16. Cutoff wavelength, double-ridged guide. (From S. Hopfer, "The Design of Ridged Waveguides," Transactions of the I.R.E., Vol. MTT-3, No. 5, 1955; Fig. 2.)

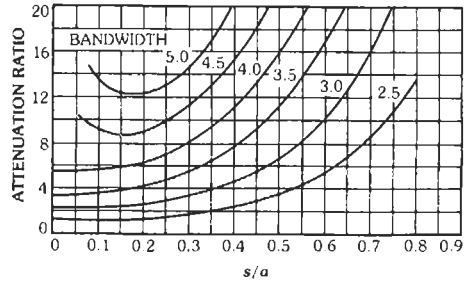


Fig. 19. Attenuation ratio, parametric in bandwidth, single-ridged guide. (From S. Hopfer, "The Design of Ridged Waveguides," Transactions of the I.R.E., Vol. MTT-3, No. 5, 1955; Fig. 11.)

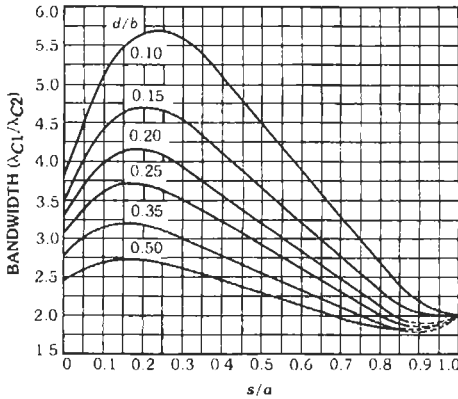


Fig. 17. Bandwidth curves, single-ridged guide. (From S. Hopfer, "The Design of Ridged Waveguides," Transactions of the I.R.E., Vol. MTT-3, No. 5, 1955; Fig. 9.)

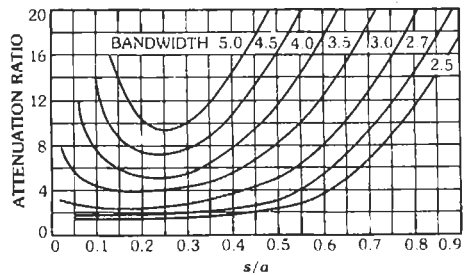


Fig. 20. Attenuation ratio, parametric in bandwidth, double-ridged guide. (From S. Hopfer, "The Design of Ridged Waveguides," Transactions of the I.R.E., Vol. MTT-3, No. 5, 1955; Fig. 12.)

TABLE 4. STANDARD WAVEGUIDES

EIA Waveguide Designation (Standard RS-261-A)	JAN Waveguide Designation (MIL-HDBK-216, 4 January 1962)	Outer Dimensions and Wall Thickness (in inches)	Frequency Range in Gigahertz for Dominant ($TE_{1,0}$) Mode	Cutoff Wavelength, λ_c , in Centimeters for $TE_{1,0}$ Mode	Cutoff Frequency, f_c , in Gigahertz for $TE_{1,0}$ Mode	Theoretical Attenuation, Lowest to Highest Frequency in dB/100 ft	Theoretical Power Rating in Megawatts for Lowest to Highest Frequency*
WR-2300	RG-290/U†	23.250×11.750×0.125	0.32-0.49	116.8	0.256	0.051-0.031	153.0-212.0
WR-2100	RG-291/U†	21.250×10.750×0.125	0.35-0.53	106.7	0.281	0.054-0.034	120.0-173.0
WR-1800	RG-201/U†	18.250×9.250×0.125	0.425-0.620	91.4	0.328	0.056-0.038	93.4-131.9
WR-1500	RG-202/U†	15.250×7.750×0.125	0.49-0.740	76.3	0.393	0.069-0.050	67.6-93.3
WR-1150	RG-203/U†	11.750×6.000×0.125	0.64-0.96	58.4	0.514	0.128-0.075	35.0-53.8
WR-975	RG-204/U†	10.000×5.125×0.125	0.75-1.12	49.6	0.605	0.137-0.095	27.0-38.5
WR-770	RG-205/U†	7.950×4.100×0.125	0.96-1.45	39.1	0.767	0.201-0.136	17.2-24.1
WR-650	RG-69/U	6.660×3.410×0.080	1.12-1.70	33.0	0.908	0.317-0.212	11.9-17.2
WR-510	—	5.260×2.710×0.080	1.45-2.20	25.9	1.16	—	—
WR-430	RG-104/U	4.460×2.310×0.080	1.70-2.60	21.8	1.375	0.588-0.385	5.2-7.5
WR-340	RG-112/U	3.560×1.860×0.080	2.20-3.30	17.3	1.735	0.877-0.572	—
WR-284	RG-48/U	3.000×1.500×0.080	2.60-3.95	14.2	2.08	1.102-0.752	2.2-3.2
WR-229	—	2.418×1.273×0.064	3.30-4.90	11.6	2.59	—	—
WR-187	RG-49/U	2.000×1.000×0.064	3.95-5.85	9.50	3.16	2.08-1.44	1.4-2.0
WR-159	—	1.718×0.923×0.064	4.90-7.05	8.09	3.71	—	—
WR-137	RG-50/U	1.500×0.750×0.064	5.85-8.20	6.98	4.29	2.87-2.30	0.56-0.71
WR-112	RG-51/U	1.250×0.625×0.064	7.05-10.00	5.70	5.26	4.12-3.21	0.35-0.46
WR-90	RG-52/U	1.000×0.500×0.050	8.20-12.40	4.57	6.56	6.45-4.48	0.20-0.29
WR-75	—	0.850×0.475×0.050	10.00-15.00	3.81	7.88	—	—
WR-62	RG-91/U	0.702×0.391×0.040	12.40-18.00	3.16	9.49	9.51-8.31	0.12-0.16
WR-51	—	0.590×0.335×0.040	15.00-22.00	2.59	11.6	—	—
WR-42	RG-53/U	0.500×0.250×0.040	18.00-26.50	2.13	14.1	20.7-14.8	0.043-0.058
WR-34	—	0.420×0.250×0.040	22.00-33.00	1.73	17.3	—	—
WR-28	RG-96/U‡	0.360×0.220×0.040	26.50-40.00	1.42	21.1	21.9-15.0	0.022-0.031

WR-22	RG-97/U‡	0.304×0.192×0.040	33.00-50.00	1.14	26.35	31.0-20.9	0.014-0.020
WR-19	—	0.268×0.174×0.040	40.00-60.00	0.955	31.4	—	—
WR-15	RG-98/U‡	0.228×0.154×0.040	50.00-75.00	0.753	39.9	52.9-39.1	0.0063-0.0090
WR-12	RG-99/U‡	0.202×0.141×0.040	60.00-90.00	0.620	48.4	93.3-52.2	0.0042-0.0060
WR-10	—	0.180×0.130×0.040	75.00-110.00	0.509	59.0	—	—
WR-8	RG-138/U§	0.140×0.100×0.030	90.00-140.00	0.406	73.84	152-99	0.0018-0.0026
WR-7	RG-136/U§	0.125×0.0925×0.030	110.00-170.00	0.330	90.84	163-137	0.0012-0.0017
WR-5	RG-135/U§	0.111×0.0855×0.030	140.00-220.00	0.259	115.75	308-193	0.00071-0.00107
WR-4	RG-137/U§	0.103×0.0815×0.030	170.00-260.00	0.218	137.52	384-254	0.00052-0.00075
WR-3	RG-139/U§	0.094×0.0770×0.030	220.00-325.00	0.173	173.28	512-348	0.00035-0.00047

* For these computations, the breakdown strength of air was taken as 15 000 volts per centimeter. A safety factor of approximately 2 at sea level has been allowed.

† Aluminum, 2.83×10^{-6} ohm-cm resistivity.‡ Silver, 1.62×10^{-6} ohm-cm resistivity.§ JAN types are silver, with a circular outer diameter of 0.156 inch and a rectangular bore matching EIA types. All other types are of a Cu-Zn alloy, 3.9×10^{-6} ohm-cm resistivity.

Note: Equivalent designations of waveguides follow.

EIA	British	IEC	EIA	British	IEC	EIA	British	IEC
WR-2300	00	-R3	WR-340	9A	-R26	WR-51	19	-R180
WR-2100	0	-R4	WR-284	10	-R32			
WR-1800	1	-R5				WR-42	20	-R220
			WR-229	11A	-R40	WR-34	21	-R260
WR-1500	2	-R6	WR-187	12	-R48	WR-28	22	-R320
WR-1150	3	-R8	WR-159	13	-R58			
WR-975	4	-R9				WR-22	23	-R400
			WR-137	14	-R70	WR-19	24	-R500
WR-770	5	-R12	WR-112	15	-R84	WR-15	25	-R620
WR-650	6	-R14	WR-90	16	-R100			
WR-510	7	-R18				WR-12	26	-R740
			WR-75	17	-R120	WR-10	27	-R900
WR-430	8	-R22	WR-62	18	-R140	WR-8	28	-R1200

TABLE 5. CHARACTERISTICS OF SINGLE-RIDGED WAVEGUIDES*

Frequency Range (GHz)	$f_{c1,0}$ (GHz)	$\lambda_{c1,0}$ (in)	$f_{c2,0}$ (GHz)	Dimensions in Inches						At $f=(3)^{1/2}f_{c1,0}$			
				a	b	d	s	t	R_1 (max)	R_2	Atten** (dB/ft)	Power Rating† (kW)	
Bandwidth 2.4:1													
0.175-0.42	0.148	79.803	0.431	28.129	12.658	5.278	4.360	—	—	1.056	0.000 24	32 870.	
0.267-0.64	0.226	52.260	0.658	18.421	8.289	3.457	2.855	—	—	0.691	0.000 45	14 100.	
0.42-1.0	0.356	33.177	1.036	11.695	5.263	2.195	1.813	0.125	0.047	0.439	0.000 87	5 682.	
0.64-1.53	0.542	21.792	1.577	7.682	3.457	1.442	1.191	0.125	0.047	0.288	0.001 64	2 451.	
0.84-2.0	0.712	16.588	2.072	5.847	2.631	1.097	0.906	0.080	0.047	0.219	0.002 48	1 421.	
1.5-3.6	1.271	9.293	3.699	3.276	1.474	0.615	0.508	0.080	0.047	0.123	0.005 91	445.8	
2.0-4.8	1.695	6.968	4.933	2.456	1.105	0.461	0.381	0.080	0.047	0.092	0.009 08	250.6	
3.5-8.2	2.966	3.982	8.632	1.404	0.632	0.264	0.218	0.064	0.031	0.053	0.021 2	81.87	
4.75-11.0	4.025	2.934	11.714	1.034	0.465	0.194	0.160	0.050	0.031	0.039	0.033 3	44.43	
7.5-18.0‡	6.356	1.858	18.498	0.655	0.295	0.123	0.1015	0.050	0.015	0.025	0.066 1	17.82	
11.0-26.5‡	9.322	1.267	27.130	0.4466	0.2010	0.0838	0.0692	0.040	0.015	0.017	0.117	8.285	
18.0-40.0‡	15.254	0.7743	44.393	0.2729	0.1228	0.0512	0.0423	0.040	0.015	0.010	0.246	3.035	
Bandwidth 3.6:1													
0.108-0.39	0.092	128.37	0.404	31.218	14.048	2.402	5.307	—	—	0.480	0.001 6	14 550.	
0.27-0.97	0.229	51.572	1.006	12.542	5.644	0.965	2.132	—	—	0.193	0.006 5	2 348.	
0.39-1.4	0.331	35.680	1.454	8.677	3.905	0.668	1.475	0.125	0.047	0.134	0.011 2	1 124.	
0.97-3.5	0.822	14.367	3.611	3.494	1.572	0.269	0.594	0.080	0.047	0.054	0.043 8	182.2	
1.4-5.0	1.186	9.958	5.210	2.422	1.090	0.186	0.412	0.080	0.047	0.037	0.075 8	87.56	
3.5-12.4	2.966	3.982	13.030	0.968	0.436	0.075	0.165	0.050	0.031	0.015	0.300	13.99	
5.0-18.0‡	4.237	2.787	18.613	0.678	0.305	0.052	0.115	0.050	0.015	0.010	0.513	6.857	
12.4-40.0‡	10.508	1.124	46.162	0.273	0.123	0.021	0.046	0.040	0.015	0.004	2.008	1.115	

* From MIL-HDBK-216, *RF Transmission Lines and Fittings*, 4 January 1962.

** Copper.

† Based on breakdown of air—15 000 volts per cm (safety factor of approx 2 at sea level). Corner radii considered.

‡ Fig. 13B in these frequency ranges only.

TABLE 6. CHARACTERISTICS OF DOUBLE-RIDGED WAVEGUIDES*

Frequency Range (GHz)	$f_{c1,0}$ (GHz)	$\lambda_{c1,0}$ (in)	$f_{c2,0}$ (GHz)	Dimensions in Inches							At $f=(3)^{1/2}f_{c1,0}$	
				a	b	d	s	t	R_1 (max)	R_2	Atten** (dB/ft)	Power Rating† (kW)
Bandwidth 2.4:1												
0.175-0.42					29.667	13.795	5.863	7.417	—	—	—	1.173
0.267-0.64					19.428	9.034	3.839	4.857	—	—	—	0.768
0.42-1.0					12.333	5.737	2.437	3.083	0.125	0.050	0.050	0.487
0.64-1.53					8.100	3.767	1.601	2.025	0.125	0.050	0.050	0.320
0.84-2.0					6.167	2.868	1.219	1.542	0.125	0.050	0.050	0.244
1.5-3.6					3.455	1.607	0.683	0.864	0.080	0.050	0.050	0.137
2.0-4.8					2.590	1.205	0.512	0.648	0.080	0.050	0.050	0.102
3.5-8.2					1.480	0.688	0.292	0.370	0.064	0.030	0.030	0.058
4.75-11.0					1.090	0.506	0.215	0.272	0.050	0.030	0.030	0.043
7.5-18.0					0.691	0.321	0.136	0.173	0.050	0.020	0.020	0.027
11.0-26.5‡					0.471	0.219	0.093	0.118	0.040	0.015	0.015	0.019
18.0-40.0‡					0.288	0.134	0.057	0.072	0.040	0.015	0.015	0.011
Bandwidth 3.6:1												
0.108-0.39	0.092	128.37	0.401	34.638	14.894	2.904	8.660	—	—	0.581	0.0014	28 830.
0.27-0.97	0.229	51.572	0.999	13.916	5.984	1.167	3.479	—	—	0.233	0.0055	4 653.
0.39-1.4	0.331	35.680	1.444	9.628	4.140	0.807	2.407	0.125	0.050	0.161	0.0097	2 227.
0.97-3.5	0.822	14.367	3.587	3.877	1.667	0.325	0.969	0.080	0.050	0.065	0.0378	361.2
1.4-5.0	1.186	9.958	5.176	2.687	1.155	0.225	0.672	0.080	0.050	0.045	0.0656	173.5
3.5-12.4	2.966	3.982	12.944	1.074	0.462	0.090	0.269	0.050	0.030	0.018	0.259	27.74
5.0-18.0	4.237	2.787	18.490	0.752	0.323	0.063	0.188	0.050	0.020	0.013	0.443	13.59
12.4-40.0‡	10.508	1.124	45.857	0.303	0.130	0.025	0.076	0.040	0.015	0.005	1.730	2.210

* From MIL-HDBK-216, *RF Transmission Lines and Fittings*, 4 January 1962.

** Copper.

† Based on breakdown of air—15 000 volts per cm (safety factor of approx 2 at sea level). Corner radii considered.

‡ Fig. 14B in these frequency ranges only.

twist, and straight section of a long waveguide run between fixed points, it is often adequate to leave a short run or bend to be filled by a flexible guide insert. Flexible sections are also used to permit thermally induced relative movement and to insulate portions of a waveguide run from shock and vibration. Flexible waveguide should not be treated as a link between a cabinet and its frequently opened doors and drawers, unless the cabinet is specified and/or designed for that type of service. Most flexible-waveguide structures are susceptible to cracking under these conditions if flexure is repeated often.

Flexible waveguide is available in many different forms. It may be made from flat ribbons wound on a rectangular mandrel with the edges convoluted or folded in and interlocked. The convoluted guide may be soldered or unsoldered, since the bending and twisting results from a flexure of each turn and not from a relative sliding as in the case of the interlocked guide. If soldered, it is more difficult to flex and essentially loses twist capability.

Corrugated flexible guide may be made by properly shaping thin-wall seamless rectangular tubing, or by bending and soldering corrugated sheet metal (with due consideration to current flow so that a low-loss joint results).

A bellows-type guide is produced from a group of radial chokes in tandem configuration and made of a flexible alloy.

Vertebral guide is made from a tandem chain of choke-cover sections contained within a neoprene or rubber jacket.

In general, all types except the seamless corrugated waveguide should be jacketed with neoprene or rubber. The unsoldered convoluted, interlocking, and vertebral guide must be jacketed to be pressurized.

Table 7 gives the properties of soldered convoluted flexible waveguide.

The wide variety of manufacturing techniques, the jacketing material and thickness, the length dependence, and other characteristics make it impossible to limit the (\pm) stretch, twist, and center-line displacement ranges of flexible waveguide. These are usually described in terms of maximum acceptable vswr or loss of the section as a function of the (\pm) stretch, twist, or displacement, and are best established on advice from the manufacturer selected.

WAVEGUIDE LOSSES

Hollow, enclosed single-conductor waveguides, propagating in the interior space, exhibit losses via dissipation in the waveguide walls and the dielectric material filling the space, leakage through the walls and connections to the guide, and localized power absorption (and heating) at the connections (flanges) because of poor contact or fabrication. The following discussion assumes that the dielectric is air, with zero loss tangent, and that the depth of penetration into the walls is very

much less than the wall thickness, so that no appreciable wall leakage occurs.

Waveguide Material and Modes

Fig. 21 shows attenuation as a function of percent conductivity for WR-112 waveguide at 9.0 GHz.* Table 8 relates material composition and percent conductivity.

To obtain the lowest attenuation where a choice of waveguide size and cross section is possible, consideration should be given to mode selection. Fig. 22 shows

* R. M. Cox and W. E. Rupp, "Fight Waveguide Losses 5 Ways," *Microwaves*, Vol. 5, No. 8, August 1966, pp. 32-40.

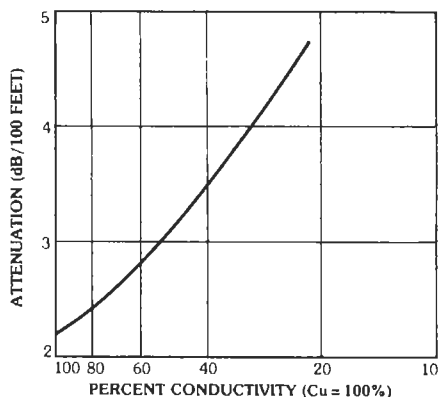


Fig. 21. Attenuation as a function of percent conductivity for WR-112 waveguide at 9.0 GHz.

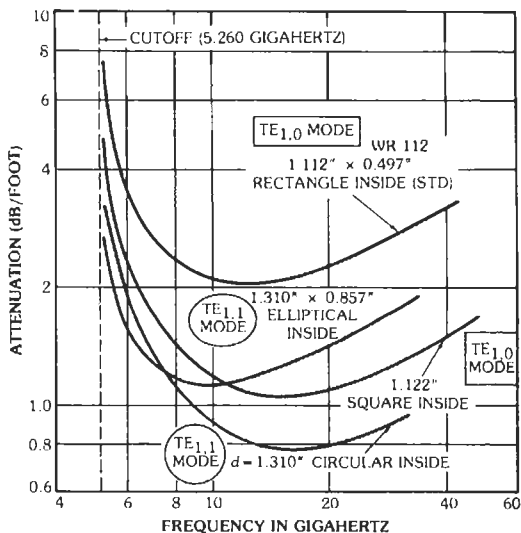


Fig. 22. Attenuation as a function of frequency for various waveguides. (From R. M. Cox and W. E. Rupp, "Fight Waveguide Losses 5 Ways," *Microwaves*, Vol. 5, No. 8, August 1966; p. 36.)

TABLE 7. PROPERTIES OF SOLDERED CONVOLUTED FLEXIBLE WAVEGUIDES*

Dimensions (inches)		Minimum Bending Radii (inches)				Equivalent Rectangular Waveguide Type	Weight (lb/ft)	Nominal Attenuation (dB/100 ft)	Nominal Power Rating (MW)	Maximum Operating Pressure (psi)
		Standard Molded Assembly		Unjacketed or Special Molded Assembly						
Inside	Outside	H Plane	E plane	H Plane	E Plane					
6.500×3.250	6.660×3.410	27	13	17	8½	RG-69/U	2.88	0.50	10	15
4.300×2.150	4.460×2.310	18	9	11½	5¾	RG-104/U	1.46	0.80	8.0	20
2.840×1.340	3.000×1.500	14	7	9	4½	RG-48/U	0.530	1.50	2.0	30
1.872×0.872	2.000×1.000	8	4	5	2½	RG-49/U	0.332	3.0	1.0	30
1.372×0.622	1.500×0.750	5	2½	3¼	1¾	RG-50/U	0.266	4.7	0.50	30
1.122×0.497	1.250×0.625	3½	1¾	2¼	1¾	RG-51/U	0.200	5.7	0.40	45
0.900×0.400	1.000×0.500	3	1½	2	1	RG-52/U	0.112	9.0	0.25	60
0.622×0.311	0.702×0.391	3	1½	2	1	RG-91/U	0.085	15.0	0.20	60
0.420×0.170	0.500×0.250	2½	1¼	1½	¾	RG-53/U	0.050	29.0	0.10	60
0.280×0.140	0.360×0.220	2½	1¼	1½	¾	RG-96/U	0.039	35.0	0.05	60

* From MIL-HDBK-216, *RF Transmission Lines and Fittings*, 4 January 1962.

TABLE 8. COMPOSITION AND CONDUCTIVITY OF WAVEGUIDE MATERIAL*

Material	Composition (%)						% Conductivity**
	Cu	Zn	P	Ag	Al	Mg	
Copper (oxygen free)	99.95† min	—	—	—	—	—	97.6 min
Copper DLP (deoxidized, low phosphorus)	99.90† min	—	0.004– 0.012	—	—	—	96.1 min
Commercial Bronze	89–91	9–11	—	—	—	—	44.2 min
Silver (fine)	0.08 max	—	—	99.90 min	—	—	100.0 min
Coin Silver	9–10.4	0.06	—	89.6–91.0	—	—	82.0 min
Aluminum 1100	0.2	0.10	—	—	99.0 min	—	59.5 min
Aluminum 6061	0.15–0.40	0.25	—	—	95	0.8–1.2	40.0 min
Magnesium	—	0.05	0.6–1.4	—	2.5–3.5	94.0	37.5‡

* From R. M. Cox and W. E. Rupp, "Fight Waveguide Losses 5 Ways," *Microwaves*, Vol. 5, No. 8, August 1966, p. 34.

** International Annealed Copper Standard.

† Any silver present is counted in the copper content.

‡ MIL-HDBK-216, Military Standardization Handbook, *RF Transmission Lines and Fittings*, 4 January 1962.

the relation among various waveguides for different modes and cross sections over a broad frequency band. Fig. 23 relates the attenuation rates of several circular and square guides (in the $TE_{0,1}$, $TE_{1,1}$, and $TE_{1,0}$ modes) over a limited frequency band.

Waveguide Flange Leakage*

It is extremely difficult to apply quantitative measurements to flange leakage without direct reference to a specific set of measuring procedures, equipment, and test environment. However, in general, measurements of flange fields, made with a probe at the flanges, have indicated that the leakage fields exhibit sharp peaks distributed around the edge of the flange connection. The levels of the peaks are of the order of -130 decibels relative to the guide power, and may be higher or lower depending on the bolt tension and RFI gaskets employed.

Flange Resistance and Bolt Torque†

The equivalent series resistance of a die-cast sealed rectangular aluminum S-band flange pair as a function of bolt torque (each bolt) is shown in Fig. 24, as an indication of the importance of proper bolt tightening at a flange connection. Table 9 gives the bolt torque for several bolt sizes to meet the recommended value of about 1000 lb/linear inch of flange connection, which is estimated to give a satisfactory waveguide seal for high-power applications.

* S. Galagan. *Electrical Characteristics of Waveguide Seals with EMI Supplement*. Unpublished report prepared for the Parker Seal Co., Culver City, Cal.; January 1964.

† Handbook Catalog No. W5460, *Waveguide Flange and EMI Sealing*, Culver City, Calif.: Parker Seal Co., Copyright 1967.

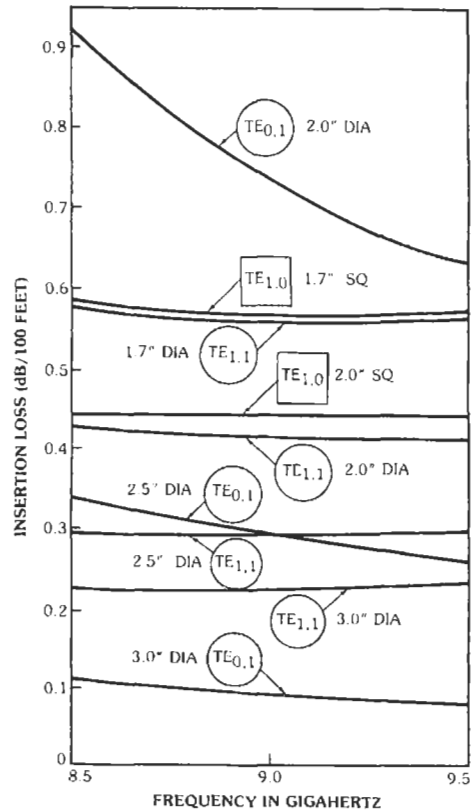


Fig. 23. Attenuation curves for various square and circular waveguides in range 8.5–9.5 GHz. (From R. M. Cox and W. E. Rupp, "Fight Waveguide Losses 5 Ways," *Microwaves*, Vol. 5, No. 8, August 1966; p. 37.)

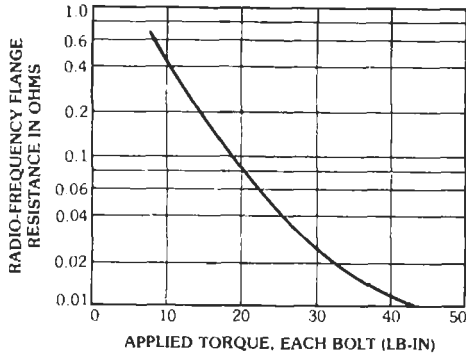


Fig. 24. Typical flange radio-frequency resistance as a function of bolt torque. (From Handbook Catalog W5460, Waveguide Flange and EMI Sealing, Culver City, Calif.: Parker Seal Co., © 1967.)

Fig. 25 shows typical flange resistance as a function of frequency for a UG-53/U plus UG-54/U choke-cover combination, with and without a mesh gasket seal. A combination of two cover flanges, without seals, may have about 10 times the choke-cover resistance.

Flange Insertion Loss**

The relationship between the flange resistance and insertion loss is

$$L(\text{decibels}) = 10 \log(1 - R_F/Z_0)$$

where R_F = radio-frequency flange resistance (measured), and

$$Z_0 = 593 b/a[1 - (f_c/f)^2]^{-1/2} \text{ (ohms)}\ddagger$$

The approximate value of flange insertion loss calculated from the above equation may be scaled to other flange sizes and frequencies by

$$R_{F1}/R_{F2} \sim (A_1/A_2)(f_1/f_2)^{1/2}$$

where A = flange area, f = frequency, and the subscripts refer to the two conditions.

Losses and Noise Temperature

In radio telescopes, satellite ground antennas, and other loss-sensitive waveguide systems, the noise-

** S. Galagan. *Electrical Characteristics of Waveguide Seals with EMI Supplement*. Unpublished report prepared for the Parker Seal Co., Culver City, Calif.; January 1964.

‡ This is the characteristic impedance, defined as the maximum transverse voltage divided by the total longitudinal current.

TABLE 9. RECOMMENDED TORQUE TABLE*

Screw Size	Threads per Inch	Torque (lb-in)	Recommended Tension** (lb)
No. 4	40	4.5	235
	80	5.5	280
No. 6	32	8.5	360
	40	10	410
No. 8	32	18	625
	36	20	685
No. 10	24	23	705
	32	32	940
1/4"	20	80	1800
	28	100	2200
	32	140	2540
3/16"	24	150	2620
	16	250	3740
3/8"	24	275	3950
	14	400	4675
7/16"	20	425	4700
	13	550	6110
1/2"	20	575	6140

* From Handbook Catalog No. W5460, *Waveguide Flange and EMI Sealing*. © 1967, Parker Seal Co., Culver City, Cal.

** Tension (lb) = torque (lb-in)/0.2 × diameter of bolt (in).

temperature contributions must be controlled and accounted for in the system design. The waveguide losses may be converted to noise temperature* by

$$T = T_R(1 - a)$$

where,

* A. J. Giger, S. Pardee, Jr., and P. R. Wickliffe, Jr., "The Ground Transmitter and Receiver." *Bell System Technical Journal*, Vol. 42, No. 4, Part 1, July 1963, p. 1096.

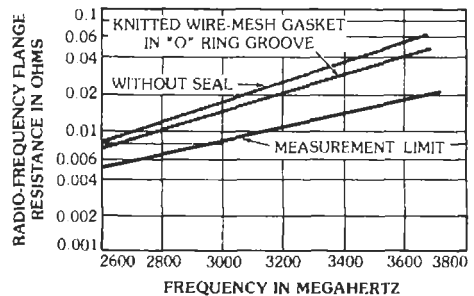


Fig. 25. Typical flange resistance as a function of frequency for choke-cover combination (UG-53/U plus UG-54/U, measured). (From Handbook Catalog W5460, Waveguide Flange and EMI Sealing, Culver City, Calif.: Parker Seal Co., © 1967.)

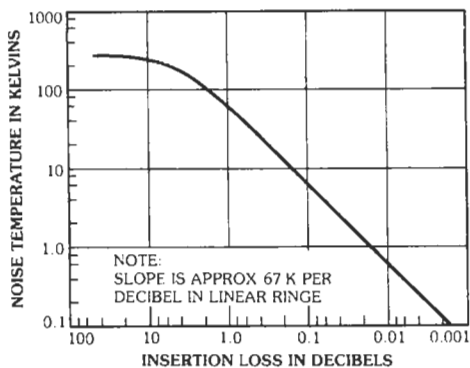


Fig. 26. Noise temperature as a function of insertion loss (290 K ambient).

T = temperature in kelvins,
 T_R = temperature of the lossy insert in kelvins,
 a = power transmission coefficient.

For A (decibels) = $-10 \log_{10} a$

$$T = T_R(1 - 10^{-A/10})$$

Fig. 26 relates the noise temperature added to a lossless waveguide system at 290 K by inserting a 290 K pad having A decibels of loss between the measuring point (or input to the receiver) and a 0 K load, as a function of insertion loss A (decibels).

Some typical values of insertion loss measured for a specific configuration are:

23-decibel directional coupler*	0.03 decibel
Flexible waveguide*	0.023 decibel
Flanges (UG-53/U plus UG-53/U)†	0.0017 decibel

WAVEGUIDE CIRCUIT ELEMENTS‡

Just as at low frequencies, it is possible to shape metallic or dielectric pieces to produce local concentrations of magnetic or electric energy within a waveguide and thus produce what are, essentially, lumped inductances or capacitances over a limited frequency bandwidth.

* A. J. Giger, S. Pardee, Jr., and P. R. Wickliffe, Jr., "The Ground Transmitter and Receiver." *Bell System Technical Journal*, Vol. 42, No. 4, Part 1, July 1963, p. 1096.

† Handbook Catalog No. W5460, *Waveguide Flange and EMI Sealing*. Culver City, Calif.: Parker Seal Co., Copyright 1967.

‡ C. G. Montgomery, R. H. Dicke, and E. M. Purcell, *Principles of Microwave Circuits*. New York: McGraw-Hill Book Co., 1948; Chapters 1 and 6. Also N. Marcuvitz, *Waveguide Handbook*. New York: McGraw-Hill Book Co., 1951.

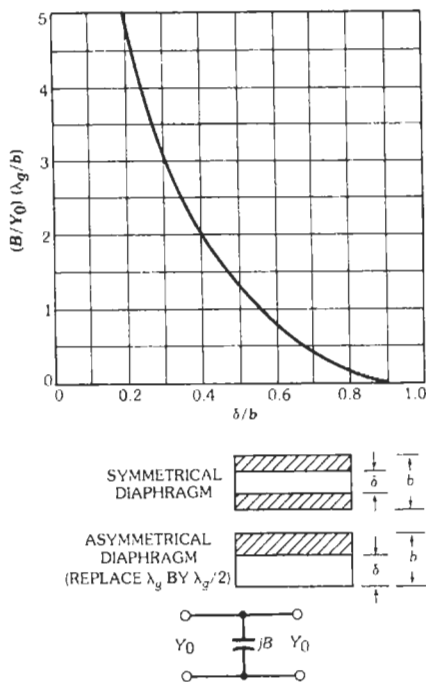


Fig. 27. Normalized susceptance of capacitive diaphragms.

This behavior as a lumped element will be evident only at some distance from the obstacle in the guide, since the fields in the immediate vicinity are disturbed.

Capacitive elements are formed from electric-field concentrating devices, such as screws or thin diaphragms inserted partially along electric-field lines. These are susceptible to breakdown under high power. Fig. 27 shows the relative susceptance, B/Y_0 , for symmetrical and asymmetrical diaphragms for small b/λ_g .

A common form of shunted lumped inductance is the diaphragm. Figs. 28 and 29 show the relative susceptance, B/Y_0 , for symmetrical and asymmetrical diaphragms in rectangular waveguides. These are computed for infinitely thin diaphragms. Finite thicknesses result in an increase in B/Y_0 .

Another form of shunt inductance that is useful because of mechanical simplicity is a round post completely across the narrow dimension of a rectangular guide (for $TE_{1,0}$ mode). Fig. 30 gives the normalized values of the elements of the equivalent four-terminal network for several post diameters.

Frequency dependence of waveguide susceptances may be given approximately as

$$\begin{aligned} \text{Inductive} &= B/Y_0 \propto \lambda_g \\ \text{Capacitive} &= B/Y_0 \propto 1/\lambda_g \text{ (distributed)} \\ &= B/Y_0 \propto \lambda_g/\lambda^2 \text{ (lumped)} \end{aligned}$$

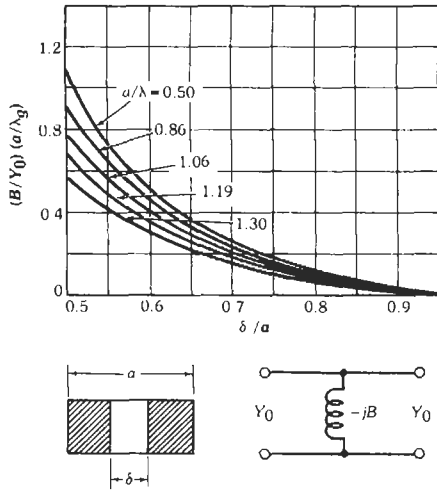


Fig. 28. Normalized susceptance of a symmetrical inductive diaphragm. (Reprinted from Microwave Transmission Circuits, by George L. Ragan, 1st ed., 1948; by permission, McGraw-Hill Book Co., N.Y.)

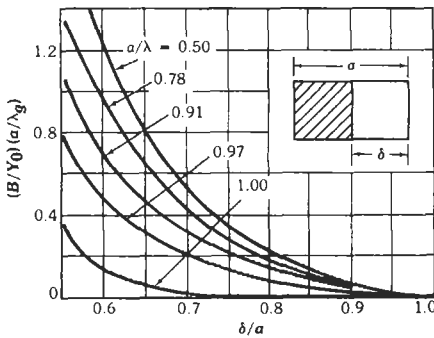


Fig. 29. Normalized susceptance of an asymmetrical inductive diaphragm. (Reprinted from Microwave Transmission Circuits, by George L. Ragan, 1st ed., 1948; by permission, McGraw-Hill Book Co., N.Y.)

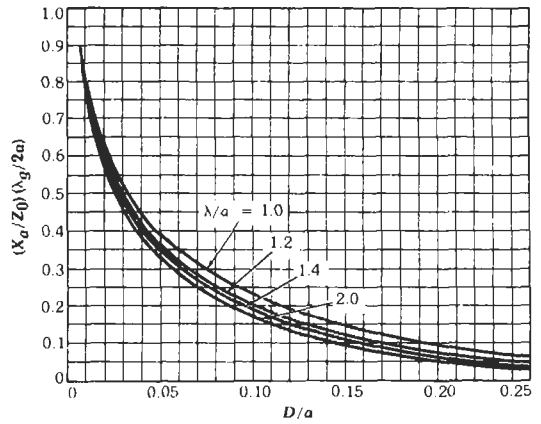
Distributed capacitances are found in junctions and slits, whereas tuning screws act as lumped capacitances.

HYBRID JUNCTIONS*

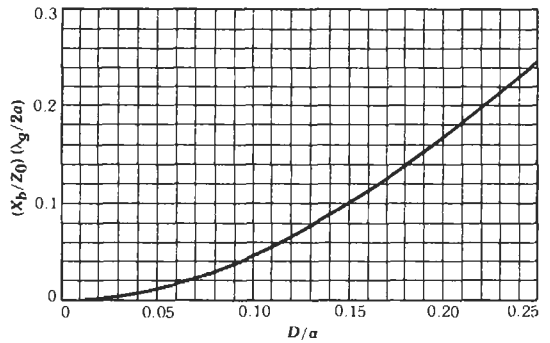
The hybrid junction is illustrated in various forms in Fig. 31. An ideal junction is characterized by the fact that there is no direct coupling between arms 1 and 4 or between 2 and 3. Power flows from 1 to 4 only by virtue of reflections in arms 2 and 3. Thus, if arm 1 is excited, the voltage arriving at arm 4 is

$$E_4 = \frac{1}{2}E_1[\Gamma_2 \exp(j2\theta_2) - \Gamma_3 \exp(j2\theta_3)]$$

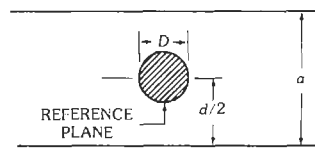
* C. G. Montgomery, R. H. Dicke, and E. M. Purcell. Principles of Microwave Circuits. New York: McGraw-Hill Book Co., 1948; Chapter 9.



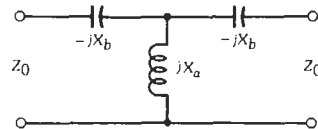
(A) Shunt reactance characteristic.



(B) Series reactance characteristic.



(C) Physical dimensions.



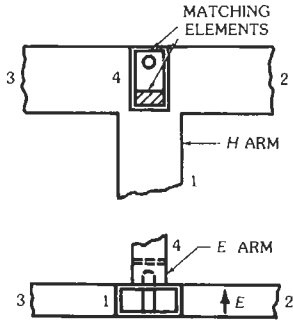
(D) Electrical equivalent.

Fig. 30. Equivalent circuit for inductive cylindrical post.

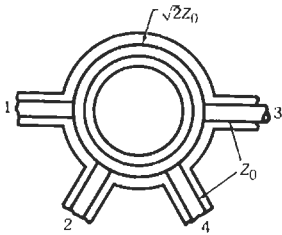
and the reflected voltage in arm 1 is

$$E_{r1} = \frac{1}{2}E_1[\Gamma_2 \exp(j2\theta_2) + \Gamma_3 \exp(j2\theta_3)]$$

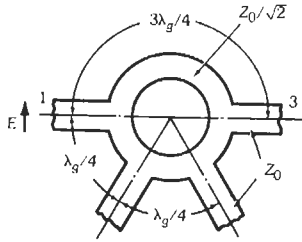
where,



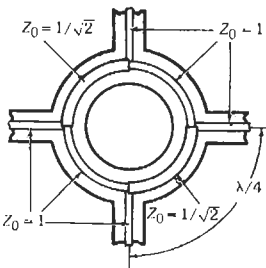
(A) Waveguide hybrid junction (magnetic T).



(B) Shunt coaxial hybrid ring.



(C) E-plane waveguide hybrid ring.



(D) Symmetrical coaxial hybrid.

Fig. 31. Hybrid junctions.

E_1 is the amplitude of the incident wave, Γ_2 and Γ_3 are the reflection coefficients of the terminations of arms 2 and 3,

θ_2 and θ_3 are the respective distances of the terminations from the junctions.

In the case of the rings, θ is the distance between the arm-and-ring junction and the termination.

If the decoupled arms of the hybrid junction are independently matched and the other arms are terminated in their characteristic impedances, then all four arms are matched at their inputs.

RESONANT CAVITIES

A cavity enclosed by metal walls has an infinite number of natural frequencies at which resonance will occur. One of the more common types of cavity resonators is a length of transmission line (coaxial or waveguide) short-circuited at both ends.

Resonance occurs when

$$2h = l(\lambda_g/2)$$

where,

- l = an integer,
- $2h$ = length of the resonator,
- λ_g = guide wavelength in resonator
 $= \lambda / [\epsilon - (\lambda/\lambda_c)^2]^{1/2}$
- λ = free-space wavelength,
- λ_c = guide cutoff wavelength,
- ϵ = relative dielectric constant of medium in cavity.

For $TE_{m,n}$ or $TM_{m,n}$ waves in a rectangular cavity with cross section a, b

$$\lambda_c = 2 / [(m/a)^2 + (n/b)^2]^{1/2}$$

where m and n are integers.

For $TE_{m,n}$ waves in a cylindrical cavity

$$\lambda_c = 2\pi a / U_{m,n}'$$

where a is the guide radius and $U_{m,n}'$ is the m th root of the equation $J_n'(U) = 0$.

For $TM_{m,n}$ waves in a cylindrical cavity

$$\lambda_c = 2\pi a / U_{m,n}$$

where a is the guide radius and $U_{m,n}$ is the m th root of the equation $J_n(U) = 0$.

For TM waves, $l = 0, 1, 2, \dots$

For TE waves, $l = 1, 2, \dots$, but not 0.

Rectangular Cavity of Dimensions $a, b, 2h$

$$\lambda = 2 / [(l/2h)^2 + (m/a)^2 + (n/b)^2]^{1/2}$$

where only one of l, m, n may be zero.

Cylindrical Cavities of Radius a and Length $2h$

$$\lambda = 1/[(1/4h)^2 + (1/\lambda_c)^2]^{1/2}$$

where λ_c is the guide cutoff wavelength.

Spherical Resonators of Radius a

$$\lambda = 2\pi a/U_{m,n} \text{ for a TE wave}$$

$$\lambda = 2\pi a/U_{m,n}' \text{ for a TM wave}$$

Values of $U_{m,n}$:

$$U_{1,1} = 4.5, U_{2,1} = 5.8, U_{1,2} = 7.64$$

Values of $U_{m,n}'$:

$$U_{1,1}' = 2.75 = \text{lowest-order root}$$

Additional Cavity Equations

Note that resonant modes are characterized by three subscripts in the mode designations of Table 10, Chart 1, and Fig. 32.

Fig. 32 is a mode chart for a right-circular-cylindrical resonator, showing the distribution of resonant modes with frequency as a function of cavity shape. With the aid of such a chart, the various possible resonances can be predicted as the length ($2h$) of the cavity is varied by a movable piston.

Effect of Temperature and Humidity on Cavity Tuning

The resonant frequency of a cavity changes with temperature and humidity (due to changes in dielectric constant of the atmosphere) and with thermal expansion of the cavity. A homogeneous cavity made of one kind of metal will have a thermal-tuning coefficient equal to the linear coefficient of expansion of the metal (Table

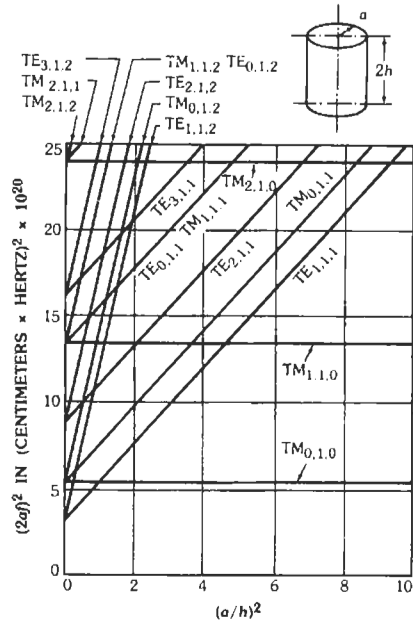


Fig. 32. Mode chart for right-circular-cylindrical cavity. (Reprinted from Techniques of Microwave Measurements, by Carol G. Montgomery, 1st ed., 1947; by permission, McGraw-Hill Book Co., N.Y.)

11), since the frequency is inversely proportional to the linear dimension of the cavity.

The relative dielectric constant of air (vacuum = 1) is given by

$$k_e = 1 + 210 \times 10^{-6}(P_a/T) + 180 \times 10^{-6}[1 + (5580/T)](P_w/T)$$

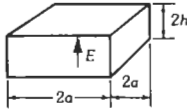
where P_a and P_w are partial pressures of air and water vapor in millimeters of mercury and T is the absolute temperature. Fig. 33 is a nomograph showing change of cavity tuning relative to conditions at 25 degrees Celsius and 60 percent relative humidity (expansion is not included).

TABLE 10. EQUATIONS FOR A RIGHT-CIRCULAR-CYLINDRICAL CAVITY

Mode	λ_0 Resonant Wavelength	Q (all dimensions in same units)
TM _{0,1,1} (E_0)	$\frac{4}{[(1/h)^2 + (2.35/a^2)]^{1/2}}$	$(\lambda_0/\delta) (a/\lambda_0)[1 + (a/2h)]^{-1}$
TE _{0,1,1} (H_0)	$\frac{4}{[(1/h)^2 + (5.93/a^2)]^{1/2}}$	$(\lambda_0/\delta)(a/\lambda_0) \left[\frac{1 + 0.168 (a/h)^2}{1 + 0.168 (a/h)^3} \right]$
TE _{1,1,1} (H_1)	$\frac{4}{[(1/h)^2 + (1.17/a^2)]^{1/2}}$	$(\lambda_0/\delta) (h/\lambda_0) \left[\frac{2.39h^2 + 1.73a^2}{3.39 (h^3/a) + 0.73ah + 1.73a^2} \right]$

CHART 1. CHARACTERISTICS OF VARIOUS TYPES OF RESONATORS

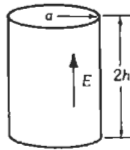
Square Prism $TE_{1,0,1}$



$$\lambda_0 = 2.83a$$

$$Q = (0.353\lambda\delta)[1 + (0.177\lambda/h)]^{-1}$$

Circular Cylinder $TM_{0,1,0}$



$$\lambda_0 = 2.61a$$

$$Q = (0.383\lambda\delta)[1 + (0.192\lambda/h)]^{-1}$$

Sphere



$$\lambda_0 = 2.28a$$

$$Q = 0.318(\lambda\delta)$$

Sphere With Cones

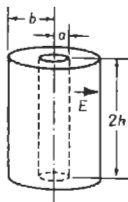


$$\lambda_0 = 4a$$

Optimum Q for $\theta = 34^\circ$

$$Q = 0.1095(\lambda\delta)$$

Coaxial TEM



$$\lambda_0 = 4h$$

Optimum Q for $(b/a) = 3.6$

($Z_0 = 77$ ohms)

$$\lambda[4\delta + 7.2(h\delta/b)]$$

Skin depth in meters = $\delta = (10^7/2\pi\omega\sigma)^{1/2}$, where σ = conductivity of wall in mhos/meter and $\omega = 2\pi \times$ frequency.

TABLE 11. LINEAR COEFFICIENTS OF EXPANSION

Metal	Linear Coefficient of Expansion/ $^\circ\text{C}$
Yellow brass	20
Copper	17.6
Mild steel	12
Invar	1.1

$\times 10^{-6}$

Coupling to Cavities and Loaded Q

Near resonance, a cavity may be represented as a simple shunt-resonant circuit, characterized by a loaded $Q = Q_l$, where $1/Q_l = (1/Q_0) + (1/Q_{ext})$, Q_0 is the unloaded Q characteristic of the cavity itself, and $1/Q_{ext}$ is the loading due to the external circuits. The variation of Q_{ext} with size of the coupling is approximately as given in Table 12.

Equations for Coupling Through a Cavity

Table 13 summarizes some of the useful relationships in a four-terminal cavity (transmission type) for three conditions of coupling: matched input (input resistance at resonance equals Z_0 of input line), equal coupling ($1/Q_{in} = 1/Q_{out}$), and matched output (resistance seen looking into output terminals at resonance equals output-load resistance). A matched generator is assumed.

In the table, g_c' is the apparent conductance of the cavity at resonance, with no output load; the transmission T is the ratio of the actual output-circuit power delivered to the available power from the matched generator. The loaded Q is Q_l and unloaded Q is Q_0 .

Cavity Coupling Techniques*

To couple power into or out of a resonant cavity, either waveguide or coaxial loops, probes, or apertures may be used.

The essentially inductive loop (a certain amount of electric-field coupling exists) is inserted in the resonator at a desired point where it can couple to a strong magnetic field. The degree of coupling may be con-

* C. G. Montgomery, R. H. Dicke, and E. M. Purcell. *Principles of Microwave Circuits*. New York: McGraw-Hill Book Co., 1948; Chapter 7.

TABLE 12. VARIATION OF Q_{ext}

Coupling	$1/Q_{ext}$ is Proportional to
Small round hole	(Diameter) ⁶
Symmetrical inductive diaphragm	(δ) ⁴ (See Fig. 28)
Small loop	(Diameter) ⁴

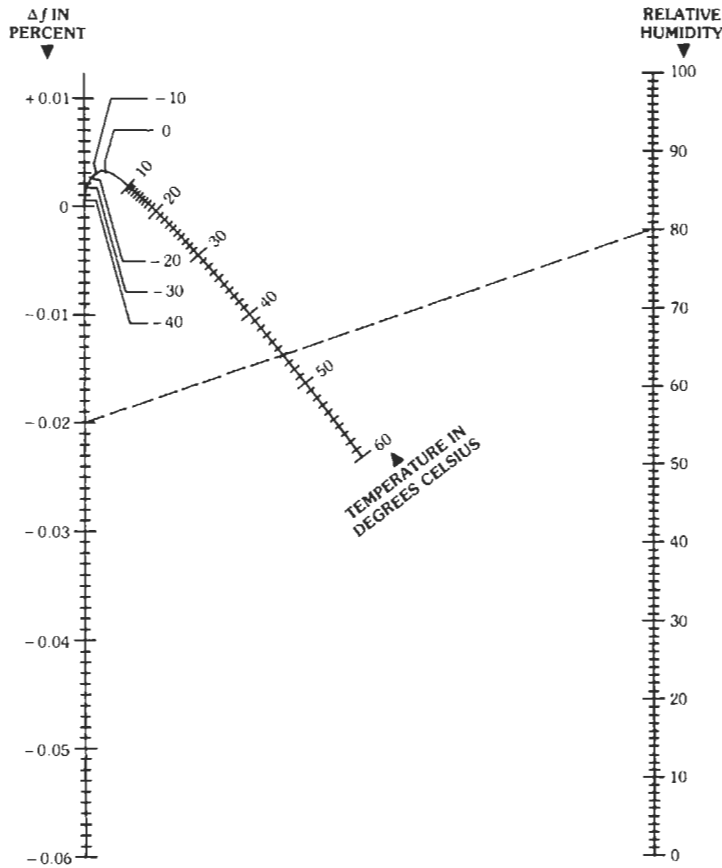


Fig. 33. Effect of temperature and humidity on cavity tuning. (Reprinted from Techniques of Microwave Measurements, by Carol G. Montgomery, 1st ed., 1947; by permission, McGraw-Hill Book Co., N.Y.)

TABLE 13. COUPLING THROUGH A CAVITY

	Matched Input	Equal Coupling	Matched Output
Input Standing-Wave Ratio	1	$1 + g_c' = 2(T^{-1/2} - 1)$	$1 + 2g_c'$
Transmission Ratio = T	$1 - g_c' = 1 - 2\rho$	$(1 + g_c'/2)^{-2} = (1 - \rho)^2$	$(1 + g_c')^{-1} = 1 - 2\rho$
$Q_l/Q_0 = \rho$	$\frac{1}{2}g_c' = \frac{1}{2}(1 - T)$	$[g_c'/(2 + g_c')] = 1 - (T^{1/2})$	$[g_c'/2(1 + g_c')] = \frac{1}{2}(1 - T)$

trolled by rotating the loop so that more or less loop area links this field. For a fixed location of the loop, the loaded Q of a loop-coupled coaxial resonator varies as the square of the effective loop area and inversely as the square of the distance of the loop center from the resonator axis of revolution.

The off-resonance input impedance of the loop is low, a feature that sometimes is helpful in series connections.

The capacitive probe is inserted in the resonator at a point where it is parallel to and can couple to strong electric fields. The degree of coupling is controlled by adjusting the length of the probe relative to the electric field.

The off-resonance input impedance of the probe-coupled resonator is high; this property is useful in parallel connections.

Aperture coupling is suitable when coupling wave-

guides to resonators or in coupling resonators together. In this case, the aperture must be located and shaped to excite the proper propagating modes.

For all means of coupling, the input impedance at resonance and the loaded Q may be adjusted by proper selection of the point of coupling and the degree of coupling.

Simple Waveguide Cavity*

A cavity may be made by enclosing a section of waveguide between a pair of large shunt susceptances, as shown in Fig. 34. Its loaded Q is given by

$$Q_l = \frac{1}{4}(\lambda_g/\lambda)^2(b^4 + 4b^2)^{1/2}\tan^{-1}(2/b)$$

and the resonant guide wavelength, λ_{g0} , is obtained from

$$2\pi l/\lambda_{g0} = \tan^{-1}(2/b)$$

Resonant Irises

Resonant irises may be used to obtain low values of loaded Q (< 30). The simplest type is shown in Fig. 35. It consists of an inductive diaphragm and a capacitive screw located in the same plane across the waveguide. For $Q_l < 50$, the losses in the resonant circuit may be ignored and

$$1/Q_l \approx 1/Q_{ext}$$

To a good approximation, the loaded Q (matched load and matched generator) is given by

* G. L. Ragan. *Microwave Transmission Circuits*. New York: McGraw-Hill Book Co., 1948; Chapter 10.

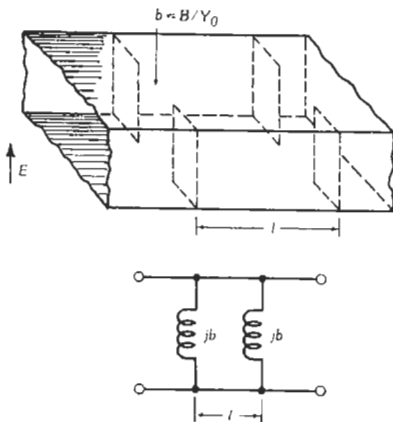


Fig. 34. Waveguide cavity and equivalent circuit.

$$Q_l = (B_l/2Y_0)(\lambda_{g0}/\lambda)^2$$

where B_l is the susceptance of the inductive diaphragm. This value may be taken from charts such as Figs. 28 and 29 as a starting point, but because of the proximity of the elements, the susceptance value is modified. Exact Q s must be obtained experimentally. Other resonant structures are given in Figs. 36 and 37. These are often designed so that the capacitive gap will break down under high power levels for use as transmit-receive (TR) switches in radar systems.

SURFACE-WAVE TRANSMISSION LINE*

The surface-wave transmission line is a single-conductor line having a relatively thick dielectric sheath

* G. Goubau. "Designing Surface-Wave Transmission Lines." *Electronics*, Vol. 27, April 1954, pp. 180-184.

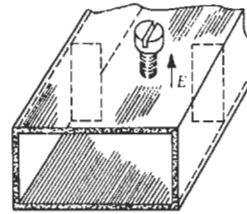


Fig. 35. Resonant iris in waveguide. Capacitive screw is tuned to resonance with inductive diaphragm.

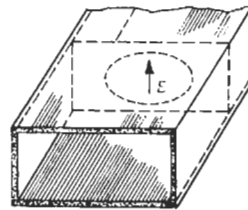


Fig. 36. Resonant element consisting of an oblong aperture in a thin transverse diaphragm.

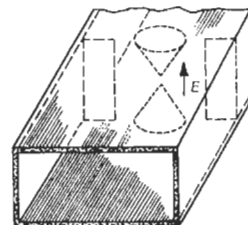


Fig. 37. Resonant structure consisting of cones with capacitive gap between apexes and thin symmetrical inductive diaphragm.

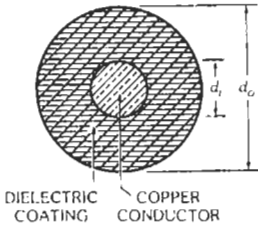


Fig. 38. Cross section of surface-wave transmission line.

(Fig. 38). The sheath diameter is often three or more times the conductor diameter. A mode of propagation that is practically nonradiating is excited on the line by means of a conical horn at each end as shown in Fig. 39. The mouth of the horn is roughly one-quarter to one wavelength in diameter. Losses are about half those of a two-wire line, but the surface-wave line has a practical lower frequency limit of about 50 megahertz. Design charts are given in Figs. 40, 41, and 42.



Fig. 39. Surface-wave transmission line with launchers at each end. (Courtesy Electronics.)

The losses in the two launchers combined vary from less than 0.5 decibel to a little more than 1.0 decibel, according to their design.

Conductor loss L_c by the equation below is 5 percent over the theoretical value for pure copper. Dielectric loss L_p for polyethylene at 100 megahertz is shown in Fig. 41. For other dielectrics and frequencies, find L_i by the equation.

$$L_c = 0.455f^{1/2}/Zd_i \text{ decibels/100 feet}$$

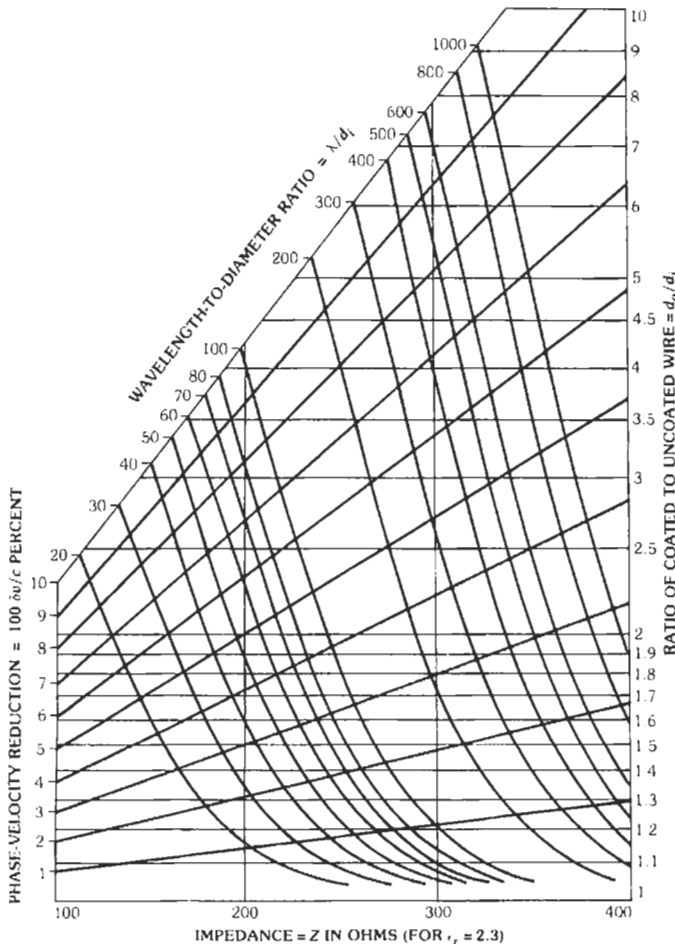


Fig. 40. Relationship among wire diameter, dielectric layer, phase-velocity reduction, and impedance (for brown polyethylene). (Courtesy Electronics.)

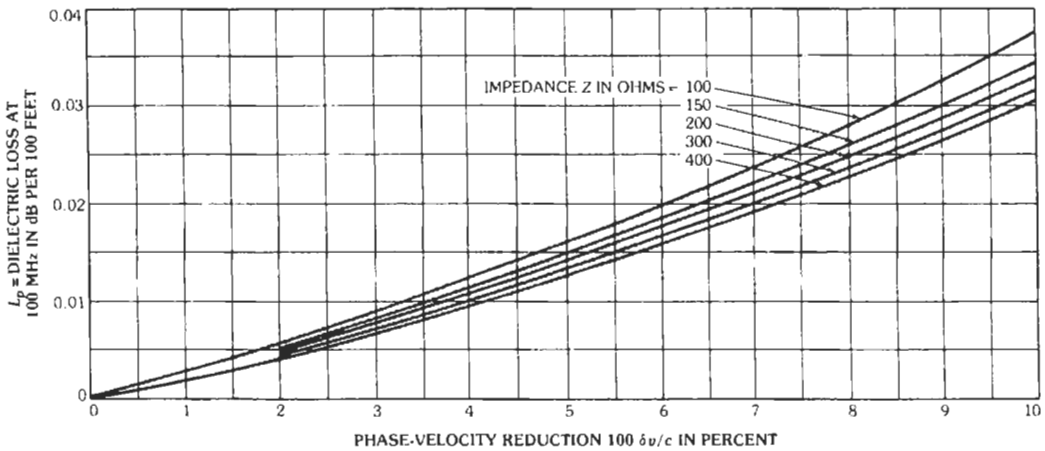


Fig. 41. Dielectric loss at 100 MHz for brown polyethylene ($\epsilon_r = 2.3$ and $F_p = 5 \times 10^{-4}$). (Courtesy Electronics.)

$$L_i = 26F_p L_p / (\epsilon_r - 1) \text{ decibels/100 feet}$$

$$L_i = L_p f / 100$$

for brown polyethylene (Fig. 41).

Symbols

- c = velocity of propagation in free space
- d_i = diameter of the conductor (inches in equation for L_c)
- d_o = outside diameter of the dielectric coating
- f = frequency in megahertz
- F_p = power factor of dielectric
- L_c = conductor loss in decibels/100 feet
- L_i = dielectric loss in decibels/100 feet
- L_p = dielectric loss shown in Fig. 41

- Z = waveguide impedance in ohms
- δv = reduction in phase velocity
- ϵ_r = dielectric constant relative to air
- λ = free-space wavelength

Example: At 900 megahertz ($\lambda = 0.333$ meter), a 200-foot line is required having a permissible loss of 1.0 decibel/100 feet (not including the launcher losses). What are its dimensions?

If 20 percent is allowed for dielectric loss, the conductor loss would be $L_c = 0.8$ decibel/100 feet. Assuming $Z = 250$ ohms as a first approximation, the formula for L_c gives $d_i = 0.068$ inch. Use No. 14 AWG wire ($d_i = 0.064$ and $\lambda/d_i = 204$). Now going to Fig. 40 and assuming that $100\delta v/c = 6$ percent is adequate, we find that $d_o/d_i = 3$ and $Z = 270$ ohms.

Recomputing, $L_c = 0.79$ decibel/100 feet. By Fig. 41, $L_p = 0.017$ at 100 megahertz for brown polyethylene. For the same material at 900 megahertz, the loss is $L_i = 0.15$ decibel/100 feet.

For 200 feet, the combined conductor and dielectric loss is 1.9 decibels, to which must be added the loss of 0.5 to 1.0 decibel total for the two launchers.

Dielectric Other Than Polyethylene

Determine Z and $\delta v/c$ for polyethylene ($\epsilon_r = 2.3$) from Fig. 40. Then use Fig. 42 to find the value of d_o/d_i required for the same performance with actual dielectric constant ϵ_r . Make a computation of the new dielectric loss, using Fig. 41 and the equation for L_i .

DIELECTRIC-ROD WAVEGUIDES

The dielectric-rod waveguide has applications in antenna structures, laser devices, fiber optics, and millimetric-wave techniques.

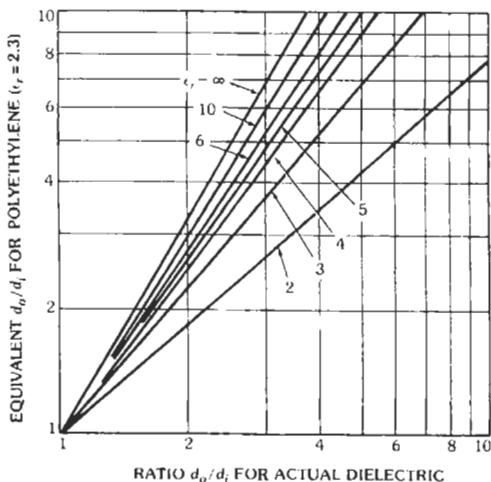


Fig. 42. Conversion chart for dielectric other than polyethylene. (Courtesy Electronics.)

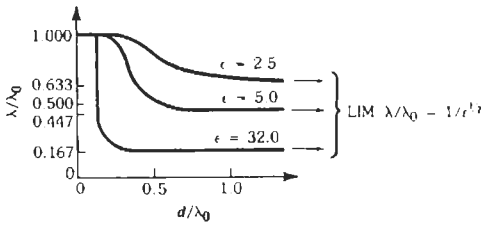


Fig. 43. Wavelength of $HE_{1,1}$ mode as a function of d/λ_0 . (From D. G. Kieley, "Dielectric Aerials," Methuen's Monographs on Physical Subjects. New York: John Wiley & Sons, Inc., 1953; p. 27.)

The field structures for the nonradiating modes fall into two classes—circularly symmetric and nonsymmetric modes. The cutoff wavelengths (λ_c) for symmetric modes are (for the $E_{0,m}$ and $H_{0,m}$ modes)*

$$\lambda_c = \pi d(\epsilon - 1)^{1/2}/j_{0,m}$$

where,

d = rod diameter,

ϵ = relative dielectric constant,

$j_{0,m}$ = m th root of $J_0(X)$.

Analysis of the field equations reveals the necessary coexistence of an E wave with an H wave to obtain a nonsymmetric field structure.†

These modes are described as HE if the H mode is predominant, and as EH if the E mode predominates. The special case of the $HE_{1,1}$ mode, referred to as the "dipole" mode because of the resemblance of the transverse-electric-field pattern to that of the electrostatic dipole, is of special interest because it has zero cutoff frequency.‡

Fig. 43 describes the relation between λ/λ_0 and d/λ_0 for rods of different ϵ (λ = operating wavelength, λ_0 = free-space wavelength, d = rod diameter), and the field structure is shown in Fig. 44.

The attenuation of the $HE_{1,1}$ mode, for material having relatively low loss,§ is found from

$$\alpha = 27.3(\epsilon/\lambda_0)R \tan\delta \text{ (dB/cm)}$$

where,

ϵ = relative dielectric constant,

$\tan\delta$ = loss tangent of dielectric,

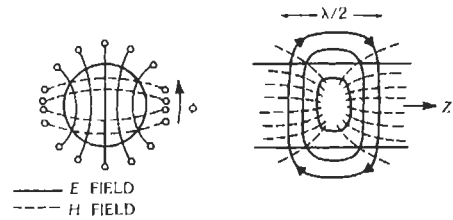


Fig. 44. $HE_{1,1}$ mode, field distribution. (From D. G. Kieley, "Dielectric Aerials," Methuen's Monographs on Physical Subjects. New York: John Wiley & Sons, Inc., 1953; p. 27.)

R = attenuation factor (dimensionless).

Note that this closely resembles the expression for TEM-mode propagation in a low-loss dielectric medium,** given by

$$\alpha = [27.3(\epsilon)^{1/2}/\lambda] \tan\delta \text{ (dB/cm)}$$

For d/λ_0 larger than 0.8, $R \approx [1/(\epsilon)^{1/2}]$. For other values see Elsasser.§

The $HE_{1,1}$ -mode waveguide is impractical for many frequencies in the UHF and lower bands because of the size of the rod and field spread adjacent to the rod outside of the dielectric. At higher microwave and millimeter-wave frequencies, lack of suitable material prevents use of this rod waveguide for transmission over long distances. Some attempts are being made to use this structure as a flexible waveguide by covering the rod with a cladding dielectric having a lower dielectric constant and further with a lossy jacket. The structure then resembles that of an optical fiber.

The optical fiber is now widely used in practice due to its extremely low propagation loss. More detailed information on optical fibers can be found in Chapter 22.

RECTANGULAR DIELECTRIC GUIDES AND IMAGE GUIDES

Rectangular dielectric waveguides (Fig. 45) find use in integrated optics and in millimeter-wave integrated circuits. In the latter case, image guides (Fig. 46) are more frequently used. Modes are classified into E_{pq}^y and E_{pq}^x . The former has a principal E field in the y direction and the latter in the x direction. The subscripts p and q indicate the number of transverse field maxima in the x and y directions.

For a well-guided mode, approximate expressions for these modes in the guide in Fig. 45 are††

* H. M. Barlow and J. Brown. *Radio Surface Waves*. Oxford at the Clarendon Press, 1962; p. 71.

† D. G. Kieley. "Dielectric Aerials," *Methuen's Monographs on Physical Subjects*. New York: John Wiley & Sons, Inc., 1953; pp. 7-29.

‡ H. M. Barlow and J. Brown. *Radio Surface Waves*. Oxford at the Clarendon Press, 1962; p. 69.

§ W. M. Elsasser. "Attenuation in a Dielectric Rod." *Journal of Applied Physics*, Vol. 20, December 1949, pp. 1193-1196.

** G. L. Ragan. *Microwave Transmission Circuits*, 1st ed. New York: McGraw-Hill Book Co., 1948; p. 29.

†† E. A. J. Marcatili. "Dielectric Rectangular Waveguide and Directional Coupler for Integrated Optics." *Bell Syst. Tech. J.*, Vol. 48, 1969, pp. 2071-2102.

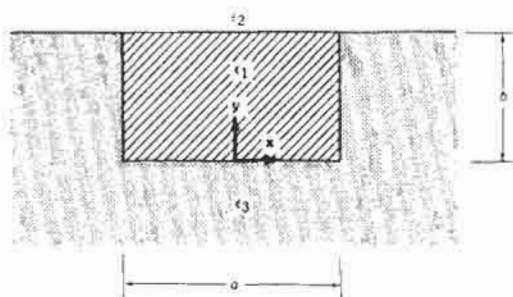


Fig. 45. Cross section of rectangular dielectric waveguide.

$$\beta = \{\epsilon_1 k^2 - (p\pi/a)^2 / (1 + 2A_3/\pi a)^2 - (q\pi/b)^2 / [1 + (\epsilon_2 A_2 + \epsilon_3 A_3) / \epsilon_1 \pi b]^2\}^{1/2}$$

for the E_{pq}^y mode; $p, q = 1, 2, \dots$

$$\beta = \{\epsilon_1 k^2 - (p\pi/a)^2 / (1 + 2\epsilon_3 A_3 / \epsilon_1 \pi a)^2 - (q\pi/b)^2 / [1 + (A_2 + A_3) / \pi b]^2\}^{1/2}$$

for the E_{pq}^x mode; $p, q = 1, 2, \dots$

where,

$$A_2 = \lambda_0 / 2 \sqrt{\epsilon_1 - \epsilon_2}$$

$$A_3 = \lambda_0 / 2 \sqrt{\epsilon_1 - \epsilon_3}$$

λ_0 is the free space wavelength,
 $\epsilon_1, \epsilon_2,$ and ϵ_3 are the relative dielectric constants
of the materials involved.

Dispersion characteristics are provided for a rectangular guide created in a glass substrate in normalized form (Fig. 47). These curves were derived with the above equations.

The above formulas are not very accurate near the cutoff frequencies. More accurate data are obtainable

** K. J. Button and J. C. Wiltse, eds. *Infrared and Millimeter Waves*, Vol. 4. New York: Academic Press, Inc., 1981; pp. 195-273.

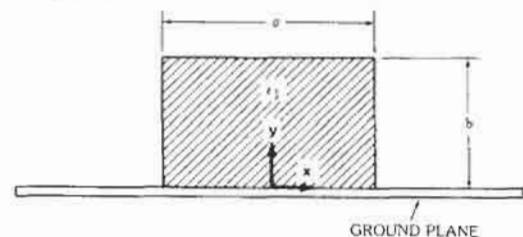


Fig. 46. Cross section of image guide.

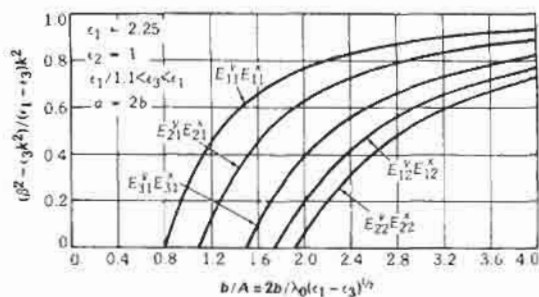


Fig. 47. Approximated dispersion characteristics of a rectangular dielectric waveguide.

by using a number of techniques, the simplest of which is the effective dielectric constant (EDC) method.**

As the image guide is usually surrounded by air and is considered one half of the dielectric rod in free space for the dominant E_{11}^y mode, the formulas given above and the EDC can be used for dispersion characteristics. Results calculated by the EDC are given in Fig. 48 for various dielectric materials when the aspect ratio is $a/b = 2$.

DIELECTRIC RESONATORS

A dielectric resonator usually consists of a high dielectric material such as barium-tetratitanate formed into a pillbox shape. The dominant resonant mode is $TE_{01\delta}$ and has magnetic fields in the radial and axial directions and an electric field in the circumferential direction. The fields do not vary in the circumferential direction for this mode.

The resonant frequencies have been calculated by several methods. Fig. 49 shows resonant frequency (in the form of the relationship of the resonant wavelength to the radius) versus the structural profile for the

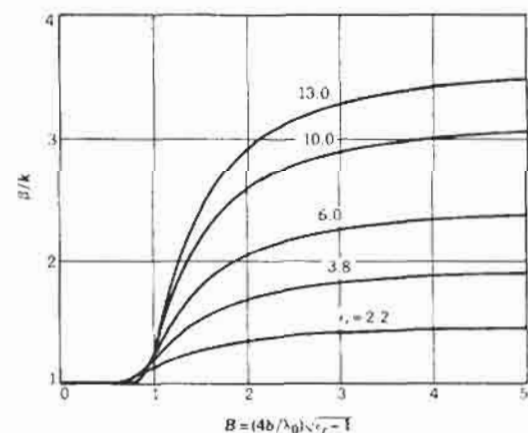


Fig. 48. Dispersion characteristics of image guide.

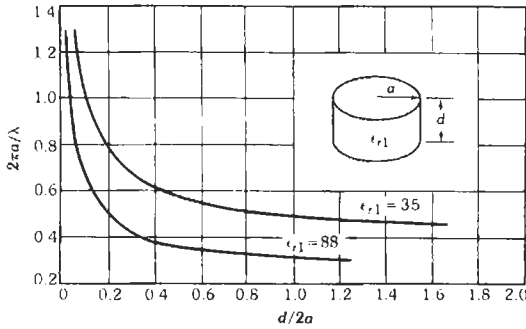


Fig. 49. Resonant frequency of the $TE_{01\delta}$ mode of a pillbox dielectric resonator. (After T. Itoh and R. Rudokus, "New Method for Computing the Resonant Frequencies of Dielectric Resonators." IEEE Trans. Microwave Theory and Techniques, Vol. MTT-25, Jan. 1977; pp. 32-54.)

resonator immersed in free space.* Unloaded Q on the order of several thousands has been obtained.†

SLOT LINES

Slot line consists of a narrow gap (or slot) in a conductive coating on a dielectric substrate, as shown in Fig. 50. The slot line normally uses a dielectric substrate of sufficiently high permittivity (e.g., $\epsilon_r = 16$) that the guide wavelength is much smaller than the free-space wavelength, and the fields are closely confined to the slot with negligible radiation loss. The nature of the slot-mode configuration is such that the electric field extends across the slot while the magnetic field is in a plane perpendicular to the slot and forms closed loops at half-wave intervals.

* Calculated by the method reported by T. Itoh and R. Rudokus in "New Method for Computing the Resonant Frequencies of Dielectric Resonators," IEEE Trans. Microwave Theory and Techniques, Vol. MTT-25, Jan. 1977, pp. 32-54.

† For the resonant frequencies and Q of the dielectric resonator placed in a microwave integrated circuit configuration, see M. Dydyk, "Apply High- Q Resonators to mm-Wave Microstrip," Microwaves, Vol. 19, Dec. 1980, pp. 62-63.

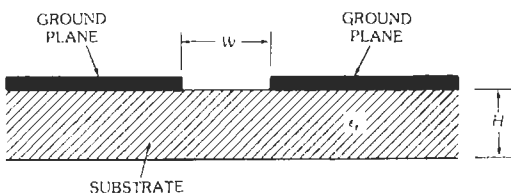
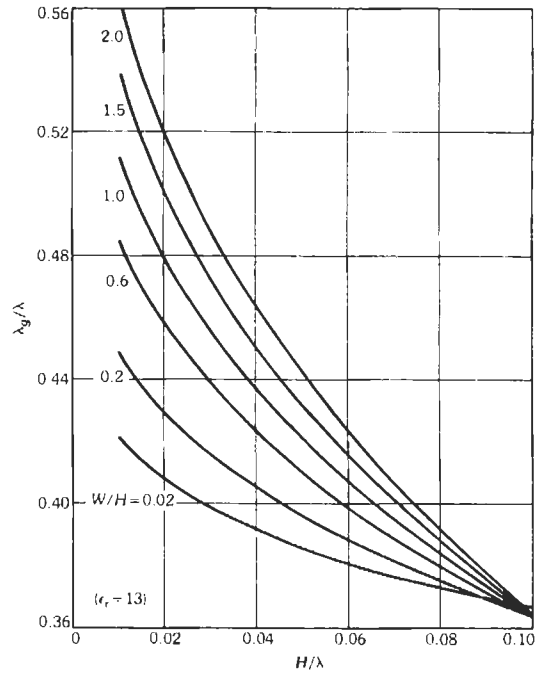
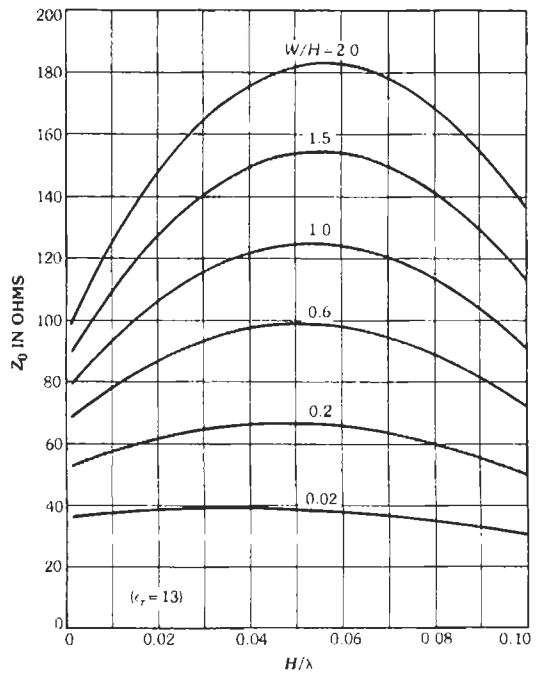


Fig. 50. Cross section of slot line.



(A) Relative guide wavelength.



(B) Characteristic Impedance.

Fig. 51. Characteristics of slot line.

Since the slot mode is a non-TEM wave, the definition of characteristic impedance is not unique. One definition could be

$$Z_0 = V^2/2P$$

where,

V is the peak voltage amplitude across the slot,
 P is the average power flow of the wave.

For the intervals given by $9.6 \leq \epsilon_r \leq 20$, $0.02 \leq W/H \leq 2.0$, and $0.015 \leq H/\lambda \leq 0.08$ and an infinitesimally thin ground plane, a closed-form approximation to λ_g/λ (where λ_g is the guide wavelength and λ is the free-space wavelength) can be expressed as*

$$\lambda_g/\lambda = f_1(\epsilon_r)\{f_2(W/h) + f_3(W/h)(H/\lambda)^{f_4(W/H)} + f_5(W/H)\}$$

* C. M. Krowne. "Approximations to Hybrid Mode Slot Line Behavior." *Electronics Letters*, Vol. 14, 13 April 1978, pp. 258-259.

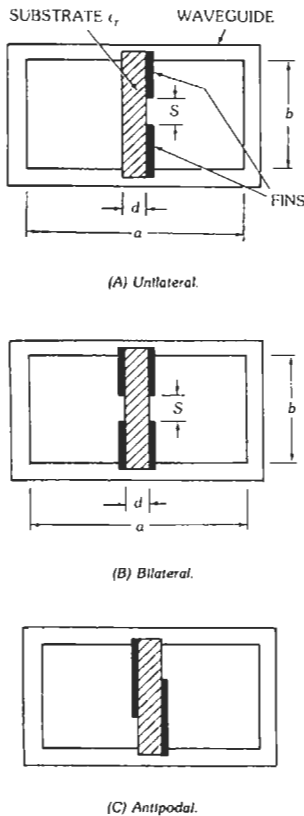
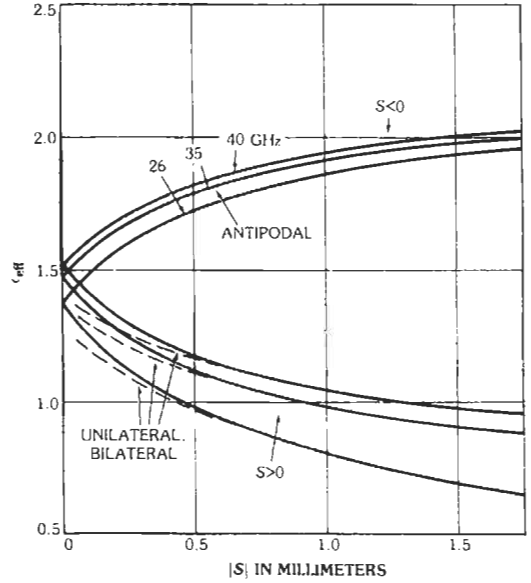
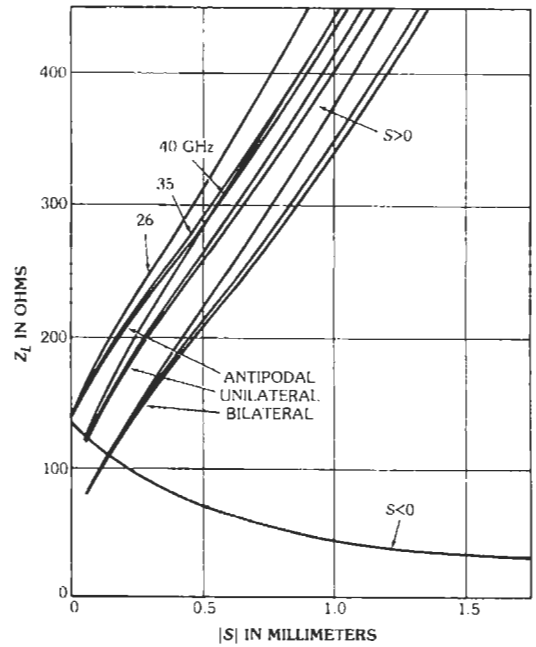


Fig. 52. Cross sections of fin-lines.



(A) Dielectric constant.



(B) Characteristic Impedance.

Fig. 53. Effective dielectric constant, $\epsilon_{eff} = (\lambda/\lambda_g)^2$, and characteristic impedance of fin-lines. (From Figs. 3 and 4 of H. Hofmann, "Calculation of Quasi-Planar Lines for mm-Wave Application," 1977 IEEE MTT-S International Microwave Symposium Digest, June 1977, pp. 381-384. © 1977 IEEE.)

$$f_1(\epsilon_r) = 3.549 \epsilon_r^{-0.56}$$

$$f_2(W/H) = 0.5632 (W/H)^{0.104(W/H)^{0.266}}$$

$$f_3(W/H) = -0.8777 (W/H)^{0.81} + 0.4233 (W/H) - 0.2492$$

$$f_4(W/H) = -1.269 \times 10^{-2} [\ln(50 W/H)]^{1.7} + 0.0674 \ln(50 W/H) + 0.20$$

$$f_5(W/H) = 1.906 \times 10^{-3} [\ln(50 W/H)]^{2.9} - 7.203 \times 10^3 \ln(50 W/H) + 0.1223$$

The accuracy of the expression is $\pm 3.7\%$
 A closed-form approximation for Z_0 is

$$Z_0 = (11/\epsilon_r)^{p(W/H, H/\lambda)} g(W/H, H/\lambda)$$

$$p(W/H, H/\lambda) = [-30.21 \ln(W/H) - 46.03](H/\lambda)^2 + [0.5073 \ln(W/H) + 3.358 (W/H) + 6.492] (H/\lambda) - 2.013 \times 10^{-2} \ln(W/H) - 0.1374(W/H) + 0.2365$$

$$g(W/H, H/\lambda) = [-1.176 \times 10^4 (W/H)^{0.502} - 6.311 \times 10^3 (W/H) - 162.7] (H/\lambda)^2 + [900.5(W/H)^{0.28} + 1262(W/H) - 123.8] (H/\lambda) + 1.637 \ln(50 W/H) + 40.99(W/H)^{0.46} + 30.96$$

Over the same parameter ranges as above, the accuracy of Z_0 is $\pm 14.5\%$. For $W/H \geq 0.2$, it is better than 4%. Data calculated by these formulas are plotted in Fig. 51.

FIN-LINES

Fin-lines (Fig. 52) consist of fins separated by a gap printed on one or both sides of a dielectric substrate that is in turn placed at the center of a rectangular waveguide along its E -plane. Therefore, fin-lines are considered printed versions of ridged waveguides and have single-mode operating bandwidth wider than the one for an enclosing waveguide itself. Fin-lines are widely used for millimeter-wave integrated circuits in the frequency range from 26.5 GHz to 100 GHz.

The effective dielectric constant, $\epsilon_{eff} = (\lambda/\lambda_g)^2$, and the characteristic impedance, Z_L , have been computed by a number of authors. Those by H. Hofmann are reproduced in Fig. 53 for K_a -band (26.5–40 GHz) application.* The characteristic impedance is defined as the ratio of the voltage across the slot to the current on the fins and is believed useful for small values of substrate thickness and gap width. Negative values of S indicate the overlap of fins in the antipodal fin-lines. The waveguides supporting the fins of Fig. 53 are WR-28 ($a = 7.112$ mm, $b = 3.556$ mm).

* H. Hofmann. "Calculation of Quasi-Planar Lines for mm-Wave Application." 1977 IEEE MTT-S International Microwave Symposium Digest, San Diego, June 1977, pp. 381–384.

31 Scattering Matrices

Georges A. Deschamps and John D. Dyson

Amplitude of a Traveling Wave	31-2
Reflection Coefficient	31-2
Definition	
Measurement	
Scattering Matrix of a Junction	31-3
Definition	
Properties	
Change of Terminal Plane	
Two-Port Junctions	31-3
Transformation Matrix	31-3
Measurement of the Scattering Matrix	31-4
Geometry of Reflection Charts	31-6
Conformal Charts	
Projective Chart	
Evaluation of Hyperbolic Distance	31-6
Problem A	
Problem B	
Problem C	
Correspondences With Current, Voltage, and Impedance	
Viewpoints	31-9
Normalized Current and Voltage	
Current and Voltage Not Normalized	
Normalized Impedance and Admittance	
Impedance and Admittance Matrix of a Junction	
Transformation Matrix	31-10

Microwave structures are characterized by dimensions that are of the order of the wavelength of the propagated signal. The notions of current, voltage, and impedance, useful at lower frequencies, have been successfully extended to these structures, but these quantities are not as directly available for measurement; there are no voltmeters or ammeters and no apparent "terminal pair" between which to connect them. The electromagnetic field itself, distributed throughout a region, becomes the relevant quantity.

Within uniform structures, which are the usual form of waveguides, the *power flow* and the *phase* of the field at a cross section are the quantities of importance. The most usual form of measurement, that of the standing-wave pattern in a slotted section, is easily interpreted in terms of *traveling waves* and gives directly the *reflection coefficient*. The scattering description of waveguide junctions was introduced* to express this point of view. It is not, however, restricted to microwaves; a low-frequency network can be considered as a "waveguide junction" between transmission lines connected to its terminal pairs, and the scattering matrix is a useful complement to the impedance and admittance descriptions.

AMPLITUDE OF A TRAVELING WAVE

In a uniform transmission line a traveling wave is characterized, for a given mode and frequency, by the electromagnetic-field distribution in a transverse cross section and by a propagation constant h . The field in any other cross section, at a distance z in the direction of propagation, has the same pattern but is multiplied by $\exp(-jhz)$. A wave propagating in the opposite direction, for the same mode and frequency, varies with z as $\exp(jhz)$. When losses are negligible, h is real.

The *amplitude* of a traveling wave, at a given cross section in the waveguide, is a complex number a defined as follows. The square $|a|^2$ of the magnitude of a is the power flow,** that is, the integral of the Poynting vector over the waveguide cross section. The phase angle of a is that of the transverse field in the cross section.†

The amplitude of a given traveling wave varies with z as $\exp(-jhz)$.

The wave amplitude has the dimensions of the square root of a power. The meter-kilogram-second unit is therefore the (watt)^{1/2}.

* C. G. Montgomery, R. H. Dicke, and E. M. Purcell, *Principles of Microwave Circuits* (New York: McGraw-Hill Book Co., 1948).

** The amplitude is sometimes defined to make the power flow equal to $\frac{1}{2}|a|^2$ rather than to $|a|^2$. This would correspond to the use of peak values instead of root-mean-square values.

† This phase is well defined for a pure mode, since the field has the same phase everywhere in the cross section.

REFLECTION COEFFICIENT

Definition

At a cross section in a waveguide, the reflection coefficient W (also often represented by Γ) is the ratio of the amplitudes of the waves traveling respectively in the negative and positive directions.

The positive direction must be specified and is usually taken as toward the load. To give a definite phase to the reflection coefficient, a convention is necessary that describes how the phases of waves traveling in opposite directions are to be compared. The usual convention is to compare in the two waves the phases of the transverse electric-field vectors.‡

For a short-circuit, produced, for instance, by a perfect conducting plane placed across the waveguide, the reflection coefficient is $W = -1$. For an open circuit, it is $W = +1$; and for a matched load, it is $W = 0$.

When the cross section is displayed by z in the positive direction, the reflection coefficient W becomes

$$W' = W \exp(2jhz) \quad (\text{Eq. 1})$$

Measurement

In a slotted waveguide equipped with a sliding voltage probe,* the position of a maximum is one where the phase of the reflection coefficient is zero.

The ratio of the maximum to the minimum (the standing-wave ratio, or *swr*) is

$$(\text{swr}) = (1 + |W|)/(1 - |W|)$$

Therefore

$$W = [(\text{swr}) - 1]/[(\text{swr}) + 1] \quad (\text{Eq. 2})$$

is the value of W at the position of a maximum. At the position of a minimum, which is easier to locate in practice, the reflection coefficient is $[1 - (\text{swr})]/[1 + (\text{swr})]$.

At any other position, the value of W is obtained by applying Eq. (1). If the reflection coefficient is wanted in some waveguide connected to the slotted section, a good match must obtain at the transition, or a correction must be applied as explained later in problems A and B (pages 31-6 to 31-8).

Reflectometers that give the reflection coefficient by direct reading, or display it on a Smith chart, are in current use.

‡ The dual convention, based on the magnetic-field vector, would give the "current" reflection coefficient, equal to minus the "voltage" reflection coefficient. The latter is used almost exclusively, and the "voltage" is implicit.

* A probe that gives a reading proportional to the electric field.

SCATTERING MATRIX OF A JUNCTION

Definition

To define accurately the waves incident on a waveguide junction and those reflected (or scattered) from it, some reference locations must be chosen in the waveguides. These locations are called the ports† of the junction. In a waveguide that can support several propagating modes, there should be as many ports as there are modes. (These ports may or may not have the same physical location in the multimode waveguide.)

At each port *i* of a junction, consider the amplitude *a_i* of the incident wave traveling toward the junction, and the amplitude *b_i* of the scattered wave, traveling away from it. As a consequence of Maxwell's equations, there exists a linear relation between the *b_i* and the *a_i*. Considering the *a_i* (where *i* varies from 1 to *n*) as the components of a vector **a**, and the *b_i* as the components of a vector **b**, this relation can be expressed by

$$\mathbf{b} = \mathbf{S}\mathbf{a}$$

where **S** = (*s_{ij}*) as an *n* × *n* matrix called the *scattering matrix* of the junction.

The *s_{ii}* is the *reflection coefficient* looking into port *i*, and *s_{ij}* is the *transmission coefficient* from *j* to *i*, all other ports being terminated in matching impedances.

Properties

For a *reciprocal* junction, the transmission coefficient from *i* to *j* equals that from *j* to *i*; the matrix **S** is symmetrical.

$$\mathbf{S} = \tilde{\mathbf{S}}$$

where $\tilde{\mathbf{S}}$ denotes the transpose of **S**.

The total power incident on the junction is

$$|\mathbf{a}|^2 = \sum_{i=1}^n |a_i|^2$$

The total power scattered is

$$|\mathbf{b}|^2 = \sum_{i=1}^n |b_i|^2$$

For a lossless junction, these two powers are equal

$$|\mathbf{a}|^2 = |\mathbf{b}|^2$$

This implies that the matrix **S** is unitary (see "Matrix Algebra" in Chapter 47).

† At lower frequencies, for a network connecting transmission lines, a port is a terminal pair.

$$\mathbf{S}^\dagger = \mathbf{S}^{-1}$$

For a *passive* junction with losses, $|\mathbf{b}|^2 < |\mathbf{a}|^2$, hence the matrix $\mathbf{I} - \mathbf{S}\mathbf{S}^\dagger$ is positive definite.

Change of Terminal Plane

If the port in arm *i* is moved away from the junction by ϕ_i electrical radians, the scattering matrix becomes

$$\mathbf{S}' = \Phi\mathbf{S}\Phi \tag{Eq. 3}$$

where

$$\Phi = \begin{bmatrix} e^{-j\phi_1} & 0 & 0 & 0 & \dots \\ 0 & e^{-j\phi_2} & 0 & 0 & \dots \\ 0 & 0 & e^{-j\phi_3} & 0 & \dots \\ \cdot & \cdot & \cdot & \cdot & \dots \\ \cdot & \cdot & \cdot & \cdot & \dots \\ \cdot & \cdot & \cdot & \cdot & \dots \end{bmatrix} \tag{Eq. 4}$$

TWO-PORT JUNCTIONS

The two-port junction includes the case of an obstacle or discontinuity placed in a waveguide as well as that of two essentially different waveguides connected to each other.

If reciprocity applies, the scattering matrix

$$\mathbf{S} = \begin{bmatrix} s_{11} & s_{12} \\ s_{21} & s_{22} \end{bmatrix} \tag{Eq. 5}$$

is symmetrical

$$s_{21} = s_{12}$$

For a lossless junction, the scattering coefficients can be expressed by

$$\begin{aligned} s_{11} &= + \tanh(u/2) \exp(-2j\alpha) \\ s_{22} &= - \tanh(u/2) \exp(-2j\beta) \\ s_{12} &= + \operatorname{sech}(u/2) \exp[-j(\alpha + \beta)] \end{aligned} \tag{Eq. 6}$$

in terms of three parameters, *u*, α , and β .

This corresponds to the representation of the junction by an ideal transformer with transformer ratio $n = \exp(-u/2)$, of hyperbolic amplitude *u*, placed between two sections of transmission line with electrical lengths α and β , respectively.

The quantity $-20 \log_{10} |s_{12}|$ is the insertion loss.

TRANSFORMATION MATRIX

To find the effect of successive obstacles in a waveguide or to combine two-port junctions placed in cascade, it is convenient to introduce the wave transformation matrix **T**. This matrix **T** relates the traveling

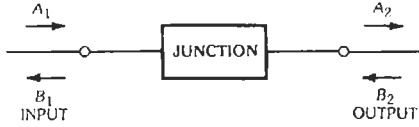


Fig. 1. Convention for wave transformation matrix T .

waves on one side of the junction to those on the other side. Using the notations of Fig. 1

$$\begin{bmatrix} A_1 \\ B_1 \end{bmatrix} = T \begin{bmatrix} A_2 \\ B_2 \end{bmatrix} \quad (\text{Eq. 7})$$

The 2×2 transformation matrix T may be deduced from the scattering matrix S

$$T = s_{21}^{-1} \begin{bmatrix} 1 & -s_{22} \\ s_n & -\det S \end{bmatrix} \quad (\text{Eq. 8})$$

Conversely, if $T = (t_{ij})$, the scattering matrix is

$$S = t_{11}^{-1} \begin{bmatrix} t_{21} & \det T \\ 1 & -t_{12} \end{bmatrix} \quad (\text{Eq. 9})$$

When reciprocity applies to the junction

$$\det T = s_{12}/s_{21} \quad (\text{Eq. 10})$$

becomes unity.

The input reflection coefficient $W' = B_1/A_1$ is related to the load reflection coefficient $W = B_2/A_2$ by

$$W' = (t_{21} + t_{22}W)/(t_{11} + t_{12}W) \quad (\text{Eq. 11})$$

$$= s_{11} + [s_{12}^2W/(1 - s_{22}W)] \quad (\text{Eq. 12})$$

When a number of junctions, 1, 2, 3, are placed in cascade (Fig. 2), the output port of each of them being the input port of the following one, the resulting junction has the transformation matrix

$$T = T_1 T_2 T_3$$

If n similar junctions with transformation matrix T are placed in cascade, the resulting transformation matrix is T^n .

Letting trace $T = t_{11} + t_{22} = 2 \cos \theta$

$$T^n = (\sin n\theta/\sin \theta) T - [\sin (n - 1) \theta/\sin \theta] \quad (\text{Eq. 13})$$

MEASUREMENT OF THE SCATTERING MATRIX

The measurement of the scattering parameters of the junction by conventional techniques is covered in Chapter 12. However, there may be occasions when it is not possible, or not desirable, to connect equipment to the

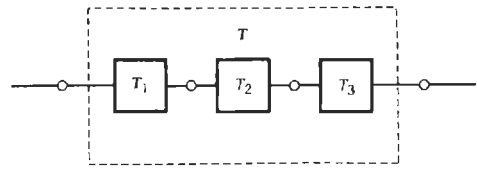


Fig. 2. Junctions in cascade.

output ports for the measurement of the transmission coefficients s_{ij} . Under these conditions, it is still possible to determine the coefficients s_{ii} and s_{ij} from measurements of the reflection coefficient at one port.

For any load with reflection coefficient W , placed on side 2 of a junction, the input reflection coefficient W' on side 1 can be measured. Coefficient W' is called the *image* of W , and they are related by the transformation

$$W' = s_{11} + [s_{12}s_{21}W/(1 - s_{22}W)] \quad (\text{Eq. 14})$$

The measurement of the images of three known loads, W_1, W_2, W_3 , will provide sufficient information to determine, either analytically or graphically, the coefficients s_{11} and s_{22} and the product $s_{12}s_{21}$.

Considering an analytical approach first, Eq. 14 can be expressed as:

$$W_i' = s_{11} + W_i W_i' s_{22} - W_i \Delta s \quad (\text{Eq. 15})$$

where,

$$\Delta s = s_{11}s_{22} - s_{12}s_{21}$$

$$i = 1, 2, 3$$

Equation 15 represents three linear equations in s_{11}, s_{22} , and Δs with solutions:

$$s_{11} = \{W_1' W_2 W_3 (W_2' - W_3') + W_2' W_3 W_1 (W_3' - W_1') + W_3' W_1 W_2 (W_1' - W_2')\}/D$$

$$s_{22} = -[W_1 (W_2' - W_3') + W_2 (W_3' - W_1') + W_3 (W_1' - W_2')]/D$$

$$\Delta s = -[W_1' W_1 (W_2' - W_3') + W_2' W_2 (W_3' - W_1') + W_3' W_3 (W_1' - W_2')]/D$$

$$D = [W_2 W_3 (W_2' - W_3') + W_3 W_1 (W_3' - W_1') + W_1 W_2 (W_1' - W_2')] \quad (\text{Eq. 16})$$

Once s_{11}, s_{22} , and $s_{12}s_{21}$ are known, the reflection coefficient W_L of an unknown termination can be determined from the measured image W_M by rewriting Eq. (14):

$$W_L = (W_M - s_{11})/[s_{12}s_{21} + s_{22}(W_M - s_{11})] \quad (\text{Eq. 17})$$

Solutions to Eq. 14 can also be obtained graphically.* Doing so provides an insight into the transformation of W into its image W' . The images of various known loads can be plotted on a reflection chart and the scattering coefficients deduced by the following procedures.

(A) With a matched load, one obtains directly s_{11} plotted as O' on Fig. 3. Point O' is called the iconocenter.

(B) With a sliding short-circuit on side 2, or any variable reactive load, the input reflection coefficient describes a circle Γ' , image of the unit circle Γ . This circle can be deduced from three or more measurements. Let C be its center and R its radius (Fig. 3). The magnitudes of the scattering coefficients result:

$$\begin{aligned} |s_{11}| &= OO' \\ |s_{22}| &= O'C/R \\ |s_{12}|^2 &= R(1 - |s_{22}|^2) \end{aligned} \quad (\text{Eq. 18})$$

The phases of these coefficients all follow from one more measurement as described in (C).

(C) The input reflection coefficient is measured with an open-circuit load placed at port 2, or for a short-circuit placed a quarter wave away from it. This may be one of the measurements taken in step (B). It gives (Fig. 4) the point P' , image of the point $P(W = +1)$.

A point P'' is constructed by projecting P' through O' onto Q on Γ' , then Q through C onto P'' on Γ' (Fig. 4). Then

$$\begin{aligned} \text{Phase of } s_{11} &= \text{angle } (OP, OO') \\ \text{Phase of } s_{22} &= \text{angle } (O'C, CP'') \\ \text{Phase of } s_{12} &= \frac{1}{2} \text{ angle } (OP, CP'') \end{aligned} \quad (\text{Eq. 19})$$

* G. A. Deschamps, "Determination of the Reflection Coefficients and Insertion Loss of a Waveguide Junction," *Journal of Applied Physics*, Vol. 24, pp. 1046-1050; August 1953; Also, *Electrical Communication*, Vol. 31, pp. 57-62; March 1954.

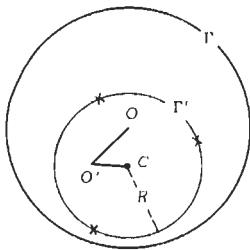


Fig. 3. Construction for the magnitudes of the scattering coefficients.

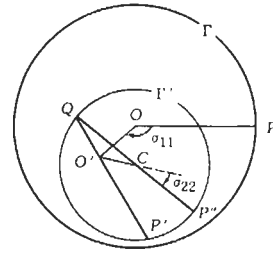


Fig. 4. Construction for the phases of the scattering coefficients.

(D) When no matched load is available, as was assumed in (A), the iconocenter O' may be obtained as in Fig. 5. Let P_1, P_2, P_3, P_4 represent the input reflection coefficients when a short circuit is placed successively at port 2 and at distances $\lambda/8, \lambda/4,$ and $3\lambda/8$ from it. These points define the circle Γ' [as in (B)], and the intersection I (the crossover point) of P_1P_3 and P_2P_4 may be used to find O' : draw perpendiculars to CI at points C and I up to their intersections with Γ' and get C' and I' ; then O' is the intersection of CI and $C'I'$.

The point P_3 is identical to P' in (C) above; hence the four measurements give the complete scattering matrix by constructing P'' and applying Eqs. 18 and 19.

(E) The construction of O' in (D) above is valid with any sliding load not necessarily reactive. Taking a load with small standing-wave ratio increases the accuracy of the construction.

(F) When exact measurements of the displacements of the sliding load are difficult to make, for instance if the wavelength is very short, the point O' may be obtained as follows. Using a reactive load, construct the circle Γ' as in (B) above; then using a sliding load as in (E) above, construct a circle Γ'' (see Fig. 6). The iconocenter O' is the hyperbolic midpoint of the diameter of Γ'' (through C) with respect to Γ' . It may be constructed by means of the hyperbolic protractor*

* G. A. Deschamps, *Hyperbolic Protractor for Microwave Impedance Measurements and Other Purposes* (New York: International Telephone and Telegraph Corp., 1953).

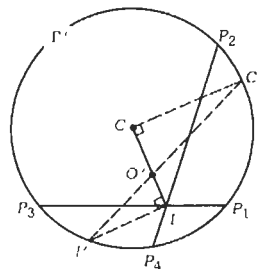


Fig. 5. Determination of O' from four measurements.

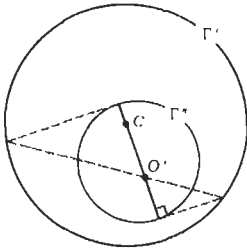


Fig. 6. Use of circles Γ'' and Γ' for determination of O' .

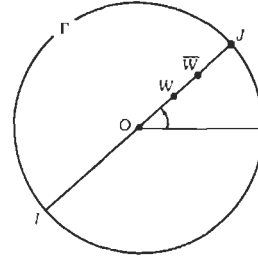


Fig. 7. Representation of a reflection coefficient by W on a Smith chart and \bar{W} on the projective chart.

(Fig. 8) or by means of the dotted-line construction (Fig. 5).

GEOMETRY OF REFLECTION CHARTS

The following brief outline is complemented by the section on hyperbolic trigonometry contained in Chapter 47.

Conformal Charts

A reflection coefficient can be represented by a point in a plane just as any complex number is represented on the Argand diagram.

The passive loads, $|W| \leq 1$, are represented by points inside a unit circle Γ . Inside this circle, the lines of constant resistance and reactance (Smith chart) or the lines of constant magnitude and phase of the impedance (Carter chart) may be drawn.

The transformation from a load reflection coefficient W to its image W' through a two-port junction is bilinear as in Eq. 11 or 12. On the reflection chart, this transformation maps circles into circles and preserves the angle between curves and the cross ratio of four points; if

$$[W_1, W_2, W_3, W_4] = \frac{W_1 - W_3}{W_1 - W_4} : \frac{W_2 - W_3}{W_2 - W_4}$$

denotes the cross ratio of four reflection coefficients W_1, W_2, W_3 , and W_4 , then

$$[W_1', W_2', W_3', W_4'] = [W_1, W_2, W_3, W_4]$$

The transformation through a lossless junction preserves also the unit circle Γ and therefore leaves invariant the *hyperbolic distance* defined in Chapter 47. The hyperbolic distance to the origin of the chart is the *mismatch*, that is, the standing-wave ratio expressed in decibels. It may be evaluated by means of the proper graduation on the radial arm of the Smith chart. For two arbitrary point W_1, W_2 , the hyperbolic distance between them may be interpreted as the mismatch that results from the load W_2 seen through a lossless network that matches W_1 to the input waveguide.

Projective Chart

The reflection coefficient W is represented by the point \bar{W} (Fig. 7) on the same radius of the circle Γ but at a distance

$$O\bar{W} = 2OW/(1 + OW^2) \quad (\text{Eq. 20})$$

from the origin.

This is equivalent to using the standing-wave ratio squared instead of the direct ratio:

$$\bar{W}J/\bar{W}I = (WJ/WI)^2 \quad (\text{Eq. 21})$$

The transformation of Eqs. 11, 12, when the junction is lossless, is represented on this chart by a projective transformation, that is, one that maps straight lines into straight lines and preserves the cross ratio of four points on a straight line. It therefore preserves the hyperbolic distance defined in Chapter 47.

EVALUATION OF HYPERBOLIC DISTANCE

On the projective chart, the hyperbolic distance $\langle AB \rangle$ between two points A and B inside the circle Γ can be evaluated by means of a hyperbolic protractor as shown in Fig. 8. The line AB is extended to its intersections I and J with Γ . The protractor is placed so that the sides OX, OY of the right angle go through I and J . (This can be done in many ways but does not affect the result.) The numbers read on the radial lines of the protractor going through A and B , respectively, are added if A and B are on opposite sides of the radial line marked O ; they are subtracted otherwise. This result divided by 2 is the distance $\langle AB \rangle$. In Fig. 8, for instance

$$\langle AB \rangle = \frac{1}{2} (12 + 4) = 8 \text{ decibels}$$

Problem A

A slotted line with 100-ohm characteristic impedance is used to make measurements on a 60-ohm

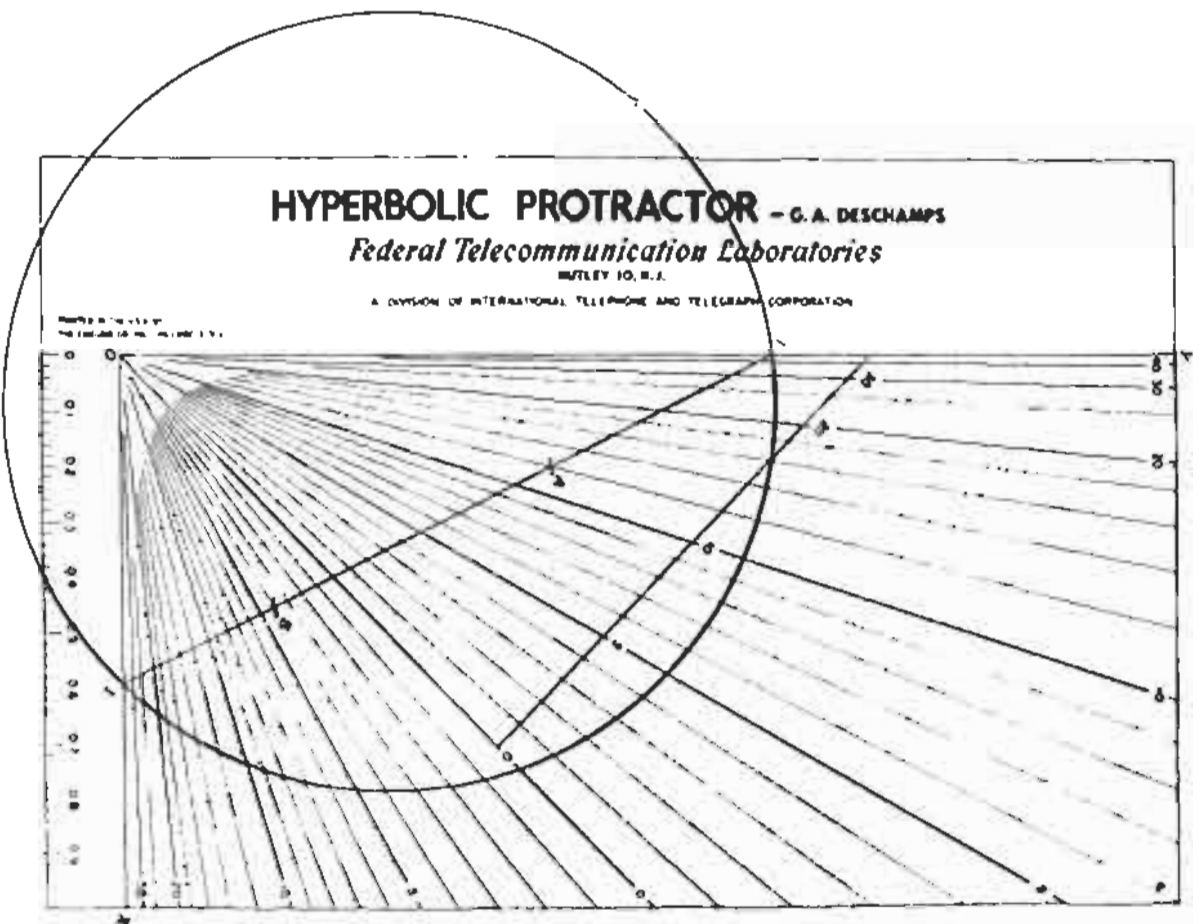


FIG. 8. Definition and evaluation of hyperbolic distance $\langle \Delta r \rangle$ using hyperbolic protractor.

coaxial line. The transition acts as an ideal transformer. Find the reflection coefficient W of an obstacle placed in the coaxial line, knowing that it produces a reflection coefficient

$$W' = 0.5 \exp(j\pi/2)$$

in the slotted line.

A match in the coaxial line appears in the slotted line as a normalized impedance of 0.6; hence the mismatch (standing-wave ratio in decibels) is 4.5 decibels. The corresponding point \bar{O}' is plotted on the projective chart as in Fig. 9 at the distance $\langle O\bar{O}' \rangle = 4.5$. (On the Smith chart drawn inside the same unit circle Γ , the point would be O' .)

The point \bar{W}' representing the unknown load is plotted at the hyperbolic distance

$$20 \log_{10}[(1 + 0.5)/(1 - 0.5)] = 9.5 \text{ decibels}$$

from the origin in the direction $+90^\circ$. The hyperbolic distance

$$\langle \bar{O}'\bar{W}' \rangle = 11 \text{ decibels}$$

is measured with the protractor. This is the mismatch produced by the obstacle in the coaxial line. It corresponds to a magnitude of the reflection coefficient of 0.56.

The phase of this reflection coefficient is the elliptic angle $\langle \bar{O}'P, \bar{O}'\bar{W}' \rangle$. It is evaluated as explained in Chapter 47: extend QO' up to R on Γ and measure the arc

$$PR = 56^\circ$$

The answer is

$$W = 0.56 \angle 56^\circ$$

Problem B

If the transition between the slotted line and the waveguide is not an ideal transformer as in problem A,

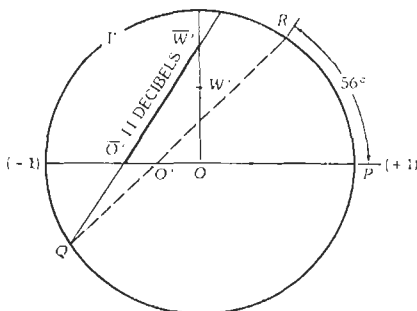


Fig. 9. Measurement of reflection coefficient with a mismatched slotted line.

its properties may be found by the method described earlier in this chapter. In particular, if the transition has no losses (the circle Γ' coincides with Γ), the point O' may be found as in (A), (D), (E), or (F), above, the point P' as in (C) or (D), above, and this completes the calibration.

For any load placed in the waveguide and producing the reflection coefficient W' in the slotted line, the corrected standing-wave ratio in decibels is the hyperbolic distance $[O'W']$. This is evaluated by constructing \bar{O}' , \bar{W}' on the projective chart and measuring $\langle \bar{O}'\bar{W}' \rangle$ with the protractor. The phase angle is the elliptic angle $\langle \bar{O}'P', \bar{O}'\bar{W}' \rangle$ (see Chapter 47).

Problem C

A section of coaxial line 90 electrical degrees in length and with 100-ohm characteristic impedance is inserted between a 50-ohm coaxial line on one side and a 70-ohm coaxial line on the other (Fig. 10). Find the transformer ratio $n = \exp(-u/2)$ and the electrical lengths α, β of the representation [Eqs. (6)].

The two discontinuities are assumed to act as ideal transformers with hyperbolic amplitudes

$$20 \log_{10} \frac{100}{50} = 6 \text{ decibels} = 0.67 \text{ neper}$$

and

$$20 \log_{10} \frac{70}{100} = -3.1 \text{ decibels} = -0.36 \text{ neper}$$

The characteristic polygon* on the projective chart is

* G. A. Deschamps, *Hyperbolic Protractor for Microwave Impedance Measurements and Other Purposes* (New York: International Telephone and Telegraph Corp., 1953; pp. 15-16 and p. 41).

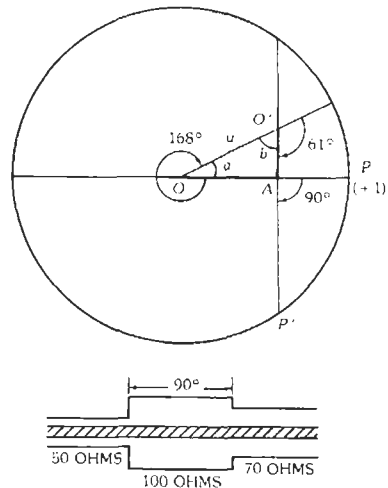


Fig. 10. Solution for transformation in transmission line.

a triangle OAO' with right angle A ; hence, $u = \langle OO' \rangle$ is given by

$$\cosh u = \cosh 0.69 \cosh 0.36$$

$$u = 0.78 \text{ neper} = 6.8 \text{ decibels}$$

$$n = \exp(-u/2) = 1/1.48$$

The lengths of line α and β can be deduced from evaluating the elliptic angles $\langle OA, OO' \rangle = a$ and

$$\langle O'A, O'O \rangle = b$$

$$\tan a = \tanh 0.36 / \sinh 0.69 = 0.46$$

$$a = 24.7^\circ$$

$$\tan b = \tanh 0.69 / \sinh 0.36 = 1.62$$

$$b = 58.4^\circ$$

$$\alpha = \frac{1}{2}(360^\circ - 24.7^\circ) = 167.6^\circ$$

$$\beta = \frac{1}{2}(180^\circ - 58.4^\circ) = 60.8^\circ$$

The resulting equivalent network is shown in Fig. 11. It could also have been obtained by geometrical evaluation of the distance $\langle OO' \rangle$ with the hyperbolic protractor and the elliptic angles a and b by constructions (page 31-7 and Chapter 47).

CORRESPONDENCES WITH CURRENT, VOLTAGE, AND IMPEDANCE VIEWPOINTS

Normalized Current and Voltage

In a waveguide, at a point where the amplitudes of the waves traveling in the positive and negative directions are, respectively, a and b , the normalized voltage v and the normalized current i are defined by the expressions

$$v = a + b$$

$$i = a - b \tag{Eq. 22}$$

The net power flow at that point in the positive direction is

$$|a|^2 - |b|^2 = \text{Re } vi^* \tag{Eq. 23}$$

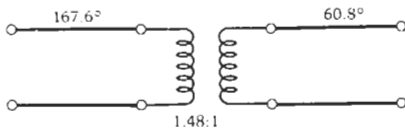


Fig. 11. Equivalent circuit for Fig. 10.

Current and Voltage Not Normalized

A more general definition for current and voltage becomes possible when a meaning has been assigned to the characteristic impedance Z_0 of the waveguide

$$V = vZ_0^{1/2}$$

$$I = iY_0^{1/2} \tag{Eq. 24}$$

where,

$Y_0 = 1/Z_0$ is the characteristic admittance, v and i are the normalized values defined above.

Conversely, if by some convention the voltage (or the current) has been defined, a characteristic impedance will result from Eqs. (24). This is the case for a two-conductor waveguide supporting the TEM mode; the characteristic impedance is the ratio of voltage to current in a traveling wave.

If V and I are the voltage and the current at a point in a waveguide of characteristic impedance $Z_0 = 1/Y_0$, the amplitudes of the waves traveling in both directions at that point are

$$a = \frac{1}{2}(V Y_0^{1/2} + I Z_0^{1/2})$$

$$b = \frac{1}{2}(V Y_0^{1/2} - I Z_0^{1/2}) \tag{Eq. 25}$$

Normalized Impedance and Admittance

At a point in a waveguide, the normalized impedance is $Z = v/i$, and the normalized admittance is the inverse, $Y = 1/Z$.

They are related to the reflection coefficient $W = b/a$ by

$$Z = (1 + W)/(1 - W)$$

$$Y = (1 - W)/(1 + W) \tag{Eq. 26}$$

hence

$$W = (1 - Y)/(1 + Y) = (Z - 1)/(Z + 1) \tag{Eq. 27}$$

Impedance and Admittance Matrix of a Junction

The Z and Y matrices of a junction are defined in terms of the scattering matrix S by

$$Y = (\mathbf{I} - \mathbf{S})(\mathbf{I} + \mathbf{S})^{-1}$$

$$Z = (\mathbf{I} + \mathbf{S})(\mathbf{I} - \mathbf{S})^{-1} \tag{Eq. 28}$$

The matrices Y and Z do not always exist, since S may have eigenvalues $+1$ or -1 , which means that $\det(\mathbf{I} - \mathbf{S})$ or $\det(\mathbf{I} + \mathbf{S})$ may be zero.

Conversely

$$S = (\mathbf{1} - \mathbf{Y})(\mathbf{1} + \mathbf{Y})^{-1} = (\mathbf{Z} - \mathbf{1})(\mathbf{Z} + \mathbf{1})^{-1} \quad (\text{Eq. 29})$$

These equations may be used as definitions for the scattering matrix of lumped-constant networks with n terminal pairs. This is equivalent to considering the network as a junction between n transmission lines of unit characteristic impedance.

If the network or the junction is reciprocal, \mathbf{Y} and \mathbf{Z} are purely imaginary.

For a two-port junction, Eq. 28 becomes

$$\begin{aligned} \mathbf{Y} &= (\mathbf{1} - \mathbf{S})(\mathbf{1} + \mathbf{S})^{-1} \\ &= [\det(\mathbf{1} + \mathbf{S})]^{-1} \\ &\times \begin{bmatrix} 1 - \det \mathbf{S} + (s_{22} - s_{11}) & -2s_{12} \\ -2s_{21} & 1 - \det \mathbf{S} - (s_{22} - s_{11}) \end{bmatrix} \end{aligned} \quad (\text{Eq. 30})$$

and

$$\begin{aligned} \mathbf{Z} &= (\mathbf{1} + \mathbf{S})(\mathbf{1} - \mathbf{S})^{-1} \\ &= [\det(\mathbf{1} - \mathbf{S})]^{-1} \\ &\times \begin{bmatrix} 1 - \det \mathbf{S} - (s_{22} - s_{11}) & 2s_{12} \\ 2s_{21} & 1 - \det \mathbf{S} + (s_{22} - s_{11}) \end{bmatrix} \end{aligned} \quad (\text{Eq. 31})$$

$$\begin{aligned} \det(\mathbf{1} + \mathbf{S}) &= 1 + \text{tr } \mathbf{S} + \det \mathbf{S} \\ &= 1 + (s_{11} + s_{22}) + (s_{11}s_{22} + s_{12}s_{21}) \end{aligned}$$

$$\begin{aligned} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} &= \mathbf{Z} \begin{bmatrix} i_1 \\ i_2 \end{bmatrix} \\ \begin{bmatrix} i_1 \\ i_2 \end{bmatrix} &= \mathbf{Y} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} \end{aligned}$$

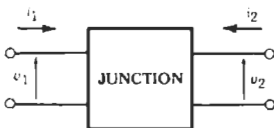


Fig. 12. Sign convention for defining the impedance and admittance of a two-port junction.

$$\begin{aligned} \det(\mathbf{1} - \mathbf{S}) &= 1 - \text{tr } \mathbf{S} + \det \mathbf{S} \\ &= 1 - (s_{11} + s_{22}) + (s_{11}s_{22} - s_{12}s_{21}) \end{aligned}$$

The matrices \mathbf{Y} and \mathbf{Z} relate normalized voltages and currents at both ports (Fig. 12) as follows:

$$\begin{aligned} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} &= \mathbf{Z} \begin{bmatrix} i_1 \\ i_2 \end{bmatrix} \\ \begin{bmatrix} i_1 \\ i_2 \end{bmatrix} &= \mathbf{Y} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} \end{aligned}$$

TRANSFORMATION MATRIX

A transformation matrix useful for composing two-port junctions in cascade relates the voltage and current on one side of the junction to the same quantities on the other side. With the notation in Fig. 13

$$\begin{bmatrix} v' \\ i' \end{bmatrix} = \mathbf{U} \begin{bmatrix} v \\ i \end{bmatrix} \quad (\text{Eq. 32})$$

The matrix \mathbf{U} , sometimes called the *ABCD* matrix, has the same properties as \mathbf{T} described earlier.

For a series element that has a normalized impedance Z

$$\mathbf{U} = \begin{bmatrix} 1 & Z \\ 0 & 1 \end{bmatrix}$$

and for a shunt element with normalized admittance Y

$$\mathbf{U} = \begin{bmatrix} 1 & 0 \\ Y & 1 \end{bmatrix}$$

A product of matrices of these types gives the transformation matrix for any ladder network.

For the shunt element Y , the scattering matrix is

$$\mathbf{S} = (2 + Y)^{-1} \begin{bmatrix} -Y & 2 \\ 2 & -Y \end{bmatrix} \quad (\text{Eq. 33})$$

hence

$$\begin{aligned} s_{11} &= s_{22} \\ s_{12} &= 1 + s_{11} \end{aligned} \quad (\text{Eq. 34})$$

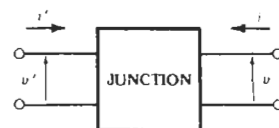


Fig. 13. Sign convention for voltages and currents related by the transformation matrix.

For the series element Z , the scattering matrix is

$$\mathbf{S} = (2 + Z)^{-1} \begin{bmatrix} Z & 2 \\ 2 & Z \end{bmatrix} \quad (\text{Eq. 35})$$

hence

$$s_{11} = s_{22}$$

$$s_{12} = 1 - s_{11} \quad (\text{Eq. 36})$$

Relations 34 and 36 are characteristic, respectively, of a shunt and a series obstacle in a waveguide.

The matrix \mathbf{T} can be deduced from \mathbf{U} and vice versa:

$$\begin{aligned} \mathbf{T} &= \frac{1}{2} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \mathbf{U} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \\ &= \frac{1}{2} \begin{bmatrix} u_{11} + u_{12} + u_{21} + u_{22} & u_{11} - u_{12} + u_{21} - u_{22} \\ u_{11} + u_{12} - u_{21} - u_{22} & u_{11} - u_{12} - u_{21} + u_{22} \end{bmatrix} \end{aligned} \quad (\text{Eq. 37})$$

A similar equation will transform \mathbf{T} into \mathbf{U} , since

$$\mathbf{U} = \frac{1}{2} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \mathbf{T} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \quad (\text{Eq. 38})$$

32 Antennas

Robert C. Hansen

General	32-3
Introduction	
Definitions	
Small Antennas	
Superconducting Antennas	
Near-Field Power Density	
Antenna Noise Temperature	
Elliptical and Circular Polarization	
Low-Gain Antennas	32-13
Half-Wave Dipole	
Radial-Wire Ground Systems	
Printed-Circuit Antennas	
Loops	
Slot Antennas	
Medium-Gain Antennas	32-22
Horns	
Helices	
Yagi-Uda Antennas	
Frequency-Independent Antennas	
Log-Periodic Antennas	
Fractal Antennas	
Arrays	32-32
General Characteristics	
Grating and Quantization Lobes	
Linear Array Feeds	

Mutual Impedance
Thinned Arrays
Tolerances
Multiple-Beam Arrays
Electronic Scanning
Adaptive Arrays

Aperture Distributions 32-47

Design of Distributions
Taylor One-Parameter Line Source Distribution
Taylor \bar{n} Line Source Distribution
Bayliss \bar{n} One-Parameter Difference Line Source
Low Sidelobe Distributions
Measurement of Low-Sidelobe Patterns
Hansen One-Parameter Circular Source Distribution
Taylor \bar{n} Circular Source Distribution

Reflectors 32-54

Parabolic Reflectors
Scanning and Multiple-Beam Reflectors

GENERAL

Introduction

The field of antennas is sufficiently broad to be beyond the scope of this chapter, even for succinct design information. Thus the intent is to be eclectic rather than inclusive. Fortunately, several excellent books on antennas are available. For more detailed design data, the reader should consult the two-volume *Handbook of Antenna Design* edited by Rudge et al.* Thorough texts on antenna theory, including extensive coverage of the powerful geometric theory of diffraction (GTD) and moment method analytical approaches, as well as exemplary computer programs, are by Balanis† and by Stutzman and Thiele.‡ A third excellent text, which has extensive coverage on fixed beam array design, is by Elliott.§ Finally, the most extensive treatment of phased arrays (electronic scanning) is *Phased Array Antennas* by Hansen;¶ another useful book is still the three-volume *Microwave Scanning Antennas* by Hansen.#

Six parts comprise this chapter. Basic antenna behavior, including definitions, fields, near-field power density, antenna noise temperature, and polarization coupling are in the general part. The next two parts cover, respectively, low- and medium-gain antennas, where dipoles, slots, loops, and microstrip patches are low gain, and horns, Yagi-Udas, helices, spirals, and log-periodics are medium-gain. Arrays of all types are covered next. Aperture distributions, because of their common importance, occupy an entire part. Finally, reflector-type antennas are the subject of the last part. Some old favorites have been left out, and this reflects somewhat the changing antenna usage.

Because of the power and ready availability of calculators and computers, tables of calculated functions have largely been omitted. Instead, key performance indices have been quantified, with design formulas given so that the designer can implement them directly.

Definitions

Directivity and gain are measures of how well energy is concentrated in a given direction. Directivity

is the ratio of power density (PD) in that direction to the power density that would be produced if the power were radiated isotropically. The reference can be linearly or circularly polarized, and directivity is often given in dBi, decibels above isotropic. Some early literature refers to gain above a dipole: this usage is deprecated, as it is confusing and unnecessary. Directivity, then, is given by

$$D = 4\pi PD / \left(\sqrt{\frac{\epsilon}{\mu}} \iint |E|^2 d\Omega \right)$$

Gain includes antenna losses; thus gain is the field intensity produced in the given direction by a fixed input power to the antenna. Gain is related to directivity by efficiency η , and is

$$G = D\eta$$

$$G = 4\pi PD/P_{in}$$

Through reciprocity, directivity is independent of transmission or reception, as is gain. Gain, as widely used in the industry, includes the impedance mismatch factor. However, the academic gain does not.

Effective area is defined by:

$$A_e = \lambda^2 G / 4\pi$$

where λ is the free-space wavelength. (All through this chapter, commonly used symbols are employed.) For an antenna matched to a load, the load power is $P_{load} = PD \cdot A_e$, where PD is the power density at the antenna in watts per square meter.

Path loss is part of the range equation, where received power is related to transmitted power as

$$P_r = P_t G_t G_r \lambda^2 / (4\pi R)^2$$

The distance between antennas is R , and the path loss is given by

$$\text{Path loss} = (4\pi R/\lambda)^2$$

Effective length relates the ability of a receiving antenna to produce open-circuit voltage. It is

$$l_e = V/E$$

where V is the open-circuit voltage for an incident field strength E . The early usage "effective height," is deprecated, as it also has a meaning for antennas over the earth. For any antenna, the preceding parameters are related through

$$30\pi l_e^2 = R_r A_e$$

Here, R_r is radiation resistance, where the radiated power of a current-driven antenna is $P_r = I^2 R_r$. Table 1 gives often used parameters for short dipoles and for half-wave dipoles. Dipole half-length and monopole lengths are h , and θ_1 is the half-power beamwidth, i.e., the width of the pattern between -3 -dB points. For loops, the diameter is d , N is the number of turns, $k = 2\pi/\lambda$, and μ_e is the effective permeability.

* Rudge, A. W., et al. (Eds.). *Handbook of Antenna Design*. London: Peter Peregrinus Ltd., 1983.

† Balanis, C. A., *Antenna Theory: Analysis and Design*. New York: John Wiley & Sons, Inc., 1997.

‡ Stutzman, W. L., and Thiele, G. A. *Antenna Theory and Design*. New York: John Wiley & Sons, Inc., 1998.

§ Elliott, R. S. *Antenna Theory and Design*. Englewood Cliffs, NJ: Prentice-Hall, Inc., 1981.

¶ Hansen, R. C. *Phased Array Antennas*. New York: John Wiley & Sons, Inc., 1998.

Hansen, R. C. *Microwave Scanning Antennas*. New York: Academic Press, Inc.; Vol. 1, 1964; Vols. 2 and 3, 1966; Los Altos, CA: Peninsula Publ., 1985.

TABLE I. LOW-GAIN-ANTENNA PARAMETERS

Type	D	l_c	A_c	θ_1	R_r
Isotropic	1	—	$\lambda^2/4\pi = 0.0796\lambda^2$	360°	—
Short Dipole	1.5	h	$3\lambda^2/8\pi = 0.1194\lambda^2$	90°	$20k^2h^2$
$\lambda/2$ Dipole	1.6409	λ/π	$30\lambda^2/\pi R_r = 0.1306\lambda^2$	78.078°	73.13
$\lambda/4$ Monopole	3.2818	λ/π	$30\lambda^2/\pi R_r = 0.2612\lambda^2$	78.078°	36.56
Small Loop	1.5	$\pi N k d^3 \mu_c/4$	$3\lambda^2/8\pi = 0.1194\lambda^2$	90°	$5\pi^2 N^2 k^4 d^4 \mu_c^2/4$

Bandwidth may be defined through pattern characteristics, efficiency, impedance, etc. The latter is often used, with the bandwidth being the range between the two frequencies where the radiated power falls to half (the 3-dB, or half-power, bandwidth) or where the VSWR reaches a fixed value, e.g., 2. Note that half-power bandwidth occurs when the input $R = |X|$, and for $VSWR = 5.828$. Since the product of antenna Q and fractional bandwidth, BW , is:

$$Q \cdot BW = (VSWR - 1)/\sqrt{VSWR}$$

the bandwidth at one VSWR is related to that at another VSWR by:

$$BW_1/BW_2 =$$

$$[(VSWR_1 - 1)/(VSWR_2 - 1)](VSWR_2/VSWR_1)^{1/2}$$

For example, the 3-dB bandwidth is 2.828 times larger than the $VSWR = 2$ bandwidth. These results include the effect of a matched load (generator); for half-power $Q \cdot BW = 2$.

Field regions must be defined carefully, as the early optical terms are inadequate for modern antennas. The terms Fraunhofer and Fresnel, which refer to specific field integral approximations in optics, are obsolete and deprecated. For example, with focused antennas, a Fraunhofer-type pattern may exist well within the usual D^2/λ near-far field boundary distance, while a Fresnel-type field may exist for smaller and for larger distances. Some antennas have no Fraunhofer-type fields anywhere, as they have no phase center, e.g., some horns and annular slots. Finally, the Fresnel approximation itself is unambiguous only in one dimension; for the more common and useful area sources, the Fresnel results vary widely with coordinate system and formulation.*

Thus, the following definitions have evolved. Space is divided into three regions as follows. That region of space immediately surrounding the antenna in which the reactive components predominate is known as the *reactive near-field region*. The size of this region varies for different antennas. For most antennas, however, the outer limit is on the order of a few wavelengths or less. For the particular case of an electrically small

dipole, the reactive field predominates to a distance of approximately $\lambda/2\pi$, where the radiating and reactive fields are equal. Beyond the reactive near-field region, the radiating field predominates. The radiating region is divided into two subregions, the *radiating near-field region* and the *far-field region*. In the radiating near-field region, the relative angular distribution of the field (the usual radiation pattern) is dependent on the distance from the antenna. The reason for this behavior is twofold: the relative phase relationship of field contributions from different elements of the antenna changes with distance, and the relative amplitudes of these field contributions also change with distance. As the observation point in space moves away from the antenna, the amplitude of the field first oscillates and then decays monotonically. This variation in the limit is given by the reciprocal of the first power of distance. Furthermore, the relative phase and amplitude relationships between the field contributions from different elements of the antenna asymptotically approach a fixed relationship, and the relative angular distribution of the field becomes independent of the distance. This occurs in the far-field region; patterns are essentially independent of distance. For most antennas, the transition distance is D^2/λ , where D is the width of the equivalent uniformly excited aperture. Precision gain measurements or measurements of nulls may require a distance of $2D^2/\lambda$ or more. Low-sidelobe antennas will require multiples of D^2/λ , depending on the sidelobe level and allowable error.† Fig. 1 shows the reactive near-field boundary at $R = \lambda$ and the D^2/λ far-field boundary for several planar apertures.

Small Antennas

Antennas that are small in wavelength are conceptually simple, and difficult to use, as will appear below. It has been observed by Wheeler‡ that all small antennas are dipoles, loops, or combinations of these two canonical types. When the dipole (or monopole over a large ground plane) is short, the current is essentially

† Hansen, R. C. *Phased Array Antennas*. New York: John Wiley & Sons, Inc., 1998.

‡ Wheeler, H. A. "Fundamental Limitations of Small Antennas." *Proc. IRE*, Vol. 35, December 1947, pp. 1479-1484. Wheeler, H. A. "Small Antennas." *Trans. IEEE*, Vol. AP-23, July 1975, pp. 462-469.

* Hansen, R. C. *Microwave Scanning Antennas*. New York: Academic Press, Inc.: 1964, Vol. 1, Chapter 1; Los Altos, CA: Peninsula Publ., 1985.

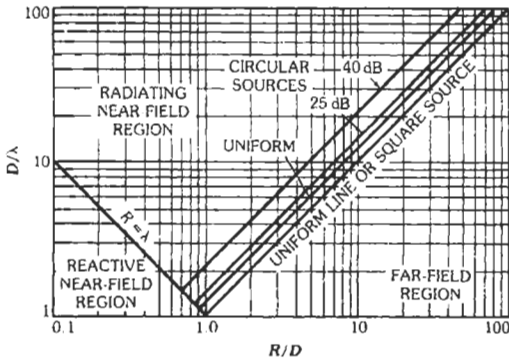


Fig. 1. Field regions.

linear from feed to end, and the fields produced, when the dipole is along the z -axis of a standard spherical coordinate system, are

$$E_{\theta} = j 60 \pi h I_0 e^{-jkr} \sin \theta / r \lambda,$$

$$H_{\phi} = -j h I_0 e^{-jkr} \sin \theta / 2 r \lambda$$

The radiation resistance is $R_r = 20k^2 h^2$, where $k = 2\pi/\lambda$ and h is the dipole half-length. Similarly, a small loop of diameter d carries essentially a constant current, and the fields with the loop axis along z are analogously

$$E_{\theta} = 30 \pi^3 d^2 I_0 e^{-jkr} \sin \theta / r \lambda^2$$

$$H_{\phi} = -\pi^2 d^2 I_0 e^{-jkr} \sin \theta / 4 r \lambda^2$$

The radiation resistance is given by:

$$R_r = (5/4) \pi^2 N^2 k^4 d^4 \mu_e^2$$

See the section on low-gain antennas for details of μ_e . When $h \ll \lambda$ or $d \ll \lambda$, the radiation resistance is very small, often smaller than the loss resistance of the conductors. Efficiency then can be small, and although the directivity is 1.5, the gain is also small. Reactance, on the other hand, is high; for short dipoles it varies as $1/kh$, whereas for small loops it varies as $\sim kd$. The Q is therefore high, and the bandwidth, which is approximately $(h/\lambda)^3$ for $Q \gg 1$, is small. It is useful to know what bandwidth can be achieved; the Wheeler papers referred to above give practical answers. A theoretical fundamental limitation was derived by Chu* and refined by Harrington† and Collin.‡ See also Hansen.§

* Chu, L. J., "Physical Limitations of Omnidirectional Antennas." *J. Appl. Phys.*, Vol. 19, December 1948, pp. 1163-1175.

† Harrington, R. F. "Effect of Antenna Size on Gain, Bandwidth, and Efficiency." *J. Res. Nat. Bur. Stand.*, Vol. 64D, January-February 1960, pp. 1-12.

‡ Collin, R. E. "Evaluation of Antenna Q." *Trans. IEEE*, Vol. AP-12, January 1964, pp. 23-27.

§ Hansen, R. C. "Fundamental Limitations in Antennas." *Proc. IEEE*, Vol. 69, February 1981, pp. 170-182.

for a discussion of fundamental limitations in antennas. Since any radiating field can be written as a sum of spherical modes, the antenna, of whatever type it happens to be, is enclosed in a sphere of radius a . The radiated power can be calculated from propagating modes within the sphere. All modes contribute to the reactive power. When the sphere is sufficiently large to support several propagating modes, this approach is of little value because the modal coefficients are difficult to calculate. With only one propagating mode, the radiated power arises primarily from that mode. The utility of the Chu work becomes apparent when the sphere is too small to allow a propagating mode; all modes are then evanescent (below cutoff), and the Q becomes large, as the evanescent modes contribute little real power. Note that, unlike the case of a closed waveguide, there is a real part of each evanescent mode. Each mode has a Q based on the ratio of stored energy to radiated energy, and the Q rises rapidly when kr drops below the mode number. For all modes well below cutoff, the result is

$$Q = (1 + k^2 a^2) / [k^3 a^3]$$

When the antenna contains loss, the Q is reduced by the efficiency, η . Fig. 2 gives fundamental limitation curves for several values of efficiency. This Q is for the

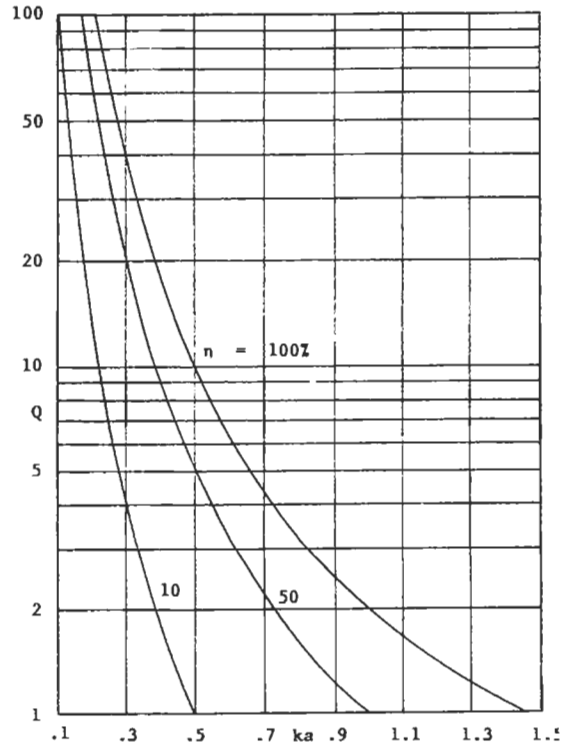


Fig. 2. Chu-Harrington fundamental limitations for single-mode antenna, various efficiencies.

lowest TM mode. When both a TM mode and a TE mode are excited, the value of Q is halved. The importance of the Chu result is that it relates the lowest achievable Q to the maximum dimension of an electrically small antenna, and this result is independent of the art that is used to construct the antenna within the hypothetical sphere, except in determining whether a pure TE or pure TM mode, or both, is excited. Since the Q grows rapidly (inverse cube) as size decreases, this indeed represents a fundamental limit that has only been approached but never equaled, much less exceeded. Bandwidth is derived from Q by assuming that the antenna equivalent is a resonant circuit with fixed values. Then the fractional bandwidth is

$$\text{Bandwidth} = \frac{f_{\text{upper}} - f_{\text{lower}}}{f_{\text{center}}} = 1/Q$$

A matched load is not included here, but use of it halves the Q . For $Q \gg 1$ this relationship is meaningful, as the fixed resonant circuit is a good approximation to the antenna. But for $Q < 2$, the representation is no longer accurate. However, the curves are still useful for low Q even though imprecise. Since most small antennas are loops or dipoles, which do not use the spherical volume efficiently, an octave antenna is large, often larger than $\lambda/2$.

Antenna bandwidth can be increased by use of one or more matching networks. The maximum improvement possible was derived by Fano, and his results have been implemented by Matthaei et al.* From their formulas, the bandwidth improvement factors for VSWR = 2 and for the half-power case are as shown in Table 2.

Electrically small antennas can be broadbanded by introducing loss as indicated,† or by utilization of a large mismatch. A small loop with a very low-impedance preamplifier at its terminals or a short monopole (dipole) with a high-impedance preamplifier at its terminals is called an "aperiodic loop" or "aperiodic

TABLE 2. BANDWIDTH INCREASE USING N MATCHING SECTIONS

N	$V = 2$	Half-Power
1	2.3094	2.0301
2	2.8596	2.4563
3	3.1435	2.6772
∞	3.8128	3.2049

* Matthaei, G. L., Young, L., and Jones, E. M. T. *Micro-wave Filters, Impedance-Matching Networks, and Coupling Structures*. McGraw-Hill, 1964.

† Kanda, M. "A Relatively Short Cylindrical Broadband Antenna with Tapered Resistive Loading for Picosecond Pulse Measurements." *Trans. IEEE*, Vol. AP-26, May 1978, pp. 439-447.

monopole." The term "integrated antenna" has been used, but it is confusing and is now obsolete. If the amplifier impedance is much larger than the monopole reactance at the lowest frequency of interest, the output signal is nearly constant over a bandwidth of one or more octaves. The large mismatch reduces the signal and external noise. Aperiodic antennas are useful at those frequencies where the external noise is so large that even after the large mismatch loss the antenna noise is still larger than the preamp noise; thus the system is external-noise limited. Let the ratio of preamp input resistance to lowest antenna reactance be $\alpha = R/X_1$, and use F_{ant} as the ratio of external antenna noise temperature to $T_0 = 290$ K, and F_N as the preamp noise figure. A key system performance factor, γ , is the ratio of actual s/n to s/n if the system were external-noise limited. That is, γ determines s/n degradation compared to a narrow-band matched antenna-preamp system at frequency f . The s/n degradation factor for a monopole that is short over the band of interest is‡

$$\gamma = \frac{(2\alpha F_{\text{ant}} f^2 / f_1^2) / (\alpha^2 + f_1^2 / f^2)}{[2\alpha (F_{\text{ant}} - 1) f^2 / f_1^2] / (\alpha^2 + f_1^2 / f^2) + 3AF_N / k_1^3 h^2}$$

In this result, $A = \ln(h/a) - 1$, and h and a are the monopole length and radius. Even with F_N in the 2-dB range, there can be significant s/n degradation in using an aperiodic monopole unless the upper frequency is close to the lower, f_1 .

Superconducting Antennas

High T_c superconducting materials offer advantages for millimeter arrays, both waveguide and microstrip types. For superdirective arrays, and for electrically small dipole antennas, efficiency improvements are significant only when the radiation resistance is such that the Q is unacceptably high, typically greater than 1000. However, these antennas have low radiation resistance and high reactance which present a very high VSWR. A transmission-line match or lumped circuit with matched loss L has an apparent loss L_a that is greatly increased due to the VSWR or V :

$$L_a = [(V + 1)^2 L^2 - (V - 1)^2] / 4VL$$

Thus, the superconducting material is needed in the matching networks. Electrically small loop antennas can realize efficiency improvements, but without the leavening effect of loss resistance the Q is often unacceptably high. For a review of superconducting antennas, see Hansen.§

‡ Radjy, A. H., and Hansen, R. C. "S/N Performance of Aperiodic Monopoles." *Trans. IEEE*, Vol. AP-27, March 1979, pp. 259-261.

§ Hansen, R. C. "Superconducting Antennas." *Trans. IEEE*, Vol. AES-26, March 1990, pp. 345-355.

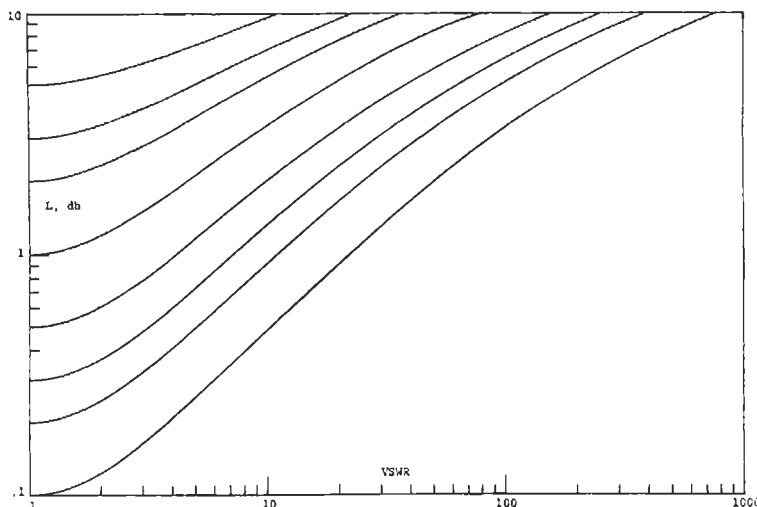


Fig. 2a. Line loss with mismatch.

Near-Field Power Density

Personnel and equipment electromagnetic radiation safety have made the understanding of radiated power density important. For points in the far-field region, the power density is given simply by

$$PD = P_r G / 4\pi r^2$$

where

- P_r is radiated power in watts,
- r is distance in meters,
- PD is in W/m^2

To convert to mW/cm^2 , divide by 10. When the point of interest is within $\lambda/2$ of a low-gain transmitting antenna, such as a single vertical tower, the near-field dipole formulas can be used for approximate results. For dipoles at greater distances, $PD = 30P_r / \pi R r^2$.

Near fields of larger antennas can be accurately estimated if the aperture (or array) distribution is known. Both line sources and planar apertures will be discussed, because the near-field behavior of these two is quite different. For a line source, the beam is usually designed to be narrow in one plane while it is broad or omnidirectional in the other plane. In the near field, the phase errors degrade the pencil beam pattern, effectively decollimating it so that $1/r$ of the field intensity variation with distance is effectively canceled. The other $1/r$ still remains, with the result that near-field power density for a line source oscillates about a $1/r$ variation. The oscillations are reduced for tapered amplitude distributions, and there is, of course, a transition region from the far-field $1/r^2$ behavior to the near-field $1/r$. Because most high-power antennas are planar sources instead of line sources, only limited information is given on line sources. For uniform excitation the near-field power density is*

$$PD = (4.05/\Delta) [C^2 (1/2\sqrt{\Delta}) + S^2 (1/2\sqrt{\Delta})]$$

Here, the distance from the antenna is normalized in terms of $2L^2/\lambda$, where $\Delta = r/(2L^2/\lambda)$; C and S are conventional Fresnel integrals. For a square uniformly excited aperture, the power density is

$$PD = 16.4 [C^2 (1/2\sqrt{\Delta}) + S^2 (1/2\sqrt{\Delta})]^2$$

Fig. 3 shows power density for both the line source and the square source.

Of more interest is the near-field behavior of circular apertures, either arrays or reflectors. The Hansen one-parameter circular space factor† is a good fit to the

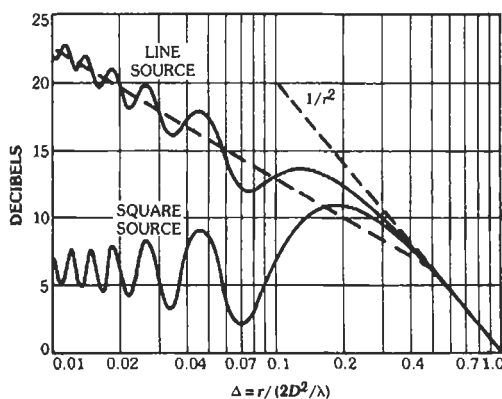


Fig. 3. On-axis power density.

* Ricardi, L. J., and Hansen, R. C. "Comparison of Line and Square Source Near Fields." *Trans. IEEE*, Vol AP-11, November 1963, pp. 711-712.

† Hansen, R. C. "A One-Parameter Circular Aperture Distribution with Narrow Beamwidth and Low Sidelobes." *Trans. IEEE*, Vol. AP-24, July 1976, pp. 477-480.

main beam and near sidelobes of most reflector antennas if the parameter is properly chosen. Similarly, it is an excellent representation for near-field power density near the axis. These robust low- Q distributions are also used for arrays. The axial power density is obtained from numerical integration with only the results given here.* Figs. 4 and 5 give axial power density for circular apertures with sidelobe ratios from 20 to 35 dB. Power density on-axis for a uniformly excited circular aperture is shown in Fig. 6.† Thus, it can be observed that the near-field phase errors do decollimate the beam in both planes, thereby producing a constant power density value around which the near field oscillates. All of these data are normalized to a value of unity at a distance of $2D^2/\lambda$, and the nominal power density at that distance is found simply from

$$PD = PG/[4\pi(2D^2/\lambda)^2]$$

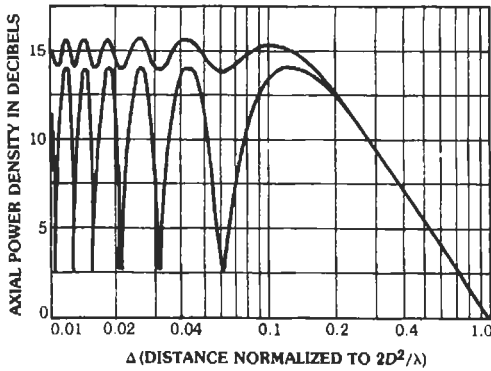


Fig. 4. Axial power density for SLR of 20 and 30 dB, one-parameter circular aperture.

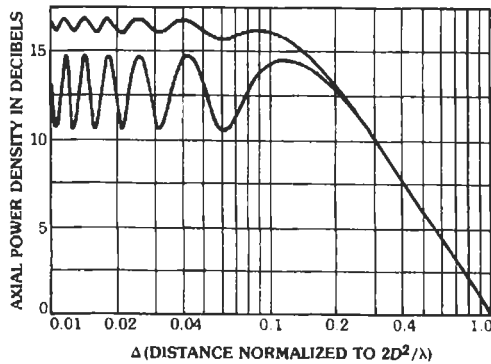


Fig. 5. Axial power density for SLR of 25 and 35 dB, one-parameter circular aperture.

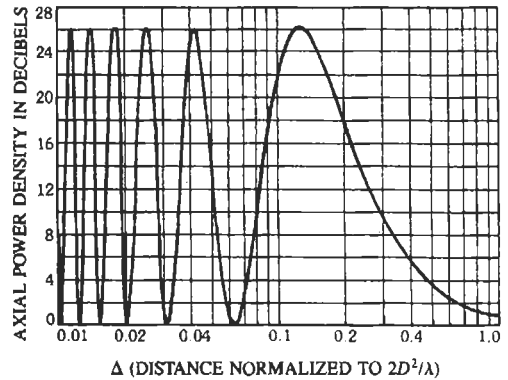


Fig. 6. Axial power density of uniform circular aperture.

where

- P is the input power,
- G is the antenna gain.

Calculation of near fields off-axis is difficult because the small-angle approximate techniques quickly become useless. At large angles, Geometric Theory of Diffraction techniques are useful.

Antenna Noise Temperature

Antenna noise contributes to system noise and thus will affect system performance. Antenna noise consists of external noise and internal noise, with the latter caused by loss in the antenna, feed, cable, etc. It is customary and convenient to refer noise measurements to the input terminals of the preamplifier or receiver. If there is loss within the antenna or in components between the antenna and the preamplifier input, this loss will change the antenna contribution to system noise. If L is the loss ratio (P_{in}/P_{out}) and T_a is the external noise temperature, the total antenna temperature at the preamplifier input is given by

$$T_{ant} = [290(L-1) + T_a]/L$$

Fig. 7 shows this in graphical form. It is assumed that the lossy components are at an ambient temperature of 290 K. For other ambient temperatures, the coefficient in the formula is, of course, changed. The loss applies only to actual dissipative loss and not to virtual losses such as aperture excitation efficiency. The systems temperature is

$$T_{sys} = [T_0(L-1) + T_a]/L + T_0(F_N-1) + T_0(F_{rec}-1)/G_{pre}$$

where the preamplifier noise figure is F_N and the receiver noise figure is F_{rec} . Note that the preamplifier noise temperature contribution involves (F_N-1) ; this is because the measurement of amplifier noise figure includes an input standard temperature of $T_0 = 290$ K. An alternate form of the system temperature is:

$$T_{sys} = T_a/L - T_0(F_N-1)/L + T_0(F_{rec}-1)/G_{pre}$$

* Hansen, R. C. "Circular-Aperture Axial Power Density," *Microwave J.*, Vol. 19, February 1976, pp. 50-52.

† Hansen, R. C. "Antenna Power Densities in the Fresnel Region," *Proc. IRE*, Vol. 47, December 1959, pp. 2119-2120.

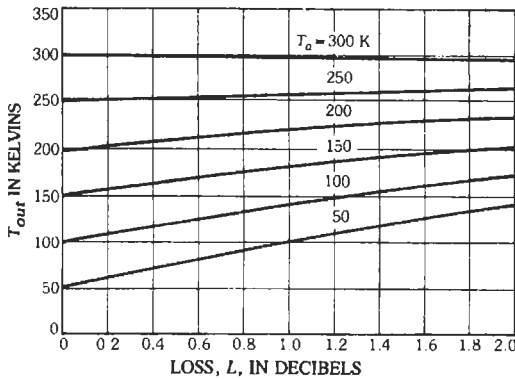


Fig. 7. External plus internal antenna noise temperature.

In these equations, the receiver noise contributions are decreased by the gain of the preamplifier. The external antenna noise temperature is a summation of all of the noises seen by the antenna, weighted by the antenna power pattern. Thus

$$T_a = \frac{\iint T(\theta, \phi)P(\theta, \phi)d\Omega}{\iint P(\theta, \phi)d\Omega}$$

In this equation, $T(\theta, \phi)$ is the spatial temperature pattern in kelvins, $P(\theta, \phi)$ is the antenna power pattern, and the integration is over all space. Next, specific external noise sources will be reviewed.

Below 10 MHz, the noise is primarily atmospheric and very large. This noise is mostly Gaussian, but there are impulsive tails at low probabilities of occurrence. Extensive data are given by the CCIR.* The highest noise generally occurs in spring in the 2000–2400 time block. The CCIR data have been computerized by Lucas and Harper.† Broadly, atmospheric noise is omnidirectional and independent of distance. A minor exception occurs in the VLF absorption notch around 4 kHz in that noise may be higher in some directions due to rapid attenuation with distance. The atmospheric noise contribution becomes negligible between 35 and 40 MHz. Galactic noise generally consists of noise from a few hot stars plus a varying background that results from many, many stars. The background results, of course, depend on how many discrete star sources have been separated out. Typical charts show measurements made by Taylor‡ at 136

* CCIR. *Characteristics and Applications of Atmospheric Radio Noise, Report 322.2*. Geneva: ITU, 1983.

† Lucas, D. L., and Harper, J. D. *A Numerical Representation of CCIR Report 322: High Frequency (3–30 Mcps) Atmospheric Radio Noise Data*, TN318, NTIS No. COM-7510374. NBS, 1965.

‡ Taylor, R. E. "13 6/400 MHz Radio Sky Maps." *Proc. IEEE*, Vol. 61, April 1973, pp. 469–472.

and 400 MHz; see Figs. 8 and 9. A composite fit to these and other measurements is:

$$T \approx (3.068 \times 10^8)/f_{\text{MHz}}^{2.3}$$

This line is shown on the composite noise plot, Fig. 10. The sun is a special case because it subtends a finite angle and is much closer. Sun noise temperature measurements by Hogg and Mumford§ show a slight oscillation about a linear frequency dependence. The latter is

$$T_s = 1.958 \times 10^8/f_{\text{MHz}}$$

For antenna beamwidths larger than 0.5° , the sun noise temperature contribution is

$$T_{\text{sun}} = T_s \theta_s^2 G/4\pi$$

where

θ_s is the subtended angle of the sun,
 G is the antenna gain in that direction.

For those portions of the antenna pattern that see the earth, the earth will contribute noise at a blackbody temperature of roughly 310 K. In the range from 20 to 600 MHz, all of these noise sources are usually dwarfed by man-made noise. This consists of power-line noise caused by arcs and corona and by user-generated pulses; noise from various transmitters including TV, FM, cellular, and push-to-talk; and noise generated by vehicles. Extensive studies and measurements of urban noise have been made, and the reader is referred to books by Skomal¶ and Herman.# The most appropriate division into cultural areas now seems to be "business," "residential," and "rural." These divisions and the data of Spaulding and Disney** are used here as seen in Fig. 10. However, it must be recognized that these urban noise envelopes apply away from vehicular traffic streams. That is to say, vehicle noise is included, but the vehicles are not closer than several hundred feet. A vehicle stream at 50 feet can generate noise 20 to 40 dB higher than these envelopes. A single "super-noisy" vehicle can generate as much noise as a stream of vehicles, and at close distances the noise levels can be extremely large. Although few data are available, vehicular noise can be significant at frequencies as high as 5–10 GHz, again for short distances. Atmospheric noise is produced by gaseous absorption, with the low-

§ Hogg, D. C., and Mumford, W. W. "The Effective Noise Temperature of the Sky." *Microwave J.*, Vol. 3, March 1960, pp. 80–84.

¶ Skomal, E. N. *Man-Made Radio Noise*. New York: Van Nostrand Reinhold Co., 1978.

Herman, J. R. *Electromagnetic Ambients and Man-Made Noise*. Gainesville, VA.: Don White Consultants, Inc., 1979.

** Spaulding, A. D., and Disney, R. T. *Man-Made Radio Noise. Part I: Estimates for Business, Residential, and Rural Areas*. OTR Report 74-38, June 1974.

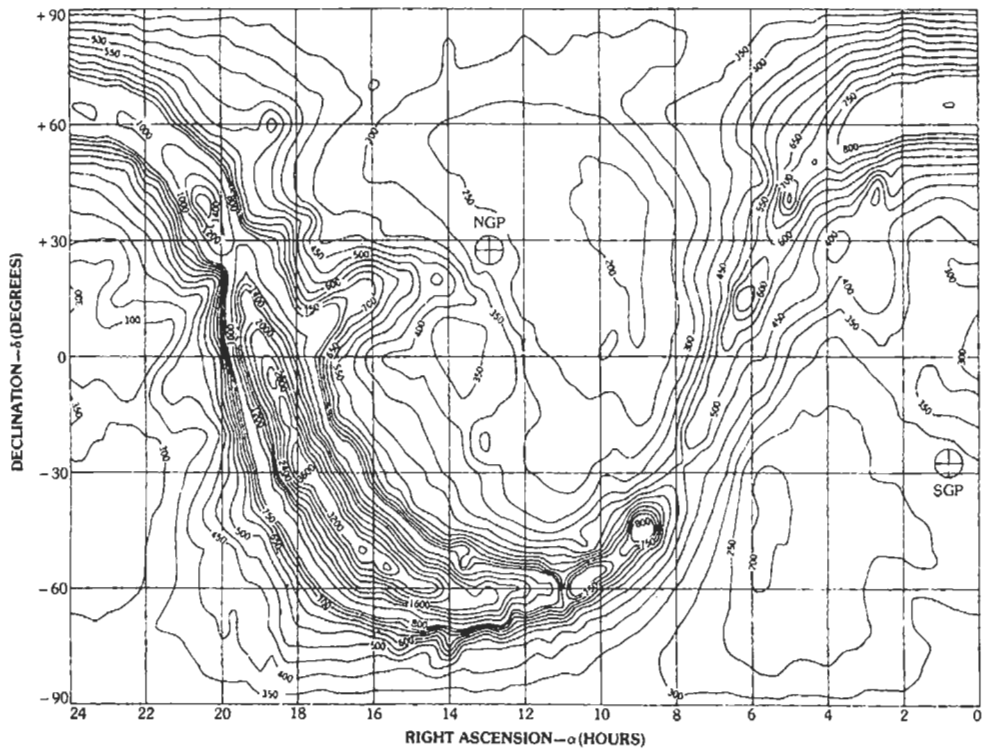


Fig. 8. Noise temperature (kelvins), 136 MHz. (After Taylor.)

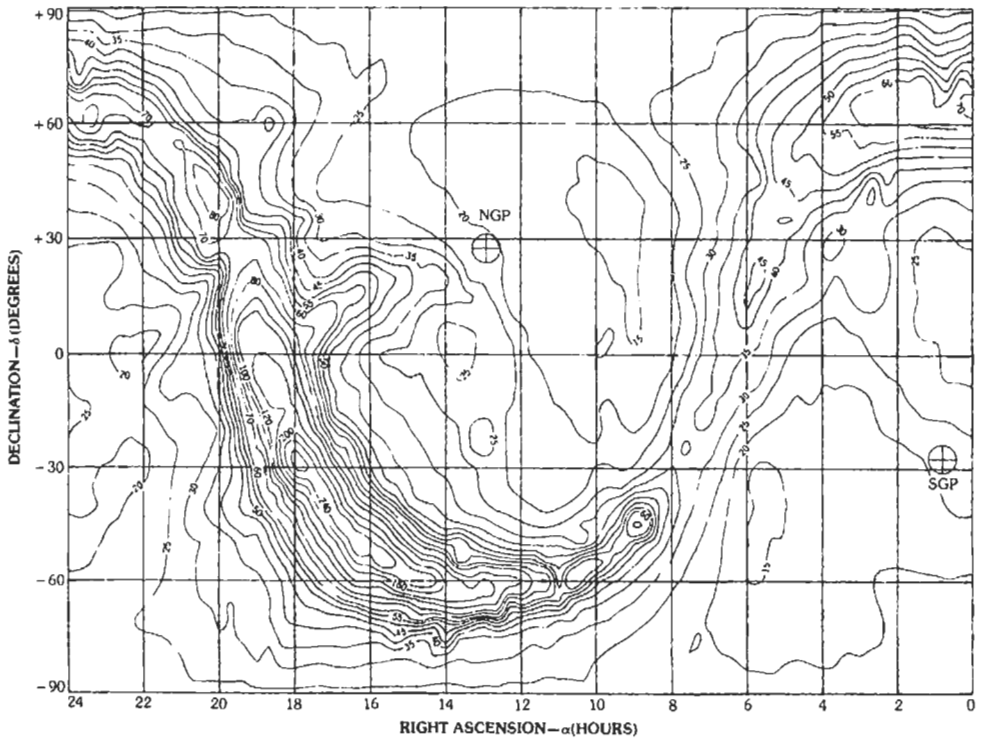


Fig. 9. Noise temperature (kelvins), 400 MHz. (After Taylor.)

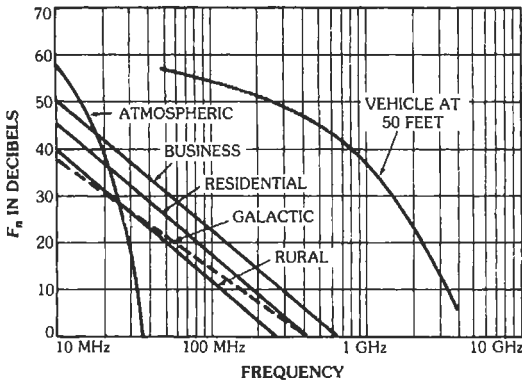


Fig. 10. Composite noise data. (After Hansen.)

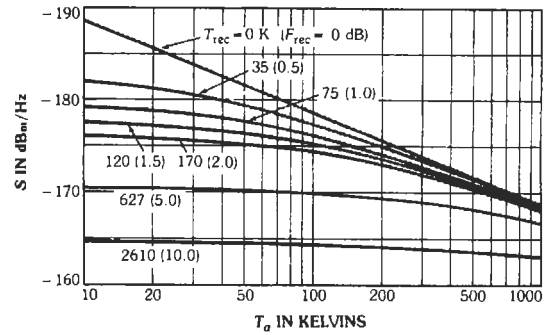


Fig. 12. Absolute sensitivity $S = K(T_a + T_{rec})$.

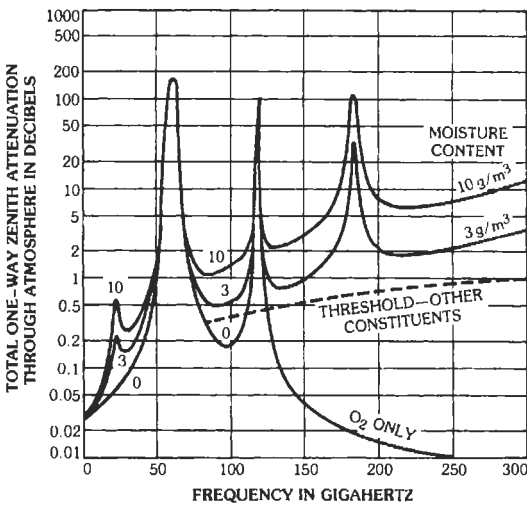


Fig. 11. Zenith atmospheric absorption.

est molecular resonance being that for water at 22.235 GHz. Fig. 11 from Smith* gives a composite of absorption data, and the transmission windows around 33, 95, 140, and 230 GHz may be observed.

From the two-dimensional antenna pattern and these noise data, the antenna noise temperature may be readily calculated. The total receiver noise power is then $N = KT_{sys}$ where $K = 198.60$ dBm. System sensitivity is shown in Fig. 12 in dBm/Hz for various receiver and antenna noise temperatures.

Elliptical and Circular Polarization

An electromagnetic wave is linearly polarized when the electric field lies wholly in one plane containing the direction of propagation. A plane electromagnetic wave, at a given frequency, is elliptically polarized when the extremity of the electric vector describes an

ellipse in a plane perpendicular to the direction of propagation, making one complete revolution during one period of the wave. If the rotation is clockwise looking in the direction of propagation, the sense is right-hand. More generally, any field vector, electric, magnetic, or other, is elliptically polarized if its extremity describes an ellipse. Two perpendicular axes, OX and OY , are chosen for reference in the plane of the polarization ellipse (Fig. 13). This plane is usually perpendicular to the direction of propagation. At a given frequency, the field components along these axes are represented by two complex numbers.

$$X = |X| \exp j\theta_1$$

$$Y = |Y| \exp j\theta_2$$

The axial ratio, r , is the ratio of the minor axis, OB , to the major axis, OA , of the polarization ellipse. The relative power coupling, η , between two elliptically polarized receiving antennas is given by

$$\eta = \frac{1}{2} \left[1 + \frac{(1-r_1^2)(1-r_2^2) \cos 2\theta \pm 4r_1r_2}{(1+r_1^2)(1+r_2^2)} \right]$$

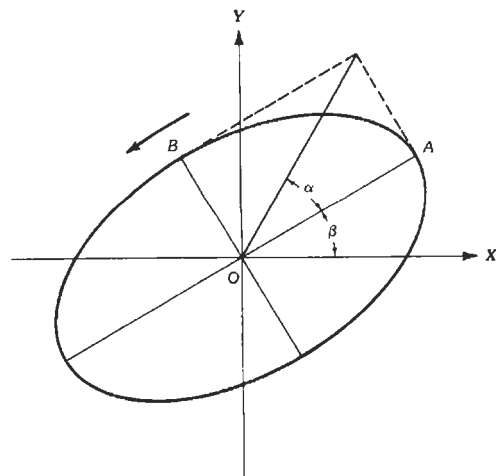


Fig. 13. Polarization ellipse.

* Smith, E. K. Private communication, 1982.

In this equation, r_1 and r_2 are the axial ratios of the antennas, and θ is the angle between the major axes of the polarization ellipses. The plus sign is used if both antennas produce the same sense of polarization, and the minus sign is used for opposite senses.* When a typical circularly polarized antenna is measured by a perfect rotating linear antenna, the circularly polarized gain is calculated from

$$G_{\text{cir}} = \frac{1}{2} [\sqrt{G_{\text{lin max}}} + \sqrt{G_{\text{lin min}}}]^2$$

Figs. 14 and 15 show transmission efficiency from a nearly linear antenna to a nearly circular antenna for both senses of polarization and for both zero and 90-degree orientation. When antenna gain is measured with a nominally linear standard-gain antenna where the orthogonal powers for two linear polarizations are summed, there is an error if the standard-gain antennas are imperfect. If the antenna under test has axial ratio r_1 and the standard-gain antenna has axial ratio r_2 , the error in gain measurement is as follows:

$$\text{Error} = 1 \pm 4 r_1 r_2 / (1 + r_1^2) (1 + r_2^2)$$

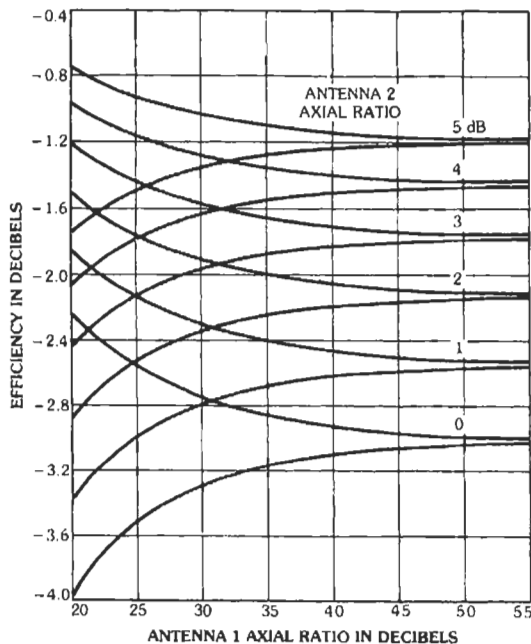


Fig. 14. Near-linear to near-circular efficiency, orientation 0. (Courtesy J. J. Murphy, Aerospace Corp.)

Figs. 16 and 17 show these errors for same sense and opposite sense of polarization.

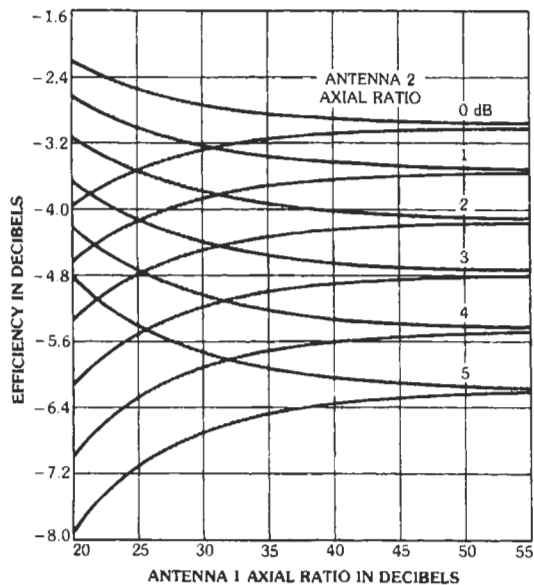


Fig. 15. Near-linear to near-circular efficiency, orientation 90. (Courtesy J. J. Murphy, Aerospace Corp.)

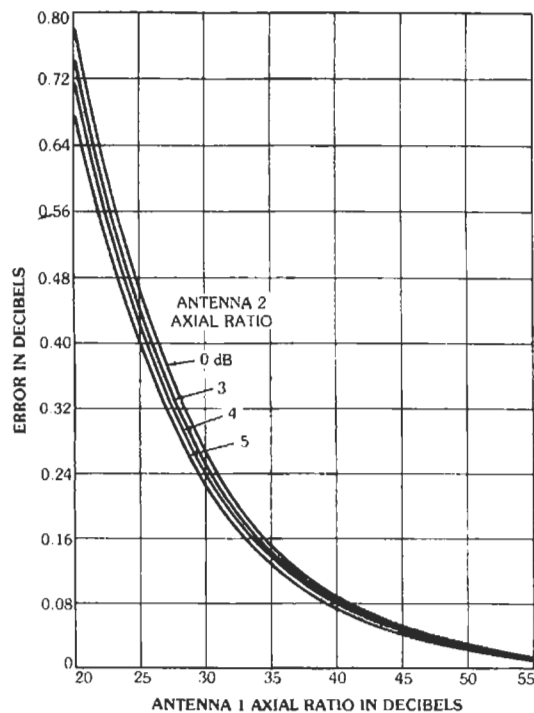


Fig. 16. Error in gain measurement, same sense. (Courtesy J. J. Murphy, Aerospace Corp.)

* Hatkin, L. "Elliptically Polarized Waves." *Proc. IRE*, Vol. 38, December 1950, p. 1455.

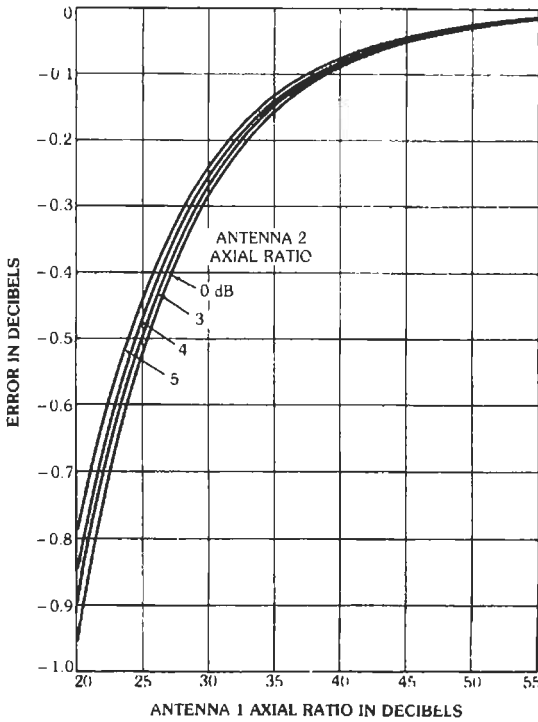


Fig. 17. Error in gain measurement, opposite sense. (Courtesy J. J. Murphy, Aerospace Corp.)

LOW-GAIN ANTENNAS

Half-Wave Dipole

The resonant dipole, with length near a half wavelength, is a principal canonical antenna. The half-wave dipole has a rotationally symmetric pattern of $\cos[(\pi/2) \cos \theta] / \sin \theta$, where $\theta=0^\circ$ gives the dipole axis. Half-power beamwidth is 78.1° , directivity is $1.641 = 2.15$ dB, effective area is $A_e = 0.131\lambda^2$, and effective length is 0.318λ . When the dipole is not oriented along a spherical-coordinate-system z axis, but instead along x (Fig. 18), the pattern becomes

$$E(\theta, \phi) = \frac{\cos\left\{\left(\frac{\pi}{2}\right) \sin \theta \cos \phi\right\}}{\sqrt{1 - \sin^2 \theta \cos^2 \phi}}$$

The far-field half-wave dipole field expressions, with I_0 the feed current, are

$$E_\theta = j \frac{60 I_0 e^{-jk r} \cos\left[\left(\frac{\pi}{2}\right) \cos \theta\right]}{r \sin \theta}$$

$$H_\phi = j \frac{I_0 e^{-jk r} \cos\left[\left(\frac{\pi}{2}\right) \cos \theta\right]}{2 \pi r \sin \theta}$$

These formulas can be used to calculate power densities from $PD = E^2 / 120 \pi$ as close as several wave-

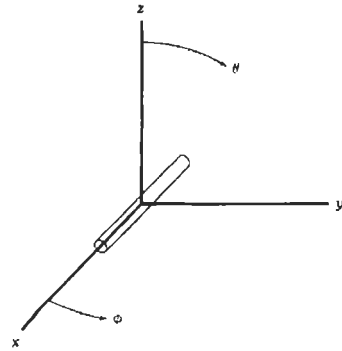


Fig. 18. Dipole along x -axis.

lengths from the antenna. At closer distances, use the near-field results given below. When the dipole length is not $\lambda/2$, the pattern becomes

$$E(\theta) = \frac{\cos(kh \cos \theta) - \cos kh}{\sin \theta}$$

Exact near-fields for thin dipoles are:

$$E_\rho = j(30 I_0 / \rho) \left[(z-h) e^{-jk r_1} / r_1 + (z+h) e^{-jk r_2} / r_2 - (2z \cos kh) e^{-jk r_0} / r_0 \right]$$

$$E_z = -j(30 I_0) \left[e^{-jk r_1} / r_1 + e^{-jk r_2} / r_2 - (2 \cos kh) e^{-jk r_0} / r_0 \right]$$

$$H_\phi = -j(I_0 / 4 \pi \rho) \left[e^{-jk r_1} + e^{-jk r_2} - (2 \cos kh) e^{-jk r_0} \right]$$

Note that this exact near-field is given in a cylindrical coordinate system (ρ, ϕ, z) with the dipole along z . Distances r_1, r_2, r_0 are from the farthest tip, the nearest tip, and the center of the dipole, respectively, to the observation point. These equations describe the discovery many years ago by Schelkunoff that a thin wire radiates a spherical wave from each tip, and from the center if the length is not a half wavelength.

Self-impedance of thin dipoles can be calculated approximately by the Carter zero-order theory, which assumes a sinusoidal current distribution. See Hansen* for an efficient computer algorithm. Fat dipoles require complex integral-equation theory, which is impractical to calculate, or moment method calculations. Because the feed-region geometry greatly affects the input impedance, a common and satisfactory procedure is to use zero-order theory for quick values, and then to measure the input impedance of the actual structure. Fig. 19 shows Z_{in} for several values of h/a and a/λ , where a is the radius. Dipoles are usually made of a

* Hansen, R. C. "Formulation of Echelon Dipole Mutual Impedance for Computer." *Trans. IEEE*, AP-20, Nov. 1972, pp. 780-781.

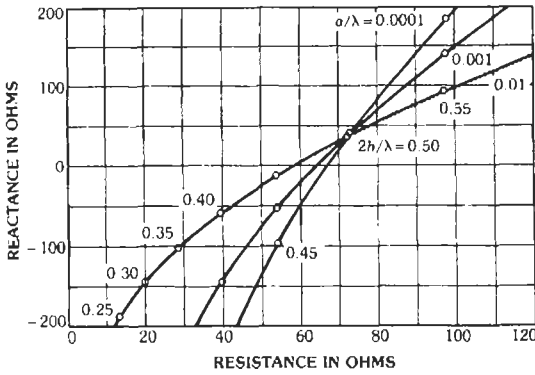


Fig. 19. Dipole self-impedance, zero order.

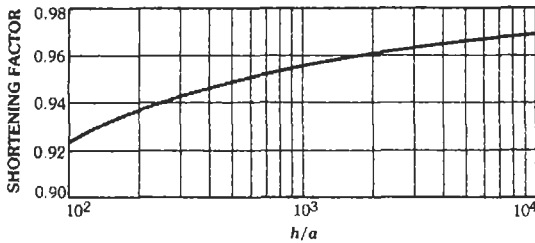


Fig. 20. Dipole shortening due to diameter.

solid or tubular cylinder; when a flat (strip) dipole is used, it is closely equivalent to a cylindrical dipole with $a = w/4$, where w is the strip width. Again, when zero-order theory is used, Fig. 20 gives shortening of resonant length versus h/a . Only a vanishingly thin dipole is resonant at halfwave. The resonant input resistance changes only slowly with a/λ , as seen in Table 3. In the limit of $a = 0$, the resonant resistance is 73.13 ohms.

Dipoles are frequently used parallel to a back screen, especially in arrays. If the dipole-screen spacing is s , the overall pattern, assuming the screen extends past the dipole by several wavelengths in each direction, is the dipole pattern multiplied by the screen array factor:

$$P(\theta) = P_{\text{dipole}}(\theta) \times \sin(ks \cos \theta)$$

TABLE 3. RESONANT RESISTANCE VS a/λ

a/λ	R_{res} (ohms)
0	73.13
0.00001	68.26
0.0001	66.79
0.001	64.11
0.01	58.16

where θ is zero normal to the screen. From this, the beamwidths in the E - and H -planes are 72.67° and 120° . Directivity is determined from the self- and mutual impedance:

$$D = 480/(R_{11} - R_{12})$$

Using zero-order mutual impedance, $R_{11} = 73.13$ ohms, and a spacing $s = \lambda/4$ gives $D = 5.603 = 7.48$ dB. Actual directivities will be different because the R_{11} for actual dipoles is not 73 ohms. In principle, use of smaller spacing should increase gain; in practice, when the dipole and its image in the screen are close together the impedance drops so that loss resistance limits the gain.* See Fig. 21.

Sometimes when space is limited, the dipole ends are bent, forming a U dipole. This affects both resonant length and resistance; these have been calculated using moment methods† for $a/\lambda = 0.01$. Fig. 22 shows the change in resonant length and the reduction in radiation resistance as more of the length is in the arms.

Monopoles (whips) are widely used for communications, e.g., on vehicles. When a quarter wavelength cannot be accommodated, a shorter whip can be loaded with a series inductance. The loading coil keeps the current distribution roughly constant from the feed to (and slightly past) the coil, thus increasing radiation resistance by the square of the current moment over the triangular current distribution of a short whip. Moment method studies‡ have shown that efficiency varies slowly with position of the coil, with a peak in the $0.3h$ to $0.5h$ range. Input resistance increases as the

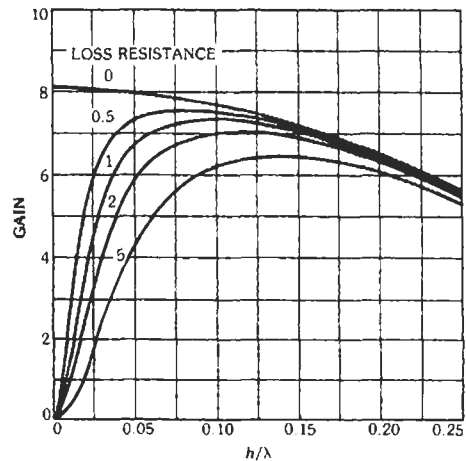


Fig. 21. Effect of loss on dipole over screen.

* Kraus, J. D. *Antennas*. 2nd ed. New York: McGraw-Hill Book Co., 1988.

† Hansen, R. C. Array chapters in *Handbook of Antenna Design*, Vol. 2. A. W. Rudge et al. (Eds.). London: Peter Pergrinus Ltd., 1983.

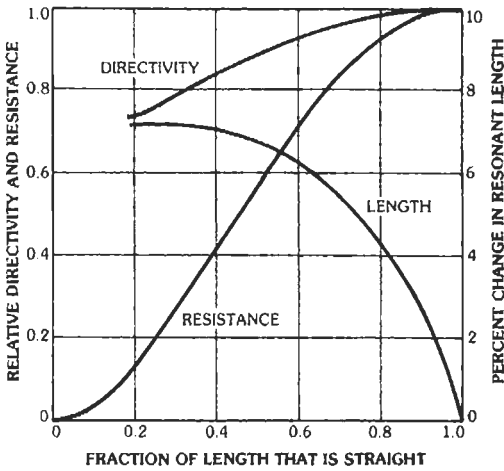


Fig. 22. Characteristics of resonant U dipoles.

coil moves toward the end, and the coil reactance also increases. For whips less than 0.1λ long, it is difficult to produce efficiencies of even 50 percent with high- Q coils. Extensive design curves are given in the reference. Short whips with high impedance are covered in the general section.

The broadest-bandwidth dipole is the open sleeve dipole, which consists of a dipole with two tubes parallel to the dipole, one on each side (Fig. 23). King and Wong* give performance data for various dimensions. The sleeves need not be tubular, but can be flat strips. When crossed dipoles are used for circular polarization, the sleeves can be metal plates or disks. Bandwidth of an octave for VSWR roughly 2 is achievable.

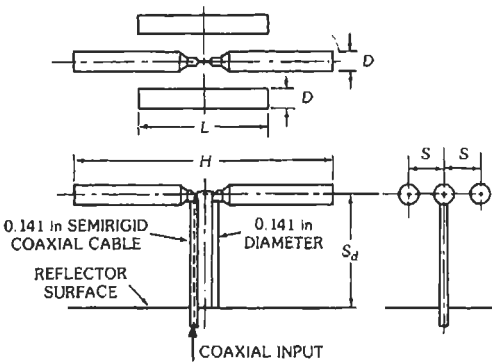


Fig. 23. Open-sleeve dipole.

‡ Hansen, R. C. "Efficiency and Matching Tradeoffs for Inductively Loaded Short Antennas." *Trans. IEEE*, Vol. COM-23, April 1975, pp. 430-435.

* King, H. E., and Wong, J. L. "An Experimental Study of a Balun-Fed Open-Sleeve Dipole in Front of a Metallic Reflector." *Trans. IEEE*, Vol. AP-20, 1972, pp. 201-204.

The behavior of the earth as a reflecting surface is considerably different for horizontal than for vertical polarization. For horizontal polarization, the earth may be considered a perfect conductor; that is, the reflected wave at all vertical angles β is substantially equal to the incident wave and 180° out of phase with it. The approximation is good for practically all types of ground.

For vertical polarization, however, the problem is much more complex, as both the relative amplitude, K , and relative phase, ϕ , change with vertical angle β and vary considerably with different types of ground. Fig. 24 is a set of curves that illustrate the problem. The subscripts to the amplitude and phase coefficients, K and ϕ , refer to the type of polarization, H for horizontal, and V for vertical.

Radial-Wire Ground Systems

A vertical tower antenna for MF and below usually requires a ground system to act as a low-loss "spreader" of current along the surface of the earth and to increase efficiency. Typically, the antenna base is surrounded by a metallic disk or closely spaced wire mesh, with radial wires emanating from this. It is not always possible or economical to make the radials sufficiently long and dense that the monopole perfor-

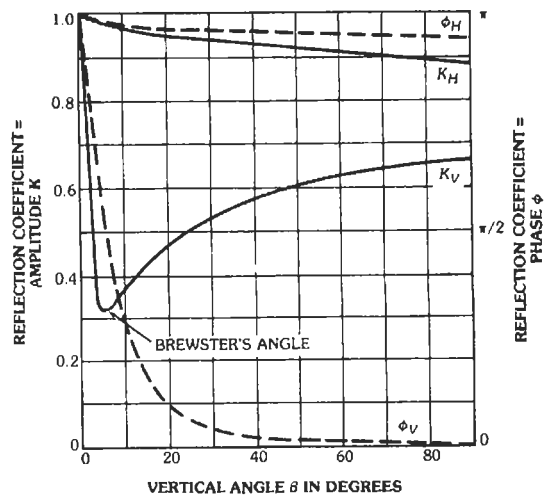
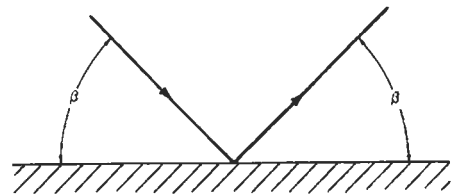


Fig. 24. Typical ground-reflection coefficients for horizontal and vertical polarizations.

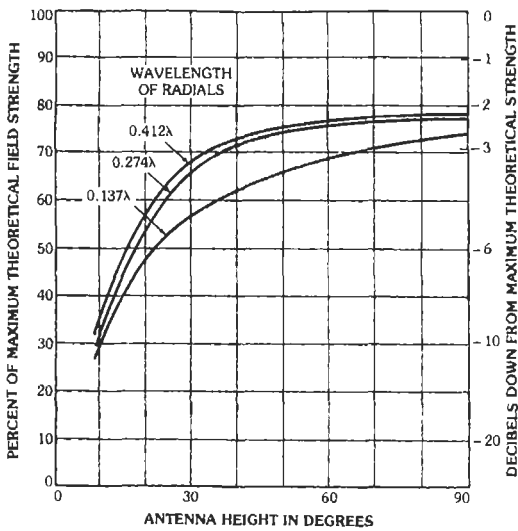
mance is independent of the earth parameters. There is some evidence that surface wire radials, if sufficiently long, give higher efficiency than buried radials. Extensive calculations have been made by Maley and King,* but comparisons with measurements are sparse.

* Maley, S. W., and King, R. J. "Impedance of a Monopole Antenna with a Radial-Wire Ground System on an Imperfectly Conducting Half Space." *Rad. Sci.*, Part 1, March-April 1962, pp. 175-180; Part 2, February 1964, pp. 159-165; Part 3, March 1964, pp. 297-301.

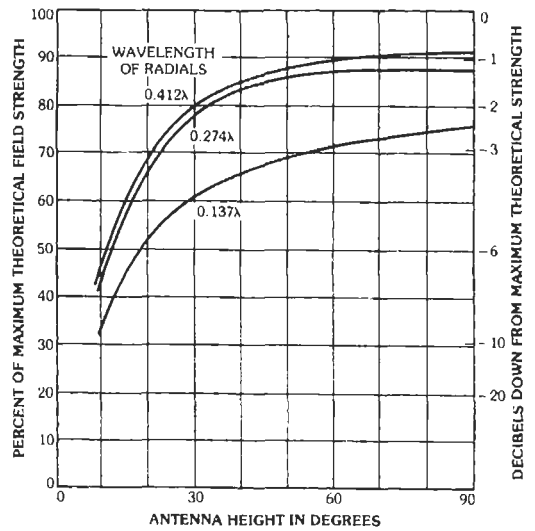
Fig. 25 shows experimental measurements at 3 MHz made by Brown et al.† and Christman.‡ Efficiency is given versus monopole height, and number and length of radials. For different earth, the absolute efficiencies will change, but the relative results are useful.

† Brown, G. H., Lewis, R. F., and Epstein, J. "Ground Systems as a Factor in Antenna Efficiency." *Proc. IRE*, Vol. 25, June 1937, pp. 753-787.

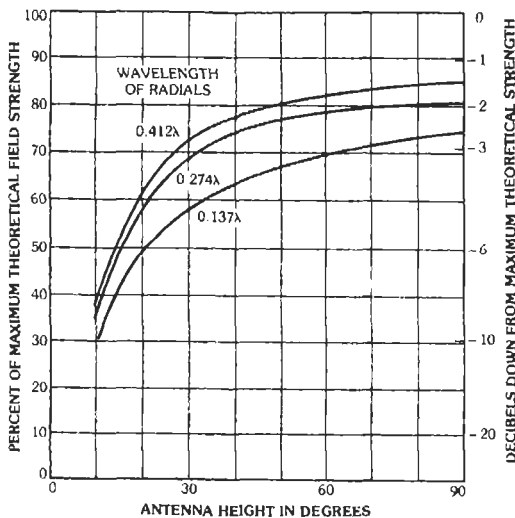
‡ Christman, A. M. "Ground Systems for Vertical Antennas." *Ham Radio*, August 1979, pp. 31-33.



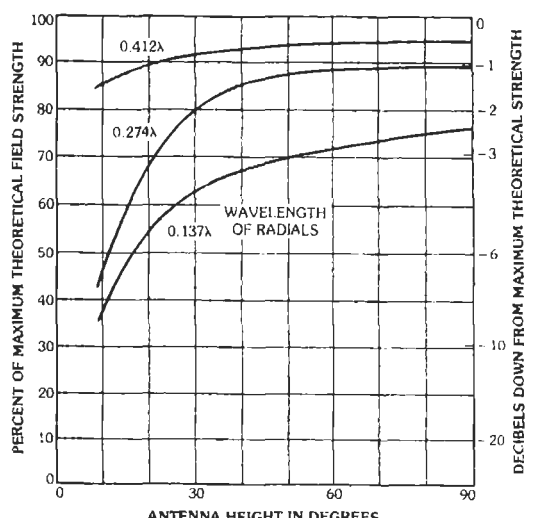
(A) 15 radials.



(C) 60 radials.



(B) 30 radials.



(D) 113 radials.

Fig. 25. Efficiency vs. monopole height.

Printed-Circuit Antennas

Printed-circuit antennas include stripline slots, printed-circuit-board dipoles, and microstrip patches. A stripline slot consists of a narrow rectangular slot cut in the top stripline ground plane, with the slot excited by proper positioning of the center conductor below the slot. A linear resonant array of collinear slots may be fed by a single center conductor, with this conductor centered under each slot at the proper angle.* Boxed stripline is used to reduce higher modes.

Dipoles may be printed for low-cost fabrication, either as an array of dipoles on a single dielectric substrate, with each dipole fed by a balun which is normal to the dipole array face, or with each dipole and balun on a separate dielectric substrate. The array comprises a stack of these sheets.

The most widely used printed-circuit antenna is the microstrip patch, which in its simplest form is a rectangular or circular patch of metal fed by the microstrip upper conductor; see Fig. 26. Thus the element and feed line, and usually other elements, power dividers, etc., can all be prepared as a single etched PC board. A patch can also be fed from below by a coaxial probe, but this in large part removes the cost advantage. Early work on microstrip antennas was done by Deschamps and Sichak.† A review is given by Pozar.‡

Since microstrip is thin, the patch and ground plane acts as a TM resonator, with the sidewalls (the gaps) acting approximately as magnetic walls. Treating the patch as a lowest mode resonator gives useful results.

More recent methods of analysis make use of a series of cavity modes to match a side wall admittance that

includes radiation conductance and discontinuity susceptance.§ More rigorous analyses take into account the effect of microstrip or probe feed." The major disadvantage of patches is their narrow bandwidth, due to the close spacing between patch and ground plane. Fig. 27 shows bandwidth versus substrate height in free space wavelengths, for VSWR < 2. #

These are zero order results based on transmission line analysis. Loss and surface wave effects are not included. The bandwidth is approximately:

$$BW \approx \frac{96t / \lambda_0}{\sqrt{2\epsilon} [4 + 17\sqrt{\epsilon}]}$$

Increasing the substrate height does increase the bandwidth, but the coupling to surface wave modes

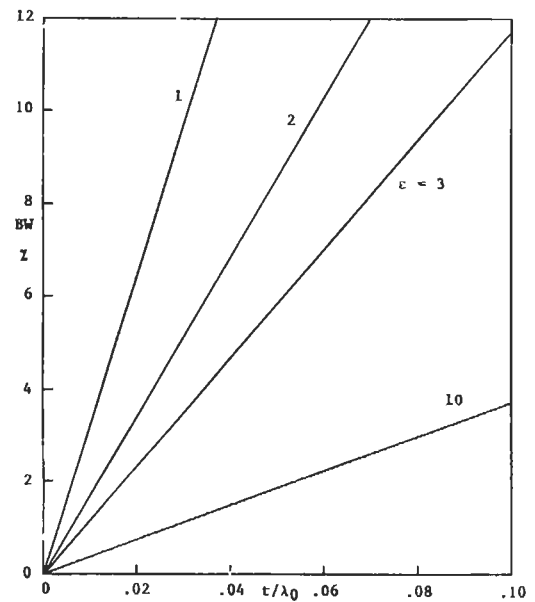


Fig. 27. Bandwidth of lossless patch, VSWR = 2.

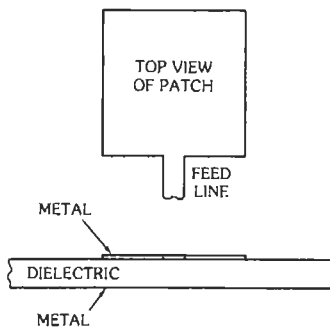


Fig. 26. Microstrip patch.

* Park, P. K., and Elliott, R. S. "Design of Collinear Longitudinal Slot Arrays Fed by Boxed Stripline." *Trans. IEEE*, Vol. AP-29, January 1981, pp. 135-140.

† Deschamps, G., and Sichak, W., "Microstrip Microwave Antenna." *Proc. 1953 Allerton Antenna Symposium*, University of Illinois.

‡ Pozar, D. M. "Microstrip Antennas." *Proc. IEEE*, Vol. 80, January 1992, pp. 79-91.

§ Lo, Y. T., et al. "Theory and Experiment on Microstrip Antennas." *Trans. IEEE*, Vol. AP-27, March 1979, pp. 137-145. Carver, K. R. "Practical Analytical Techniques for the Microstrip Antenna." *Printed Circuit Antenna Technology Workshop*, New Mexico State University, 1979. Hammerstad, E.O. "Equations for Microstrip Circuit Design." *Proc. European Microwave Conf.*, 1975, pp. 268-272.

Pozar, D. M., and Voda, S. M. "A Rigorous Analysis of a Microstripline Fed Patch Antenna." *Trans. IEEE*, Vol. AP-35, December 1987, pp. 1343-1350. Aberle, J. T., and Pozar, D. M. "Analysis of Infinite Arrays of One- and Two-Probe-Fed Circular Patches." *Trans. IEEE*, Vol. AP-38, April 1990, pp. 421-432.

Hansen, R. C. and Burke, M. "Antennas with Magneto-dielectrics." *Micro. Opt. Tech. Ltrs.*, Vol. 26, July 20, 2000, pp. 75-78.

and the cross-polarized radiation also increase, with undesirable pattern effects. Surface waves may be a problem for larger arrays or large substrates. An easy technique to obviate these is to make the dielectric anisotropic, by drilling holes perpendicular to the slab. Such a slab is sometimes called a Photonic Band Gap material in analogy to quantum mechanics, but this nomenclature is at best misleading. Most techniques for extending the bandwidth of patches have done so at the expense of efficiency, i.e., through use of matching networks with high standing waves. However, the parasitic patch* avoids these problems by keeping the substrate thickness low, with a parasitic patch above the driven patch increasing the effective radiation height. Bandwidth can be doubled, with the parasitic patch dimensions and height above the driven patch adjusted to give a symmetrical impedance curve. Fig. 28 sketches the parasitic patch configuration. These may be arrayed as are ordinary patches. In practice, the parasitic patches are printed on a thin dielectric substrate, with a foam layer used to support this substrate above the microstrip.

Slot coupling can also increase bandwidth. Here coupling is through a slot in the patch groundplane; the latter is usually a stripline plane.† Bandwidth can be roughly doubled as the electromagnetic coupler serves as a matching circuit.

Rectangular or circular patches, as described, are mostly linearly polarized radiators, with patch widths roughly a third of a wavelength. Accordingly, the patterns are between those of a short dipole and a half-wave dipole. Although a square patch could be fed on two adjacent sides with the proper phases to obtain circularly polarized radiation, simpler circularly polarized patches have been developed.‡ The key is to modify dimensions to allow the two cross-polarized

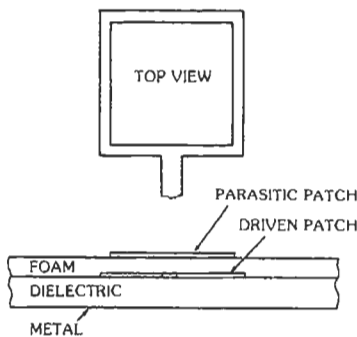


Fig. 28. Parasitic microwave patch.

* Hall, P. S., et al. "Wide Bandwidth Microstrip Antennas for Circuit Integration." *Elec. Ltrs.*, 15, 1979, pp. 458-460.

† Sullivan, P. L., and Schaubert, D. H. "Analysis of Aperture Coupled Microstrip Antenna." *Trans. IEEE*, Vol. AP-34, Aug. 1986, pp. 977-984.

modes to be of equal amplitude and 90 degrees out of phase. Fig. 29 sketches four ways of accomplishing this. A simple analysis has been given by Lo and Richards.§ Separate feeds provide wider bandwidth.

More sophisticated patch cutouts, such as ell, cross, etc., have been used to produce CP, or even dual frequency operation. All of these, however, tend to reduce the bandwidth.

Loops

Loops small in wavelengths are equivalent to short magnetic dipoles (see the general section at the beginning of this chapter). When the circumference of a loop is small in wavelengths and the area is A, the patterns are given by

$$E_{\theta} = (120\pi^2 I_0 A \mu_e \sin \theta) / r \lambda^2$$

$$H_{\theta} = (\pi I_0 A \mu_e \sin \theta) / r \lambda^2$$

where

I_0 is the loop current,

θ is measured from the loop axis.

Radiation resistance is

$$R_{rad} = 20N^2 k^4 A^2 \mu_e^2$$

where μ_e is the effective permeability of a magnetic core. Effective length is $l_e = NkA\mu_e$.

The effective area and directivity are those of a short dipole: $3\lambda^2/8\pi$ and 1.5. Magnetically cored loops almost universally use ferrite cores with the highest permeability available for the frequencies of interest.

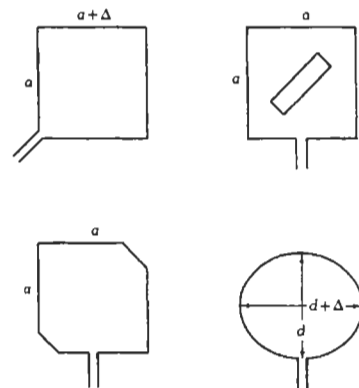


Fig. 29. Circularly polarized patches.

‡ Kerr, J. L. "Microstrip Antenna Developments." *Proc. Printed Circuit Antenna Technology Workshop*, October 1979, New Mexico State University.

§ Lo, Y. T., and Richards, W. F. "Perturbation Approach to Design of Circularly Polarized Microstrip Antennas." *Elec. Ltrs.*, Vol. 17, May 28, 1981, pp. 383-385.

Typically, the coil diameter is small in wavelengths, and the core is long compared with the diameter. Whether the coil is distributed over the core or concentrated near the center is important mostly for practical factors such as distributed capacitance. Tape cores can be made of Permalloy type alloys or of rapidly cooled amorphous alloys such as Metglas; these can be used up to a few megahertz with permeabilities of many tens of thousands (see Chapter 4); ferrite cores extend into several hundred megahertz, with permeabilities of thousands at lower frequencies and on the order of 20 at upper frequencies. Effective permeability can be calculated by assuming the core is a prolate spheroid. With the formulas of Wait,* calculations of μ_e , in terms of L/d , the core length/diameter ratio, have been made. Fig. 30 shows these data, which come from:

$$\mu_e/\mu = 1/[(1 + \mu - 1)(\zeta^2 - 1) Q_1]$$

where

$$\zeta = (L/d)/\sqrt{L^2/d^2 - 1}$$

and

$$Q_1 = (\zeta/2) \ln [(\zeta + 1)/(\zeta - 1)] - 1$$

Lower permeabilities require larger L/d to realize μ_e as a significant fraction of μ . The core need not be solid, as much of the magnetic field is concentrated

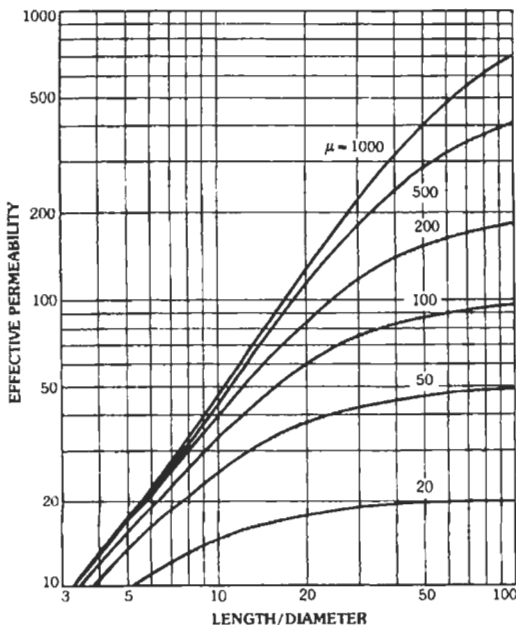


Fig. 30. Effective permeability of spheroidal core.

* Wait, J. R. "Receiving Properties of a Wire Loop with a Spheroidal Core." *Can. J. Tech.*, Vol. 31, January 1953, pp. 9-14.

near the surface. Again, a hollow cylindrical core can be approximated by two prolate spheroids, and a formula for μ_e was developed by Wait.† Figs. 31, 32, and 33 show hollow-core effective permeability for $\mu = 20$, 100, and 500, respectively.‡ It is apparent that a significant part of the solid core may be eliminated. Such a hollow core can be assembled from parallel rods or strips around a foam block.

Slot Antennas

The properties of many slot antennas can be deduced from the properties of the complementary metallic antenna as enunciated by Booker in his exten-

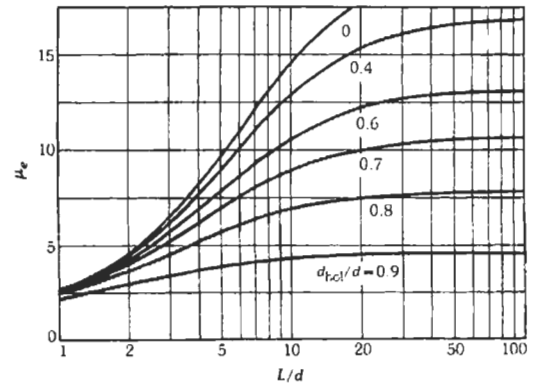


Fig. 31. Effective permeability of hollow spheroidal core, for $\mu = 20$.

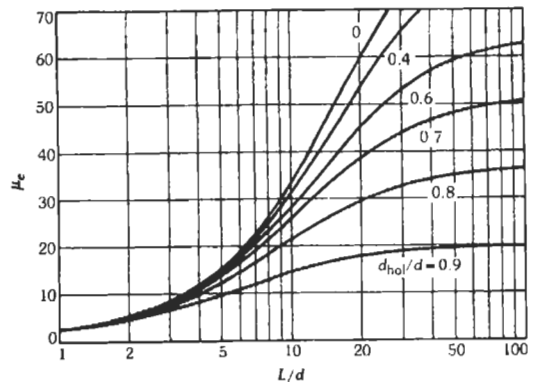


Fig. 32. Effective permeability of hollow spheroidal core, for $\mu = 100$.

† Wait, J. R. "The Receiving Loop with a Hollow Prolate Spheroidal Core." *Can. J. Tech.*, Vol. 31, June 1953, pp. 132-137.

‡ Hansen, R. C. In *Microwave Engineers Handbook*, Vol. 2, T. S. Saad (Ed.). Dedham, MA: Horizon House, 1971; pp. 44-46 and 89-96.

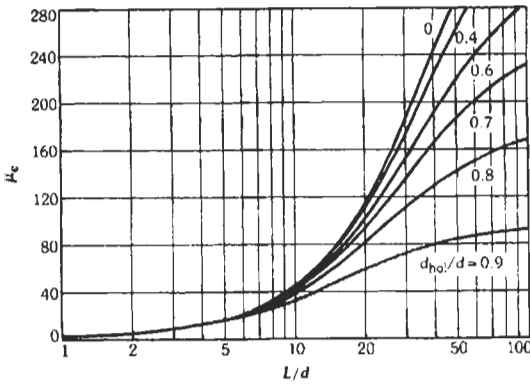


Fig. 33. Effective permeability of hollow spheroidal core, for $\mu = 500$.

sion of the optical work of Babinet. The admittance, Y_s , of the slot antenna is related to the impedance, Z_m , of the metallic antenna by

$$Z_m/Y_s = (120\pi)^2 / 4$$

The magnitude of the electric field, E_s , produced by the slot is proportional to the magnitude of the magnetic field, H_m , of the metallic antenna, and H_s is proportional to E_m . The electric- and magnetic-plane patterns of the slot are similar to the magnetic- and electric-plane patterns, respectively, of the metallic antenna. In Fig. 34, a rectangular slot antenna in an infinite metallic plane is shown with the complementary strip dipole. If the slot radiates on only one side, the relationship is

$$Z_m/Y_s = (120\pi)^2 / 2$$

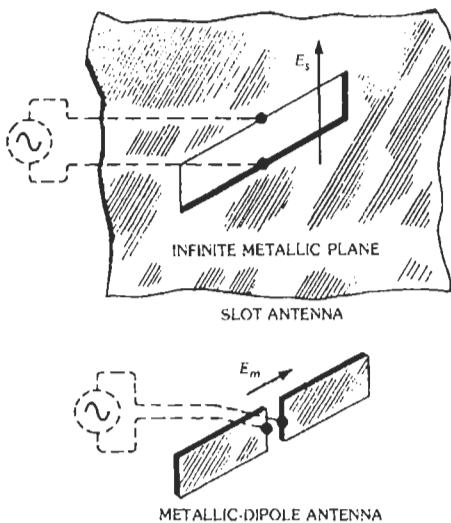


Fig. 34. Slot antenna and strip dipole.

The E -plane pattern of the slot and H -plane pattern of the dipole are omnidirectional, while the slot H -plane pattern is the same as the dipole E -plane pattern.

Due to the ease of feeding, most slots are located in waveguide. Fig. 35 shows displaced longitudinal, inclined, and edge slots; the first and last are widely used for linear and planar arrays. The pattern and mutual admittance of slots in an array are close to those of the Booker complementary strip dipole array, and dipole mutual impedances are used in the design of such arrays. The slot self-admittance, however, contains a term involving reactive energy in the waveguide in the vicinity of the slot. Longitudinal slots are covered here because they are simpler and easier to use. Radiation conductance obtained from the Stevenson formula and for a sinusoidal slot field distribution is*

$$\frac{G_{rad}}{Y_0} = \frac{480(a/b)[\sin^2(\pi x/a)](\cos \beta l - \cos kl)^2}{\pi(\beta/k)R_0}$$

where

a and b are waveguide width and height,
 x and l are the slot displacement and length,
 R_0 is the Booker equivalent strip dipole resistance,
 β is the guide wavenumber.

Susceptance is difficult to calculate, with both moment methods† and variational methods‡ used. In any high performance array, it is necessary to make careful measurements of single-slot admittance for a family of slot offsets and lengths. Fig. 36 shows the Stevenson resonant conductance, which is a function only of offset. Slot resonant length involves three variables: x/a , a/b , and β/k . Fig. 37 gives resonant length for several frequency-guide combinations. These two

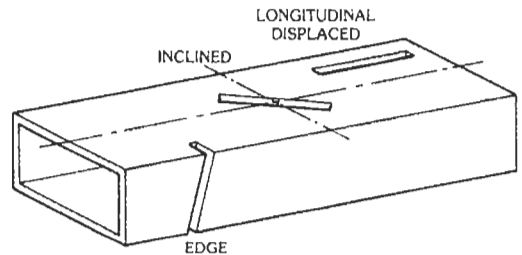


Fig. 35. Waveguide slots.

* Elliott, R. S., and Kurtz, L. A. "The Design of Small Slot Arrays." *Trans. IEEE*, Vol. AP-26, March 1978, pp. 214-219.

† Khac, T. V., and Carson, C. T. "Impedance Properties of a Longitudinal Slot Antenna in the Broad Face of a Rectangular Waveguide." *Trans. IEEE*, Vol. AP-21, Sept. 1973, pp. 708-710.

‡ Yec, H. Y. "Impedance of a Narrow Longitudinal Shunt Slot in a Slotted Waveguide Array." *Trans. IEEE*, Vol. AP-22, 1974, pp. 589-592. Hansen, R. C. *Phased Array Antennas*. New York: John Wiley & Sons, Inc., 1998.

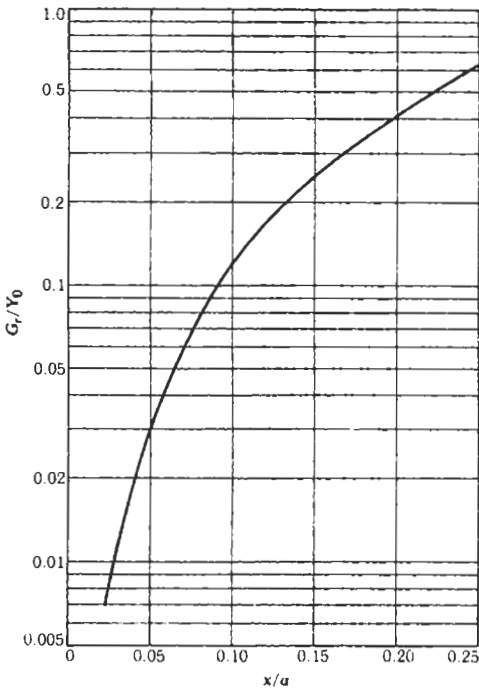


Fig. 36. Resonant conductance of longitudinal slot.

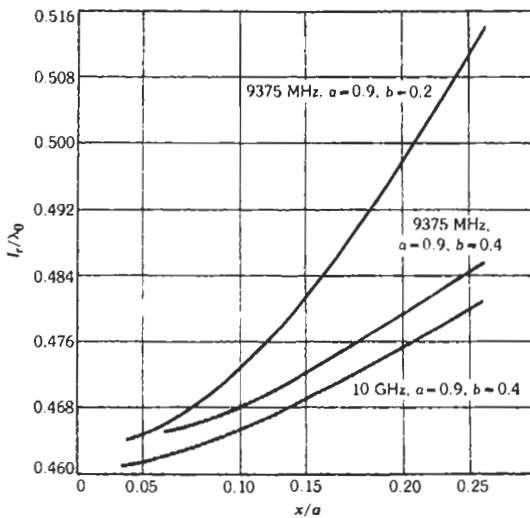


Fig. 37. Resonant length of longitudinal slot.

parameters, resonant length and conductance, are important because they allow slot admittances versus frequency to be plotted in the Kaminow and Stegen "universal" form.* Slot length is given in terms of resonant length, and both conductance and susceptance are given in terms of resonant conductance. The result

was originally thought to be a single symmetric curve for G/Y_0 and a single asymmetric curve for B/Y_0 . As Fig. 38 shows, it now appears that although all longitudinal slots fit these curves approximately, there are differences that vary with x/a , a/b , and β/k . The precise control of aperture distribution allowed by waveguide slots and the use of computer design have allowed high-performance low-sidelobe or tailored-sidelobe arrays to be constructed.†

Slots are sometimes used on cylinders, in either an axial or circumferential configuration. A circumferential half-wave slot has a smooth pattern, with the cylinder reducing the back lobes. Axial half-wave slots also exhibit shadowing, with pattern ripples around $\phi = \pi$. As ka increases, the number of ripples increases, but their depth and angular extent decrease. For both types of slots, larger cylinders give more shadowing and lower field at $\phi = \pi$. Fig. 39 shows typical axial and circumferential slot patterns in the azimuth plane for $ka = 8$. These were calculated with the formulas of Wait:‡

Axial

$$E_\phi = \frac{1}{\pi ka} \sum_{n=0}^{\infty} \frac{\epsilon_n j^n \cos n\phi}{H_n^{(2)'}(ka)}$$

Circumferential

$$E_\theta = \frac{ka}{\pi} \sum_{n=0}^{\infty} \frac{\epsilon_n j^n [\cos(n\pi/2ka)] \cos n\phi}{(k^2 a^2 - n^2) H_n^{(2)}(ka)}$$

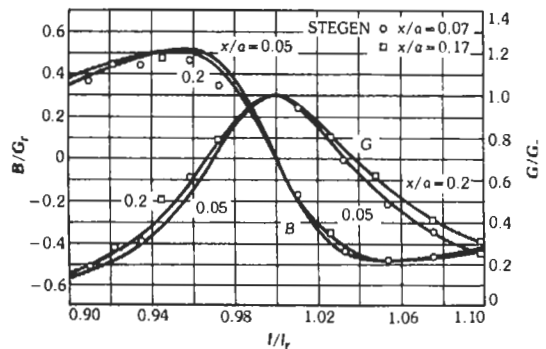


Fig. 38. Longitudinal slot admittance, $a = 0.9$, $b = 0.4$, frequency = 9375 MHz.

* Blass, J. "Slot Antennas." Chapter 8 of *Antenna Engineering Handbook*. R. C. Johnson and H. Jasik (Eds.). New York: McGraw Hill Book Co., 1984.

† Elliott, R. S. *Antenna Theory and Design*. Englewood Cliffs, NJ: Prentice-Hall, Inc., 1981. Hansen, R. C. *Phased Array Antennas*. New York: John Wiley & Sons, Inc., 1998.

‡ Wait, J. R. *Electromagnetic Radiation from Cylindrical Structures*. Elmsford, NY: Pergamon Press, Inc., 1959.

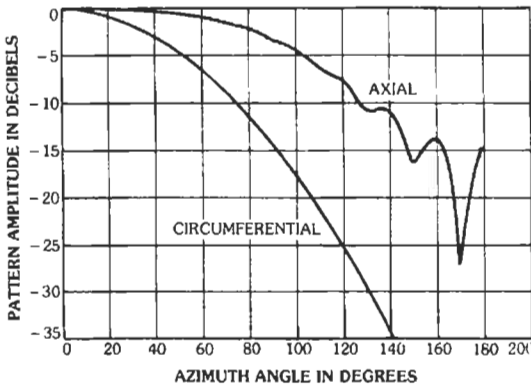


Fig. 39. Half-wave slot on cylinder, $ka = 8$.

Here $\epsilon_n = 1$ for $n = 0$, $\epsilon_n = 2$ for $n > 0$, and $H_n^{(2)}$ is the Hankel function of the second kind. Three-dimensional patterns can be calculated by replacing ka by $ka \sin \theta$, except that finite length cylinders will have end effects that require other calculation methods such as geometric theory of diffraction.*

MEDIUM-GAIN ANTENNAS

Horns

Horns are often used as gain standards, as feeds for dish antennas, and as microwave antennas directly. The simplest type of horn is a flared waveguide. For a flare only in the E -plane or only in the H -plane, directivity can be accurately and simply calculated:†

$$G_H = (64 a R_e / \pi \lambda b) [C^2(w) + S^2(w)] = 32 ab F_E / \pi \lambda^2$$

where

$$w = b / \sqrt{2\lambda R_e}$$

and

$$G_H = (4\pi b R_m / \lambda a) \{ [C(u) - C(v)]^2 + [S(u) - S(v)]^2 \} = 32 ab F_H / \pi \lambda^2$$

where

$$\sqrt{2}u = \sqrt{\lambda R_m} / a + a / \sqrt{\lambda R_m}$$

$$\sqrt{2}v = \sqrt{\lambda R_m} / a - a / \sqrt{\lambda R_m}$$

In the formulas above, R_e and R_m are the slant lengths, and C and S are conventional Fresnel integrals.‡ The formulas are also written in terms of directivity reduction factors F_E and F_H , which are shown in

* Hansen, R. C. (Ed.). *Geometric Theory of Diffraction*. IEEE Press/Wiley, 1981.

† Schelkunoff, S. A., and Friis, H. T. (Eds.). *Antennas, Theory and Practice*. New York: John Wiley & Sons, Inc., 1952.

Fig. 40. Figs. 41 and 42 show normalized directivity for each of these, and it may be seen that for each slant length there is an optimum width or height. More commonly, a pyramidal horn is used with dimensions selected to be optimum in both planes. However, this may not have coincident phase centers; see below. The

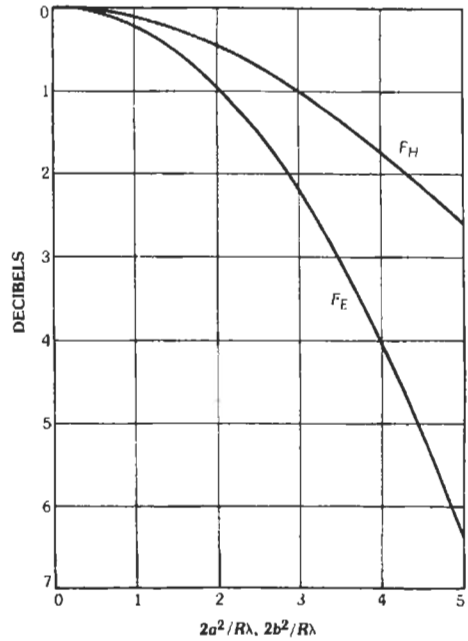


Fig. 40. Sectoral-horn directivity factors.

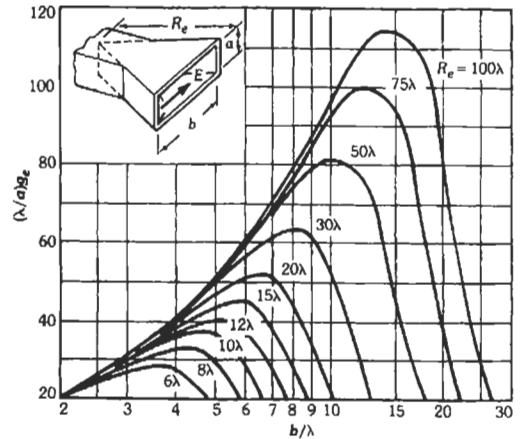


Fig. 41. Directivity of a large horn flared in the electric plane (After Schelkunoff.)

‡ Abramowitz, M. and Stegun, I. A. (eds.). *Handbook of Mathematical Functions*. NBS Applied Math Series, 1970.

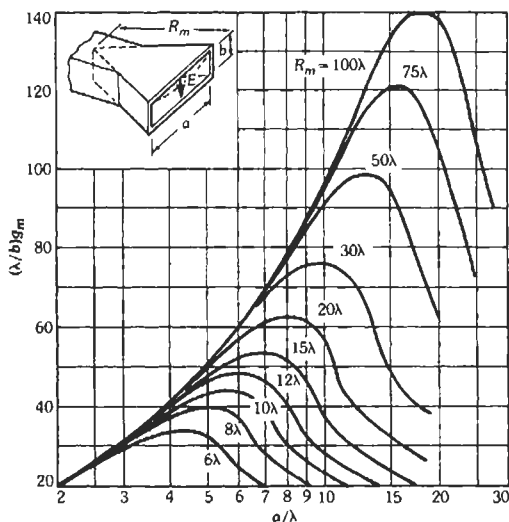


Fig. 42. Directivity of a large horn flared in the magnetic plane (After Schelkunoff.)

pyramidal horn directivity can be calculated from the constituent sectoral horn directivities by

$$G_{\text{pyramidal}} = \pi\lambda^2 G_E G_H / 32 ab = 32 ab F_E F_H / \pi\lambda^2$$

This can also be written in terms of *E*- and *H*-plane gain reduction factors as shown in the formula above. Half-power beamwidths are given approximately by $\theta_3 = 51 \lambda/b$ for the *E*-plane and $70\lambda/a$ for the *H*-plane. The beamwidths in each plane are essentially independent of dimensions in the other plane. Fig. 43 shows experimental patterns in both planes for a variety of radial lengths and flare angles.* Wide-angle patterns may be computed by using geometric theory of diffraction; see reprint book by Hansen.† Phase centers of pyramidal horns may also be calculated by using Fresnel integrals.‡ The phase center is generally different for each plane, and it is located toward, but not at, the horn apex. The optimum rectangular pyramidal horn has slightly different dimensions from one designed with coincident phase centers. For some applications this is important.

Pyramidal horns are usually rectangular, whereas corrugated horns are usually circular. The corrugated horn offers a pattern that is nearly symmetric; i.e., the *E*- and *H*-plane beamwidths are nearly equal. In addition, the sidelobes can be better controlled. Further discussion of corrugated horns is outside the scope

* Rhodes, D. R. "An Experimental Investigation of the Radiation Patterns of Electromagnetic Horn Antennas." *Proc. IRE*, Vol. 36, September 1948, pp. 1101–1105.

† Hansen, R. C. (Ed.). *Geometric Theory of Diffraction*. IEEE Press/Wiley, 1981.

‡ Muehldorf, E. I. "The Phase Center of Horn Antennas." *Trans. IEEE*, Vol. AP-18, November 1970, pp. 753–760.

here; refer to papers in the reprint book by Love§ and Olver et al."

In the TEM horn, developed by Kerr,# elimination of the sidewalls allows much greater bandwidth. Best performance obtains when both the conductor spacing and width increase roughly exponentially from throat to mouth. These horns have found wide use for short-pulse and transient applications. Printed circuit versions are also available.** These horns have been integrated with solid-state devices for use in the millimeter and submillimeter bands.††

Helices

A cylindrical helix can radiate a broadside pattern (normal mode) or an end-fire beam (axial mode). Only the latter is discussed here. The circumference must be on the order of a wavelength, and to utilize the end-fire properties effectively, a length on the order of several wavelengths is usually used. Fig. 44 shows beamwidth versus length for several circumferences. The axial-mode helix is broadband but is limited to roughly a 2.5:1 frequency ratio.‡‡ A wider bandwidth can be covered by use of a long multifilar helix such as a quadhelix.§§ In this configuration, the helix consists of four windings spaced 90 degrees with the wires fed with a 90-degree phase progression at the base. Although the feed is complicated, this type of helix does not require a ground plane, although a ground plane may improve performance. A short resonant version of the quadhelix that radiates an omnidirectional circularly polarized pattern is also available.¶¶ This quadhelix utilizes a fractional turn but operates with-

§ Love, A. W. (Ed.). *Electromagnetic Horn Antennas*. IEEE Press/Wiley, 1976.

¶ Olver, A. D., et al. *Microwave Horns and Feeds*. IEE Press, 1994.

Kerr, J. L. "Short Axial Length Broad-Band Horns." *Trans. IEEE*, Vol. AP-21, September 1973, pp. 710–714.

** Janaswamy, R., and Schaubert, D. H. "Analysis of the Tapered Slot Antenna." *Trans. IEEE*, Vol. AP-35, September 1987, pp. 1058–1065; Chio, T. H., and Schaubert, D. H. "Parameter Study and Design of Wide-Band Widescan Dual-Polarized Tapered Slot Antenna Arrays." *Trans. IEEE*, Vol AP-48, June 2000, pp. 879–886.

†† Rebeiz, G. M., et al. "Monolithic Millimeter-Wave Two-Dimensional Horn Imaging Arrays." *Trans. IEEE*, Vol. AP-38, September 1990, pp. 1473–1482. Schwering, F. K. "Millimeter Wave Antennas." *Proc. IEEE*, Vol. 80, January 1992, pp. 92–102.

‡‡ Kraus, J. D. *Antennas*. 2nd. ed. New York: McGraw-Hill Book Co., 1988.

§§ Adams, A. T., Greenough, R. K., Wallenberg, R. F., Mendelovicz, A., and Lumjiak, C. "The Quadrifilar Helix Antenna." *Trans. IEEE*, Vol. AP-22, March 1974, pp. 173–178.

¶¶ Kilgus, C. C. "Shaped-Conical Radiation Pattern Performance of the Backfire Quadrifilar Helix." *Trans. IEEE*, Vol. AP-23, May 1975, pp. 392–397.

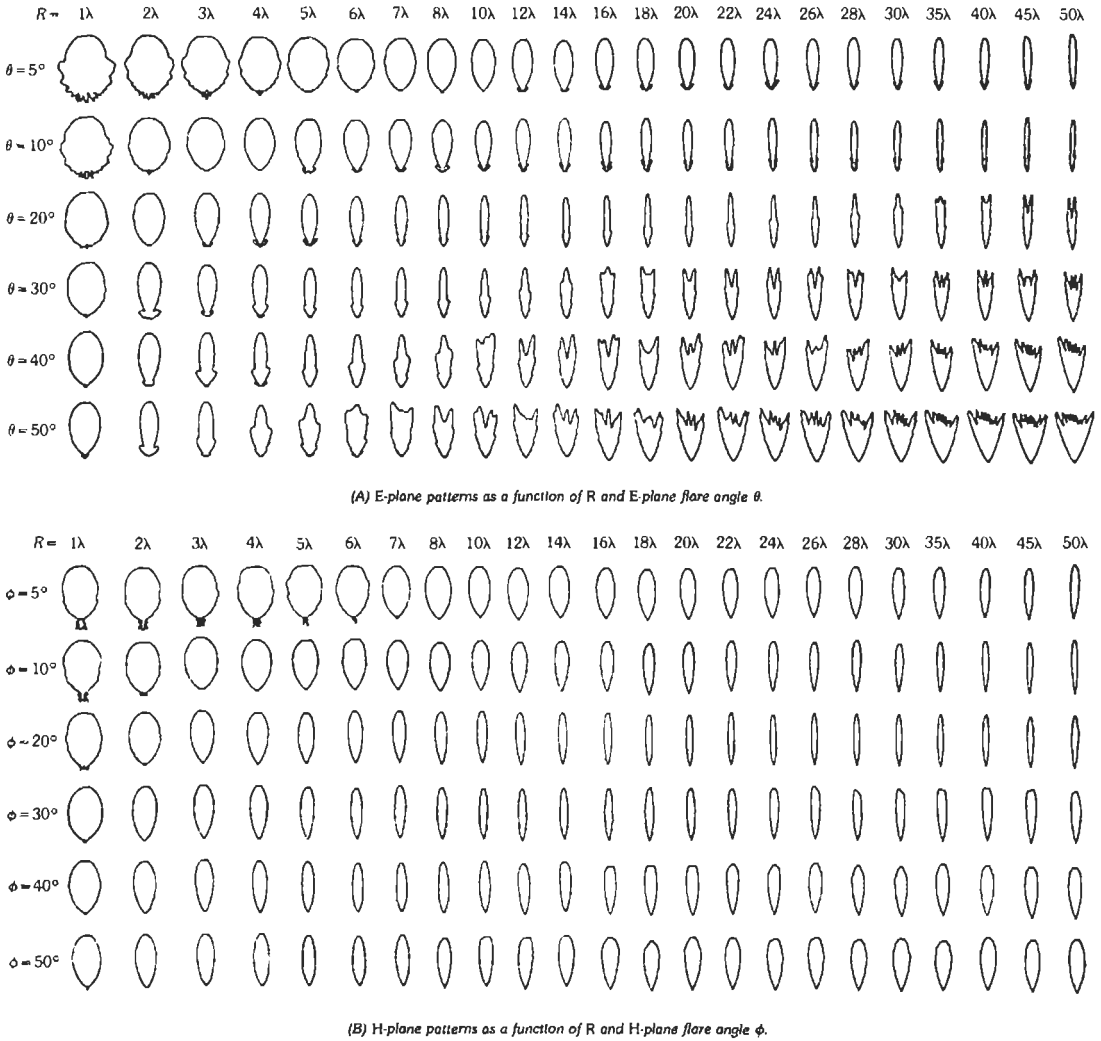


Fig. 43. Patterns of a rectangular horn antenna with radial horn length R . (From Rhodes, D. R. "An Experimental Investigation of the Radiation Patterns of Electromagnetic Horn Antennas." *Proc. IRE*, Vol. 36, Sept. 1948, pp. 1101–1105.)

out ground plane and can be small in size. A principal application is for GPS.*

Yagi-Uda Antennas

The Yagi-Uda antenna is an end-fire array constructed usually of a single driven dipole with a reflector dipole behind and one or more parasitic director dipoles in front. Fig. 45 sketches a three-element Yagi-Uda antenna. Extensive design curves for three- and

four-element Yagi-Uda antennas are given by Uda and Mushiake.† Fig. 46 is extracted from that work. One of the difficulties of designing Yagi-Uda antennas of many elements is the large number of variables. Typically, the spacing between directors is kept constant, and the lengths are adjusted. Table 4 from Stutzman and Thiele‡ shows performance of equally spaced Yagi-Uda antennas with up to seven elements. When the number of elements becomes large, the directors act

* Tranquilia, J. M., and Best, S. R. "A Study of the Quadrifilar Helix Antenna for Global Positioning System (GPS) Applications." *Trans. IEEE*, Vol. AP-38, Oct. 1990, pp. 1545–1550.

† Uda, S., and Mushiake, Y. *Yagi-Uda Antenna*. Research Inst. of Electrical Communication, Tohoku Univ., Sendai, Japan, 1954.

‡ Stutzman, W. L., and Thiele, G. A. *Antenna Theory and Design*. New York: John Wiley & Sons, Inc., 1998.

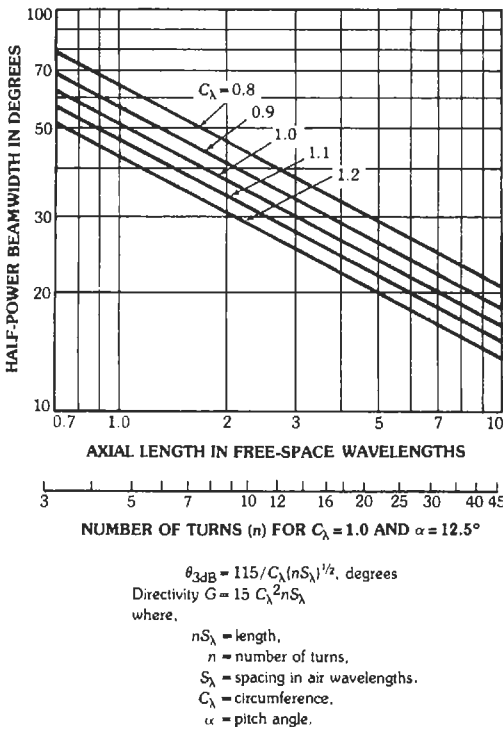


Fig. 44. Axial-mode helical-antenna beamwidth vs length. (From J. D. Kraus, *Antennas*, Fig. 7-21, ©1988, McGraw-Hill Book Co.)

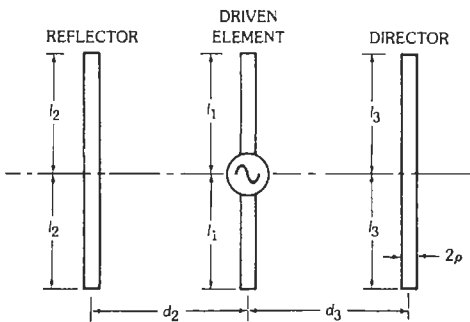


Fig. 45. Three-element Yagi-Uda antenna.

as a slow-wave transmission line.* and the directivity can be approximately calculated from this.† Unfortunately, as the number of elements increases, the bandwidth decreases significantly so that arrays of more than three or four elements are quite narrow-band.

* Serracchioli, F., and Levis, C. A. "The Calculated Phase Velocity of Long End-Fire Uniform Dipole Arrays." *Trans. IEEE*, Vol. AP-7, December 1959, pp. S424-S434.

† Ehrenspeck, H. W., and Pochler, H. "A New Method of Obtaining Maximum Gain from Yagi Antennas." *Trans. IEEE*, Vol. AP-7, October 1959, pp. 379-386.

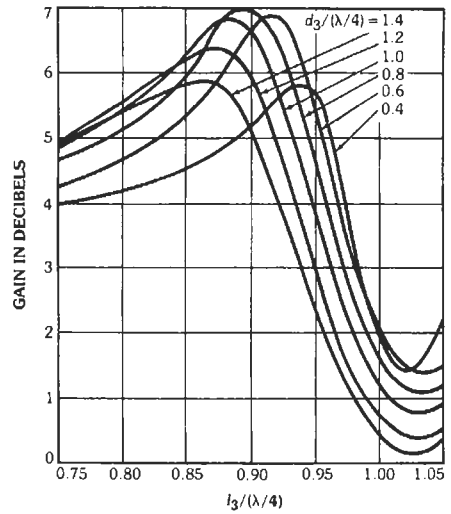


Fig. 46. Calculated gain of three-element Yagi-Uda antenna for indicated values of $d_3/(\lambda/4)$: $l_1 = l_2 = d_2 = \lambda/4$ and $\rho = \lambda/200$. (From Uda, S. and Mushiake, Y. *Yagi-Uda Antenna*, Fig. 9-3, 1954, Sasaki Printing and Publishing Co.)

Even for three- and four-element arrays, if the parameters are adjusted for maximum directivity at one frequency, the bandwidth will be narrow; however, one can adjust the parameters to have slightly less gain over a reasonable bandwidth. A Yagi-Uda antenna can be completely designed by use of array impedance matrix techniques. Here, the set of antenna currents is found from the inverse of the array impedance matrix, with the individual mutual impedance terms calculated by using the zero-order theory of Carter. Extensive calculations have been made by Lawson in a series of articles,‡ using this technique. Fig. 47 shows gain and front-to-back ratio for various three-element Yagi-Uda antennas. Table 5 gives the parasitic lengths and relative resonances for the six curves in each graph. It should be noted that although a single Yagi-Uda can give excellent performance, arraying of Yagi-Udas where the booms are parallel is only partially satisfactory. Surface-wave antennas do not array well because their near fields interact in an undesirable way, and thus the gain realized from an array of Yagi-Uda arrays is always significantly less than the product of the gain of a single Yagi-Uda array times the number of arrays.

A related antenna is the short backfire antenna of Ehrenspeck.§ As shown in Fig. 48, it consists of a flat

‡ Lawson, J. L. "Yagi Antenna Design." *Ham Radio*, Jan., Feb., May, June, July, 1980.

§ Ehrenspeck, H. W. "A New Class of Medium-Size High-Efficiency Reflector." *Trans. IEEE*, Vol. AP-22, March 1974, pp. 329-332. Large, A. C. "Short Backfire Antennas with Waveguide and Linear Fields." *Microwave J.*, Vol. 19, August 1976, pp. 49-52. Kumar, A. "Backfire Antennas Aim at Direct Broadcast TV." *Microwaves Magazine*, April 1978, pp. 106-112.

TABLE 4. CHARACTERISTICS OF EQUALLY SPACED YAGI-UDA ANTENNAS*

N , No. of Elements	Spacing (wavelengths)	Element Lengths			Gain (dB)	Front- to-Back Ratio(dB)	Input Impedance (ohms)	H -Plane		E -Plane	
		Reflector, L_R (wavelengths)	Driver, L (wavelengths)	Directors, L_D (wavelengths)				HP_H (degrees)	SLL_H (dB)	HP_E (degrees)	SLL_E (dB)
3	0.25	0.479	0.453	0.451	9.4	5.6	$22.3 + j15.0$	84	-11.0	66	-34.5
4	0.15	0.486	0.459	0.453	9.7	8.2	$36.7 + j9.6$	84	-11.6	66	-22.8
4	0.20	0.503	0.474	0.463	9.3	7.5	$5.6 + j20.7$	64	-5.2	54	-25.4
4	0.25	0.486	0.463	0.456	10.4	6.0	$10.3 + j23.5$	60	-5.8	52	-15.8
4	0.30	0.475	0.453	0.446	10.7	5.2	$25.8 + j23.2$	64	-7.3	56	-18.5
5	0.15	0.505	0.476	0.456	10.0	13.1	$9.6 + j13.0$	76	-8.9	62	-23.2
5	0.20	0.486	0.462	0.449	11.0	9.4	$18.4 + j17.6$	68	-8.4	58	-18.7
5	0.25	0.477	0.451	0.442	11.0	7.4	$53.3 + j6.2$	66	-8.1	58	-19.1
5	0.30	0.482	0.459	0.451	9.3	2.9	$19.3 + j39.4$	42	-3.3	40	-9.5
6	0.20	0.482	0.456	0.437	11.2	9.2	$51 - j1.9$	68	-9.0	58	-20.0
6	0.25	0.484	0.459	0.446	11.9	9.4	$23.2 + j21.0$	56	-7.1	50	-13.8
6	0.30	0.472	0.449	0.437	11.6	6.7	$61.2 + j7.7$	56	-7.4	52	-14.8
7	0.20	0.489	0.463	0.444	11.8	12.6	$20.6 + j16.8$	58	-7.4	52	-14.1
7	0.25	0.477	0.454	0.434	12.0	8.7	$57.2 + j1.9$	58	-8.1	52	-15.4
7	0.30	0.475	0.455	0.439	12.7	8.7	$35.9 + j21.7$	50	-7.3	46	-12.6

* Conductor diameter = 0.005λ .From Stutzman, W. L., and Thiele, G. A. *Antenna Theory and Design*. New York: John Wiley & Sons, Inc., 1981.

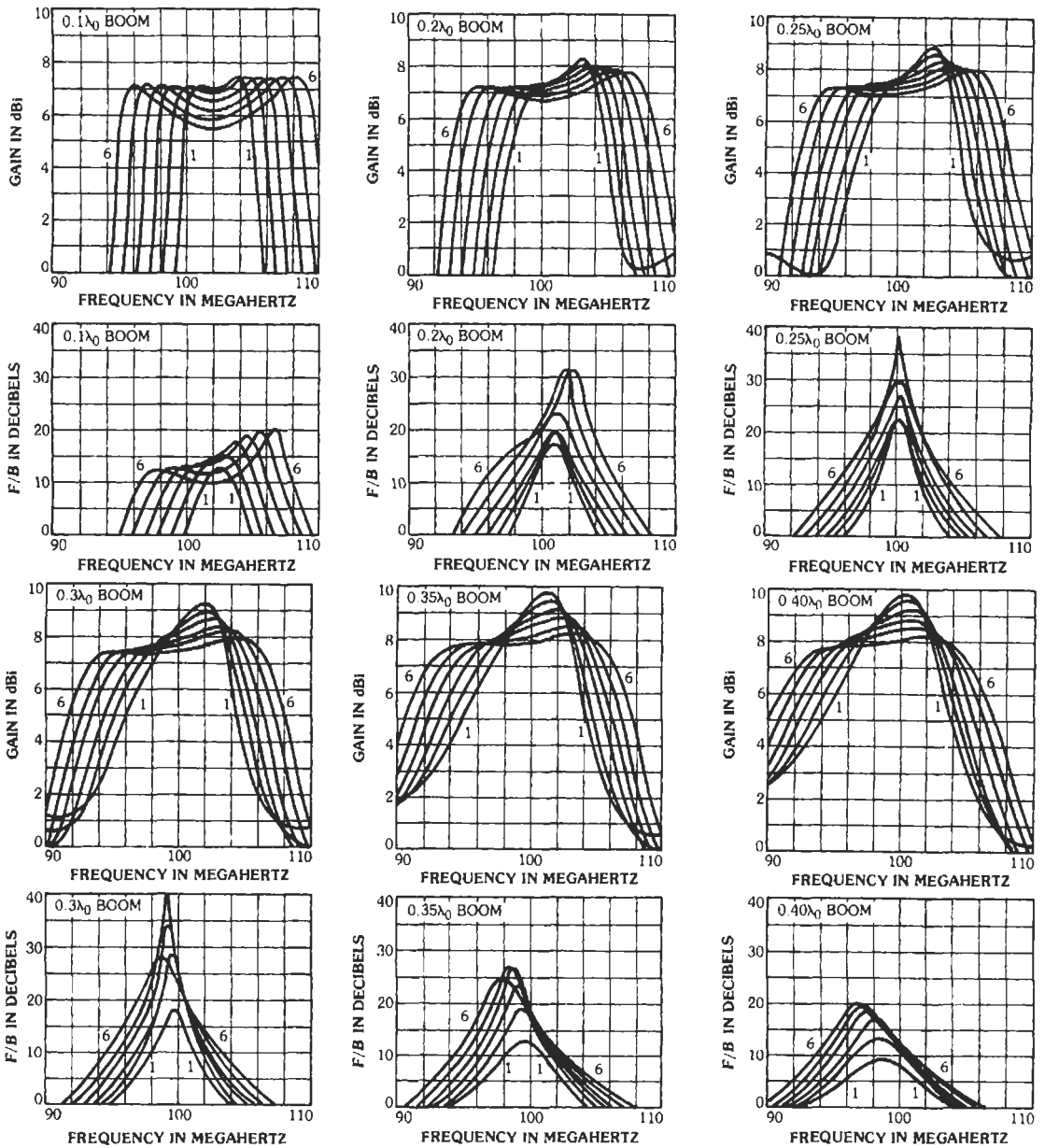


Fig. 47. Gain and front-to-back (F/B) ratio for three-element Yagi beams with varying boom lengths and for changing reflector and director lengths. (From Lawson, J. L. "Yagi Antenna Design." *Ham Radio*, Jan., Feb., May, June, July, 1980.)

plate with a dipole in front of the plate, and with a cylindrical rim around the plate. The dipole feed typically has a dipole and small plate reflector. Performance of the antenna is affected by both the plate diameter and the depth of the rim. A single feed can produce gains in the region of 10 dB. Several of these may be arrayed, all over a flat reflector with a common

rim. However, when more than four feed elements are used, a parabolic reflector may be competitive. Because the antenna has dimensions that are in the resonance region, analyses have only been partially successful. Precise analysis will probably require patch moment methods. Information on bandwidth is sketchy, but the aperture efficiencies tend to be well

TABLE 5. PARASITIC LENGTHS AND RELATIVE RESONANT FREQUENCIES

Curve	Reflector		Director	
	Length/ λ	Resonance	Length/ λ	Resonance
1	0.49150	0.98	0.47223	1.02
2	0.49657	0.97	0.46764	1.03
3	0.50174	0.96	0.46314	1.04
4	0.50702	0.95	0.45873	1.05
5	0.51241	0.94	0.45441	1.06
6	0.51792	0.93	0.45016	1.07

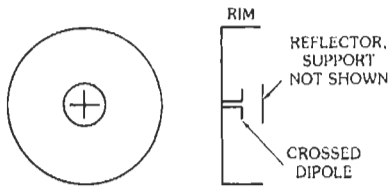


Fig. 48. Backfire antenna.

above those for conventional parabolic-reflector antennas. Higher gain may be realized by feeding the reflecting plate by a feed dipole and dipole reflector, with a number of director dipoles between the feed dipole and the plate. In this scheme, the Yagi-Uda performance is combined with that of the backfire.

Frequency-Independent Antennas

The principle of frequency-independent antennas was established circa 1957 through the recognition that conventional antenna bandwidth limitations occur because critical dimensions change in wavelengths. An antenna specified only in angles should then be frequency independent. Of course, all antennas are of finite size, so an absolute low-frequency cutoff must exist, but within the maximum size the angular prescription can be followed.* A commonly used frequency-independent antenna is the spiral, which exists in two forms. The equiangular spiral is shown in Fig. 49. The conductor edges are formed in polar coordinates by

$$r_1 = r_0 e^{i\phi}, r_2 = cr_1$$

Another version of the spiral is the Archimedean shown in Fig. 50. Although this spiral is not prescribed

* Rumsey, V. H. *Frequency Independent Antennas*. New York: Academic Press, Inc., 1966. Mayes, P. E. "Frequency-Independent Antennas and Broad-Band Derivatives Thereof." *Proc. IEEE*, Vol. 80, Jan. 1992, pp. 103-112.

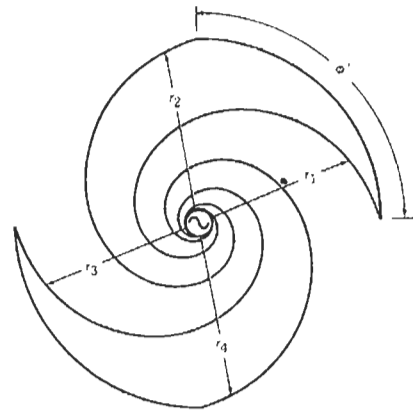


Fig. 49. Two-arm equiangular spiral. (Courtesy E. A. Wolff. *Antenna Analysis*. John Wiley & Sons, 1967)

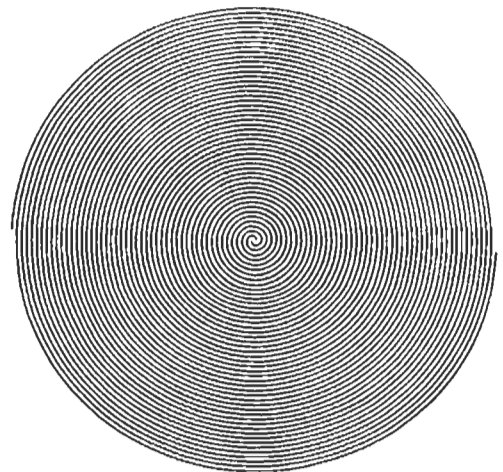


Fig. 50. Archimedean spiral. (Courtesy E. A. Wolff.)

only by angles, it gives excellent performance and is widely used. Both spirals operate through an "active region" from which the radiation takes place. At the

highest frequency, the active region is contiguous with the feed, and as the frequency decreases, the active region moves out to the edge. Corresponding with the active-region movement is a rotation of the pattern. The pattern is circularly polarized corresponding to the direction of the spiral arms, and the pattern exhibits a peak on the axis. Typically, the pattern is broad and the gain slightly more than that of a dipole. Since the circumference of the active region is a wavelength, the spiral must have a diameter of more than λ/π at the lowest frequency. The spiral may be used for azimuth direction finding by utilizing four, six, or eight arms. A four-arm spiral, for example, would have four interleaved arms fed with 90-degree phase rotation by a simple network of hybrids. The higher mode utilized in the four-arm spiral requires a diameter twice as large as that of the two-arm spiral so that the diameter now approaches a wavelength. Spirals with more arms require even larger diameters. The spiral by itself radiates on both sides, and a cavity is often used to allow surface mounting. Unfortunately, the cavity usually significantly reduces the bandwidth, and in order to achieve several octaves, it is necessary to accept some gain degradation via the use of absorbing materials and other techniques. Further information may be found in several references.*

The conical spiral consists of a two-arm spiral deployed on a cone with the feed point at the truncated apex as sketched in Fig. 51. This type of antenna gives

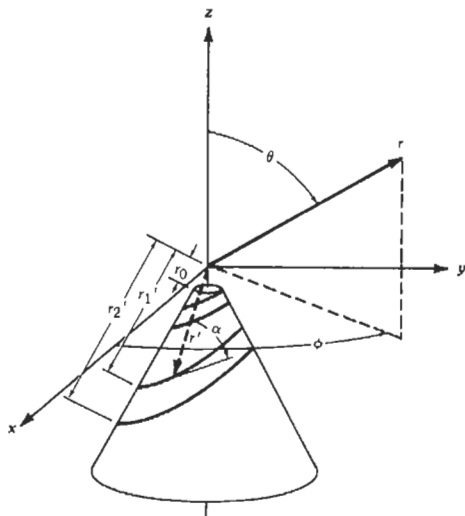


Fig. 51. Conical spiral. (Courtesy E. A. Wolff.)

* Sivan-Sussman, R. "Various Modes of the Equiangular Spiral Antenna." *Trans. IEEE*, Vol. AP-11, Sept. 1963, pp. 533-539. Kaiser, J. A. "The Archimedean Two-Wire Spiral Antenna." *Trans. IEEE*, Vol. AP-8, May 1960, pp. 312-323. Curtis, W. L. "Spiral Antennas." *Trans. IEEE*, Vol. AP-8, May 1960, pp. 298-306. Dyson, J. D. "The Equiangular Spiral Antenna." *Trans. IEEE*, Vol. AP-7, April 1959, pp. 181-187.

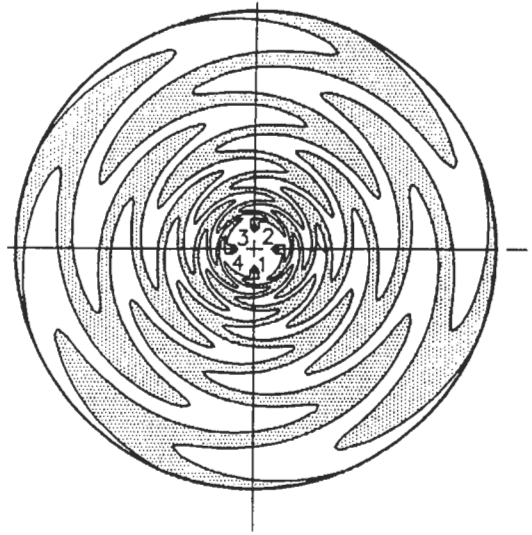


Fig. 51a. Sinuous antenna. (Courtesy R. H. DuHamel.)

a unidirectional pattern with the beamwidth controlled by the cone angle and the size. Further information may be found in a number of references.†

DuHamel‡ achieved dual circular polarization by configuring a spiral in a zig-zag manner: this is the sinuous antenna§ shown in Fig. 51a. Four arms are required for dual CP; six or eight arms allow angle-of-arrival modes to be used.

Another way of achieving a frequency-independent antenna is to design an active region that moves not continuously as in the case of the spiral, but in a discrete fashion. The antenna must now have a performance that is repetitive over a number of suitably related frequency subbands, and this, of course, is the log-periodic antenna discussed next.

† Dyson, J. D. "Measuring the Capacitance per Unit Length of Two Infinite Cones of Arbitrary Cross Section." *Trans. IEEE*, Vol. AP-7, January 1959, pp. 102-103. Dyson, J. D. "Characteristics and Design of the Conical Log-Spiral Antenna." *Trans. IEEE*, Vol. AP-13, July 1965, pp. 488-499. Dyson, J. D. and Mayes, P. E. "New Circularly Polarized Frequency-Independent Antennas With Conical Beam or Omnidirectional Patterns." *Trans. IEEE*, Vol. AP-9, July 1961, pp. 334-342. Tang, C. H. "A Class of Modified Log-Spiral Antennas." *Trans. IEEE*, Vol. AP-11, July 1963, pp. 422-427. Atia, A. E. and Mei, K. K. "Analysis of Multiple-Arm Conical Log-Spiral Antennas." *Trans. IEEE*, Vol. AP-19, May 1971, pp. 320-331.

‡ DuHamel, R. H., and Scherer, J. P. "Frequency Independent Antennas." *Antenna Engineering Handbook*. R. C. Johnson (Ed.). New York: McGraw-Hill, Inc., 1993.

Log-Periodic Antennas

Although the term "log-periodic" can be applied to any antenna designed with a structure that is periodic in the logarithm of some normalized dimensions, almost all log-periodic antennas in use are of the dipole-array type or of the trapezoidal-tooth type. Common usage now calls these simply "log-periodic antennas." The log-periodic array is shown in Fig. 52, where the critical parameters σ and τ are described. It can be seen that the dipoles fit within an envelope of fixed angle and that the ratio of spacing is constant. The log-periodic antenna operates through an "active

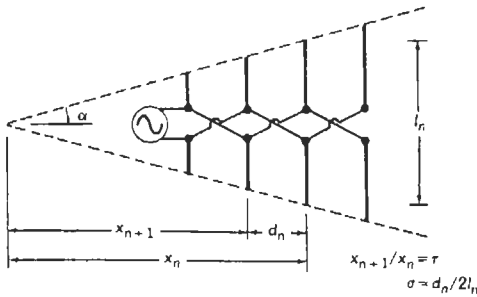


Fig. 52. Log-periodic antenna. (Courtesy W. L. Weeks. *Electromagnetic Theory for Engineering Applications*. John Wiley & Sons, 1964)

region" from which the radiation takes place. This active region consists of those elements that are around $\lambda/2$ in length. Obviously, the closer τ is to unity, the larger the active region and the higher the gain. The low-frequency cutoff occurs when the active region has reached the large end of the antenna. Conversely, the high-frequency cutoff occurs with the active region at the small end. In principle, the log-periodic antenna may be designed to cover many octaves, but, in practice, the upper frequency is limited by the precision required in constructing the small elements, feed lines, and support structure, and by the feeder loss, since the antenna is fed from the small end. Fig. 53 from Carrel* shows directivity versus σ and τ . Fig. 54 shows typical construction using a two-wire feed line. The dipole diameter is not critical, although it is generally desirable to use smaller diameter dipoles at the small end than at the large end when the bandwidth is large. Log-periodic antennas may be printed on a dielectric substrate, but care must be exercised in the design because substrate thickness will limit performance at the upper frequencies. In principle, one can taper the substrate thickness to maintain scaling along the log periodic, but this, of course, is not practical. The length of the

* Carrel, R. L. "The Design of Log-Periodic Dipole Antennas." *1961 IRE Conv. Rec.*, Part 1, pp. 61-75.

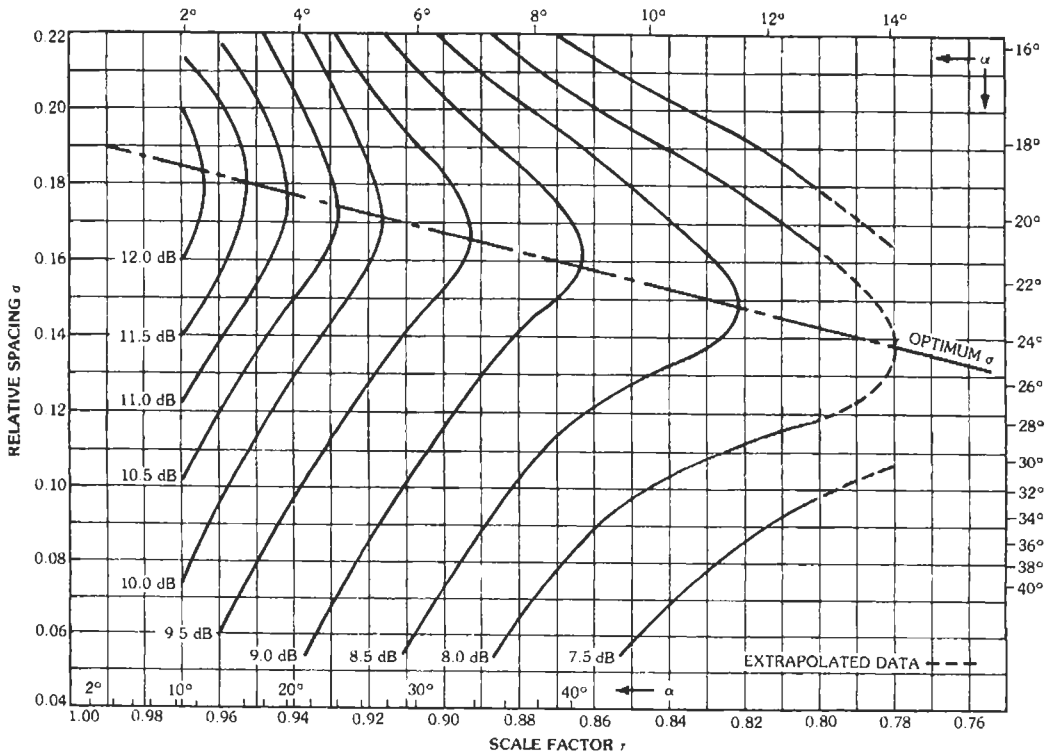


Fig. 53. Constant-directivity contours. (After R. L. Carrel.)

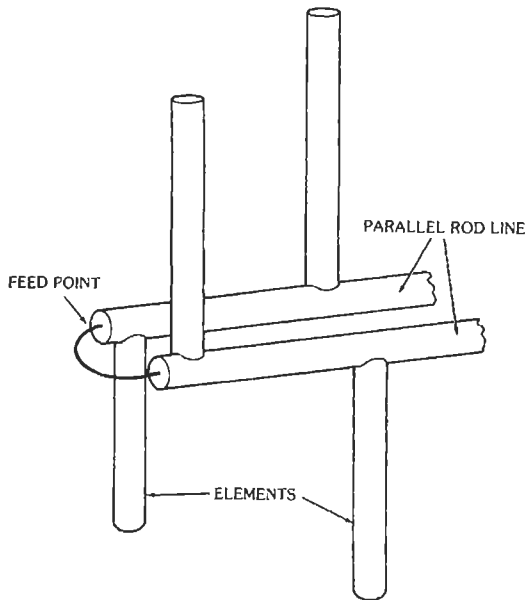


Fig. 54. Log-periodic construction. (Courtesy W. L. Weeks.)

dipoles can be reduced somewhat by loading techniques such as bending with a modest reduction in performance. Reduction of the spacing between elements does not appear practical. Care must be exercised to ensure that resonances are not set up in the feed line; swept impedance and gain measurements are essential for any new log-periodic design. There are several useful log-periodic references.*

The trapezoidal-toothed log-periodic can be in planar form as shown in Fig. 55, but it is more commonly in pyramidal form where the feed is at the apex of a rectangular pyramid and the toothed arms form the sides of the pyramid. Fig. 56 shows a wire version, which, again, is usually used in pyramidal form. In the pyramidal form, the beam is unidirectional. Cross polarized levels are sometimes high.†

* Isbell, D. E. "Log-Periodic Dipole Arrays." *Trans. IEEE*, Vol. AP-8, May 1960, pp. 260-267. Carrel, R. L. "The Design of Log-Periodic Dipole Antennas." *TR 52*, Antenna Lab. University of Illinois, 1961. Smith, C. E. (Ed.). *Log-Periodic Antenna Design Handbook*. Cleveland, Ohio: Smith Electronics, Inc., 1966. DeVito, G., and Stracca, G. B. "Comments on the Design of Log-Periodic Dipole Antennas." *Trans. IEEE*, Vol. AP-21, May 1973, pp. 303-308; see also Butson, P. C., and Thompson, G. T. "A Note on the Calculation of the Gain of Log-Periodic Dipole Antennas." *Trans. IEEE*, Vol. AP-24, January 1976, pp. 105-106. Hilbert, M., et al. "Resonance Phenomena of Log-Periodic Antennas: Characteristic-Mode Analysis." *Trans. IEEE*, Vol. AP-37, October 1989, pp. 1224-1234.

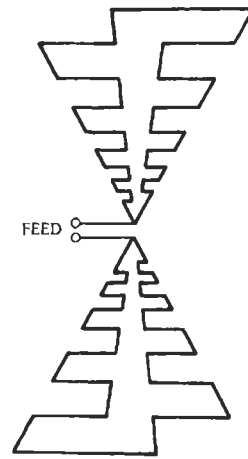


Fig. 55. Planar trapezoidal toothed log-periodic antenna.

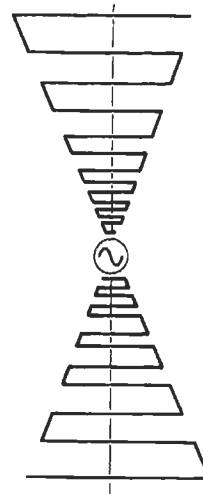


Fig. 56. Log-periodic trapezoid wire antenna (wires are connected to booms).

The log-periodic antenna can be analyzed and designed with the impedance matrix techniques described in the section on arrays, except that the transmission feed line plays an inseparable role in antenna performance. Thus, the transmission line is modeled as series of cascaded short sections with each section loaded by the impedance of the dipole at that

† DuHamel, R. H., and Ore, F. R. "Logarithmically Periodic Antenna Design." *IRE Convention Record*, Part 1, 1958, pp. 139-151. DuHamel, R. H., and Isbell, D. E. "Broadband Logarithmically Periodic Antenna Structures." *IRE Convention Record*, Part 1, 1957, pp. 119-128.

point on the feeder. Of course, the dipole impedance includes all mutual impedances. This results in a transmission-line transformation matrix that is solved in conjunction with the dipole-array impedance matrix. It has been found that although the current distribution on the active-region elements is nearly sinusoidal, elements on either side may have significantly different current distributions. Thus, for accurate results, it is necessary to use either moment method segmentation on the dipoles or to use the three-term current theory of Chang.* The latter is simpler and gives quite satisfactory results.

Fractal Antennas

Fractal theory has been applied to antennas but the results have not been promising. Many parts of nature are well modeled by fractals, e.g. certain ferns; unfortunately Maxwell's equations are not fractal. To see this, consider frequency independent antennas such as the planar spiral and the log-periodic dipole array. High frequency active regions are small and located at (near) the feed; low frequency active regions are large and are located farthest from the feed. Thus, feed currents can activate large active regions with only minor effects in passing through the intervening smaller regions. In contrast, the classic Mandelbrot fractal diagram† has one (largest) blob where the feed would be, with more smaller blobs on the large blob periphery; each smaller blob has its own cluster of still smaller blobs, etc. This topology is just the inverse of the frequency independent antenna topology. Generally, a conventional antenna in the same area (volume) performs better.

ARRAYS

General Characteristics

Arrays of antenna elements are almost always regularly spaced on a rectangular or triangular lattice for planar arrays, equally spaced for linear arrays. Non-uniformly spaced (thinned) arrays will be discussed later. Since fields add, the pattern of an array is the product of the element pattern and the array factor, the latter being that of an array of isotropic elements. Thus, the element and array behavior can be separated for pattern purposes. The array factor for N elements can be written:

$$E(\theta, \phi) = \sum_{n=1}^N A_n e^{-j2n(u-u_0)}$$

* Chang, D. C., Lee, S. W., and Rispin, L. "Simple Formula for Current on a Cylindrical Receiving Antenna." *Trans. IEEE*, Vol. AP-26, September 1978, pp. 683-690.

† Peritgen, H. O., and Saupe, D. *The Science of Fractals*. New York: John Wiley & Sons, 1997.

where

$$u = \pi(d/\lambda) \sin \theta \cos \phi$$

$$u_0 = \pi(d/\lambda) \sin \theta_0 \cos \phi_0$$

The interelement phase shift, u_0 , positions the beam at θ_0, ϕ_0 . Element spacing is d , and the excitation coefficients are A_n . Symmetric arrays can be written as a real series with terms $A_n \cos n(u - u_0)$. Factoring the array expression displays the zeros in u ; these are of critical importance in controlling array-factor behavior. Uniform excitation gives a $\sin \pi u / \pi u$ or $\text{sinc } \pi u$ pattern, and this is both a constituent of a shaped beam pattern (using Woodward-Lawson synthesis‡) and a prototype for low-sidelobe designs. Zeros of $\text{sinc } \pi u$ occur at $u = \pm n$; shifting these zeros controls the sidelobe envelope. Although in the early years of antennas, array and aperture excitation functions were chosen for simplicity and integrability, they are now designed through zero placement to yield optimum performance given the requirements. Computer codes are then employed to furnish needed details. Large arrays, roughly N of 20 or more, are usually designed by sampling a continuous distribution. See the section on aperture distributions. For small arrays, the zeros of a polynomial, such as the Chebyshev, are adjusted; the order of the polynomial matches the number of elements in the array.§ For all types of arrays, design using physical principles of zero placement should be used to obtain optimum performance." This also obviates comparing various well established but meretricious array/aperture distributions.

Examples of two-element arrays for various spacings and phasings have been given by Southworth as shown in Fig. 57.

The 3-dB beamwidth is given by

$$\theta_3 = \arcsin(\sin \theta_0 + 0.443\lambda/Nd) - \arcsin(\sin \theta_0 - 0.443\lambda/Nd)$$

For $N \geq 6$ this is

$$\theta_3 \approx 0.886\lambda / (Nd \cos \theta_0)$$

These formulas show the beamwidth broadening with scan angle θ_0 . With a linear array factor, the pattern is rotationally symmetric about the array axis, with a disk-shaped pattern at broadside. As the beam is moved toward end-fire, the disk becomes a cone, and the inner -3 dB points coalesce along the axis for an angle of:

$$\theta_0 = \arcsin(1 - 0.443\lambda/Nd)$$

‡ Balanis, C. A. *Antenna Theory, Analysis and Design*. New York: John Wiley & Sons, 1997.

§ Villeneuve, A. T. "Taylor Patterns for Discrete Arrays." *Trans. IEEE*, Vol. AP-32, October 1984, pp. 1089-1093.

" Hansen, R. C. *Phased Array Antennas*. New York: John Wiley & Sons, Inc., 1998.

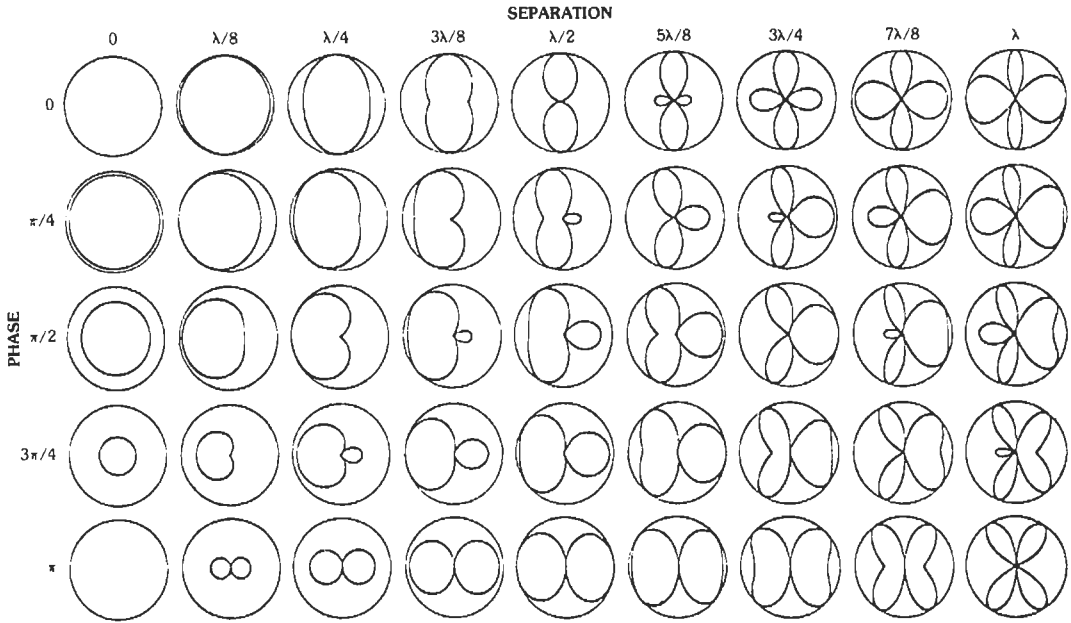


Fig. 57. Patterns of two-element arrays. (From Southworth, G. C. "Certain Factors Affecting the Gain of Directive Antenna Arrays." *Proc. IRE*, Vol. 18, Sept. 1930, pp. 1502-1536.)

For beam angles beyond this angle, a single end-fire beam exists. At end-fire, the beamwidth is larger than broadside beamwidth by the factor:

$$\theta_{3 \text{ endfire}} / \theta_{3 \text{ broadside}} \approx 2 \sqrt{0.886\lambda / (Nd)}$$

When analog phasers are used to produce the inter-element phase shift, u_0 , required to scan the beam, any beam position may be reached if the phasers can be set precisely. Digital phasers, however, have a least count of phase shift, and the array positioning will similarly have a least position change. An m -bit phaser has a least count phase shift of

$$\phi_{lc} = 2\pi / 2^m$$

and this produces a least beam shift of $kd \sin \theta_{lc}$, $k = 2\pi / \lambda$. These are related by

$$kd \sin \theta_{lc} = \pi / 2^{m-1}$$

For the common $\lambda/2$ -spaced array, the relationship is $\sin \theta_{lc} = 2^{m-1}$. Table 6 shows the fineness of beam steering available versus the number of bits per phaser. Pseudo-randomization can reduce the pointing error.*

* Smith, M. S., and Guo, Y. C. "A Comparison of Methods for Randomizing Phase Quantization Errors in Phased Arrays." *Trans. IEEE*, Vol. AP-31, November 1983, pp. 821-827.

TABLE 6. FINENESS OF BEAM STEERING

m	θ_{lc} (degrees)
2	30
3	14.48
4	7.18
5	3.58
6	1.79

Positioning the beam of an array at one frequency requires only a phaser per element. Positioning the beam over a large bandwidth requires the proper time delay, i.e., line length, at each element. The bandwidth allowed with the use of phasers depends on array size; bandwidth is taken as beam movement to the -3-dB point:

$$BW = (f_2 - f_1) / f_0 \approx (f_2 - f_1) / f_2 = (\sin \theta_2 - \sin \theta_1) / \sin \theta_1 \approx \lambda / 2L \sin \theta_1$$

Assumed here is a not small array length. For a beam angle of 30 degrees, for example, the fractional bandwidth is

$$BW \approx \lambda / L$$

Thus, long arrays have less bandwidth in terms of beam shift at the band edges.

Directivity of arrays, unlike fields, cannot be found by combining element and array-factor directivity. The latter for a linear array is given by:

$$G = N \sqrt{1 + (2/N) \sum_{n=1}^{N-1} (N-n) \text{sinc } nkd}$$

Fig. 58 shows this versus spacing. Linear arrays of half-wave dipoles are often used; directivity is obtained from mutual resistances (for an efficient computer algorithm, see Hansen*) with the proper values used for parallel or collinear dipoles:

$$G = 120N \sqrt{R_0 + (2/N) \sum_{n=1}^{N-1} (N-n) R_n}$$

Here R_0 is the self-resistance and R_n is the mutual resistance between the first and n th dipoles. See Figs. 59 and 60. These directivity curves are for a uniform array with a broadside beam.

The directivity drops at a spacing somewhat less than a wavelength due to the emergence of a grating lobe. (See subsection on grating lobes.) Each time the

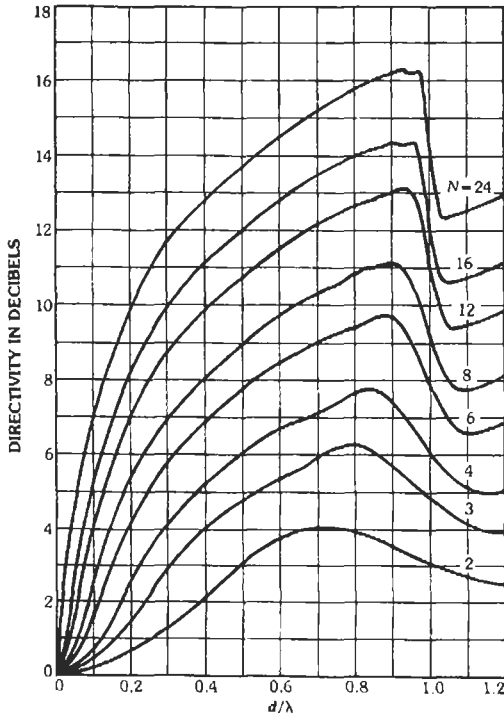


Fig. 58. Array factor directivity.

* Hansen, R. C. "Formulation of Echelon Dipole Mutual Impedance for Computer." *Trans. IEEE*, Vol. AP-20, Nov. 1972, pp. 780-781.

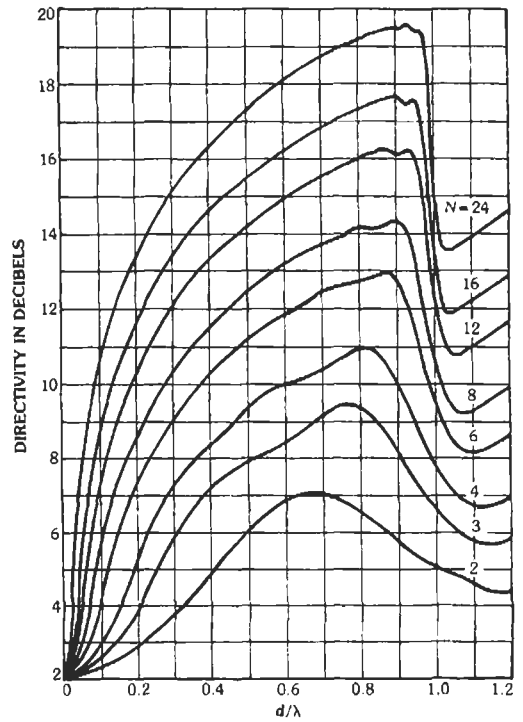


Fig. 59. Directivity of parallel dipole array.

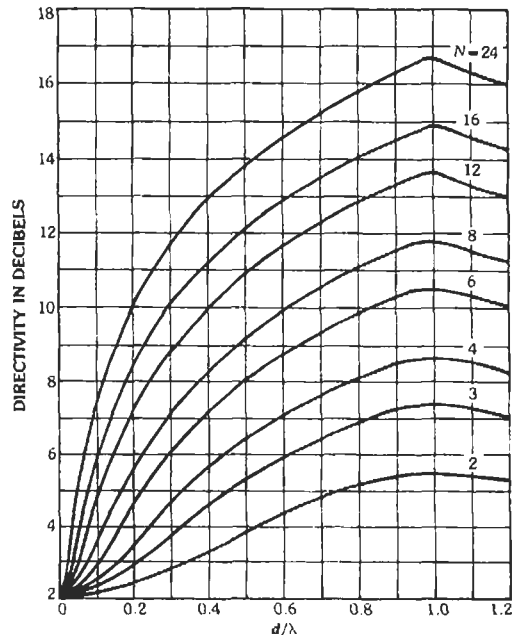


Fig. 60. Directivity of collinear dipole array.

spacing is increased enough to admit another grating lobe, the gain drops proportionately. Directivity of planar arrays can be approximated in lieu of exact calculations, by computing the gain of constituent x - and y -direction linear arrays of the exact elements, multiplying these directivities together, and then adding a correction factor.* Since dipole directivity of 1.604 is included twice, it must be divided out.

Grating and Quantization Lobes

The uniform array equation indicates that a maximum of unity occurs whenever $u = n$. Grating lobes for a small array are shown in Fig. 61. For another array, the pattern is as shown in Fig. 62, which is for a spacing of 0.707λ . For a broadside main beam, a spacing of one wavelength produces a grating lobe at -90° with a symmetric lobe at $+90^\circ$. Similarly, a main beam scanned to 90° will produce a grating lobe at -90° for half-wave spacing. Fig. 63 shows the grating-lobe angles for various spacings and main beam angles. Here $\theta_0 = 0$ represents broadside. These curves allow the designer to select a spacing to minimize grating-lobe effects. However, they are for appearance of the grating lobe at $\pm 90^\circ$ so that half of the grating lobe is

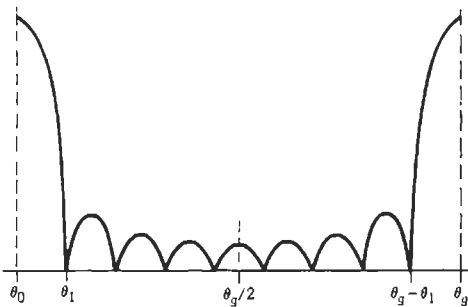


Fig. 61. Array pattern with grating lobe.

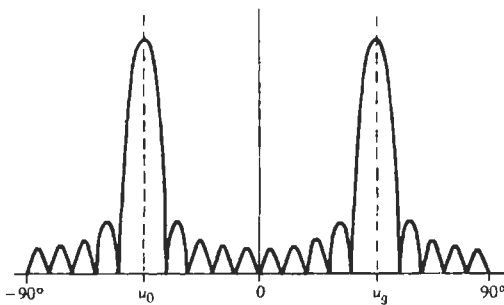


Fig. 62. Array pattern scanned to 45° .

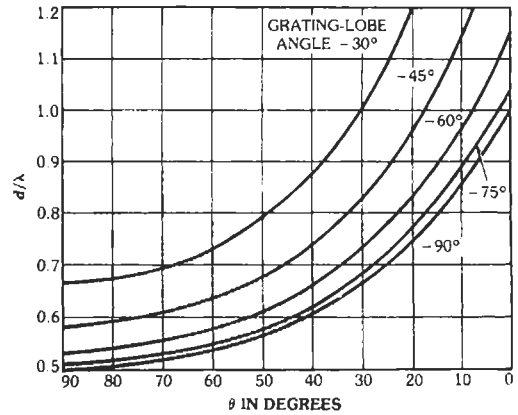


Fig. 63. Element spacing vs grating-lobe angle.

visible. In some cases, it is desirable to reduce the spacing such that the entire grating lobe is precluded, with the grating-lobe null at $\pm 90^\circ$. This spacing reduction has been calculated for the general Taylor one-parameter line source, which includes the uniform as a special case.† The spacing reduction factor is

$$d/d_{gl} = \left(N - \sqrt{1 + B^2} \right) / N$$

B is the Taylor parameter; see the section on aperture distributions. Grating lobes for arrays whose elements are on a rectangular lattice can be analyzed by using two linear-array cases. Arrays using hexagonal or triangular lattices are more complex.‡

Phaser quantization, as mentioned, affects the beam position least count, and it also distorts the pattern. The desired phase versus position along the array is a straight line whose slope is related to the beam angle, but digital phasers produce a stair-step approximation to this. The effects of the sawtooth error depend on the number of elements that fit within a stair step. If there are two or more elements per step (two element sub-array), one or more discrete lobes called quantization lobes will be produced. These have a fixed height that is less than the main beam height, with the height and position both calculable from the element spacing and scan angle. Because of these factors, quantization lobes should not be called grating lobes. Table 7 gives the height of the quantization lobe in terms of number of phaser bits. Where there are fewer than two elements per sawtooth, the errors approach randomness, and the sidelobe level is raised. Fig. 64 shows rms sidelobe level due to phaser quantization for a broad-

* Hansen, R. C. "Comparison of Square Array Directivity Formulas." *Trans. IEEE*, Vol. AP-20, 1972, pp. 100-102.

† Hansen, R. C. Array chapters in *Handbook of Antenna Design*. Vol. 2. A. W. Rudge, et al. (Eds.). London: Peter Peregrinus Ltd., 1983.

‡ Hansen, R. C. *Phased Array Antennas*. New York: John Wiley & Sons, Inc., 1998.

TABLE 7. PHASER QUANTIZATION LOBE HEIGHT

Phaser Bits	QL (dB)
2	-10.5
3	-17.1
4	-23.6
5	-29.8
6	-36.0

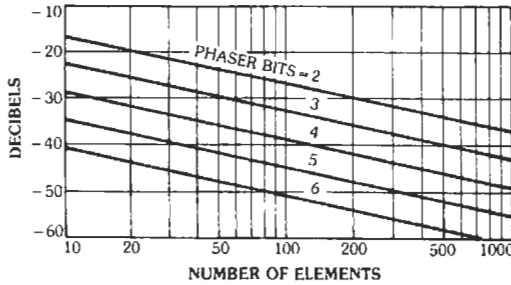


Fig. 64. RMS sidelobe level due to random phaser quantization error, uniform array at broadside.

side main beam. When a quantization lobe exists, it can be suppressed by introducing a pseudorandomization into the phaser control bits.* However, the energy in the quantization lobe is dispersed, thereby raising the sidelobe level, and the random curves of Fig. 64 apply.

Linear Array Feeds

Resonant arrays have elements spaced along a transmission line feed at half-guide wavelength intervals. The transmission line ends are shorted so that a standing wave is set up in the feed. At center frequency, each element is resonant, and the input conductance is the sum of the element conductances. The bandwidth of resonant arrays is narrow, because at frequencies off resonance the conductances are out of phase. The feed admittance can be calculated by using a cascade of A-B-C-D matrices with the results shown in Fig. 65. The abscissa is the number of elements times percent bandwidth. Within the range shown, the VSWR is monotonic, but for larger frequency excursions, VSWR oscillates about a value of two for the uniformly illuminated array. In these results, mutual coupling has not been included; thus the bandwidth indicated is optimistic. A design procedure for resonant arrays including the effects of mutual impedance has been developed by Elliott.† To avoid splitting of the main

* Hansen, R.C., and Charlton, G. C. "Subarray Quantization Lobe Decollimation." *Trans. IEEE*, Vol. AP-47, August 1999, pp. 1237-1239.

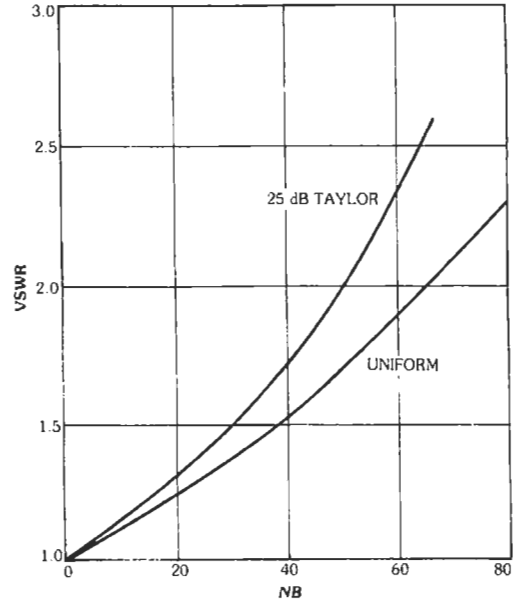


Fig. 65. VSWR of resonant array vs number of elements times percent bandwidth.

beam, arrays should be designed with a VSWR well below 2.0.‡

A traveling-wave array consists of a transmission line with equally spaced elements coupled to it, with spacings that are not any multiple of half-guide wavelength. Power is fed into the feed at one end, and a load is placed at the other end. A portion of the feed energy is coupled to each element, with the power remaining after the last element dissipated in the load. This array is thus nonresonant, and the beam angle will be either forward or backward and will change with frequency. The broadside resonant condition must be avoided. For traveling-wave arrays of many elements, the coupling of each element is small. The element admittances then add nearly randomly, which makes the array reasonably well matched at all points. From transmission line formulations, the conductance of each element may be calculated, including loss in the transmission line. If the loss factor $s = \exp(-2\alpha)$, the conductance for the n th element is given by

$$G_n = F_n s^{-n} \left/ \left[(1 - L s^{1-N})^{-1} \sum_{m=1}^N F_m s^{-m} - \sum_{m=1}^n F_m s^{-m} \right] \right.$$

‡ Elliott, R. S. *Antenna Theory and Design*. Englewood Cliffs, NJ: Prentice-Hall, Inc., 1981.

‡ Kummer, W. H. "Feeding and Phase Scanning." *Micro-wave Scanning Antennas*, Vol. 111. R. C. Hansen (Ed.). New York: Academic Press, Inc., 1966; Los Altos, CA: Peninsula Publ., 1985. Hansen, R. C. *Phased Array Antennas*. New York: John Wiley & Sons, Inc., 1998, Chapter 6.

The F_n excitation coefficients are usually those of the Taylor one-parameter or Taylor \bar{n} distributions described elsewhere in this chapter. A continuous version of the conductance formula is sometimes used, but unless the number of elements is very large, the discrete form above is preferable. Fig. 66 gives slot conductance values for a 29-element array with uniform excitation for several fractions of power dissipated in the load. As expected, the element conductances increase along the feed as there is less power to which to couple. Fig. 67 gives similar curves

for a Taylor one-parameter distribution, and Fig. 68 is for a Taylor \bar{n} distribution.

The beam position of a traveling-wave array is easily determined from the phase equation:

$$kd \sin \theta_0 + 2n\pi = \beta d - \pi$$

where

θ is the beam angle measured from broadside,

d is the element spacing,

$k = 2\pi/\lambda$,

β is the waveguide wave number.

For a given waveguide velocity ratio, a curve such as that in Fig. 69 can be plotted; this curve is important in showing the number of beams versus beam angles.

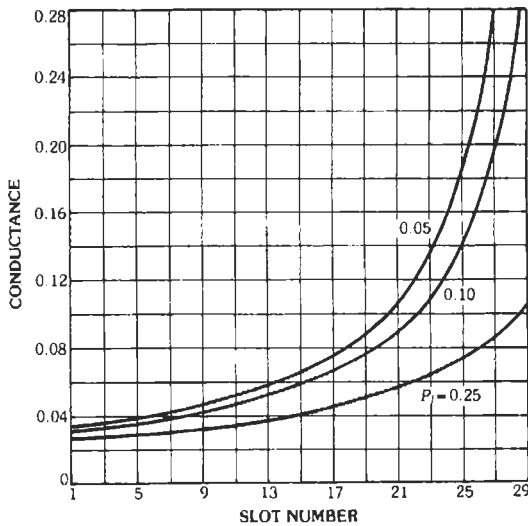


Fig. 66. Slot conductance values, uniform excitation, $N = 29$.

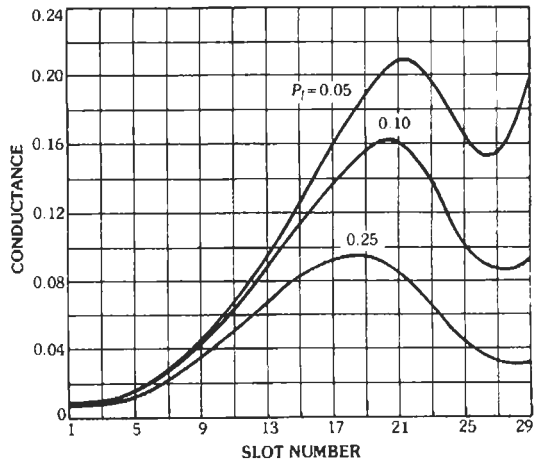


Fig. 68. Slot conductance values, Taylor $\bar{n} = 5$, SLR 25 dB, $N = 29$.

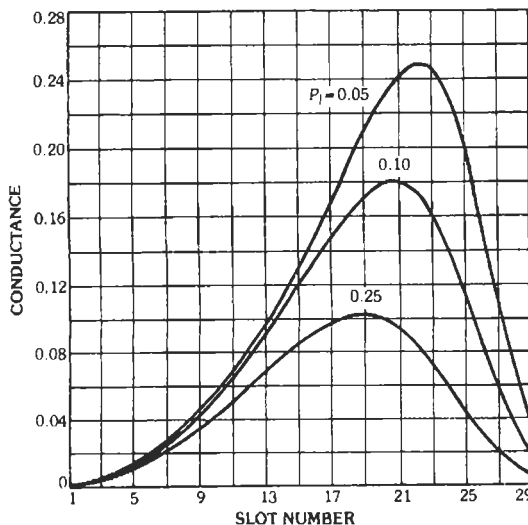


Fig. 67. Slot conductance values, Taylor one-parameter, SLR 25 dB, $N = 29$.

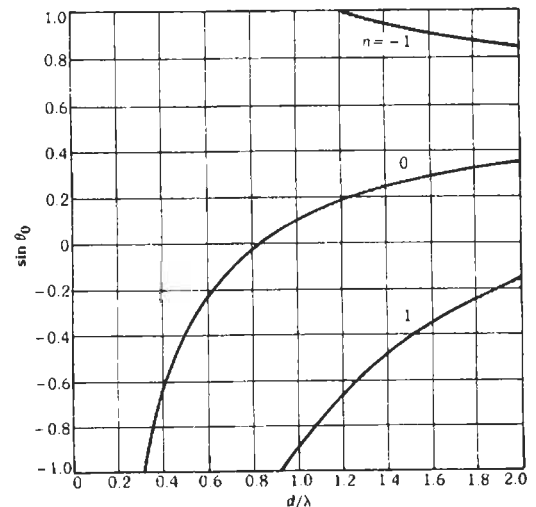


Fig. 69. TW array beam angle vs spacing, $\beta/k = 0.6$.

For example, for spacings smaller than 0.9375λ , only one beam exists, and for most spacings it is a backward-directed beam. Multiple beams are, of course, undesirable; this type of plot enables the single-beam range to be quickly discovered. It is necessary to inspect beam behavior over the entire frequency band of interest, which can be done with other types of plots.* Change of beam angle with frequency is important and can readily be written in a normalized form:

$$f[(d/df)(\sin \theta_0)] = 1/(\beta/k) - \sin \theta_0$$

It can be seen that the normalized slope depends simply on β/k and the beam angle.

Corporate feeds can be used with a linear array; power dividers or hybrid junctions can usually provide significantly greater bandwidth than the array elements. Thus, for wideband arrays corporate feeds are recommended. However, corporate feeds are generally bulky and expensive except for implementation in stripline or microstrip. Thus, printed-circuit arrays are often built with corporate feed networks.

Mutual Impedance

Mutual impedance between elements in any array affects both fixed beam and scanning arrays. In the former, the mutual impedance effects can be accommodated in the design of the array.† Arrays of medium gain elements such as horns, spirals, or dishes may exhibit only minor mutual coupling effects. However, scanning with such large virtual subarrays is severely limited due to quantization lobes (see the subsection on electronic scanning in this chapter). Mutual impedance between low-gain elements is typified by that of halfwave dipoles. Figs. 70, 71, and 72 show complex-plane spiral plots of mutual impedance between parallel

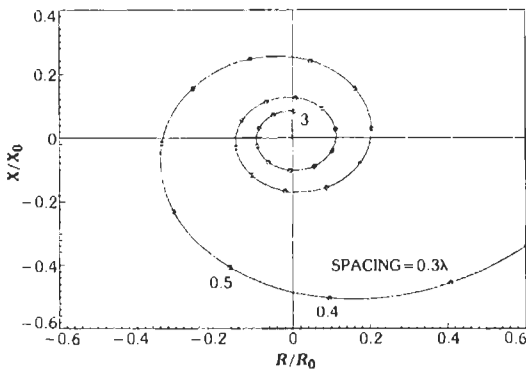


Fig. 70. Parallel dipole mutual impedance/ R_0 .

* Hansen, R. C. Array Chapters in *Handbook of Antenna Design*. Vol. 2. A. W. Rudge et al. (Eds.). London: Peter Peregrinus Ltd., 1983.

† Elliott, R. S. *Antenna Theory and Design*. Englewood Cliffs, NJ: Prentice-Hall, Inc., 1981.

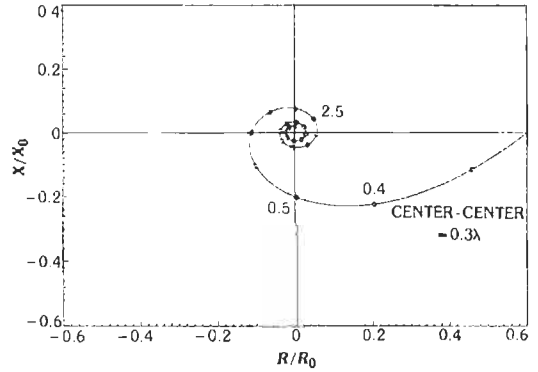


Fig. 71. 45° echelon dipole mutual impedance/ R_0 .

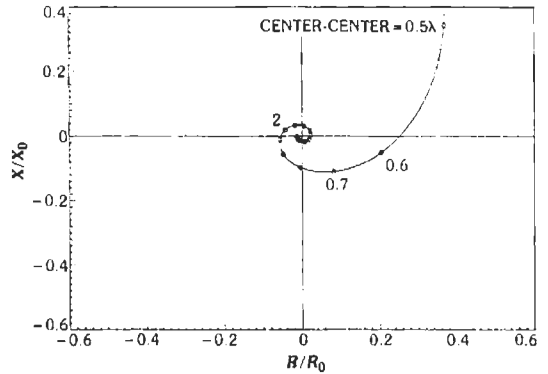


Fig. 72. Collinear dipole mutual impedance/ R_0 .

el, collinear, and 45° echelon dipoles. The circles on the graph represent different values of dipole center-to-center spacing. It is important to note that the H -plane coupling is strong while the E -plane coupling is weak. Accurate and efficient computer codes in Fortran for the calculation of mutual impedance have been given by Hansen‡ for equal lengths and Hansen and Brunner§ for unequal dipoles. These results and codes can be used for slots by replacing the dipole with a Booker equivalent Babinet dipole. The slot admittance and dipole impedance are related by $2Z_{12} = \eta^2 Y_{12}$ where $\eta = 120\pi$. If the slot radiates on both sides, the factor 2 is replaced by 4.

The strong mutual coupling could be reduced if the element power pattern could be transformed into the ideal symmetric $\cos \theta$ pattern. One way of approxi-

‡ Hansen, R. C. "Formulation of Echelon Dipole Mutual Impedance for Computer." *Trans. IEEE*, AP-20, Nov. 1972, pp. 780-781.

§ Hansen, R. C., and Brunner, G. "Dipole Mutual Impedance for Design of Slot Arrays." *Microwave J.*, Vol. 22, December 1979, pp. 54-56.

mating this is to use a round or square open-end waveguide radiator with dielectric plugs, etc., to produce equal TE and TM modes.* A simpler way is to use a slot straddled by a pair of monopoles, sometimes called a Clavin slot. The monopole length and spacing can be adjusted to produce a slot *E*-plane pattern similar to the narrow *H*-plane pattern. Fig. 73 shows a Clavin slot. These slots were empirically designed by Clavin;† an analytical design was given by Papierz et al.‡ Mutual coupling calculations have also been developed by Elliott.§ Figs. 74 and 75 show complex-plane plots of mutual coupling between Clavin slots for parallel and collinear geometries. It can be noted that these couplings are both small and are more alike. Thus, this type of element should give improved scanning performance in electronically scanned arrays.

Mutual coupling effects in electronic scanning are difficult to calculate and outside the scope here. Only a few general remarks will be made. An array with a feed network that provides fixed drive voltages (using dipoles as an example) can be simply analyzed from an impedance matrix equation. That is, the dipole currents are found by multiplying the array voltage vector by the inverse of the mutual impedance matrix. Since the scan angles are contained in the drive vector, it is

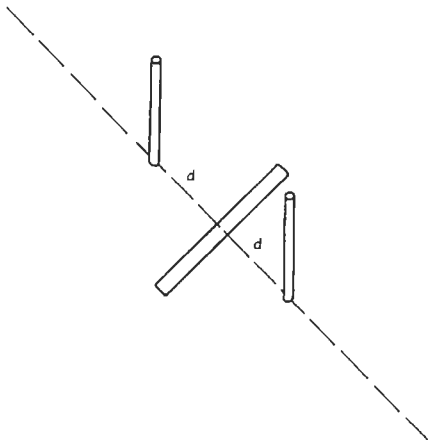


Fig. 73. Slot with parasitic monopoles (Clavin slot).

* Hansen, R. C. *Phased Array Antennas*. New York: John Wiley & Sons, Inc., Chapter 7.

† Clavin, A., Huebner, D. A., and Kilburg, F. J. "An Improved Element for Use in Array Antennas." *Trans. IEEE*, Vol. AP-22, July 1974, pp. 521-526.

‡ Papierz, M. S., et al. "Analysis of Antenna Structure with Equal E- and H-Plane Patterns." *IEE Proc.*, Vol. 124, Jan. 1977, pp. 25-30.

§ Elliott, R. S. "On the Mutual Admittance between Clavin Elements." *Trans. IEEE*, Vol. AP-28, Nov. 1980, pp. 864-870.

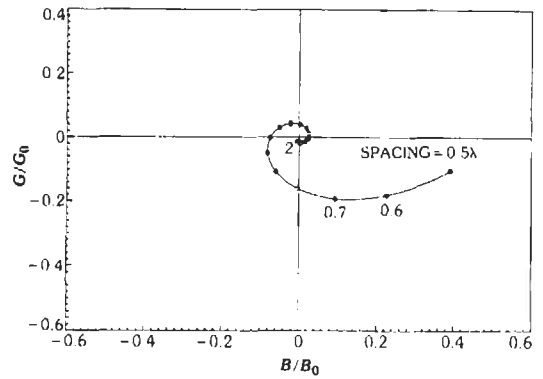


Fig. 74. Parallel Clavin-slot mutual admittance/ G_0 .

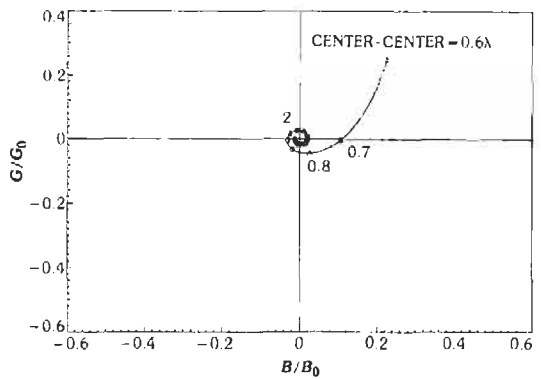


Fig. 75. Collinear Clavin-slot mutual admittance/ G_0 .

simple to find the Scan Impedance of each element and the resulting array pattern for a given scan angle. Many arrays do not have a constant voltage feed network but are of the constant-available-power type. Here an element impedance mismatch reduces the applied voltage. This approach is more suited to a scattering equation. For details, see Hansen." Lengthy formulas have been developed for Scan Impedance of dipole arrays,# and these can be translated to slot arrays. The Scan Impedance is the element input impedance as modified by mutual coupling from neighboring elements, and this Scan Impedance, of course, varies with scan angle. The array gain versus scan can be calculated by using the isolated element pattern times the nominal array pattern factor times one minus the scan power reflection coefficient. The scan reflection coefficient is, of course, determined

" Hansen, R. C. Array chapters in *Handbook of Antenna Design*. Vol. 2. A. W. Rudge et al. (Eds.). London: Peter Peregrinus Ltd., 1983.

Ibid.

from the Scan Impedance. The array gain can also be calculated from a Scan Element Pattern times the nominal array pattern factor. Here the Scan Element Pattern is the element pattern as affected by mutual coupling. All of these comments are useful only for large arrays, i.e., arrays with many elements in the coordinate under consideration. Small arrays must, in essence, be analyzed by adding the contribution of each element, including mutual coupling or mutual scattering. See the subsection on electronic scanning.

Edge effects in finite arrays have been elucidated by an exact analysis of finite-by-infinite dipole arrays. Oscillations in Scan Impedance occur over the finite dimension of the array, even for larger arrays. These are modeled by a Gibbs type phenomenon, where the oscillation amplitudes increase toward the edges.* The period of the oscillation is scan related and is proportional to $1/(1 - \sin \theta_0)$ where θ_0 is the polar scan angle.

Thinned Arrays

Thinned arrays were originally conceived to offer the resolution of a large array without the large number of elements required for a filled array. However, the sidelobe levels have proved to be unattractively high except when the number of elements used is large. A probabilistic analysis was given by Agrawal and Lo† with results that compare closely with Monte Carlo simulations for both short and long arrays and for low and high probabilities. The results are in terms of array length L/λ and a parameter, α , which is number of elements times power sidelobe level: $\alpha = N \times \text{SLL}$. The probability is given by

$$P = (1 - e^{-\alpha}) e^{-(2L/\lambda) \sqrt{\pi\alpha/3} \exp(-\alpha)}$$

Fig. 76 shows the probability of achieving a sidelobe level versus parameter α for different length arrays. For example, a 300-wavelength aperture with 150 elements, giving an average element spacing of 2 wavelengths, has only a 40% probability of achieving a -13 dB sidelobe level. Table 8 shows the α values for several different lengths of arrays for a 90% probability. Also shown is the limit of sidelobe ratio that can be achieved when the number of elements is sufficient to fill the array. Of course, if this number of elements were rearranged with constant half-wave spacing, the sidelobe level would be -13.26 dB. The filling ratio necessary to achieve a 90% probability of sidelobes below -10 dB is also shown. From this, it can be seen that either high filling ratios or large arrays are needed to achieve even modest sidelobe levels. A thinned array designed for one frequency and no scan, frequently provides unacceptably high sidelobes at other

* Hansen, R. C. *Phased Array Antennas*. John Wiley & Sons, Inc., 1998, Section 8.3.

† Agrawal, V. D., and Lo, Y. T. "Mutual Coupling in Phased Arrays of Randomly Spaced Antennas." *Trans. IEEE*, Vol. AP-20, May 1972, pp. 288-295.

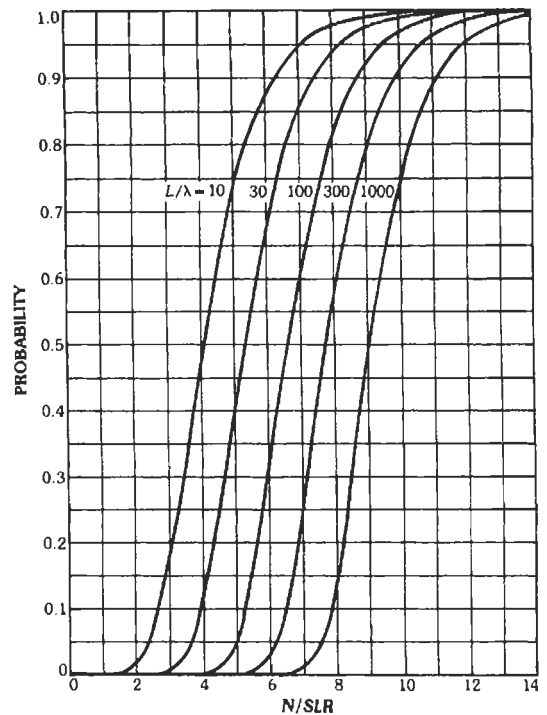


Fig. 76. Probability of random array sidelobe ratio greater than SLR.

TABLE 8. 90% SIDELOBE PROBABILITY

L/λ	$\alpha = N \cdot \text{SLL}$	Maximum SLR (dB)	Filling Ratio for 10 dB SLL
10	6.201	5.1	—
30	7.373	9.1	—
100	8.652	13.6	43%
300	9.813	17.9	16%

frequencies, or at other scan angles. A special type of thinned array is the space-tapered array. Here, the spacing of uniformly excited elements is tapered to match a low sidelobe distribution of the type covered in the section on aperture distribution. This scheme is attractive for arrays that have distributed transmitter/receiver modules because it allows all modules and elements to be alike. Experience has shown that the fall-off of sidelobes is not quite as optimistic as the theory predicts. With any space-tapered array, thorough calculations are needed to ensure that no high sidelobes exist, and it may be necessary to adjust the filling ratio and spacing to ensure sidelobe control over the bandwidth and range of scan angles. For

information on degree of thinning versus sidelobe level, see Hansen.*

A new thinned array design† has roughly constant sidelobe envelopes over a wideband and over scan; typically the peak sidelobe envelope is -15dB and the RMS envelope is -19dB, both with respect to the main beam peak. Gain is approximately equal to the number of elements.

Tolerances

Random errors in element position, element orientation, and element excitation will affect array performance. Since mechanical tolerances can be controlled much more tightly than excitation, only phase and amplitude excitation errors are of concern. They are assumed to have normal (Gaussian) distribution with zero mean and variance σ^2 . These errors reduce the gain a small amount, and on the average add a small constant amount to the sidelobes.‡ The effect on sidelobes is less for large arrays and more for low sidelobes. The gain, with respect to error-free gain, is given by

$$G/G_0 \approx \exp(-\sigma_{amp1}^2 - \sigma_{ph}^2) \approx 1/(1 + \sigma_{amp1}^2 + \sigma_{ph}^2)$$

Thus, errors reduce gain or directivity only to second order, and gain reduction can generally be neglected. For sidelobes, the effect, however, is first order. The mean sidelobe level is related to the error-free sidelobe level by

$$SL = SL_0 \sqrt{1 + \sigma^2 / G \cdot SL_0^2}$$

where

- G is the antenna gain,
- σ^2 is the total error variance,
- SL₀ is the voltage sidelobe level.

Fig. 77 shows mean sidelobe level versus the universal factor σ / \sqrt{G} . From this, it is apparent that sidelobe designs from 13 to 30 dB are quite robust, whereas low sidelobe designs have sensitive tolerances. Similarly, these curves clearly show that for a given sidelobe level and sidelobe degradation, larger arrays allow larger errors. The probability density function of the pattern with random errors is a modified Rayleigh, and the probability can be written in terms of the Marcum Q function. Curves of probability of certain sidelobes exceeding specified values are beyond the scope here, and are given in Hansen.§

* Hansen, R. C. *Phased Array Antennas*. New York: John Wiley & Sons, Inc., 1998.

† Loepr, D. G. "Isosporic Arrays - Massively Thinned Phased Arrays with Well-Controlled Sidelobes." *Trans. IEEE*, Vol. AP-47, December 1999, pp. 1825-1835.

‡ Ruze, J. "Antenna Tolerance Theory - A Review." *Proc. IEEE*. Vol. 54, April 1966; pp. 633-640.

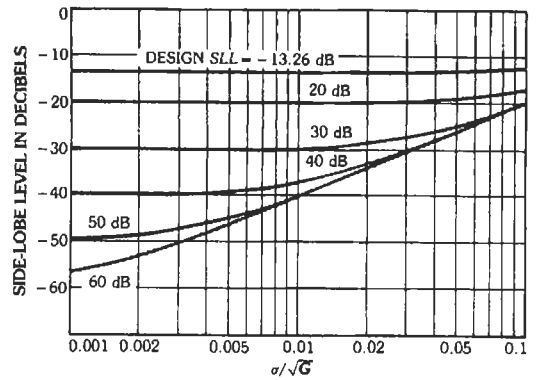


Fig. 77. Mean sidelobe level vs universal factor.

Multiple-Beam Arrays

Multiple-beam arrays usually utilize a beam-forming network (BFN) such as the Butler matrix or Rotman-Gent lens. The BFN, in essence, takes the discrete Fourier transform of the array distribution and produces the transform (space factor) at each beam port. The Butler matrix is, as has been pointed out innumerable times, the microwave equivalent of the FFT. An eight-element, eight-beam Butler BFN is shown in Fig. 78. The beams produced by this network are sketched in Fig. 79. The beam position for any spacing is

$$\sin \theta_i = \pm(2i - 1)\lambda/(2Nd), \quad i = 1, 2, \dots N/2$$

The crossover level is independent of spacing; when the element spacing is increased, the beamwidths become narrower and the beams move closer together. Beam coverage from the center of the leftmost beam to the center of the rightmost beam is

$$\theta_{cov} = 2 \arcsin (N - 1)\lambda/(2Nd)$$

When $d = \lambda/2$, the space -90° to +90° is just filled with beams, as indicated in Fig. 79. Larger spacing moves the beams closer together so that the coverage angle of the N beams is less. However, the space will be filled with beams, and thus there is aliasing and directivity loss due to the extra beams, which are essentially grating lobes. The foregoing assumes isotropic or low-gain elements. When moderate-gain elements such as horns are used with larger spacing, the grating lobes are replaced by subarray quantization lobes; i.e., the grating-lobe amplitude is reduced by the element pattern. The sidelobe ratio and crossover level of a Butler BFN approach those of a uniform line source for large N. Table 9 shows these parameters for

§ Hansen, R. C. *Phased Array Antennas*. New York: John Wiley & Sons, Inc., 1998.

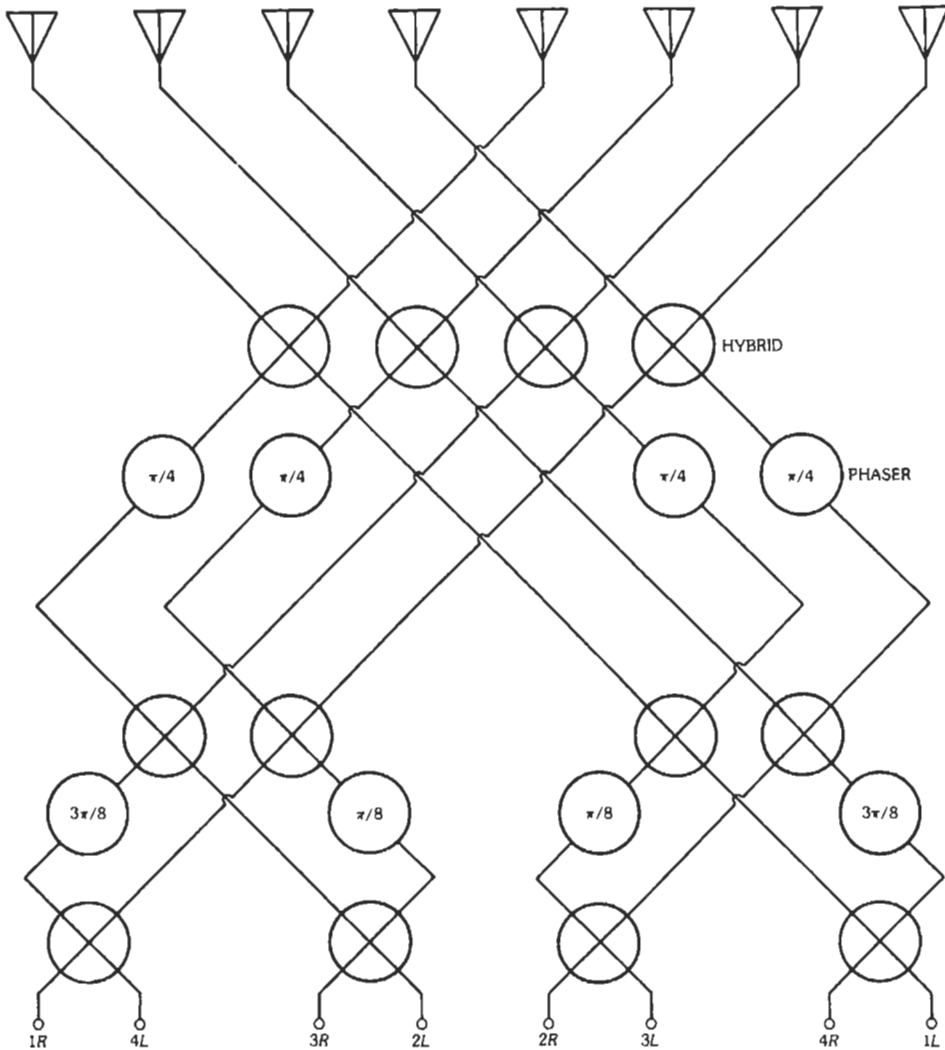


Fig. 78. Butler matrix BFN.

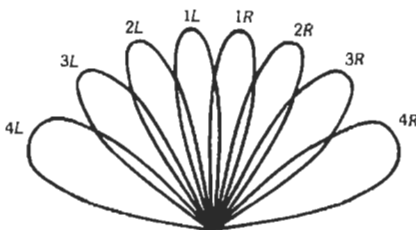


Fig. 79. Butler BFN beams.

popular-size Butlers. The performance of a Butler BFN over a frequency band corresponds to that which is the result of change in spacing. At frequencies above

TABLE 9. BUTLER ARRAY PARAMETERS

N	SLR	Crossover Level
4	11.30 dB	-3.70 dB
8	12.80 dB	-3.87 dB
16	13.15 dB	-3.91 dB
32	13.23 dB	-3.92 dB

the design frequency, the beams become narrow, and additional beams are added at each side. Conversely, at lower frequency, the beams spread out, and one or more beams on each side may disappear.

The Rotman-Gent BFN or lens was originally designed with circular arcs for mechanical scanning. However, they are now usually implemented in microstrip with fixed feeds, and the input-output faces are now allowed to follow the optimum curves. See Fig. 80. When the lens is symmetric, i.e., the input and output curves are the same, the lens can be specified with a single parameter.* Nonsymmetric lenses offer more flexibility and are almost always used.† This type of lens admits a three-point correction, where the center feed and two other symmetrically located feeds can produce a perfectly collimated beam. Note that the bootlace lines connecting the output face to the radiating arrays are an integral part of the lens and are specified by the lens design. The Rotman-Gent lens, being an optical device, has a fixed set of beam positions controlled by the design of the lens. Accordingly, the crossover level at center frequency is adjustable and is a tradeoff with feed mutual coupling and other design factors. At frequencies above or below the design frequency, the beams become narrower or broader, and the crossover level changes accordingly. Rotman-Gent lenses tend to be lossy, and information has not been published on the minimum loss of which this design is capable.

It is known that a multiple-beam antenna designed such that the beams are orthogonal‡ has minimal cross talk between beams. When beams are not orthogonal, the cross talk, which affects both the radiated pattern and the feed reflected power, is produced by the mutual coupling between feed ports and between array ports. When the spacing, d , and excitation are such as

to provide orthogonal beams, mutual coupling still exists, of course, but its effects cancel out in the proper directions. Unfortunately, although many configurations have approximately orthogonal beams, only a uniform line source and line sources with a cosine to an integral power distribution are exactly orthogonal. Contrary to older literature, a cosine-on-a-pedestal does not produce orthogonal beams. A discrete array of N isotropic elements can approach orthogonality if N is large, but for small N the condition is not satisfied. The addition of practical elements such as dipoles or slots destroys the orthogonality condition. In spite of all this, arrays with as few as eight elements show excellent sidelobes for collinear dipoles and modest degradation for parallel dipoles, both at half-wave spacing. Table 10 shows the drawbacks of using purely orthogonal distributions in that only the uniform line source has acceptable efficiency and crossover level.

The effects of nonorthogonality can be reduced by those techniques used to reduce mutual coupling between elements in an array. For example, the use of monopole pins astride a slot reduces E -plane coupling.

Low sidelobes in multiple beams can be produced in several ways. However, the only way that is efficient and simple involves the superposition of beams. For example, the Taylor \bar{n} distribution (see the section on aperture distributions) has been written as a sum over \bar{n} sinc beams. Since a Butler matrix or Rotman-Gent lens produces sinc beams if uniformly excited, several beam ports can be combined with suitable weighting to produce a single beam with low sidelobes. Using this technique, Thomas§ achieved -28-dB sidelobes over a 10% bandwidth with a Rotman-Gent lens. All of the above techniques and devices have a one-dimensional radiating array. A two-dimensional array may be handled by cascading rows of Butler BFNs or of Rotman-Gent BFNs, but this requires much hardware. There is no rotationally symmetrical three-dimensional equivalent of the Rotman-Gent lens; the three-dimensional lens has four perfect points but is not symmetric."

There is, however, a Butler matrix for two-dimensional arrays. This applies to an array disposed on a

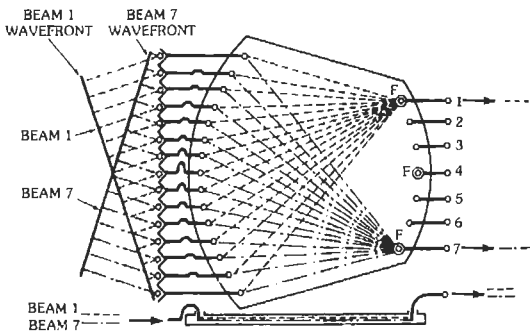


Fig. 80. Rotman BFN

* Shelton, J. P. "Focusing Characteristics of Symmetrically Configured Bootlace Lenses." *Trans. IEEE*, Vol. AP-26, July 1978, pp. 513-518.

† Hansen, R. C. "Design Trades for Rotman Lenses." *Trans. IEEE*, Vol. AP-39, April 1991, pp. 464-472.

$$\ddagger \int_{-1/2}^{1/2} F_i(\pi v) F_j(\pi v) dv = \begin{cases} 1 & \text{for } i = j \\ 0 & \text{for } i \neq j \end{cases}$$

where $v = (d/\lambda) \sin \theta$.

TABLE 10. ORTHOGONAL SPACE FACTORS

Distribution	SLR	Efficiency	Crossover
Uniform	13.26 dB	1	-3.92 dB
cos	23.00 dB	0.8106	-9.54 dB
cos ²	31.47 dB	0.6667	-15.40 dB

§ Thomas, D. T. "Multiple Beam Synthesis of Low Sidelobe Patterns in Lens Fed Arrays." *Trans. IEEE*, Vol. AP-26, November 1978, pp. 883-886.

" Rao, J. B. L. "Multifocal Three-Dimensional Bootlace Lenses." *Trans. IEEE*, Vol. AP-30, November 1982, pp. 1050-1056.

regular hexagonal lattice. The resulting BFN is a three-dimensional microwave network with multiple-arm hybrid junctions but is simpler than cascading two-dimensional Butler BFNs.*

Electronic Scanning

Some of the design tradeoffs for electronic scanning, such as grating lobes, have been addressed above. A large electronic scanning array, or phased array, can be characterized in terms of Scan Impedance and Scan Element Pattern; the modifier "scan" implies a quantity that varies with scan angles. In the Scan Impedance approach, the pattern is written as the product of three factors. First is the array factor (of isotropic elements), which is h_N ; the second is the isolated element pattern, g_i ; the last term involves the Scan Impedance. These combine to give a gain pattern of:

$$G(\theta)/(h_N g_i R_i / R_a) [1 - \Gamma_a]^2$$

where

R_i is the isolated element resistance,

R_a is the scan element resistance,

Γ_a is the scan reflection coefficient.

Scan Impedance can be measured only when all elements are radiating with the proper amplitude and phase. The Scan Element Pattern can be measured with a distant transmitter and all elements terminated in a matched load. It shows relative behavior versus scan angles. If one element is excited, relative angle behavior is useful but the absolute value is not. Although these formulas show only θ , they apply in general to two-dimensional scanning over θ , ϕ . For such large arrays, a powerful analysis was developed by Oliner and colleagues using the unit cell concept. Here, one element in a large array is seen to be closely equivalent to an element in an infinite array; in the latter, the periodic nature of the array allows a single cell to contain the entire characteristics of the array. This cell, called the unit cell, consists of a virtual waveguide with suitably chosen impedance walls with the waveguide symmetrically located about the element (typically a slot) in the element lattice. The wall boundary conditions are derived from the scan angles. These unit cells are normal to the array face and are contiguous. Interior modes are LSE and/or LSM. If no grating lobes exist, the conductance is given by the single propagating mode in the unit cell waveguide, and thus a closed form result occurs. All modes contribute susceptance, and thus many series terms must

be computed. This Floquet method has proved to be powerful and perceptive in the understanding of scanning behavior of arrays. For more detailed information, refer to Oliner and Malech† and Hansen‡. Applications to arrays of microstrip patches have also been made.§

A blind spot can be produced in an array at a particular angle; the array radiation is zero at this angle. This occurs when the dominant mode is canceled by a higher mode, and can occur because of external or internal structure. External structure often used for matching or for protective purposes includes dielectric sheets over or near the array face and dielectric plugs protruding from waveguide-type elements. Internal structure includes dielectric loaded waveguide elements, monopole loaded slots, etc. Blind spots can be severe for large arrays but usually are negligible for small arrays. In case a blind spot is experienced, the array design can usually be changed to remove it. Typical measures are reducing lattice dimensions, changing design to improve VSWR, and altering the periodicity of the external structure. Simulators have also proved useful. For more detailed information, see Oliner and Malech† and Hansen‡.

Small arrays are usually analyzed with the element-by-element approach. Here, a matrix equation is set up relating the drive voltage to the element currents times the impedance matrix, using dipoles as an example. These simultaneous equations are solved for the currents, from which the impedance of each element is immediately available. Note that solution of complex simultaneous equations is several times faster than matrix inversion. Since the scan direction appears in the drive vector, Scan Impedance is easily obtained. This approach allows the variation of impedance from center to edge elements to be determined. When the elements are slots or dipoles, a convenient and efficient computer algorithm for mutual impedance is useful." The effects of Scan Impedance change over a range of scan angles can be reduced by designing the element to have a pattern close to the "ideal" element pattern.

† Oliner, A. A., and Malech, R. G. Chapters 3 and 4 in *Microwave Scanning Antennas*. Vol. 11. R. C. Hansen (Ed.). New York: Academic Press, Inc., 1966; Los Altos, CA: Peninsula Publ., 1985.

‡ Hansen, R. C. Array chapters in *Handbook of Antenna Design*. Vol. 2. A. W. Rudge et al. (Eds.). London: Peter Peregrinus Ltd., 1983; *Phased Array Antennas*. New York: John Wiley & Sons, Inc., 1998.

§ Pozar, D. M. "General Relations for a Phased Array of Printed Antennas Derived from Infinite Current Sheets." *Trans. IEEE*, Vol. AP-33, May 1985, pp. 498-504. Lubin, Y., and Hessel, A. "Wide-Band, Wide-Angle Microstrip Stacked-Patch-Element Phased Arrays." *Trans. IEEE*, Vol. AP-39, Aug. 1991, pp. 1062-1070.

" Hansen, R. C. "Formulation of Echelon Dipole Mutual Impedance for Computer." *Trans. IEEE*, Vol. AP-20, Nov. 1972, pp. 780-781.

* Williams, W. F., and Schroeder, K. G. "Performance Analysis of Planar Hybrid Matrix Arrays." *Trans. IEEE*, Vol. AP-27, July 1969, pp. 526-528. McFarland, J. L. "The RN² Multiple Beam-Array Family and Beam-Forming Matrix." *IEEE APS Symp. Digest*, 1979, pp. 728-731.

This is a power pattern of $\cos \theta$ (for half-wavelength spacing); it is a rotationally symmetric cosine pattern.* Compensation always involves establishing both TE and TM modes (to make the pattern equal in both planes), and perhaps additional modes for pattern shape control. This can be done with open-ended waveguides with appropriately chosen dielectric plugs and slabs.† A simpler but less effective technique uses monopoles astride a slot, the "Clavin pins."‡

Use of sub-arrays, where two or more elements are connected to one phaser, produces quantization lobes (QL) similar to those mentioned under Grating Lobes. These QL occur at the grating lobe angles (for the sub-array spacing) but are reduced by the sub-array patterns. It has recently been discovered that these sub-array QL can be partially decollimated by using a pseudo-random phase.§ Fig. 80a shows the suppressed QL for an array of 2-element sub-arrays and 30-degree scan. This technique is most useful for 2-element (or 2x2) sub-arrays, due to the pattern effects of larger sub-arrays.

Adaptive Arrays

Adaptive arrays are sometimes used to suppress interfering sources; the array is automatically adjusted to provide a pattern null in the direction of the inter-

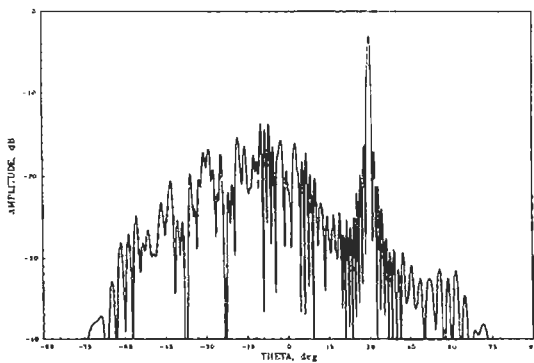


Fig. 80a. Linear array, 100 elements, 50 sub-arrays, $\theta = 30^\circ$.

* Wheeler, H. A. "Simple Relations Derived from a Phased-Array Antenna Made of an Infinite Current Sheet." *Trans. IEEE*, Vol. AP-13, July 1965, pp. 506-514.

† Knittel, G. H. "Wide-Angle Impedance Matching of Phased-Array Antennas: A Survey of Theory and Practice." In Oliner, A. A., and Knittel, G. H. (Eds.). *Phased Array Antennas*. Dedham, MA: Artech House, Inc., 1972.

‡ Elliott, R. S. "On the Mutual Admittance Between Clavin Elements." *Trans. IEEE*, Vol. AP-28, November 1980, pp. 864-870.

§ Hansen, R. C., and Charlton, G. C. "Subarray Quantization Lobe Decollimation." *Trans. IEEE*, Vol. AP-47, August 1999, pp. 1237-1239.

ferer. The simplest type is the sidelobe canceller in which an auxiliary element is used with the main antenna to provide a single null. A feedback loop is used to adjust the phase and amplitude of the single element to produce a subtraction of the interfering signal. A more powerful configuration is the adaptive array in which each element is controlled. There are three general types of adaptive arrays, which differ in the reference signal used for driving the control loops that provide the weights (amplitude and phase) for each element. The first type uses an externally provided steering vector that indicates the desired signal direction. This type of array is called Applebaum-Howells (Fig. 81); it is sensitive to errors." The second type uses a signal replica that correlates with the desired signal. This replica may be from a subcarrier oscillator, a binary sequence of the proper bit rate, etc. This type of adaptive array is associated with Widrow and is frequently called an LMS array. Fig. 82 shows a simplified configuration. The control loops produce a least-mean-square error between the desired signal and the replica, hence the LMS name. The third type of adaptive array does not have an externally supplied reference signal. This array is called a power inversion array and has been associated with Compton. It will form a null on the strongest signal with additional nulls formed on weaker signals.

All of these adaptive arrays have $N - 1$ degrees of freedom, where N is the number of elements in one plane, and $N - 1$ nulls can, in general, be placed in this plane. If there are more degrees of freedom than interferers, the power inversion array tends to null out the desired signal as well, since it cannot tell the desired signal from interfering signals. The improvement available, essentially independent of the type of adaptation, is related to the interference-to-noise ratio. To illustrate, a single interferer is assumed. The ratio of signal to interference plus noise before and after adaptation indicates the improvement (Fig. 83). For large I/N , the improvement in $S/(I + N)$ is essentially just I/N , but for small values the improvement rapidly decreases. Convergence depends on the ratio of the strongest signal to the desired signal; strong interferers require reduced loop gain to avoid instability, and this slows convergence. Methods of accelerating convergence can utilize orthogonalization of Eigenvalues, or a modified control law. The control loops and weighting can be performed either with analog or digital hardware. Thus, there are several configurations possible. Analog control loops are simple, but some control laws such as accelerated convergence are difficult to implement other than digitally. Adaptation in beam space, after a beamformer, allows preservation of desired signals and better rejection control. # Adaptation of wideband

" Compton, R. T. *Adaptive Antennas*. Englewood Cliffs, NJ: Prentice Hall, 1988.

Gabriel, W. F. "Adaptive Processing Array Systems." *Proc. IEEE*, Vol. 80, January 1992, pp. 152-161.

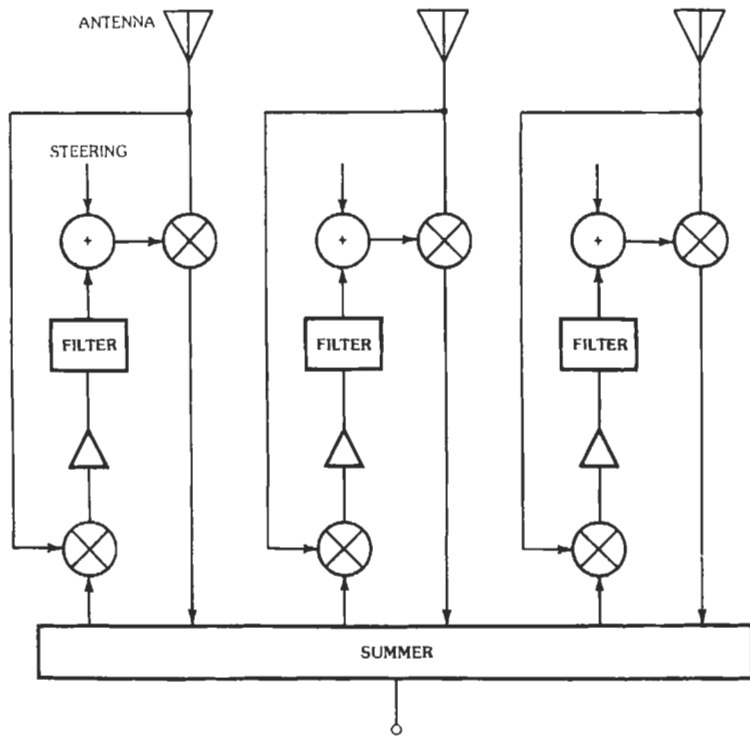


Fig. 81. Applebaum-Howells adaptive array.

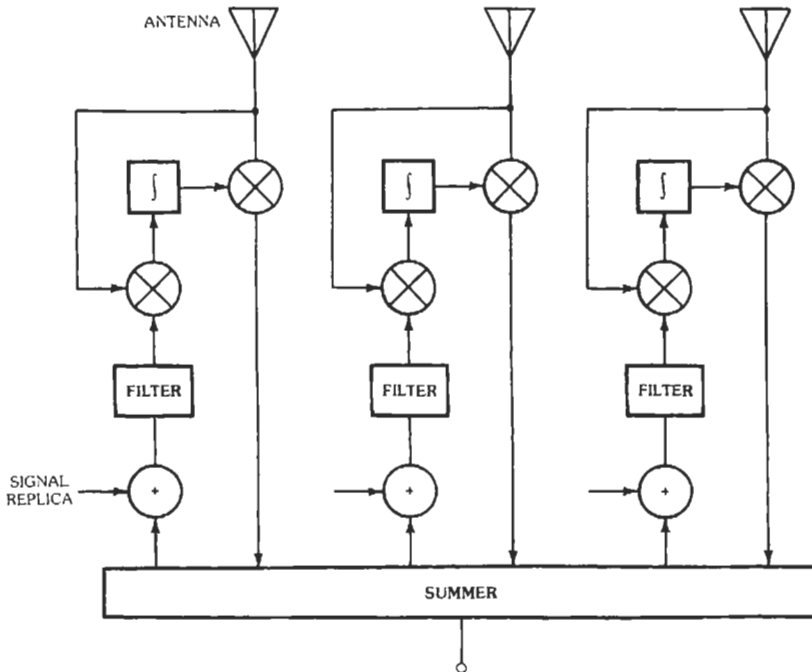


Fig. 82. Widrow LMS adaptive array.

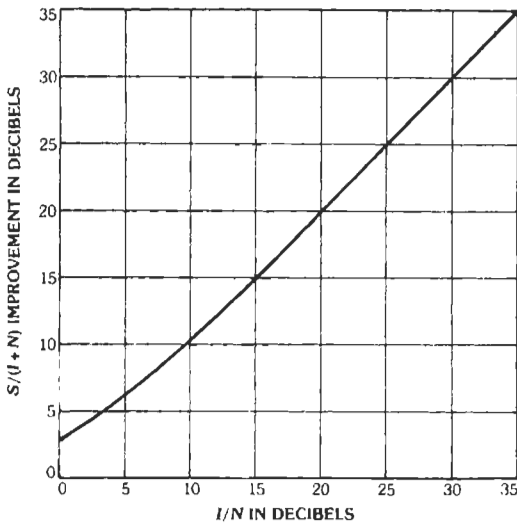


Fig. 83. Improvement in $S/(I + N)$ due to adaption.

signals may require a tapped (and adjustable) delay line at each element. Further information is in several sources.*

APERTURE DISTRIBUTIONS

Design of Distributions

Aperture distributions of concern here are used to produce patterns with one or two narrow main beams with the rest of space occupied by sidelobes. These distributions are often sampled to yield excitation coefficients for arrays of many elements. For arrays of a few elements, adjustment of zeros of a polynomial is used.† In earlier years, aperture distributions were chosen for easy integrability. With the advent of computers, it is possible to design aperture distributions based on physical principles without concern about the form of the result. These physical principles were

* Widrow, B., et al. "Stationary and Nonstationary Learning Characteristics of the LMS Adaptive Filter." *Proc. IEEE*, Vol. 64, August 1976, pp. 1151-1162. Gabriel, W. F. "Adaptive Arrays—An Introduction." *Proc. IEEE*, Vol. 64, February 1976, pp. 239-272. Gabriel, W. F. (Ed.). "Special Issue on Adaptive Antennas." *Trans. IEEE*, Vol. AP-24, September 1976. Monzingo, R. A., and Miller, T. W. *Introduction to Adaptive Arrays*. New York: John Wiley & Sons, Inc., 1980. Hudson, J. E. *Adaptive Array Principles*. London: Peter Peregrinus Ltd., 1981. Widrow, B., and Stearns, S. D. *Adaptive Signal Processing*. Englewood Cliffs, NJ: Prentice Hall, 1985.

† Villeneuve, A. T. "Taylor Patterns for Discrete Arrays." *Trans. IEEE*, Vol. AP-32, October 1984; pp. 1089-1093.

enunciated primarily by Taylor:‡ symmetric amplitude distributions are more efficient; array polynomial zeros should be real; far-out zeros should be separated by unity; and a distribution with a pedestal gives a $1/\lambda$ far-out sidelobe envelope and is more efficient than a distribution that is zero at the ends. In all that follows, the variable $u = (L/\lambda) \sin \theta$ is used. In general, then, a narrow beam distribution is designed by making the spacing of zeros for large u approximately unity, and then by adjusting the close-in zeros to give the desired first sidelobe level. The achievement of a $1/u$ envelope assures a robust, low- Q distribution. It is convenient to use sidelobe ratio (SLR) as the ratio of beam peak to first sidelobe; SLR, of course, is the inverse of sidelobe level.

The distributions that will be discussed here differ in the shape of the close-in sidelobe envelope, and whether they are for line (or rectangular) or circular sources. These distributions are now widely used and have superseded such distributions as cosineⁿ-on-a-pedestal, Gaussian, etc. The latter distributions not only have inefficient zero placement, but they involve two or more parameters. For example, a practical Gaussian must be truncated, and two parameters must be specified to determine performance indices. Incidentally, when the Gaussian is truncated, its advantages of a simple transform disappear. For these distributions, it is difficult to perform optimization: for example, determination of parameters to give maximum efficiency for a given sidelobe level is difficult. In contrast, the Taylor one-parameter distribution yields all performance indices as a function of that single parameter, while the Taylor \bar{n} distribution is easily optimized because of the physical meaning of the parameters. See Hansen.§

Taylor One-Parameter Line Source Distribution

In this distribution, the far-out zeros are placed at integer values of u , while close-in zeros are shifted to reduce close-in sidelobes to suitable values. Taylor accomplished this by using zeros

$$u = \sqrt{n^2 + B^2}$$

where B is the single parameter. The pattern is written in two forms, depending on whether u is less or greater than B . These are

$$F(u) = (\sinh \pi \sqrt{B^2 - u^2}) / \pi \sqrt{B^2 - u^2}, \quad u \leq B$$

$$F(u) = (\sin \pi \sqrt{u^2 - B^2}) / \pi \sqrt{u^2 - B^2}, \quad u \geq B$$

‡ Taylor, T. T. "Design of Line-Source Antennas for Narrow Beamwidth and Low Sidelobes." *IRE Trans.*, AP-3, January, 1955, pp. 16-28.

§ Hansen, R. C. "Array Pattern Control and Synthesis." *Proc. IEEE*, Vol. 80, January 1992, pp. 141-151.

The sidelobe ratio is given by

$$SLR = 20 \log (\sinh \pi B / \pi B) + 13.26 \text{ dB}$$

Taking the inverse transform of the pattern gives the Taylor one-parameter line source distribution:

$$g(p) = I_0 (\pi B \sqrt{1 - p^2})$$

where

p is 0 at the aperture center and 1 at each end,
 I_0 is the modified Bessel function.

A special case occurs for $B = 0$. This is the uniformly excited line source (constant amplitude), which has a pattern of simply $\sin \pi u / \pi u$. The Taylor pattern is a modified $\sin \pi u / \pi u$ pattern, with a transition from that to the hyperbolic form at $u = B$, on the side of the main beam. The hyperbolic form provides the central part of the main beam as indicated in the formula above for SLR. Fig. 84 shows a typical one-parameter pattern or space factor. Typical aperture distributions are shown in Fig. 85. Only half of a symmetric distribution is shown. Not unexpectedly, the pattern with lower sidelobes has a smaller pedestal at the end of the aperture. The integral for aperture efficiency can be reduced to a tabulated function, \bar{I}_0 .*

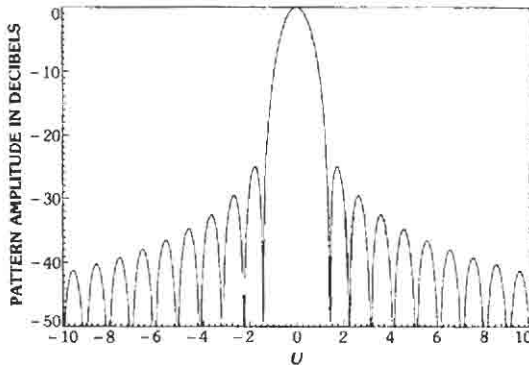


Fig. 84. Taylor one-parameter line source pattern for SLR 25 dB.

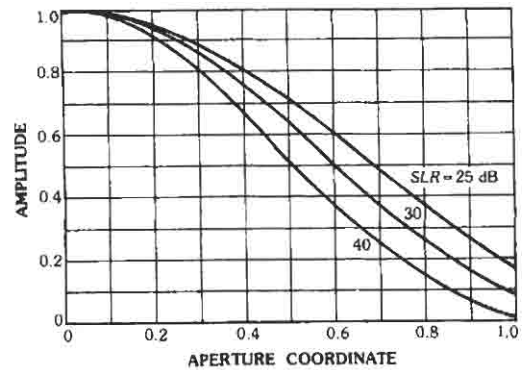


Fig. 85. Taylor one-parameter aperture distributions.

$$\eta = (2 \sinh^2 \pi B) / [\pi B \bar{I}_0 (2 \pi B)]$$

Although all of the performance indices can be given in terms of a table of equally spaced values of B , it is more convenient to specify sidelobe ratio. Table 11 gives beamwidth, efficiency, and beam efficiency. The beamwidth is normalized and is converted to an actual angle by

$$\theta_3 / (\lambda / L) = (2L / \lambda) \arcsin [u_3 / (L / \lambda)] \approx 2u_3$$

The efficiency relates the directivity to that of a uniformly excited source. Beam efficiency is the fraction of radiated energy contained in the main beam, null to null.† Given a desired performance index, the formulas allow B to be found, and from B all other indices can be determined. Taylor described this work in a Hughes Aircraft report; a review and tables are contained in Hansen.‡

* Abramowitz, M., and Stegun, L. *Handbook of Mathematical Functions*. NBS, 1970.

† Hansen, R. C. *Phased Array Antennas*. New York: John Wiley & Sons, Inc., 1998.

‡ Hansen, R. C. *Microwave Scanning Antennas*. Vol. 1. New York: Academic Press, Inc., 1964; Chapter 1; Los Altos, CA: Peninsula Publ., 1985.

TABLE 11. TAYLOR ONE-PARAMETER LINE SOURCE CHARACTERISTICS

SLR (dB)	B	U_3 (rad)	η	η_0
13.26	0	0.4429	1	0.9028
20	0.7386	0.5119	0.9330	0.9820
25	1.0229	0.5580	0.8626	0.9950
30	1.2762	0.6002	0.8014	0.9986
35	1.5136	0.6391	0.7509	0.9996
40	1.7415	0.6752	0.7090	0.9999
45	1.9628	0.7091	0.6740	1.0
50	2.1793	0.7411	0.6451	1.0

Taylor \bar{n} Line Source Distribution

A modest improvement in efficiency can be obtained by making the first few sidelobes at equal level, with a transition from the equal-level envelope to the $1/u$ envelope. This offers a compromise between the one-parameter distribution and the Chebyshev distribution. In the latter, the sidelobes are all of equal height, the aperture distribution is singular at the ends, the energy storage is high, and the distribution is sensitive to errors. For \bar{n} small, the distribution gives improved efficiency and narrower beamwidth without significantly degrading its robust nature. The distribution is a modification of the continuous equivalent of a Chebyshev distribution, with a dilation factor used to modify the first \bar{n} zeros. The pattern is given by a canonical product on zeros.*

$$F(u) = \left[\frac{(\sin \pi u) / \pi u}{\prod_{n=1}^{\bar{n}-1} (1 - u^2 / z_n^2)} \right] / (1 - u^2 / \bar{n}^2)$$

$$= \sum_{n=-(\bar{n}-1)}^{n=\bar{n}-1} F(n, A, \bar{n}) \operatorname{sinc} \pi(u+n)$$

In the second form, the pattern is a superposition of \bar{n} sinc beams. The pattern coefficients are

$$F(n, A, \bar{n}) = \frac{[(\bar{n}-1)!]^2}{(\bar{n}-1+n)!(\bar{n}-1-n)!} \prod_{m=1}^{\bar{n}-1} (1 - n^2 / z_m^2)$$

with $F(0, A, \bar{n}) = 1$. Here the zeros are:

$$Z_n = \pm \sigma \sqrt{A^2 + (n-1/2)^2} \quad 1 \leq n \leq \bar{n}$$

$$Z_n = \pm n \quad \bar{n} \leq n$$

where

A is the sidelobe parameter,
 σ is the dilation index.

* Taylor, T. T. "Design of Line-Source Antennas for Narrow Beamwidth and Low Sidelobes." *IRE Trans.*, Vol. AP-3, Jan. 1955, pp. 16-28.

The dilation index is

$$\sigma = \bar{n} / \sqrt{A^2 + (\bar{n}-1/2)^2}$$

The corresponding aperture distribution is

$$g(p) = 1 + 2 \sum_{n=1}^{\bar{n}-1} F(n, A, \bar{n}) \cos n \pi p$$

Tables of aperture distribution are not given because of the ease of calculating these with a computer. Values for check purposes can be obtained from Hansen.† Table 12 gives parameter A and beamwidth factor u_3 for several values of SLR. Also shown are σ values for \bar{n} . The beamwidth in u is approximately $2\pi u_3$; this may be converted to angle by using the formula in the previous subsection. Fig. 86 shows a typical Taylor \bar{n} pattern; Fig. 87 shows typical aperture distributions. The value of \bar{n} must be carefully chosen. If \bar{n} is too small, the zero spacing in the transition region will oscillate, thereby producing oscillations in the sidelobe envelope. If \bar{n} is large, the aperture distribution will not be monotonic; for suitably large \bar{n} the end value

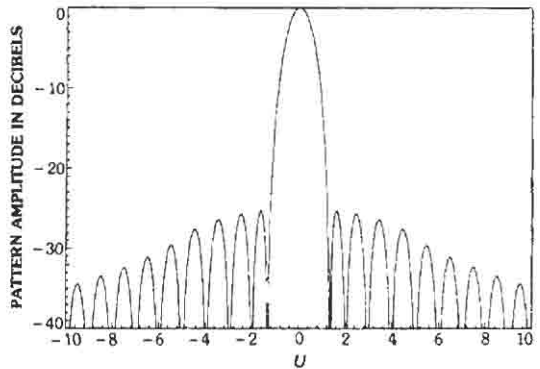


Fig. 86. Taylor line source pattern for SLR = 25 dB, $\bar{n} = 5$.

† Hansen, R. C. *Microwave Scanning Antennas*. Vol. 1. New York: Academic Press, Inc., 1964; Chapter 1; Los Altos, CA: Peninsula Pub., 1985.

TABLE 12. TAYLOR \bar{n} LINE SOURCE CHARACTERISTICS

SLR (dB)	A	u_3	σ				
			$\bar{n} = 2$	4	6	8	10
20	0.9528	0.4465	1.1255	1.1027	1.0749	1.0582	1.0474
25	1.1366	0.4890		1.0870	1.0683	1.0546	1.0452
30	1.3200	0.5284		1.0693	1.0608	1.0505	1.0426
35	1.5032	0.5653			1.0523	1.0459	1.0397
40	1.6865	0.6000			1.0430	1.0407	1.0364
45	1.8697	0.6328				1.0350	1.0328
50	2.0530	0.6639					1.0289

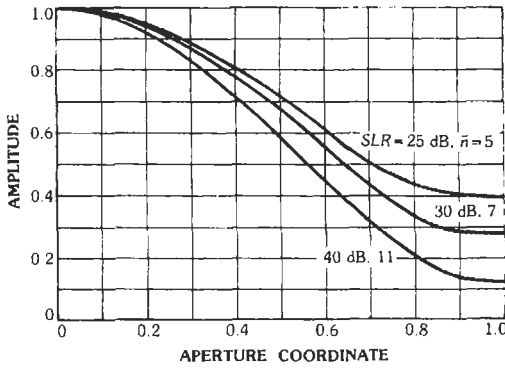


Fig. 87. Taylor \bar{n} aperture distributions.

may be larger than the center value, which is clearly disastrous. There is a value of \bar{n} that gives maximum efficiency. Table 13 gives these values and the largest \bar{n} that maintains a monotonic distribution. Note that the maximum efficiency values require an aperture distribution with peaks at the ends and so are generally not desirable.

Aperture efficiency is easily calculated from

$$\eta = 1 / \left[1 + 2 \sum_{n=1}^{\bar{n}-1} F^2(n, A, \bar{n}) \right]$$

The difference in efficiency between the maximum \bar{n} case and the monotonic \bar{n} case is roughly 1%, so the

TABLE 13. TAYLOR \bar{n} EFFICIENCIES

SLR	Max η Values		Monotonic \bar{n}	
	\bar{n}	η	\bar{n}	η
25	12	0.9252	5	0.9105
30	23	0.8787	7	0.8619
35	44	0.8326	9	0.8151
40	81	0.7899	11	0.7729

TABLE 14. BAYLISS LINE SOURCE PARAMETERS

	SLR, Decibels					
	15	20	25	30	35	40
A	1.00790	1.22474	1.43546	1.64126	1.84308	2.04154
z_1	1.51240	1.69626	1.88266	2.07086	2.26025	2.45039
z_2	2.25610	2.36980	2.49432	2.62754	2.76748	2.91231
z_3	3.16932	3.24729	3.33506	3.43144	3.53521	3.64518
z_4	4.12639	4.18544	4.25273	4.32758	4.40934	4.49734
z_0	0.66291	0.71194	0.75693	0.79884	0.83847	0.87649

monotonic \bar{n} should be used as an upper limit. The best \bar{n} can be selected by comparing computer runs of patterns for various values of \bar{n} . For additional information, see Hansen.*

Bayliss \bar{n} One-Parameter Difference Line Source

The Bayliss space factor is a difference pattern (for tracking purposes) constructed to have the same features as the Taylor \bar{n} space factor. That is, it has \bar{n} roughly equal-level sidelobes adjacent to each difference beam, with a $1/u$ envelope beyond. Taylor started with an "ideal" pattern, $\cos \pi \sqrt{(u^2 - A^2)}$, which has equal-level sidelobes. Then a space factor was constructed with \bar{n} close-in zeros matching those of the "ideal" and with remaining zeros to give a $1/u$ envelope. Finally, a dilation factor was used to make a smooth transition around $u = \bar{n}$. A good starting point to get a difference pattern is to differentiate the "ideal" sum pattern. This gives $u \text{ sinc } \pi \sqrt{(u^2 - A^2)}$. Unfortunately, not all sidelobes are of equal level; the first several are tapered. Bayliss used an iterative procedure to adjust these zeros to yield equal-level sidelobes. It was necessary to adjust only four. Since these four zeros (and A) depend on the sidelobe ratio, results for each were given in terms of fourth-order polynomials in SLR. Table 14 gives the results, for various values of SLR, for z_1, z_2, z_3, z_4, A , and u_0 . The latter is the value of u at the difference peak. These zeros are used for both Bayliss line sources and circular sources.† For $n > \bar{n}$, $z_n = \sqrt{A^2 + n^2}$. Next, the Taylor procedure is followed, where the envelope is made to approach $1/u$

* Hansen, R. C. Array chapters in *Handbook of Antenna Design*. Vol. 2. A. W. Rudge et al. (Eds.). London: Peter Peregrinus Ltd., 1983; *Phased Array Antennas* New York: John Wiley & Sons, Inc., 1998.

† Bayliss, E. T. "Design of Monopulse Antenna Difference Patterns With Low Sidelobes." *BSTJ*, Vol. 47, May-June 1968, pp. 623-650. Hansen, R. C. "Array Pattern Control and Synthesis." *Proc. IEEE*, Vol. 80, January 1992, pp. 141-151.

for large u , and with a smooth transition. The difference pattern is

$$F(u) = u \cos \pi u = \frac{\prod_{n=1}^{\bar{n}-1} (1 - u^2 / \sigma^2 z_n^2)}{\prod_{n=0}^{\bar{n}-1} (1 - u^2 / (n+1/2)^2)}$$

As before,

$$\sigma = (\bar{n} + 1/2) / \sqrt{A^2 + \bar{n}^2}$$

Fig. 88 shows a Bayliss line source pattern for SLR 25 dB, $\bar{n} = 5$. The aperture distribution for $-1 \leq p \leq 1$ is

$$g(p) = \sum_{n=0}^{\bar{n}-1} B_n \sin \pi p (n+1/2)$$

where

$$B_m = -(-1)^m (m-1/2)^2 \frac{\prod_{n=1}^{\bar{n}-1} [1 - (m+1/2)^2 / \sigma^2 z_n^2]}{\prod_{\substack{n=0 \\ n \neq m}}^{\bar{n}-1} [1 - (m+1/2)^2 / (n+1/2)^2]}$$

These formulas allow rapid calculation of Bayliss space factors and aperture distributions. Fig. 89 gives the aperture distribution corresponding to the space factor of Fig. 88. Fig. 89 shows only amplitude; each half of the aperture has constant phase, with the two halves out of phase. Excitation efficiency, which is the ratio of directivity (at one difference peak) to that of a uniformly excited line source, is given by

$$\eta = \frac{2u_0^2 \cos^2 \pi u_0 \sum_{n=0}^{\bar{n}-1} (-1)^n B_n / [u_0^2 - (n+1/2)^2]}{\sum_{n=0}^{\bar{n}-1} B_n^2}$$

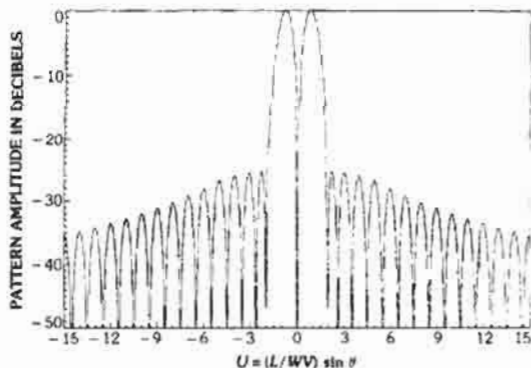


Fig. 88. Bayliss space factor for SLR = 25 dB, $\bar{n} = 5$.

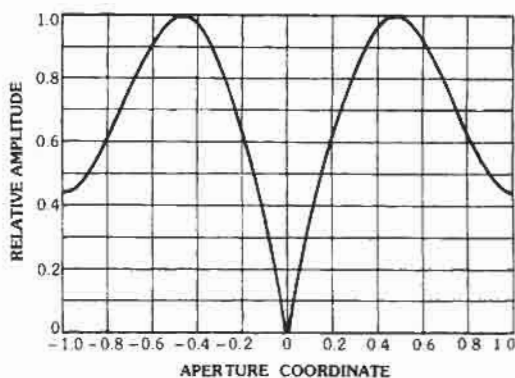


Fig. 89. Bayliss distribution for SLR = 25 dB, $\bar{n} = 5$.

The slope at $u = 0$ is given by

$$S = (2/\pi) \sum_{n=0}^{\bar{n}-1} (-1)^n B_n / (n+1/2)^2$$

This value is 0.6366 over the range of SLR in Table 15. Table 15 gives efficiency and normalized slope values (normalized for maximum linear phase slope of unity) for Bayliss space factors.

Low Sidelobe Distributions

High sidelobe ratios, in the range from 30 to 60 dB, have recently attracted more interest. The one-parameter space factors have very low pedestals for this range of SLR, and thus the Taylor \bar{n} space factors are more appropriate. Again, efficient design requires properly spaced space-factor zeros. In contrast, most of the popular distributions were coined before the computer allowed performance evaluations to be easily performed. The Taylor \bar{n} distribution is, in essence, perfect for low sidelobes because it is low- Q , which eases the tolerance problem, and because the zeros are monotonically and smoothly spaced. This can be made

TABLE 15. BAYLISS EFFICIENCY AND PATTERN SLOPE

SLR (dB)	\bar{n}	Efficiency η	Normalized Slope S
15	4	0.5959	0.9567
20	4	0.5846	0.8974
25	5	0.5633	0.8427
30	6	0.5393	0.7912
35	7	0.5162	0.7448
40	8	0.4951	0.7037

evident by comparing the Taylor with the popular Hamming distribution. The latter* is

$$g(p) = 0.54 + 0.46 \cos \pi p = a + b \cos \pi p$$

The excitation efficiency of the Hamming is

$$\eta = 2a^2 / (2a^2 + b^2) = 0.7338$$

The space factor has zeros at $u = 2, 3, 4, \dots$ and, in addition, a zero at $u = \sqrt{a/(a-b)} = 2.5981$. This close spacing of the first three zeros produces an uneven sidelobe envelope with the fourth sidelobe the highest at -42.7 dB. Thus, the Hamming is compared with a Taylor \bar{n} distribution with SLR = 42.7 dB. Table 16 shows the zeros, and it can be observed that both of the Taylor space factors have a smoothly increasing zero spacing from the first zero out to the transition point, beyond which the spacing is unity. The Hamming, on the other hand, has a first spacing of 0.598, a second spacing of 0.402, and the remaining spacings all unity. This is obviously not as good a design, as the following will show.

Table 17 gives efficiency and normalized beamwidth for the three Taylor cases and for the Hamming. It may be seen that the Hamming beamwidth is

TABLE 16. ZEROS OF DISTRIBUTIONS FOR SLR = 42.7 DB

n	Taylor		Hamming	
	One-Parameter	Taylor		
		$\bar{n} = 6$	$\bar{n} = 10$	
1	2.112	1.894	1.897	2
2	2.732	2.398	2.396	2.598
3	3.550	3.173	3.166	3
4	4.412	4.069	4.056	4
5	5.335	5.020	5.002	5
6	6.282	6	5.978	6
7	7.243	7	6.970	7
8	8.214	8	7.974	8
9	9.190	9	8.984	9
10	10.172	10	10	10

TABLE 17. COMPARISON OF DISTRIBUTIONS FOR SLR = 42.7 DB

	Taylor		Hamming	
	One-Parameter	Taylor		
		$\bar{n} = 6$	$\bar{n} = 10$	
u_3	0.694	0.637	0.635	0.651
η	0.690	0.754	0.755	0.734

roughly 3% broader, while the efficiency is roughly 3% lower. Thus, by use of appropriate space factors such as the Taylor \bar{n} , improved performance can be obtained at no additional cost in complexity, hardware, or tolerances. For rectangular arrays with separable aperture distributions, Taylor distributions may, of course, be used along each coordinate. For circular disk apertures, the circular Taylor \bar{n} distribution is similarly excellent.

Measurement of Low-Sidelobe Patterns

In addition to the obvious problem of providing a pattern range with a low background level, the measurement of low sidelobe patterns imposes restrictions on the measurement distance. A distance of $2D^2/\lambda$ is adequate for patterns with modest sidelobe ratios (25 dB or less), but larger distances are needed to measure lower sidelobes accurately, due to the quadratic phase error produced by a finite distance. As the measurement distance is reduced, the first sidelobe and first null fills in, and then the sidelobe becomes a shoulder on the main beam. (The second sidelobe is only slightly raised at this point.) At closer distances, the main beam broadens, and the second sidelobe and second null rise, and again the main beam eventually absorbs the sidelobe. Calculations have been made of the change in sidelobe level versus measurement distance for the universal Taylor \bar{n} linear space factor.† Fig. 90 gives these results for a uniform line source,

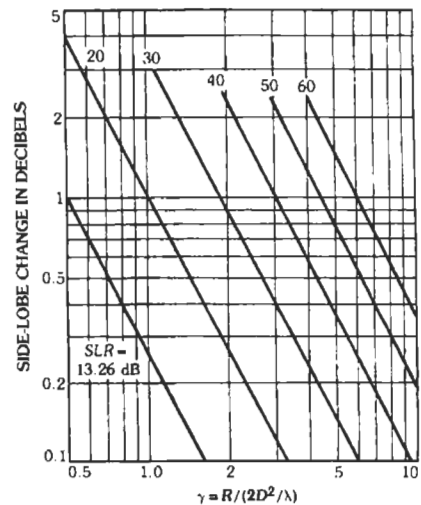


Fig. 90. Sidelobe increase due to measurement distance.

* Blackman, R. B., and Tukey, J. W. *Measurement of Power Spectra*. New York: Dover Publications, Inc., 1958.

† Hansen, R. C. "Measurement Distance Effects on Low Sidelobe Patterns." *Trans. IEEE*, Vol. AP-32, June 1984; pp. 591-594.

and for 20 (10) 60-dB Taylor \bar{n} line sources. Distance is normalized to $2D^2/\lambda$, where

$$\gamma = R/(2D^2/\lambda)$$

Each curve terminates when the sidelobe becomes a shoulder on the main beam, without any dip.

Optimum low sidelobe difference patterns are also of interest. Calculations have been made for Bayliss space factors for a uniform line source and for sidelobe ratios of 20 (10) 50 dB.* See Fig. 91.

Hansen One-Parameter Circular Source Distribution

A symmetric distribution for circular disk apertures analogous to the Taylor one-parameter line source distribution has been developed by Hansen.† This is a modified $2J_1(\pi u)/\pi u$ pattern just as the Taylor one-parameter was a modified sinc πu . The close-in zeros are shifted to produce the desired SLR. The single parameter is H , and the space factor is again written in two forms depending on whether u is less or greater than H . These are

$$F(u) = 2J_1(\pi\sqrt{u^2 - H^2})/\pi\sqrt{u^2 - H^2} \quad u \geq H$$

$$F(u) = 2I_1(\pi\sqrt{H^2 - u^2})/\pi\sqrt{H^2 - u^2} \quad u \leq H$$

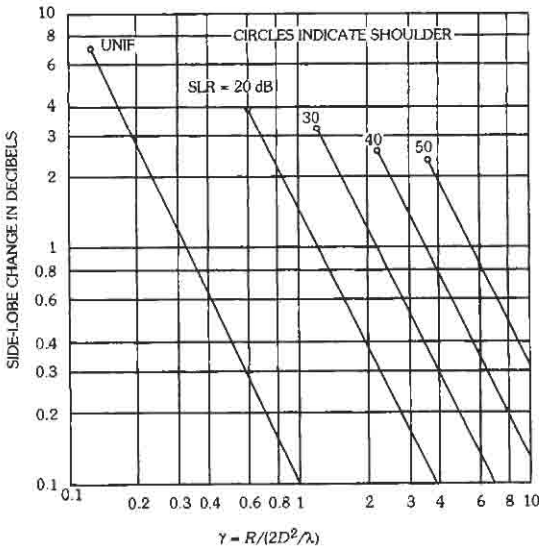


Fig. 91. Sidelobe change versus normalized measurement distance for Bayliss line source.

* Hansen, R. C. "Measurement Distance Effects on Bayliss Difference Patterns." *Trans. IEEE*, Vol. AP-40, 1992, pp. 1211-1214.

† Hansen, R. C. "A One-Parameter Circular Aperture Distribution with Narrow Beamwidth and Low Sidelobes." *Trans. IEEE*, Vol. AP-24, July 1976, pp. 477-480.

The sidelobe ratio is given by

$$SLR = 20 \log [2I_1(\pi H)/\pi H] + 17.57 \text{ dB}$$

In these formulas, J_1 and I_1 are the usual Bessel and modified Bessel functions of first kind and order one. Again, the top portion of the main beam is provided by the modified Bessel form, while the remainder of the main beam and the sidelobe structure are provided by the Bessel form. The aperture distribution is

$$g(\rho) = I_0(H\sqrt{\pi^2 - \rho^2})$$

where ρ is zero at the center of the aperture and the aperture radius is π . The patterns and aperture distributions are much like those of the Taylor \bar{n} line source and so have not been included. Table 18 gives parameter H , beamwidth u_3 , the aperture excitation efficiency, and the beam efficiency. The latter is the fraction of energy that is contained in the main beam, null to null. The actual beamwidth is approximately $2u_3$. A more exact result for apertures small in wavelengths is found in a previous section for the Taylor one-parameter line source. The excitation efficiency is given by the expression

$$\eta = 4I_1^2(\pi H)/\pi^2 H^2 [I_0^2(\pi H) - I_1^2(\pi H)]$$

The I_0 used above is the modified Bessel function of order zero. This distribution, like the Taylor distribution, is low- Q and robust, and has been found to fit well to many reflector antenna patterns, for the main beam and close-in sidelobes.

Taylor \bar{n} Circular Source Distribution

The Taylor \bar{n} circular source distribution offers an improvement in efficiency and beamwidth over the Hansen one-parameter circular source distribution, just as the Taylor \bar{n} line source distribution shows an improvement over the Taylor one-parameter source distribution. The circular aperture is a disk with rotationally symmetric excitation, and, again, the starting point is the uniform $2J_1(\pi u)/\pi u$ space factor. On each side of the main beam, \bar{n} zeros are modified by mov-

TABLE 18. CHARACTERISTICS OF HANSEN ONE-PARAMETER DISTRIBUTION

SLR (dB)	H	u_3 (rad)	η	η_b
17.57	0	0.5145	1	—
25	0.8899	0.5869	0.8711	0.9745
30	1.1977	0.6304	0.7595	0.9930
35	1.4708	0.6701	0.6683	0.9981
40	1.7254	0.7070	0.5964	0.9994
45	1.9681	0.7413	0.5390	0.9998
50	2.2026	0.7737	0.4923	1.0000

ing them to produce the desired sidelobe ratio.* Again, a dilation factor, σ , is used to provide a smooth transition between the roughly equal-level sidelobes and the tapered-envelope sidelobes. The pattern is given by a canonical product on zeros:

$$F(u) = [2J_1(\pi u)/\pi u] \prod_{n=1}^{\bar{n}-1} (1 - u^2/u_n^2)(1 - u^2/\mu_n^2)$$

where μ_n are the zeros of $J_1(\pi u)$. The close-in pattern zeros are given by

$$u_n = \pm \sigma \sqrt{A^2 + (n-1/2)^2} \quad 1 \leq n \leq \bar{n}$$

while

$$\sigma = \mu_{\bar{n}} \sqrt{A^2 + (\bar{n}-1/2)^2}$$

Table 19 gives the sidelobe parameter, A , the beamwidth, u_3 , and values of σ for various values of \bar{n} . The actual beamwidth for large apertures is $2\sigma u_3$. The patterns and aperture distributions are much like those of the Taylor \bar{n} line source and so have not been included. The aperture distribution is

$$g(p) = (2/\pi^2) \sum_{n=0}^{\bar{n}-1} F_n J_0(p\mu_n) / [J_0(\pi\mu_n)]^2$$

where

$$p = 2\pi p/D$$

* Taylor, T. T. "Design of Circular Apertures for Narrow Beamwidths and Low Sidelobes." *IRE Trans.*, AP-8, 1960 pp. 17-22.

and

$$F_m = -J_0(\pi\mu_m) \frac{\prod_{l=1}^{\bar{n}-1} (1 - \mu^2_m / u^2_l)}{\prod_{l=1, l \neq m}^{\bar{n}-1} (1 - \mu^2_m / \mu^2_l)}$$

All of these formulas are easily computed, so no tables have been included. For checking purposes, reference may be made to tables of Hansen.† The aperture excitation efficiency is given by

$$\eta = 1 / \left[1 + \sum_{n=1}^{\bar{n}-1} F_n^2 / J_0^2(\pi\mu_n) \right]$$

Table 20 gives the aperture excitation efficiency for several combinations of SLR and \bar{n} .

REFLECTORS

Parabolic Reflectors

The parabolic reflector commonly exists in both focal feed and Cassegrain form (Figs. 92 and 93). Offset reflectors are covered later. For a front-fed reflector, the reflector f/D must be matched to the feed pattern. Reflectors with pattern sidelobes roughly -25 dB and

† Hansen, R. C. "Tables of Taylor Distributions for Circular Aperture Antennas." *IRE Trans.*, Vol. AP-8, January 1960, pp. 23-26. Hansen, R. C. *Phased Array Antennas*. New York: John Wiley & Sons, Inc., 1998.

TABLE 19. TAYLOR \bar{n} CIRCULAR SOURCE CHARACTERISTICS

SLR (dB)	A	u ₃	η						
			$\bar{n} = 4$	5	6	7	8	9	10
20	0.9528	0.4465	1.1692	1.1398	1.1186	1.1028	1.0906	1.0810	1.0732
25	1.1366	0.4890	1.1525	1.1296	1.1118	1.0979	1.0870	1.0782	1.0708
30	1.3200	0.5284	1.1338	1.1180	1.1039	1.0923	1.0827	1.0749	1.0683
35	1.5032	0.5653	1.1134	1.1050	1.0951	1.0859	1.0779	1.0711	1.0653
40	1.6865	0.6000	1.0916	1.0910	1.0854	1.0789	1.0726	1.0670	1.0620

TABLE 20. EXCITATION EFFICIENCY VERSUS \bar{n}

SLR (dB)	η				
	$\bar{n} = 4$	5	6	8	10
20	0.9723	0.9356	0.8808	0.7506	0.6238
25	0.9324	0.9404	0.9379	0.9064	0.8526
30	0.8482	0.8623	0.8735	0.8838	0.8804
35	0.7708	0.7779	0.7880	0.8048	0.8153
40	0.7056	0.7063	0.7119	0.7252	0.7365

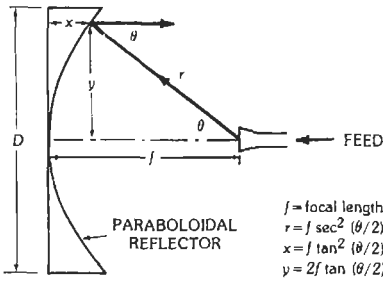


Fig. 92. Paraboloidal-reflector design.

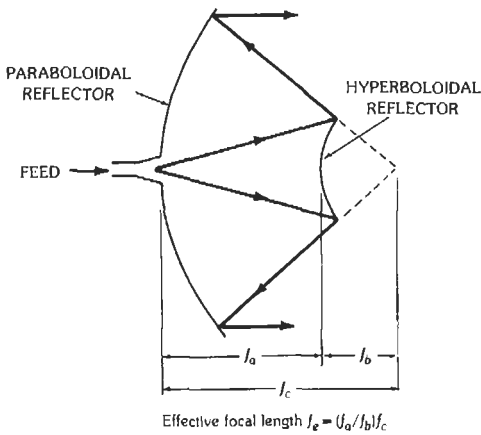


Fig. 93. Cassegrain reflector system.

good efficiency typically have illumination edge tapers of -10 to -11 dB. For lower sidelobes, the edge taper must be lower. The tradeoff between edge taper and sidelobe level can be accurately determined from the circular one-parameter distributions in the section on aperture distributions. Fig. 94 shows the total reflector subtended angle versus f/D , where

$$\theta = 2 \arctan [(8f/D)/(16f^2/D^2 - 1)]$$

Fig. 95 gives the differential path loss between the edge ray and the apex ray, also as a function of f/D :

$$R_{cdgc}/f = 1 + 1/(16f^2/D^2)$$

The edge illumination is, of course, the feed pattern value at the edge angle plus the differential path loss. Fig. 96 allows horn beamwidth at any level to be converted to the -10 dB beamwidth; this curve is an excellent empirical fit to a large amount of experimental data:

$$E_{dB} = 10 (\theta/\theta_{10})^2$$

Aperture blockage is a major limitation of front-fed reflectors, especially for low design sidelobes. Fig. 97 shows deterioration of sidelobe level versus blockage

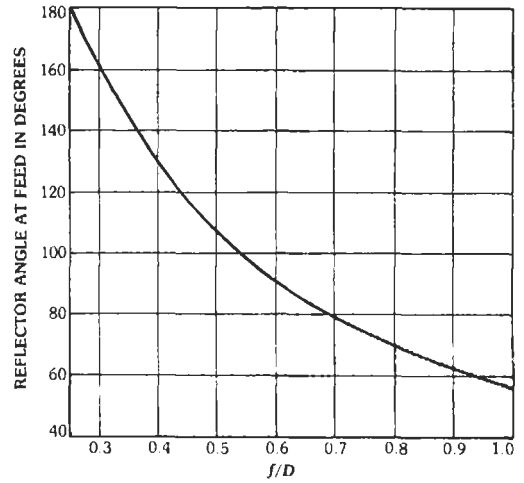


Fig. 94. Angle of reflector at feed.

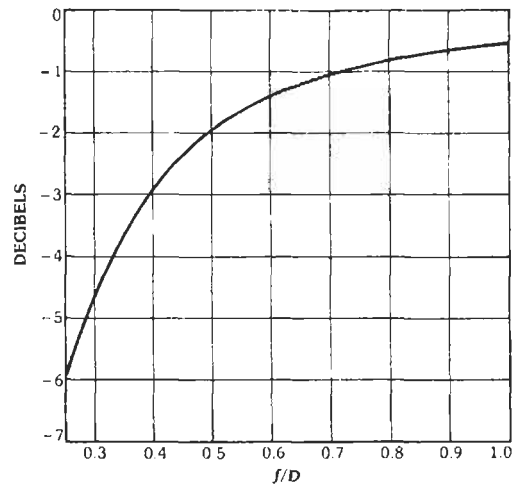


Fig. 95. Edge-center space loss.

diameter ratio for the Hansen universal one-parameter circular distribution. The reflector curvature produces cross polarization off the axis with a resulting gain loss. Determination of aperture efficiency or gain is complex and involves calculation of aperture taper efficiency, feed spillover, feed cross polarization loss, blockage loss, and reflector cross polarization loss. For more information, refer to Rusch et al.* Random

* Rusch, W. V. T., et al. "Quasi-Optical Antenna Design and Applications." Chapter 3 in *Handbook of Antenna Design*. Vol. 1. A. W. Rudge et al. (eds.). London: Peter Peregrinus Ltd., 1983. Rusch, W. V. T. "The Current State of the Reflector Antenna Art-Entering the 1990's." *Proc. IEEE*, Vol. 80, January 1992, pp. 113-126.

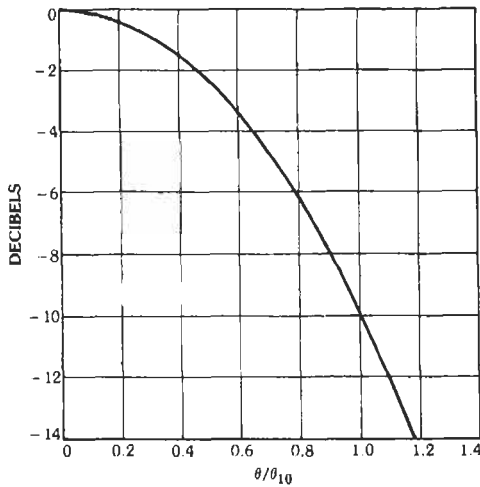


Fig. 96. Universal horn beamwidth conversion.

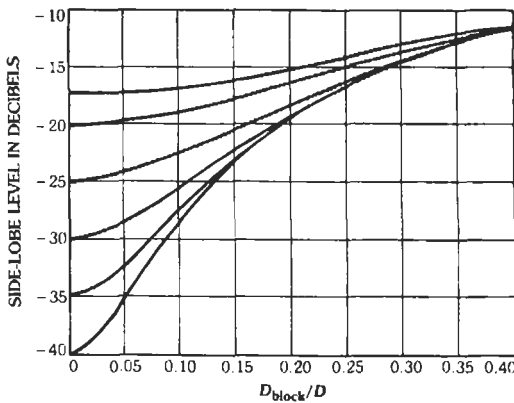


Fig. 97. Sideline level vs the blockage diameter ratio for Hansen one-parameter circular distribution.

errors in the reflector surface raise the sidelobes and reduce the gain. The effect on sidelobes can be estimated by using the universal curve in Fig. 97, where sidelobe degradation is plotted against σ/\sqrt{G} . Note, however, that the σ measured for a reflector surface must be doubled before being used in that figure, due to the reflective operation of the surface. Directivity loss has been calculated by Ruze* with a modification added by Wested† to take into account f/D . His result, which is

$$G/G_0 = (A + 1)/[A + \exp(4\pi\epsilon/\lambda)^2]$$

* Ruze, J. "Antenna Tolerance Theory—A Review." *Proc. IEEE*, Vol. 54, April 1966, pp. 633–640.

where

$$A = 1/[16f^2/D^2 \ln(1 + 1/16f^2/D^2)] - 1$$

is shown in Fig. 98, which gives gain decrease against one sigma tolerance in wavelengths (ϵ/λ) for several values of f/D . It can be seen that for loss of more than a few tenths of a decibel, f/D needs to be included. Most narrow-band reflector antennas have aperture efficiency (gain compared to that of a uniformly excited area) in the 55–65% range. Wideband antennas using log periodic or conical spiral feeds may have efficiencies anywhere from 20 to 50%.

The Cassegrain antenna, through use of a hyperboloidal subreflector, allows the feed to protrude through a hole in the dish. Thus, low-noise and/or high-power components can be conveniently located behind the dish. The minimum blockage condition exists when the feed blockage equals the subreflector blockage. Fig. 99 shows gain decrease for equal-blockage Cassegrain reflectors versus beamwidth. From this and previous material, it is apparent that front-fed reflectors are attractive for moderate beamwidths, whereas Cassegrains are attractive for beamwidths of several degrees or less. Cross-polarization effects of a Cassegrain reflector are shown in Fig. 100. A high-efficiency Cassegrain may be designed by shaping the two reflector surfaces; only small changes on the order of a wavelength at the edge arc needed.‡ This accomplishes an aperture distribution closer to uniform;

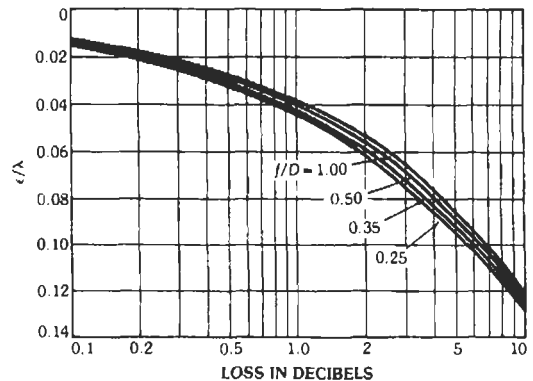


Fig. 98. Reflector directivity loss vs RMS tolerance.

† Wested, J. H. "Effect of Deviation from the Ideal Paraboloidal Shape of Large Antenna Reflectors." *Proc. IEE Conf. Large Steerable Antennas*. (London, England), 1966, pp. 115–119.

‡ Galindo, V. "Design of Dual-Reflector Antennas with Arbitrary Phase and Amplitude Distribution." *Trans. IEEE*, Vol. AP-12, July 1964, pp. 403–408. Williams, W. F. "High Efficiency Antenna Reflector." *Microwave J.*, Vol. 8, July 1965, pp. 79–82. Galindo-Israel, V., et al. "On the Theory of Synthesis of Single and Dual Offset Shaped Reflector Antennas." *Trans. IEEE*, Vol. AP-35, August 1987, pp. 887–896.

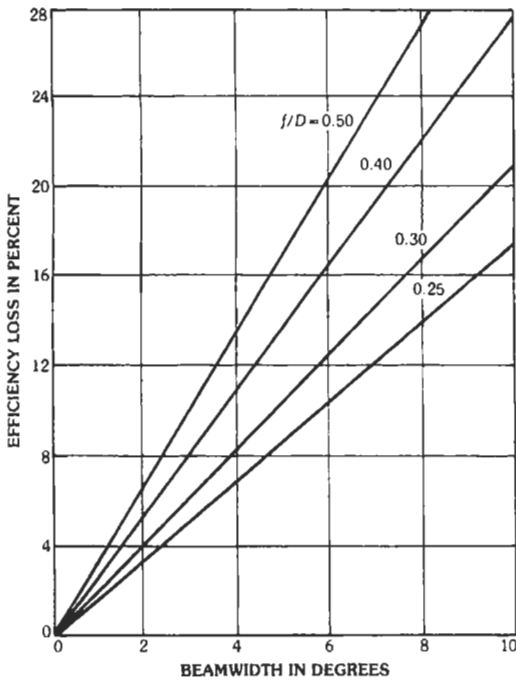


Fig. 99. Minimum blockage loss vs antenna beamwidth. (From W. V. T. Rusch, et al., "Quasi-Optical Antenna Design and Applications." *Handbook of Antenna Design*, Vol. 1. A. W. Rudge et al. (Eds.). London: Peter Peregrinus Ltd., 1983.)

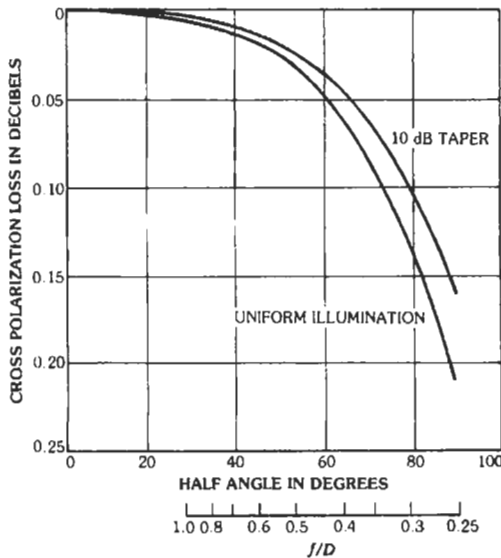


Fig. 100. Cross polarization vs subtended half angle due to reflector curvature. (From W. V. T. Rusch, et al., "Quasi-Optical Antenna Design and Applications." *Handbook of Antenna Design*, Vol. 1. A. W. Rudge et al. (Eds.). London: Peter Peregrinus Ltd., 1983.)

overall efficiencies as high as 80% have been reported. However, the (corrugated) horn required to provide the requisite feed pattern tends to provide only limited bandwidth.

Scanning and Multiple-Beam Reflectors

When the feed is moved laterally, the beam moves in the opposite direction, but through a slightly smaller angle. The ratio of beam movement to feed movement is called "beam deviation factor." It is shown in Fig. 101 versus f/D for the universal one-parameter circular distribution.

There is a concomitant loss of directivity due to coma, and this was calculated by Ruze* for several distributions. A recalculation using the circular one-parameter distribution is shown in Fig. 102. These cal-

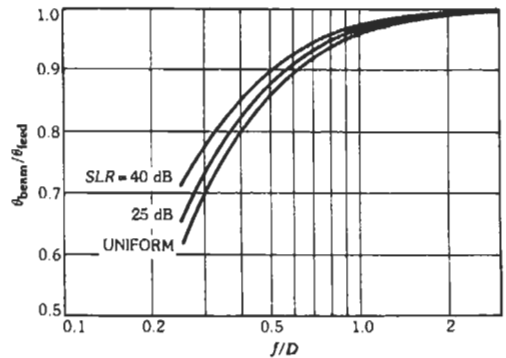


Fig. 101. Beam deviation factor of reflector with Hansen one-parameter distribution.

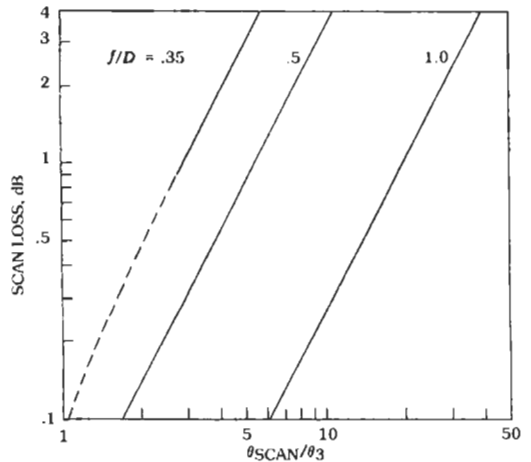


Fig. 102. Approximate parabolic reflector scan loss.

* Ruze, J. "Lateral Feed Displacement in a Paraboloid." *Trans. IEEE*, Vol. AP-13, September 1965, pp. 660-665.

culations are only approximate as they are based on the pattern peak value, instead of integrating over the entire pattern.

When the coma loss is unacceptably large, an array feed can be used in the transverse focal region. One coma correction scheme sums the array element outputs with weighting that is the conjugate of the focal field. As the scan angle becomes larger the array feed must be larger, to cover the larger coma distribution. Since the coma is all on one side of boresight, the feed array need only cover this area, but as the azimuth scan angle changes, the feed array must be rotated to match.

Another method requires only phase correction and utilizes a beam-forming network connected to the feed array. The phases at the BFN output are then corrected with line lengths (in transform space) and summed.* When a single horn is moved off axis, the surface giving the smallest diameter circle of confusion is not plane and is called the Petzval surface. The maximum-

gain surface also is not plane and deviates slightly from the Petzval surface. Generally the precise focal surface used is not important. Multiple-beam antennas typically employ a cluster of horns with each horn producing a separate beam. These configurations are usually offset to reduce the large blockage that would occur with a front-fed reflector. For further information, see Rusch et al.†

* Rudge, A. W., and Davies, D. E. N. "Electronically Controllable Primary Feed for Profile-Error Compensation of Large Parabolic Reflectors." *Proc. IEE*, Vol. 117, February 1970, pp. 351-357.

† Rusch, W. V. T., et al. "Quasi-Optical Antenna Design and Applications." *Handbook of Antenna Design*, Vol. 1. A. W. Rudge et al. (Eds.). London: Peter Peregrinus Ltd., 1983.

33

Electromagnetic-Wave Propagation

*Revised by
Douglass D. Crombie*

Very-Low Frequencies—Up to 30 Kilohertz	33-3
Low and Medium Frequencies—30 to 3000 Kilohertz	33-3
High Frequencies—3 to 30 Megahertz	33-6
Angles of Departure and Arrival	
Forecasts of High-Frequency Propagation	
Bandwidth Limitations	
Diversity	
Great-Circle Calculations	
Effect of Nuclear Explosions on Radio Propagation	33-12
Ionospheric Scatter Propagation	33-12
Meteor-Burst Propagation	33-13
Propagation Above 30 Megahertz, Line-of-Sight Conditions	33-13
Radio Refraction	
Path Plotting and Profile-Chart Construction	
Fresnel Zones	
Required Path Clearance	
Interference Between Direct and Reflected Rays	
Space-Diversity Reception	
Variation of Field Strength With Distance	
Fading and Diversity	
Atmospheric Absorption	

Free-Space Transmission Equations
Effective Areas of Typical Antennas
Antenna Gain Relative to Hypothetical Isotropic Antennas
Antenna Beam Width
Transmitter Power for a Required Output Signal/Noise Ratio

Knife-Edge Diffraction Propagation 33-29

Tropospheric Scatter Propagation 33-29

Earth-Space Communication 33-33

Radio waves may be propagated* from the transmitting antenna to the receiving antenna through or along the surface of the earth, through the atmosphere, or by reflection or scattering from natural or artificial reflectors. The conductivity and dielectric constant of the ground vary considerably from those of the atmosphere. At very-low frequencies, ground waves may be satisfactorily propagated for distances of several thousand kilometers. At high frequencies, however, the losses are so great that signals can be propagated for only a few hundred kilometers by ground wave. Propagation in the medium- and high-frequency bands is chiefly by ground wave and by reflection from the ionosphere, and severe fading is caused in these frequency bands by the interference between ground and ionospheric waves.

The refractive index of the atmosphere is an important factor in radio propagation. At frequencies between about 100 and 8000 megahertz, scattering of radio waves by inhomogeneities in the electromagnetic characteristics of the atmosphere can be used to provide satisfactory wideband communication up to several times the line-of-sight distance. New techniques are being developed for generating coherent high-power waves in the optical spectrum. Atmospheric absorption at these frequencies is high, but the large bandwidths and small antenna beam widths may make such frequencies practical for certain applications.

VERY-LOW FREQUENCIES— UP TO 30 KILOHERTZ

The propagation of long radio waves is of considerable importance in reliable communication, long-range navigation, and the detection of nuclear explosions. Considerable progress has been made in recent years in understanding the propagation of such waves in the earth-ionosphere waveguide.†

At short distances from a transmitter, the received signal is chiefly by a ground or surface wave, and at very-low frequencies its intensity is essentially inversely proportional to distance. At greater distances, the field intensity falls at a higher rate because of losses in the ground and because of the curvature of the earth. These losses increase with frequency. At sufficiently great distances, the received level is chiefly due to sky waves reflected from the ionosphere. At intermediate distances, the field is a combination of sky waves and ground waves that result in an interference pattern. The total field at the receiver may be obtained in two distinct ways. The first method, which leads to the geometric-optics theory, directly sums the contributions at the receiver from the primary source and each of its images.

* CCIR XVth Plenary Assembly, Geneva, 1982; Vol. V, *Propagation in Non-Ionized Media*, and Vol. VI, *Propagation in Ionized Media*.

† Watt, A. D. *VLF Radio Engineering*, Vol. 14, International Series of Monographs in Electromagnetic Waves. New York: Pergamon Press, Inc., 1967.

The second method treats the source and its images as self-illuminating diffraction gratings, one above the earth and one below, and leads to the waveguide mode theory. The advantages of mode theory are restricted to very-low frequencies, where relatively few modes can be supported in the earth-ionosphere waveguide. For example, when the height of the ionosphere is 80 kilometers and the wavelength is 20 kilometers, only the first eight modes can be supported. When the wavelength is 2 kilometers, however, all modes up to the 80th can be supported.

Thus at very-low frequencies, and for distances greater than say 3000 kilometers, it is simpler to use the mode of lowest order to obtain the received field. At distances less than about 1000 kilometers, it is simpler to use ray theory. The above is based on an idealized condition, since it has been assumed that the earth is flat, that the ionosphere is sharply bounded, and that the effect of the magnetic field of the earth can be ignored. Even with these simplifying assumptions, the results are useful. A full treatment of the general case can be obtained by reference to several sources.*

The results of calculations made by Wait and Spies, taking into account the curvature of the earth and the conductivities of the earth and the ionosphere, are shown in Figs. 1 and 2, where σ_g = ground conductivity, ω_r = ionospheric conductivity parameter, and n = mode order. It is seen that the lowest attenuation of mode 1 in the daytime ($h = 70$ km) is at about 18 kilohertz and at night ($h = 90$ km) it is at about 15 kilohertz.

The phase velocity ratio of mode 1 in the daytime is greater than 1 for frequencies less than about 13 kilohertz and less than 1 for higher frequencies. At night the crossover frequency is about 9 kilohertz.

In daylight, the attenuation of mode 1 is always less than that of mode 2. At night, the attenuation for the two modes may be of the same order.

LOW AND MEDIUM FREQUENCIES—30 TO 3000 KILOHERTZ†

For low and medium frequencies of approximately 30 to 3000 kilohertz with a short vertical antenna over perfectly reflecting ground

$$E = 186.4(P_r)^{1/2} \text{ millivolts rms/meter at 1 mile}$$

* Wait, J. R., *Electromagnetic Waves in Stratified Media*, New York: Pergamon Press, Inc., 1962. Budden, K. G., *The Wave-Guide Mode Theory of Wave Propagation*, New York: Prentice-Hall, Inc., 1962. Johler, J. R., "Propagation of the Low-Frequency Radio Signal," *Proceedings of the IRE*, Vol. 50, No. 4, 1962; pp. 404-427. CCIR XVth Plenary Assembly, Geneva, 1982, Vol. VI, Report 895.

† *Radio Spectrum Utilization*. New York: Joint Technical Advisory Committee (IEEE and EIA), 1964. CCIR XVth Plenary Assembly, Geneva, 1982, Vol. VI.

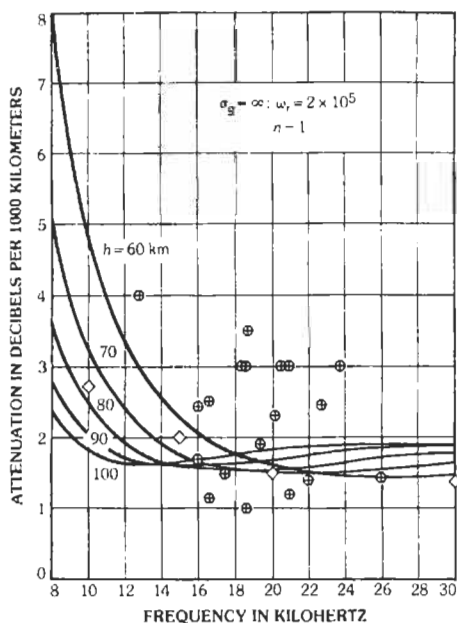


Fig. 1. Relation among attenuation, frequency, and height of the ionosphere. The diamonds represent some experimental observations by Taylor on the average daytime attenuation for west-to-east propagation over sea water. Attenuation rate in the opposite direction is greater by about 1 decibel per 1000 kilometers. (From Radio Spectrum Utilization, Joint Technical Advisory Committee (IEEE and EIA), IEEE, New York, 1964; p. 104)

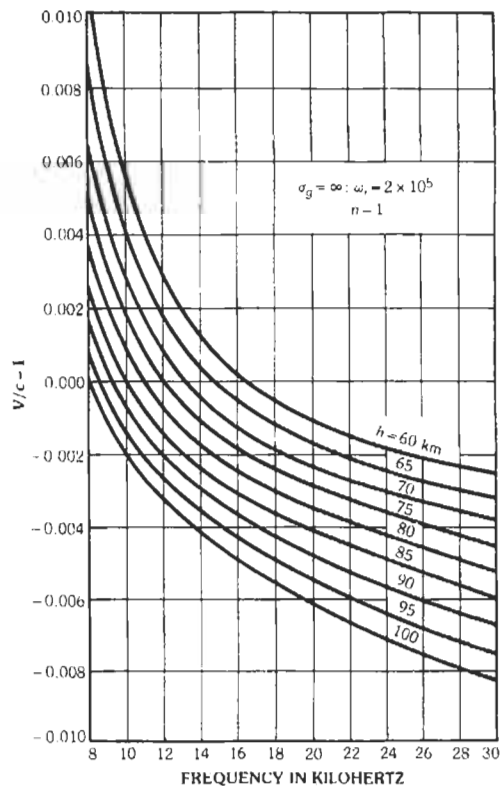


Fig. 2. Phase velocity V as a function of ionospheric height h and frequency, relative to the velocity in free space, c . (From Radio Spectrum Utilization, Joint Technical Advisory Committee (IEEE and EIA), IEEE, New York, 1964; p. 105.)

$$E = 300(P_t)^{1/2} \text{ millivolts rms/meter at 1 kilometer}$$

where P_t = radiated power in kilowatts.

Actual inverse-distance fields at 1 kilometer for a given transmitter output power depend on the height and power radiation efficiency of the antenna and associated circuit losses.

Typical values found in practice for well-designed stations are:

Small L or T antennas as on ships:

$$40(P_t)^{1/2} \text{ millivolts/meter at 1 kilometer}$$

Vertical radiators 0.15 to 0.25 λ high:

$$290(P_t)^{1/2} \text{ millivolts/meter at 1 kilometer}$$

Vertical radiators 0.25 to 0.40 λ high:

$$322(P_t)^{1/2} \text{ millivolts/meter at 1 kilometer}$$

Vertical radiators 0.40 to 0.60 λ high or top-loaded vertical radiators:

$$386(P_t)^{1/2} \text{ millivolts/meter at 1 kilometer}$$

where P_t = transmitter output power in kilowatts. These values can be increased by directive antenna systems.

It has been found that the concept of basic transmission loss, also called path loss, is convenient for the analysis of radio communication systems. Basic transmission loss is the dimensionless ratio P_R/P_A , where P_R is the power radiated from a lossless, isotropic transmitting antenna and P_A is the power available from a lossless, isotropic receiving antenna in a matched load. The isotropic antennas are at the same physical locations and operate in the same band of frequencies as the actual antennas.

Surface-wave (commonly called ground-wave) basic transmission loss is plotted in Fig. 3 for vertically polarized propagation over land having a representative conductivity and dielectric constant and in Fig. 4 for vertically polarized propagation over sea water. Both antennas are 30 feet above the surface in both figures.

In the low-frequency and medium-frequency ranges, propagation losses for horizontally polarized transmission between antennas on the surface of the earth are impractically high. Ground constants typical of various terrain types are listed in Table 1.

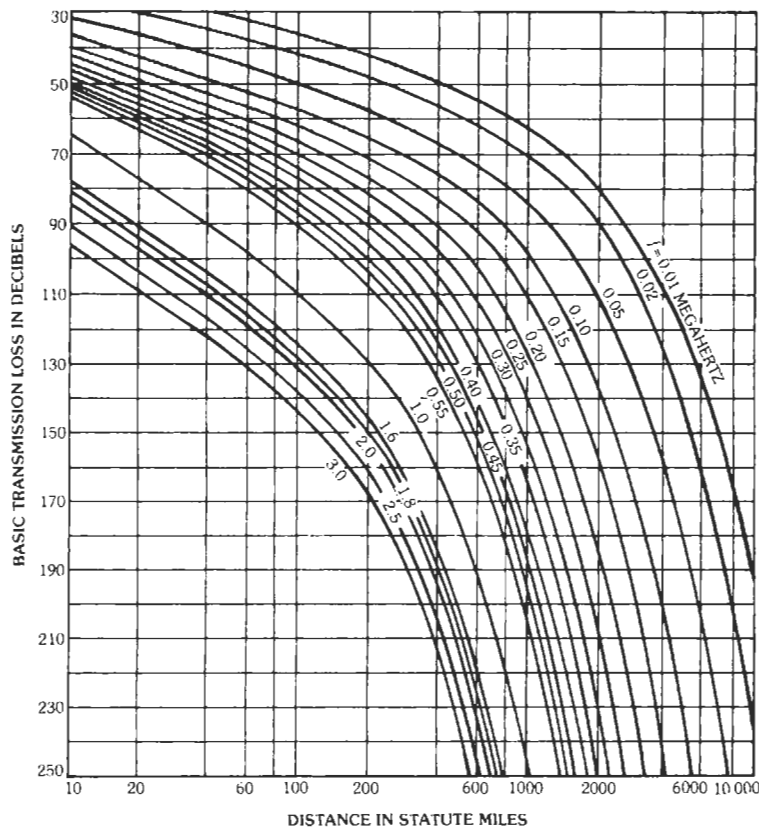


Fig. 3. Basic transmission loss expected for surface waves propagated over a smooth spherical earth. Over land: $\sigma = 0.005$ mho/meter, $\epsilon = 15$. Lossless isotropic antennas 30 feet above the surface. Vertical polarization. (Adapted from K. A. Norton, Transmission Loss in Radio Propagation: II, National Bureau of Standards Technical Note 12, June 1959; Fig. 7.)

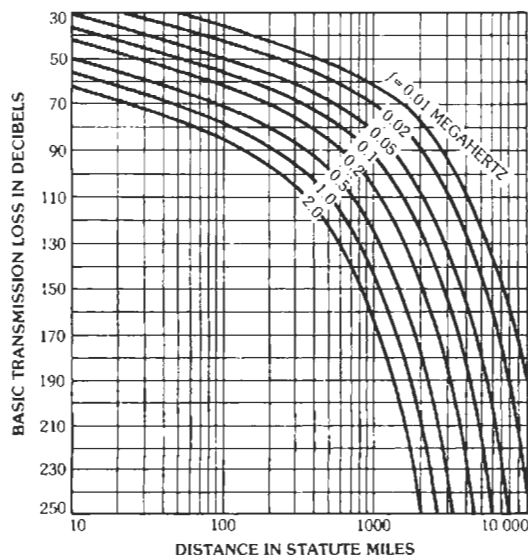


Fig. 4. Basic transmission loss expected for surface waves propagated over a smooth spherical earth. Over sea water: $\sigma = 5$ mhos/meter, $\epsilon = 80$. Lossless isotropic antennas 30 feet above the surface. Vertical polarization. (Adapted from K. A. Norton, Transmission Loss in Radio Propagation: II, National Bureau of Standards Technical Note 12, June 1959; Fig. 8.)

Under the conditions used in Figs. 3 and 4, the surface of the earth behaves like a nearly perfect reflector for the isotropic antennas that are only a small fraction of a wavelength from it and that are used to calculate the basic transmission loss. As a result, each isotropic antenna together with the surface of the earth has a gain of nearly 3.01 dB in the general direction of the horizon. By contrast, a lossless quarter-wave monopole erected over a good ground screen would have a gain of 5.16 dB, and a short lossless monopole would have a gain of 4.77 dB. It follows that the transmission loss between quarter-wave monopoles on the surface of the earth is very nearly $2 \times (5.16 - 3.01) = 4.30$ dB less than the basic transmission loss given in Figs. 3 and 4, and the transmission loss between short monopoles is $2 \times (4.77 - 3.01) = 3.52$ dB less.

Figs. 3 and 4 do not include the effect of sky waves reflected from the ionosphere. Sky waves cause fading at medium distances and produce higher field strengths than the surface wave at longer distances, particularly at night. Sky-wave field strength is subject to diurnal, seasonal, and irregular variations due to changing properties of the ionosphere.

Fig. 5 shows a family of propagation curves for F_0 computed from

$$F_0 = 80.2 - 10 \log D - 0.00176f^{0.26}D$$

TABLE 1. GROUND CONDUCTIVITY AND DIELECTRIC CONSTANT FOR MEDIUM-AND LONG-WAVE PROPAGATION TO BE USED WITH NORTON, BURROWS, BREMMER, OR OTHER DEVELOPMENTS OF SOMMERFELD PROPAGATION EQUATION

Terrain	Conductivity σ (mhos/meter)	Dielectric Constant ϵ (esu)
Sea water	5	80
Fresh water	8×10^{-3}	80
Dry, sandy, flat coastal land	2×10^{-3}	10
Marshy, forested flat land	8×10^{-3}	12
Rich agricultural land, low hills	1×10^{-2}	15
Pastoral land, medium hills and forestation	5×10^{-3}	13
Rocky land, steep hills	2×10^{-3}	10
Mountainous (hills up to 3000 feet)	1×10^{-3}	5
Cities, residential areas	2×10^{-3}	5
Cities, industrial areas	1×10^{-4}	3

In this equation, D = distance in kilometers, f = frequency in kilohertz, and F_0 is the annual median received field strength in decibels above 1 microvolt per meter that would be produced by a short, vertical transmitting dipole at or near the surface of the earth and radiating 1 kilowatt. The empirical equation is based on measured data. Fig. 5 therefore includes the effects of both sky-wave and surface-wave propagation. More recent work has been described by the CCIR.*

The extent to which the lower strata influence the effective ground constants depends on the depth of penetration of the radio energy. This in turn depends on the value of the constants and the frequency. If the depth of penetration is defined as that depth in which

* CCIR XVth Plenary Assembly, Geneva, 1982, Vol. VI, Reports 435-4, 431-3, 432-1, and 575-2.

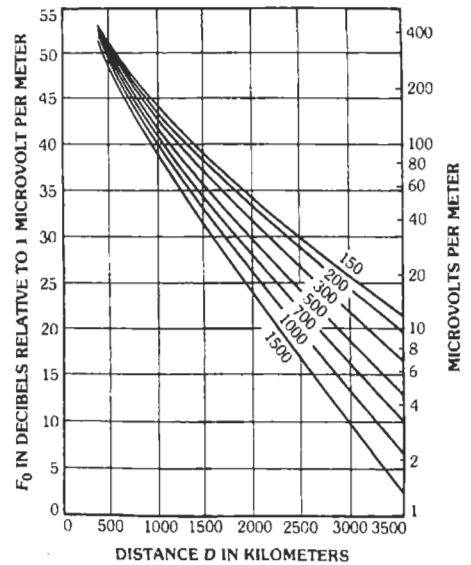


Fig. 5. Family of basic curves of F_0 to be used to determine the annual median value of the field strength for the frequencies (in kilohertz) on the curves. (From CCIR XIIIth Plenary Assembly, Geneva, 1974, Vol. VI, Report 264-3; p. 108.)

the wave has been attenuated to $1/e$ (37%) of its value at the surface, then over the frequency range from 10 kHz to 10 MHz, δ has the values shown in Table 2. It will be seen that, at frequencies of 10 MHz and above, only the surface of the ground need be considered, but at lower frequencies, strata down to a depth of 100 meters or more must be taken into account. It is particularly important to take account of the lower strata when the upper strata are of lower conductivity, since more energy penetrates to the lower levels than happens with an upper layer of higher conductivity.

HIGH FREQUENCIES—3 TO 30 MEGAHERTZ*

At frequencies between about 3 and 25 megahertz and distances greater than about 100 miles, transmis-

* Davies, K., *Ionospheric Radio*. London: Peter Peregrinus Press, 1989. CCIR XVth Plenary Assembly, Geneva, 1982, Vol. VI.

TABLE 2. DEPTH OF PENETRATION OF WAVES INTO THE GROUND

Frequency	Depth δ (m)		
	$\sigma = 4$ mho/m $\epsilon = 80$	$\sigma = 10^{-2}$ mho/m $\epsilon = 10$	$\sigma = 10^{-3}$ mho/m $\epsilon = 5$
10 kHz	2.5	50	150
100 kHz	0.80	15	50
3 MHz	0.14	5	17
10 MHz	0.08	2	9

sion depends chiefly on sky waves reflected from the ionosphere. This is a region high above the surface of the earth where the rarefied air is sufficiently ionized (primarily by ultraviolet sunlight) to reflect or absorb radio waves, such effects being controlled almost exclusively by the free-electron density. The ionosphere is usually considered as consisting of the following layers.

D Layer: At heights from about 50 to 90 kilometers, it exists only during daylight hours, and ionization density corresponds with the elevation of the sun.

This layer reflects very-low- and low-frequency waves, absorbs medium-frequency waves, and weakens high-frequency waves through partial absorption.

E Layer: At a height of about 110 kilometers, this layer is important for high-frequency daytime propagation at distances less than 1000 miles, and for medium-frequency nighttime propagation at distances in excess of about 100 miles. Ionization density corresponds closely with the elevation of the sun. Irregular cloud-like areas of unusually high ionization, called sporadic *E*, may occur up to more than 50 percent of the time on certain days or nights. Sporadic *E* occasionally prevents frequencies that normally penetrate the *E* layer from reaching higher layers and also causes occasional long-distance transmission at very-high frequencies. Some portion (perhaps the major part) of the sporadic-*E* ionization is ascribable to visible- and subvisible-wavelength bombardment of the atmosphere.

F₁ Layer: At heights of about 175 to 250 kilometers, it exists only during daylight. This layer occasionally is the reflecting region for high-frequency transmission, but usually oblique-incidence waves that penetrate the *E* layer also penetrate the *F₁* layer and are reflected by the *F₂* layer. The *F₁* layer introduces additional absorption of such waves.

F₂ Layer: At heights of about 250 to 400 kilometers, *F₂* is the principal reflecting region for long-distance high-frequency communication. Height and ionization density vary diurnally, seasonally, and over the sunspot cycle. Ionization does not follow the elevation of the sun in any fashion, since (at such extremely low air densities and molecular-collision rates) the medium can store received solar energy for many hours, and, by energy transformation, can even detach electrons during the night. At night, the *F₁* layer merges with the *F₂* layer at a height of about 300 kilometers. The absence of the *F₁* layer, and reduction in absorption of the *E* layer, causes nighttime field intensities and noise to be generally higher than during daylight.

As indicated to the right on Fig. 6, these layers are

contained in a thick region throughout which ionization generally increases with height. The layers are said to exist where the ionization gradient is capable of refracting waves back to earth. Obliquely incident waves follow a curved path through the ionosphere due to gradual refraction or bending of the wave front. When attention need be given only to the end result, the process can be assimilated to a reflection.

Depending on the ionization density at each layer, there is a critical or highest frequency, f_c , at which the layer reflects a vertically incident wave. Frequencies higher than f_c pass through the layer at vertical incidence. At oblique incidence, and distances such that the curvature of the earth and ionosphere can be neglected, the maximum usable frequency is given by

$$muf = f_c \sec \phi$$

where,

muf is the maximum usable frequency for the particular layer and distance,

ϕ is the angle of incidence at the reflecting layer.

At greater distances, curvature is taken into account by the modification

$$muf = kf_c \sec \phi$$

where k is a correction factor that is a function of distance and vertical distribution of ionization.

Both f_c and height, and hence ϕ for a given distance, vary for each layer with local time of day, season, latitude, and throughout the 11-year sunspot cycle. The various layers change in different ways with these parameters. In addition, ionization is subject to frequent abnormal variations.

Ionospheric losses are a minimum near the maximum usable frequency and increase rapidly for lower frequencies during daylight.

High frequencies travel from the transmitter to the receiver by reflection from the ionosphere and earth in one or more hops as indicated in Figs. 6 and 7. Additional reflections may occur along the path between the bottom edge of a higher layer and the top edge of a lower layer, the wave finally returning to earth near the receiver.

Fig. 6 indicates transmission on a common frequency, (1) single-hop via *E* layer, Denver to Chicago, and (2) single-hop via *F₂*, Denver to Washington, with (3)

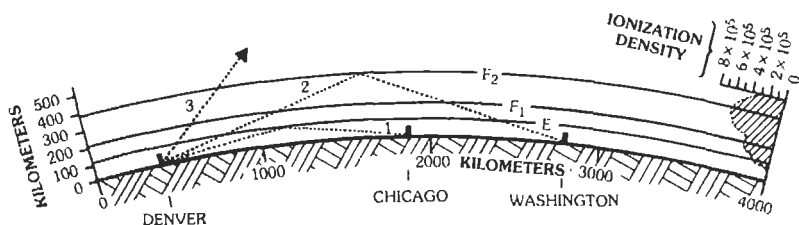
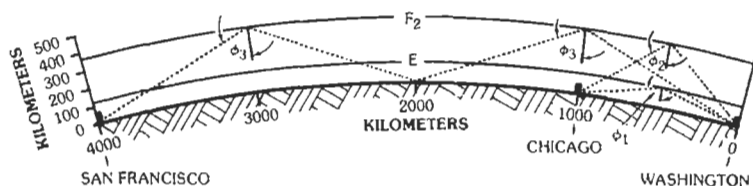


Fig. 6. Schematic explanation of skip-signal zones.

Fig. 7. Single-hop and two-hop transmission paths due to E and F_2 layers.



the wave failing to reflect at higher angles, thus producing a skip region of no signal between Denver and Chicago. Fig. 7 illustrates single-hop transmission, Washington to Chicago, via the E layer (ϕ_1). At higher frequencies over the same distance, single-hop transmission would be obtained via the F_2 layer (ϕ_2). Fig. 7 also shows two-hop transmission, Washington to San Francisco, via the F_2 layer (ϕ_3).

Actual transmission over long distances is more complex than indicated by Figs. 6 and 7, because the layer heights and critical frequencies differ with time (and hence longitude) and with latitude. Further, scattered reflections occur at the various surfaces.

Typical values of critical frequency for Washington, D.C., are shown in Fig. 8.

Preferably, operating frequencies should be selected from a specific frequency band that is bounded above and below by limits that are systematically determinable for the transmission path under consideration. The recommended upper limit is called the *optimum working frequency* (f_{ot}) and is selected below the muf to provide some margin for ionospheric irregularities and turbulence, as well as for the statistical deviation of day-to-day ionospheric characteristics from the predicted monthly median value. So far as may be consistent with available frequency assignments, operation in reasonable proximity to the upper frequency limit is preferable, in order to reduce absorption loss.

The lower limit of the normally available band of frequencies is called the *lowest useful high frequency*

(luf). Below this limit, ionospheric absorption and radio noise levels are likely to be such that radiated-power requirements become uneconomical. For a given path, season, and time, the luf may be predicted by a systematic graphic procedure. Unlike the muf , the predicted luf must be corrected by a series of factors dependent on radiated power, directivity of transmitting and receiving antennas in azimuth and elevation, class of service, and presence of local noise sources. Available data include atmospheric-noise maps, transmission-loss charts, antenna diagrams, and nomograms facilitating the computation. The procedure is formidable but worthwhile.

The upper and lower frequency limits change continuously throughout the day, whereas it is ordinarily impracticable to change operating frequencies correspondingly. Each operating frequency, therefore, should be selected to fall within the above limits for a substantial portion of the daily operating period.

Angles of Departure and Arrival

Angles of departure and arrival are of importance in the design of high-frequency antenna systems. These angles, for single-hop transmission, are obtained from the geometry of a triangular path over a curved earth with the apex of the triangle placed at the virtual height assumed for the altitude of the reflection. Fig. 9 is a family of curves showing radiation angle for different distances.

- D = great-circle distance in statute miles
- H = virtual height of ionosphere layer in kilometers
- Δ = radiation angle in degrees
- ϕ = semiangle of reflection at ionosphere

Forecasts of High-Frequency Propagation

The CCIR publishes "Basic Indices for Ionospheric Propagation" predictions several months in advance in the *Telecommunications Journal* (ITU, Geneva). A list of organizations concerned with issuing forecasts of propagation conditions is published by the CCIR.*

In designing a high-frequency communication cir-

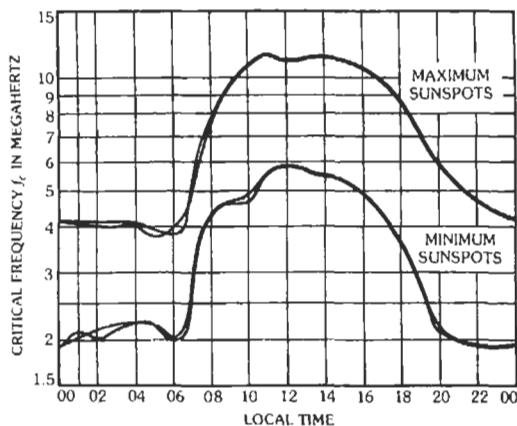


Fig. 8. Critical frequency for Washington, D.C. (From National Bureau of Standards Circular 462.)

* CCIR XVth Plenary Assembly, Geneva, 1982, Vol. VI, Report 313-4.

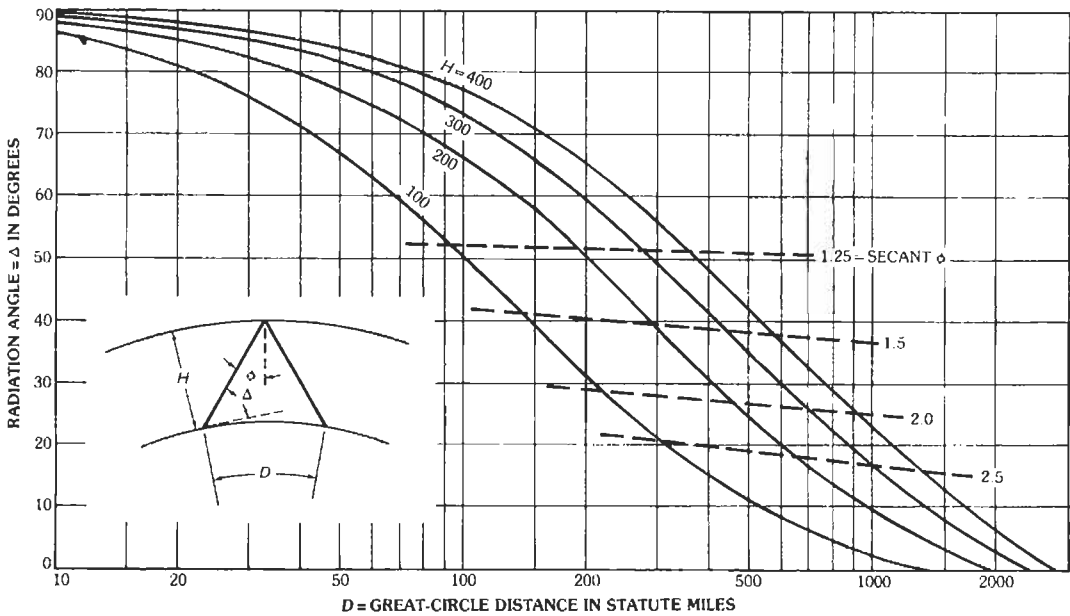


Fig. 9. Single-reflection radiation angle and great-circle distance.

cuit, it is necessary to determine the optimum traffic frequencies, system loss, signal-to-noise ratio, angle of arrival, and circuit reliability. Manual methods for calculating the values of these factors have been described,* as has the use of electronic computers for predicting the performance of high-frequency sky-wave communication circuits. †

Table 3 is a typical performance prediction prepared by computer. A general description of the circuit parameters used in the calculations is shown in the heading of the computer printout. Starting at the top of the page and reading from left to right, the heading may be described as follows.

The first line contains the month, the solar activity level in 12-month moving average Zurich sunspot number, and a circuit identification number. The second and third lines contain the transmitter and receiver locations, the bearings, and the distance. The fourth and fifth lines contain the antenna system for each terminal and their orientation relative to the great-circle path. The minimum angle indicates the lowest vertical angle considered in the mode selection.

The sixth line is the power delivered to the transmitting antenna, the man-made noise level assumed for the area in dBW in a 1-hertz bandwidth at 3 megahertz, and

the hourly median signal-to-noise ratio required to provide the service requested. The signal is in the same units as the transmitter power, and the noise is in a 1-hertz bandwidth. The seventh line contains the heading for the operating frequencies, which are given in megahertz in the eighth line. In addition to the operating frequencies, the eighth line also contains the time heading (GMT) and the classically defined maximum-usable-frequency heading (muf), i.e., frequency which has a 50% probability of having a sky-wave path.

For each time and operating frequency the body of the tabulation contains: (A) the mode having the greatest probability (MODE), (B) the median vertical angle associated with this mode (ANGLE), (C) the propagation time in tenths of milliseconds (DELAY), (D) the percentage of days that any sky-wave mode is expected to exist, circuit probability (C.PROB.), (E) the median of the hourly median signal-to-noise ratios for the days sky-wave modes exist (S/N..DB), and (F) the percentage of days within the month that the median required signal-to-noise ratio is expected to be equalled or exceeded (REL.).

Bandwidth Limitations*

In high-frequency transmission, the communication bandwidth is limited by multipath propagation. The greatest limitation occurs when two or more paths exist with a different number of hops. The bandwidth may

* Davies, K. *Ionospheric Radio Propagation*. National Bureau of Standards Monograph 80, 1 April 1965.

† Lucas, D. L., and Haydon, G. W. *Predicting Statistical Performance Indexes for High-Frequency Ionospheric Telecommunication Systems*. ESSA Technical Report IER1-ITSA, 1 August 1966.

* *Multipath Propagation Over High-Frequency Radio Circuits*. CCIR, Geneva, 1982, Vol. III, Report 203-1.

TABLE 3. SYSTEM PERFORMANCE PREDICTIONS

1		JAN		SSN = 20		CH 5.029		AZIMUTHS		N.MILES						
LONDON		DERRY		TO		CHELTENHAM										
55.00N — 7.31W		168L		70 DEG.		38.75N — 76.85W		280.5 46.3		2880.3				ANT = 0 DB		
RHOMBIC		50H		0 DEG.		MIN. ANGLE = 0 DEG.		OFF AZIMUTH		0 DEG.						
OFF AZIMUTH		0 DEG.		3 MHZ		MAN. NOISE = -154 DBW		REQ. S/N = 61 DB								
PWR = 200.00 KW																
OPERATING FREQUENCIES																
GMT	MUF	3	4	5	6	7	8	10	12	15	17	20	22	25	27	
2	8.6															MODE
		2F	2F	2F	2F	2F	2F	2F	—	—	—	—	—	—	—	ANGLE
		7	5	5	5	5	6	7	—	—	—	—	—	—	—	DELAY
		187	185	185	185	185	186	187	—	—	—	—	—	—	—	C.PROB.
		50	99	99	97	87	66	20	—	—	—	—	—	—	—	S/N..DB
		89	59	71	78	83	86	94	—	—	—	—	—	—	—	REL.
		50	42	81	91	85	66	20	—	—	—	—	—	—	—	
4	8.4															MODE
		2F	2F	2F	2F	2F	2F	2F	—	—	—	—	—	—	—	ANGLE
		8	5	5	5	6	7	7	—	—	—	—	—	—	—	DELAY
		188	185	185	185	186	187	188	—	—	—	—	—	—	—	C.PROB.
		50	99	99	96	84	61	16	—	—	—	—	—	—	—	S/N..DB
		89	58	69	77	82	87	95	—	—	—	—	—	—	—	REL.
		50	34	77	89	82	60	16	—	—	—	—	—	—	—	
6	8.3															MODE
		2F	2F	2F	2F	2F	2F	2F	—	—	—	—	—	—	—	ANGLE
		7	5	5	5	6	7	7	—	—	—	—	—	—	—	DELAY
		187	185	185	185	186	187	187	—	—	—	—	—	—	—	C.PROB.
		50	99	99	94	81	57	10	—	—	—	—	—	—	—	S/N..DB
		89	56	68	76	82	87	97	—	—	—	—	—	—	—	REL.
		50	28	76	89	80	57	10	—	—	—	—	—	—	—	
8	7.3															MODE
		2F	2F	2F	2F	2F	2F	—	—	—	—	—	—	—	—	ANGLE
		7	5	5	5	6	7	—	—	—	—	—	—	—	—	DELAY
		187	185	185	185	186	187	—	—	—	—	—	—	—	—	C.PROB.
		50	99	97	85	59	27	—	—	—	—	—	—	—	—	S/N..DB
		84	56	69	78	83	87	—	—	—	—	—	—	—	—	REL.
		48	29	71	76	56	26	—	—	—	—	—	—	—	—	
10	7.9															MODE
		2F	2F	2F	2X	2X	2F	—	—	—	—	—	—	—	—	ANGLE
		6	4	4	1	1	6	—	—	—	—	—	—	—	—	DELAY
		186	184	184	181	182	186	—	—	—	—	—	—	—	—	C.PROB.
		50	99	99	98	87	47	—	—	—	—	—	—	—	—	S/N..DB
		87	56	70	77	84	87	—	—	—	—	—	—	—	—	REL.
		49	30	75	86	84	46	—	—	—	—	—	—	—	—	
12	14.5															MODE
		2F	3E	4F	3F	2F	2F	2F	2F	2F	2F	—	—	—	—	ANGLE
		4	3	15	10	4	3	3	3	4	4	—	—	—	—	DELAY
		184	181	190	187	184	183	183	183	184	184	—	—	—	—	C.PROB.
		50	99	99	99	99	99	99	90	39	8	—	—	—	—	S/N..DB
		101	5	53	67	75	81	91	93	101	103	—	—	—	—	REL.
		50	0	14	67	84	93	98	90	39	8	—	—	—	—	

then be as small as 100 hertz, but such multipath may be minimized by operating near the muf. Operation at a frequency within approximately 10% of the muf is necessary for paths less than about 600 kilometers to obtain bandwidths greater than, say, 1 kilohertz. The

multipath reduction factor (mrf) is defined as the smallest ratio of muf to operating frequency for which the range of multipath propagation time difference is less than a specified value. The mrf thus defines the frequency above which a specified minimum protection

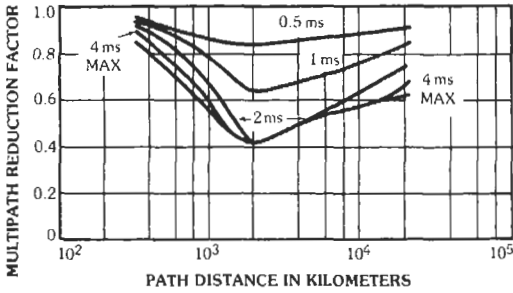


Fig. 10. Multipath reduction factor as a function of path distance. (From Salaman, R. K. "A New Ionospheric Multipath Reduction Factor (MRF)." © 1962, Institute of Radio Engineers.)

against multipath is provided. Fig. 10 shows the mrf for various lengths of path.*

Diversity†

It has been shown that if two or more high-frequency radio channels are sufficiently separated in space, frequency, angle of arrival, time, or polarization, the fading on the various channels is more or less independent. Diversity systems make use of this fact to improve the overall performance, combining or selecting separate radio channels on a single high-frequency circuit.

Satisfactory diversity improvement can be obtained if the correlation coefficient of the fading on the various channels does not exceed about 0.6, and experiments

have indicated that a frequency separation of the order of 400 hertz gives satisfactory diversity performance on long high-frequency paths. Spacing between antennas at right angles to the direction of propagation should be about 10 wavelengths. Polarization diversity has been found to be about equivalent to space diversity in the high-frequency band. Measurements have indicated that times varying from 0.05 to 95 seconds may be necessary to obtain fading correlation coefficients as low as 0.6 in high-frequency time-diversity systems. Angle-of-arrival diversity requires the use of large antennas so as to obtain the required vertical directive characteristics. Differences in the angle of arrival of 2° have been shown to give satisfactory diversity improvement on high-frequency circuits.

Great-Circle Calculations

With reference to Fig. 11, A and B are two places on the surface of the earth the latitudes and longitudes of which are known. In the figure, B = place of greater latitude (nearer the pole), L_A = latitude of A, L_B = latitude of B, and C = difference of longitude between A and B. Angles X and Y at A and B of the great circle passing through the two places and the distance, Z, between A and B along the great circle can be calculated as follows:

$$\tan \frac{1}{2} (Y - X) = \cot \frac{1}{2} C \frac{\sin \frac{1}{2} (L_B - L_A)}{\cos \frac{1}{2} (L_B + L_A)}$$

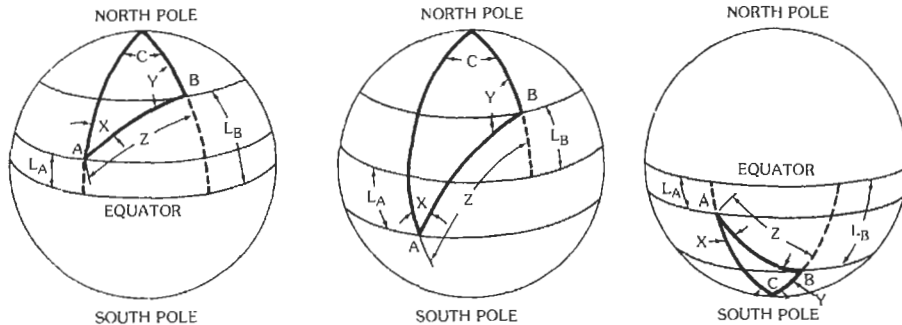
and

$$\tan \frac{1}{2} (Y + X) = \cot \frac{1}{2} C \frac{\cos \frac{1}{2} (L_B - L_A)}{\sin \frac{1}{2} (L_B + L_A)}$$

give the values of $\frac{1}{2} (Y - X)$ and $\frac{1}{2} (Y + X)$ from which

$$\frac{1}{2} (Y + X) + \frac{1}{2} (Y - X) = Y$$

and



(A) Both points located in northern hemisphere.

(B) In opposite hemispheres.

(C) Both points located in southern hemisphere.

Fig. 11. Three globes representing points A and B.

* Salaman, R. K. "A New Ionospheric Multipath Reduction Factor (MRF)." *IRE Transactions on Communication Systems*, Vol. CS-10, June 1962; pp. 220-222.

† "Bandwidth and Signal-to-Noise Ratios in Complete Systems," CCIR, Geneva, 1982, Vol. III, Report 195. Grisdale, G. L., Morriss, J. H., and Palmer, D. S. "Fading of Long-Distance Radio Signals and a Comparison of Space and Polarization Diversity Reception in the 6-18 Mc Range," *Proceedings of the IEE*, Part B, No. 13, Jan. 1957, pp. 39-51.

$$\frac{1}{2}(Y + X) - \frac{1}{2}(Y - X) = X$$

In the above equations, north latitudes are taken as positive and south latitudes as negative. For example, if B is latitude 60° N and A is latitude 20° S,

$$\frac{L_B + L_A}{2} = \frac{60 + (-20)}{2} = \frac{60 - 20}{2} = \frac{40}{2} = 20^\circ$$

$$\frac{L_B - L_A}{2} = \frac{60 - (-20)}{2} = \frac{60 + 20}{2} = \frac{80}{2} = 40^\circ$$

If both places are in the southern hemisphere and $L_B + L_A$ is negative, it is simpler to call the place of greater south latitude B and to use the above method for calculating bearings from true south and to convert the results afterward to bearings east of north.

The distance Z (in degrees) along the great circle between A and B is given by the following:

$$\tan \frac{1}{2} Z = \tan \frac{1}{2} (L_B - L_A) \times [\sin \frac{1}{2} (Y + X)] / [\sin \frac{1}{2} (Y - X)]$$

The angular distance Z (in degrees) between A and B may be converted to linear distance as follows.

Z (in degrees) $\times 111.12 =$ kilometers

Z (in degrees) $\times 69.05 =$ statute miles

Z (in degrees) $\times 60.00 =$ nautical miles

In multiplying, the minutes and seconds of arc must be expressed in decimals of a degree. For example, $Z = 37^\circ 45' 36''$ becomes 37.755° .

Example: Find the great-circle bearings at Brentwood, Long Island, longitude $73^\circ 15' 10''$ W, latitude $40^\circ 48' 40''$ N, and at Rio de Janeiro, Brazil, longitude $43^\circ 22' 07''$ W, latitude $22^\circ 57' 09''$ S; and the great-circle distance in statute miles between the two points. Refer to Chart 1.

Great-circle initial courses and distances are conveniently determined by means of navigation tables such as:

- (A) Navigation Tables for Navigators and Aviators—HO No. 206.
- (B) Dead-Reckoning Altitude and Azimuth Table—HO No. 211.
- (C) Large Great-Circle Charts:

HO Chart No.

- 1280—North Atlantic Ocean
- 1281—South Atlantic Ocean
- 1282—North Pacific Ocean
- 1283—South Pacific Ocean
- 1284—Indian Ocean

The above tables and charts may be obtained at a nominal charge from the United States Navy Department Hydrographic Office, Washington, D.C.

EFFECT OF NUCLEAR EXPLOSIONS ON RADIO PROPAGATION*

Nuclear explosions below an altitude of about 15 kilometers have little effect on radio transmission. However, a detonation occurring at an altitude between 15 and 60 kilometers can produce blackout in the low-frequency, medium-frequency, and high-frequency bands over a radius of several hundred kilometers. This effect lasts only for a few minutes except in an area close to the site of the explosion. In general, it can be said that the effect of nuclear explosions is greatest near the site of the detonation, but the effects of ionization and shock waves do not last longer than a few minutes at distances greater than a few hundred kilometers from the site of the explosion.

High-altitude nuclear explosions (at altitudes greater than about 150 km) increase the electron density in the upper atmosphere by orders of magnitude. As a result, the propagation of signals at frequencies as high as 100 GHz between satellites and from satellites to earth is impaired. The impairments take the form of large attenuation and time delays together with rapid variation of signal strength and frequency-selective fading which can limit the available bandwidth and data rates.

IONOSPHERIC SCATTER PROPAGATION†

This type of transmission permits communication in the frequency range from approximately 30 to 60 megahertz and over distances from about 1000 to 2000 kilometers. It is believed that this type of propagation is due to scattering from the lower D region of the ionosphere and that the useful bandwidth is restricted to less than 10 kilohertz. The greatest use for this type of transmission has been for printing-telegraph channels, particularly in the auroral regions where conventional high-frequency ionospheric transmission is often unreliable.

The median attenuation over paths between 800 and 1000 miles in length is about 80 decibels greater than the free-space path attenuation at 30 megahertz and about 90 decibels greater than the free-space value at 50 megahertz.

* Glasstone, S. *The Effects of Nuclear Weapons*. Washington: US Government Printing Office, 1962.

Middlestead, R. W., et al. "Satellite Crosslink Vulnerability in a Nuclear Environment." *IEEE Journal* on selected areas of communications, Vol. SAC-5, pp. 138-145, February 1987.

Mohanty, N., ed. *Space Communication and Nuclear Scintillation*. New York: Van Nostrand Reinhold, 1991.

† "Ionospheric Scatter Transmission," *Proceedings of the IRE*, Vol. 48, No. 1, 1960; pp. 5-29. CCIR XVth Plenary Assembly, Geneva, 1982, Vol. VI, Report 260-3, and Vol. III, Report 109-2.

CHART 1. EXAMPLE OF GREAT-CIRCLE CALCULATIONS

	Longitude	Latitude	
Brentwood	73°15'10"W	40°48'40"N	L_B
Rio de Janeiro	43°22'07"W	(-)22°57'09"S	L_A
	29°53'03"		C
		17°51'31"	$L_B + L_A$
		63°45'49"	$L_B - L_A$
$\frac{1}{2} C = 14°56'31"$	$\frac{1}{2} (L_B + L_A) = 8°55'45"$	$\frac{1}{2} (L_B - L_A) = 31°52'54"$	
$\log \cot 14°56'31" = 10.57371$		$\log \cot 14°56'31" = 10.57371$	
plus $\log \cos 31°52'54" = 9.92898$		plus $\log \sin 31°52'54" = 9.72277$	
	0.50269		0.29648
minus $\log \sin 8°55'45" = 9.19093$		minus $\log \cos 8°55'45" = 9.99471$	
$\log \tan \frac{1}{2} (Y + X) = 1.31176$		$\log \tan \frac{1}{2} (Y - X) = 0.30177$	
$\frac{1}{2} (Y + X) = 87°12'26"$		$\frac{1}{2} (Y - X) = 63°28'26"$	
Bearing at Brentwood = $\frac{1}{2} (Y + X) + \frac{1}{2} (Y - X) = 150°40'52"$ East of North			
Bearing at Rio de Janeiro = $\frac{1}{2} (Y + X) - \frac{1}{2} (Y - X) = X = 23°44'00"$ West of North			
$\frac{1}{2} (L_B - L_A) = 31°52'54"$		$\log \tan 31°52'54" = 9.79379$	
$\frac{1}{2} (Y + X) = 87°12'26"$		plus $\log \sin 87°12'26" = 9.99948$	
			9.79327
$\frac{1}{2} (Y - X) = 63°28'26"$		minus $\log \sin 63°28'26" = 9.95170$	
		$\log \tan \frac{1}{2} Z = 9.84157$	
		$\frac{1}{2} Z = 34°46'24"$	$Z = 69°32'48"$

$$69°32'48" = 69.547°$$

$$\text{Linear distance} = 69.547 \times 69.05 = 4802 \text{ statute miles}$$

METEOR-BURST PROPAGATION*

Frequencies in the very-high- and ultrahigh-frequency bands may be propagated by reflection from columns of ionization produced by meteors entering the lower E region. Experimental single-channel two-way telegraph circuits have been operated in the frequency range from 30 to 40 megahertz over distances of 600 to 1300 kilometers with transmitter powers of 1 to 3 kilowatts. One-way transmission of voice and facsimile have also been made with transmitter powers of 1 kilowatt and 20 kilowatts, respectively.

The frequency range from about 50 to 80 megahertz has been found best suited for meteor-burst transmission.

* "Communication by Meteor-Burst Propagation," CCIR XVth Plenary Assembly, Geneva, 1982, Vol. VI, Report 251-3. Oetting, J. D., "An Analysis of Meteor Burst Communications for Military Applications," *IEEE Transactions on Communications*, Vol. Com-28, Sept. 1980.

PROPAGATION ABOVE 30 MEGAHERTZ, LINE-OF-SIGHT CONDITIONS*

Radio Refraction†

Under normal propagation conditions, the refractive index of the atmosphere decreases with height so that radio rays travel more slowly near the ground than at higher altitudes. This variation in velocity with height

* Bullington, K., "Radio Propagation at Frequencies Above 30 Megacycles," *Proceedings of the IRE*, Vol. 35, October 1947, pp. 1122-1136. Kerr, D. E., *Propagation of Short Radio Waves*. New York: McGraw-Hill Book Co., 1951. CCIR XVth Plenary Assembly, Geneva, 1982, Vol. V.

† Bean, B. R., and Dutton, E. J. *Radio Meteorology*. Monograph 92, Institute for Telecommunication Sciences and Aeronomy, Environmental Science Services Administration (ESSA). Washington: Superintendent of Documents, 1966.

results in bending of the radio rays. Uniform bending may be represented by straight-line propagation, but with the radius of the earth modified so that the relative curvature between the ray and the earth remains unchanged. The new radius of the earth is known as the effective earth radius, and the ratio of the effective earth radius to true earth radius is usually denoted by K . The average value of K in temperate climates is about 1.33; however, values from about 0.6 to 5.0 are to be expected.

The decrease in the refractive index with height may at times be so great that the ray is bent down with a radius equal to that of the earth so that the earth may then be considered to be flat. A further increase in the refractive-index gradient results in the radio ray being bent down sufficiently to be reflected from the earth. The ray then appears to be trapped in a duct between the earth and the maximum height of the radio path.

Under certain atmospheric conditions, the refractive index may increase with height, causing the radio rays to bend upward. Such inverse bending results in a decrease in path clearance on line-of-sight paths.

The distance to the radio horizon over smooth earth, when the height, h , is very small compared with the radius of the earth, is given with a good approximation by the expression

$$d = (3Kh/2)^{1/2}$$

where,

- h = height in feet above the earth,
- d = distance to radio horizon in miles,
- K = ratio of the effective to the true radius of the earth.

Over a smooth earth, a transmitter antenna at height h_t (feet) and a receiving antenna at height h_r (feet) are in radio line-of-sight provided the spacing in miles is less than $(2h_t)^{1/2} + (2h_r)^{1/2}$ (assuming $K = 1.33$).

The nomogram in Fig. 12 gives the radio-horizon distance between a transmitter at height h_t and a receiver at height h_r . Fig. 13 extends the first nomogram to give the maximum radio-path length between

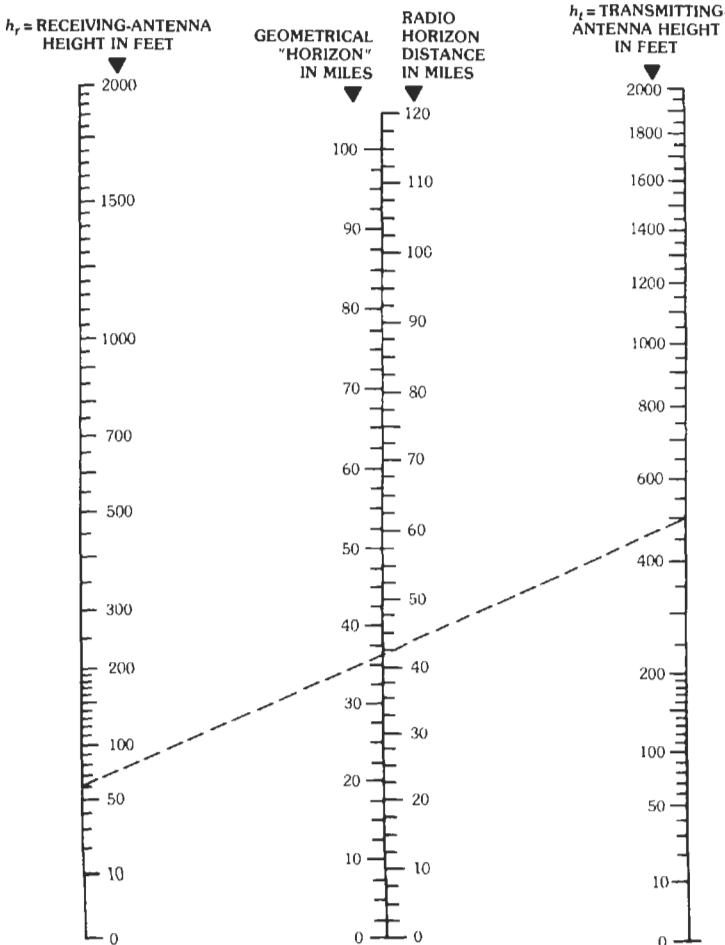


Fig. 12. Nomogram giving radio-horizon distance in miles when h_t and h_r are known. Example shown: Height of receiving antenna 60 feet; height of transmitting antenna 500 feet; maximum radio-path length = 41.5 miles. ($K = 1.33$)

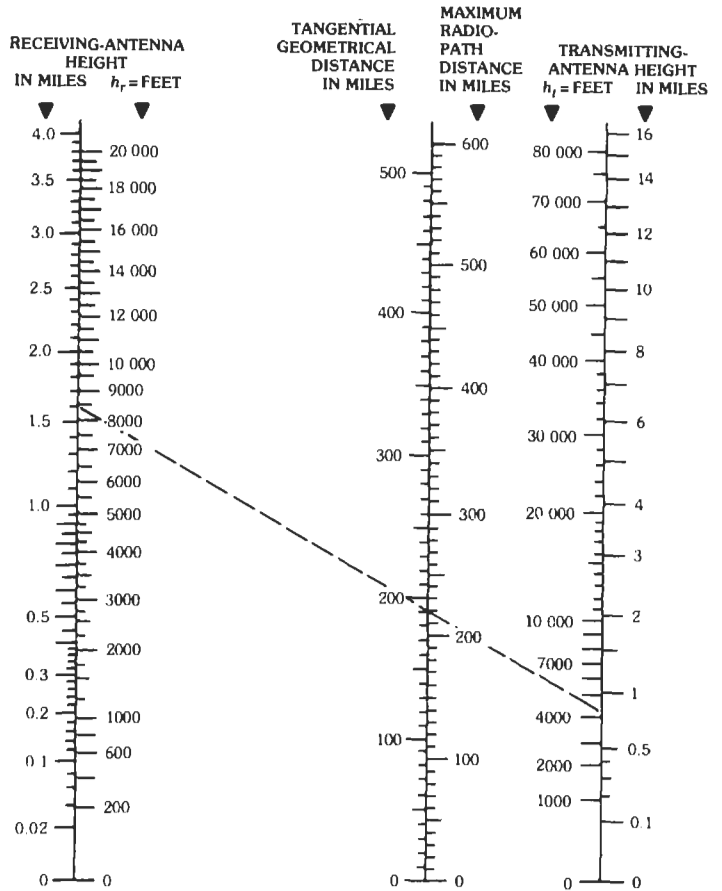


Fig. 13. Nomogram giving radio-path length and tangential distance for transmission between two airplanes at heights h_r and h_t . Example shown: Height of receiving-airplane 8500 feet (1.6 miles); height of transmitting-airplane 4250 feet (0.8 mile); maximum radio-path distance = 220 miles. ($K = 1.33$)

two airplanes whose altitudes are known. Both figures assume a value of $K = 1.33$.

Path Plotting and Profile-Chart Construction

Path Plotting—When laying out a microwave system, it is usually convenient to plot the path on a profile chart. Such charts are scaled to indicate the departure of the curvature of the earth from a straight line. With reference to Fig. 14,

$$D^2 + R^2 = (h + R)^2 = h^2 + 2Rh + R^2$$

$$D^2 = h^2 + 2Rh$$

where,

- D = distance,
- R = radius of earth (3960 miles),
- h = altitude.

Since $h \ll R$, $D = (2Rh)^{1/2}$, and inserting the true earth radius with R and D in statute miles and h in feet

$$D = \left(\frac{2 \times 3960}{5280} h \right)^{1/2}$$

$$D = [(3/2)h]^{1/2}$$

$$h = (2/3)D^2$$

for true earth, where D is in miles and h in feet. For a value of $K = 1.33$

$$D = [(3/2)h]^{1/2} (4/3)^{1/2} = (2h)^{1/2}$$

$$h = D^2/2$$

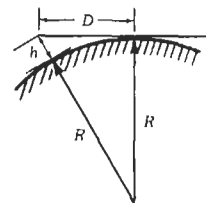


Fig. 14. Straight line tangent to surface of earth.

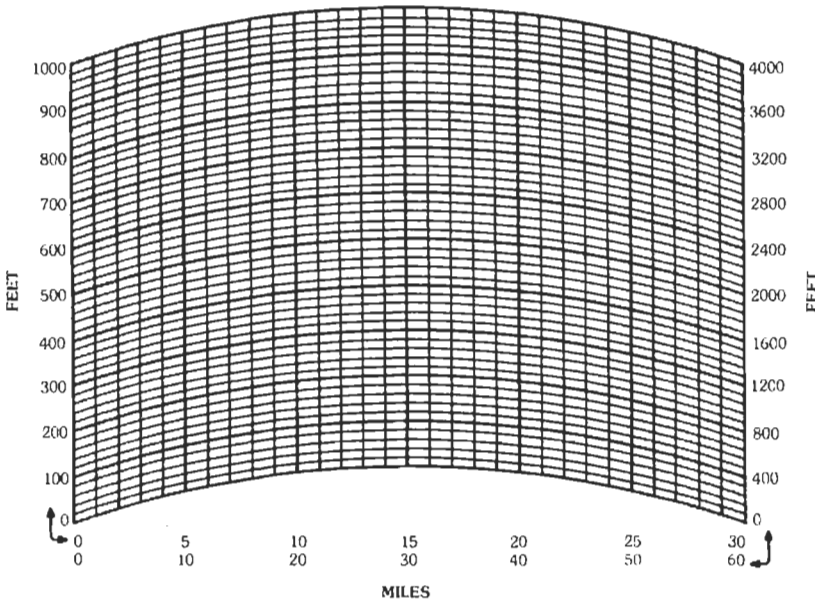


Fig. 15. Typical 4/3-earth profile paper, 1000-foot scale.

Or more generally

$$h = 2D^2/3K$$

Profile Paper—If a 4/3 effective-radius factor is used, the departure from a horizontal tangent line is

$$h = D^2/2$$

where symbols are as above. By using this equation, a template can be made for convenient drawing of profile paper (Fig. 15). For instance, if the horizontal scale is 10 miles/inch, the vertical scale 100 feet/inch, and a width corresponding to 40 miles is desired, the points in Table 4 may be plotted.

A typical example of a template constructed according to these figures is given in Fig. 16. If a different scale is desired than is provided on available profile-chart paper (for example, if a 50-mile hop is to be plotted on 30-mile paper), then the scale of miles may be doubled to extend the range of the paper to 60 miles. The vertical scale in feet must then be quadrupled; i.e.,

100-foot divisions become 400-foot divisions, as on the right in Fig. 15.

Fresnel Zones

The Fresnel-Kirchhoff theory was originally developed to account for the diffraction of light when obstructed by diaphragms, and when transmitting through apertures of various shapes and sizes. This theory may be applied to radio and sound waves and is based on the concept that any small element of space in the path of a wave may be considered as the source of a secondary wavelet, and that the radiated field can be built up by the superposition of all these wavelets (Huygens principle).

Consider a transparent screen between a distant transmitter, *T*, and a receiver, *R*, with the distance from screen to transmitter being at least 10 times the distance from screen to receiver, and with the plane of the screen perpendicular to direction *T-R*. Concentric circles may be drawn on this screen, with the centers at the point where line *T-R* intersects the screen at *O*, the radius of the first circle being such that the difference in length between path *O-R* and the path from the circumference of this circle to *R* is $\frac{1}{2} \lambda$. The radii of the other circles are such that the corresponding path-length differences are integral multiples of $\frac{1}{2} \lambda$. The radius of the first circle is $(d\lambda)^{1/2}$, where *d* is the distance from *O* to *R*, and the radius of the second circle is $(2d\lambda)^{1/2}$, of the third $(3d\lambda)^{1/2}$, etc. The area within the first circle is called the first Fresnel zone, and the other ring-shaped areas are the second, third, etc., Fresnel zones. The fields from the odd-number zones are in phase at *R*, and the fields from the even-number zones are also in phase at *R* but are opposite in phase to the fields from the

TABLE 4. POINTS FOR CONSTRUCTING TEMPLATE

Distance From Center (Horizontal)	Distance From Horizontal (Vertical)
0 miles = 0 inches	0 feet = 0 inches
5 miles = 1/2 inch	12 1/2 feet = 1/8 inch
10 miles = 1 inch	50 feet = 1/2 inch
15 miles = 1 1/2 inches	112 1/2 feet = 1 1/8 inches
20 miles = 2 inches	200 feet = 2 inches

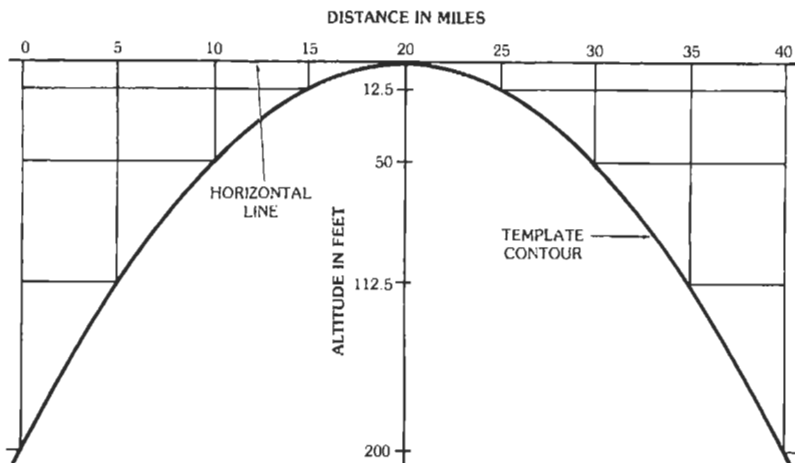


Fig. 16. Construction of a template for profile charts.

odd-number zones. It can be shown that the effect at R of each zone is nearly equal. If an infinitely absorbing screen is provided with an aperture of the same diameter as the first Fresnel zone, it will be found that the field at R is twice as great as the unobstructed or free-space field. If the aperture is increased to include the second zone, the field at R will then be nearly zero, since the fields from zones 1 and 2 are nearly equal in amplitude and opposite in phase. With a continued increase in the diameter of the aperture, further maxima and minima appear; the amplitude of these oscillations decreases very gradually until eventually the field at R approaches the free-space value, which is half that due to the first Fresnel zone. If the distance from the screen to the transmitter is d_1 , and from the screen to the receiver is d_2 , then the general expression for the radius of the n th Fresnel zone is

$$\{n\lambda[(d_1 \times d_2)/(d_1 + d_2)]\}^{1/2}$$

Required Path Clearance

A criterion to determine whether the earth is sufficiently removed from the radio line-of-sight ray to allow mean free-space propagation conditions to apply is to have the first Fresnel zone clear all obstacles in the path of the rays. This first zone is bounded by points for which the transmission path from transmitter to receiver is greater by one-half wavelength than the direct path. Let d be the length of the direct path and d_1 and d_2 be the distances to the transmitter and receiver from a point, P . The radius of the first Fresnel zone at P is approximately given by

$$R_1^2 = \lambda(d_1 d_2 / d)$$

where all quantities are expressed in the same units.

If d is in miles and frequency F in megahertz, the first Fresnel-zone radius in feet at P is given by

$$R_{1m} = 2280(d_1 d_2 / Fd)^{1/2}$$

The maximum occurs when $d_1 = d_2$ and is equal to

$$R_{1m} = 1140(d/F)^{1/2}$$

While a fictitious earth of $4/3$ of true earth radius is generally accepted for determining first Fresnel-zone clearance under normal refraction conditions, unusual conditions that occur in the atmosphere may make it desirable to allow first Fresnel clearance of an effective earth radius of 0.7 to 0.5 of the true radius.

Fig. 17 shows the effect of path clearance on radio transmission.*

Interference Between Direct and Reflected Rays

Where there is one reflected ray combining with the direct ray at the receiving point (Fig. 18), the resulting field strength (neglecting the difference in angles of arrival, and assuming perfect reflection at T) is related to the free-space intensity, irrespective of the polarization, by

$$E = 2E_d \sin 2\pi(\delta/2\lambda)$$

where,

- E = resulting field strength, } same
- E_d = direct-ray field strength, } units
- δ = geometrical length difference between direct and reflected paths, which is given to a close approximation by $\delta = 2h_{at}h_{ar}/d$,

where h_{at} and h_{ar} are the heights of the antennas above a

* Bullington, K. "Radio Propagation Fundamentals," *Bell System Tech. J.*, Vol. 36, No. 3, 1957, pp. 593-626.

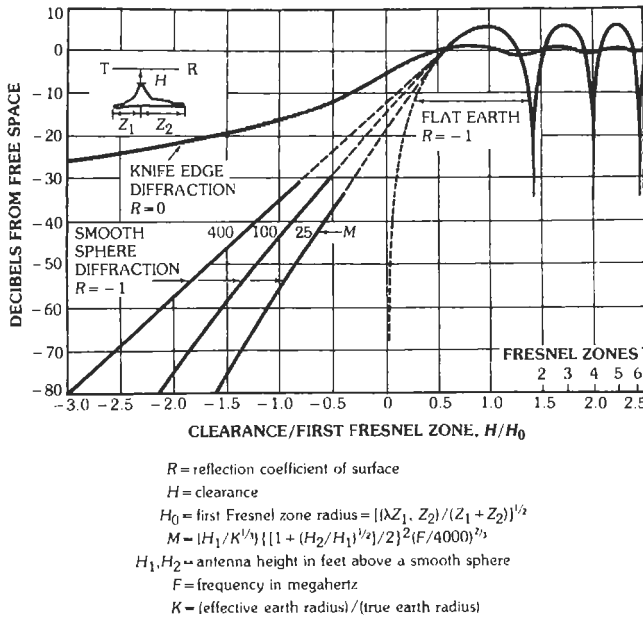


Fig. 17. Effect of path clearance on radio transmission. (From Bullington, K. "Radio Propagation Fundamentals." Bell System Tech. J., Vol. 36, No. 3, Fig. 8, ©1957 American Telephone and Telegraph Co.)

reflecting plane tangent to the effective earth. (See Fig. 18.)

The following cases are of interest:

$$E = 0 \quad \text{for } h_{at} h_{ar} = d\lambda/2$$

$$E = 2E_d \quad \text{for } h_{at} h_{ar} = d\lambda/4$$

$$E = E_d \quad \text{for } h_{at} h_{ar} = d\lambda/12$$

In case $h_{at} = h_{ar} = h$,

$$E = 0 \quad \text{for } h = (d\lambda/2)^{1/2}$$

$$E = 2E_d \quad \text{for } h = (d\lambda/4)^{1/2}$$

$$E = E_d \quad \text{for } h = (d\lambda/12)^{1/2}$$

All these equations are written with the same units for all quantities.

Space-Diversity Reception

When h_{ar} is varied, the field strength at the receiver varies approximately according to the preceding equa-

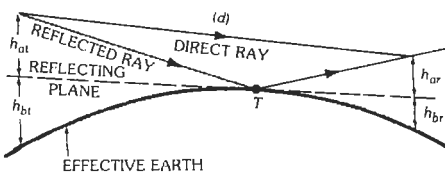


Fig. 18. Interference between direct and reflected rays.

tion. The use of two antennas at different heights provides a means of compensating to a certain extent for changes in electrical-path differences between direct and reflected rays (space-diversity reception).

The antenna spacing at the receiver should be approximately such as to give a $\lambda/2$ variation between geometrical-path differences in the two cases. An approximate value of the spacing is given by $\lambda d / 4h_{at}$ when all quantities are in the same units.

The spacing in feet for d in miles, h_{at} in feet, λ in centimeters, and f in megahertz is given by

$$\text{spacing} = 43.4\lambda d / h_{at}$$

$$= 1.3 \times 10^6 d / f h_{at}$$

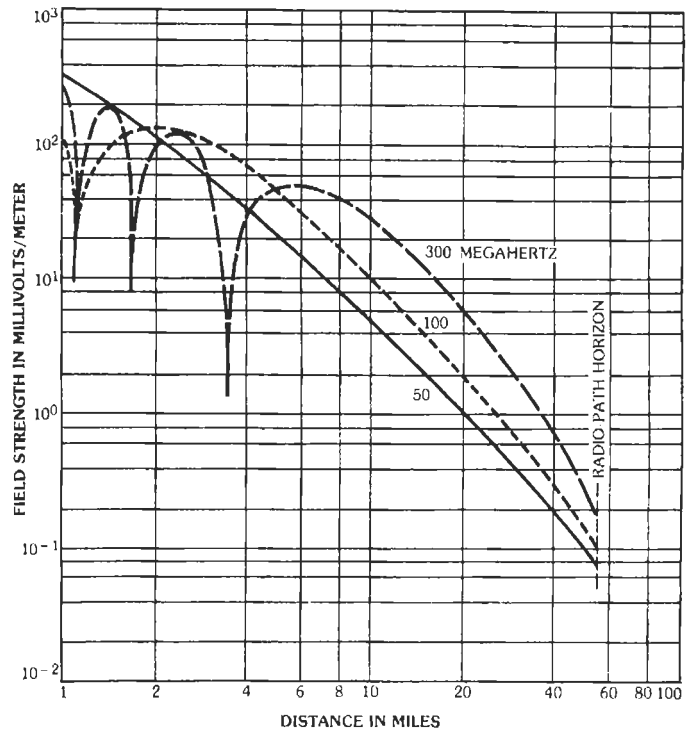
Example: $\lambda = 3$ centimeters, $d = 20$ miles, and $h_{at} = 50$ feet; therefore spacing = 52 feet.

Variation of Field Strength With Distance

Fig. 19 shows the variation of resultant field strength with distance and frequency; this effect is due to interference between the free-space wave and the ground-reflected wave as these two components arrive in or out of phase.

To compute the field accurately under these conditions, it is necessary to calculate the two components separately and to add them in the correct phase relationship. The phase and amplitude of the reflected ray are determined by the geometry of the path and the change in magnitude and phase at ground reflection. For horizontally polarized waves, the reflection coefficient can be taken as approximately 1, and the phase shift at

Fig. 19. Variation of resultant field strength with distance and frequency. Antenna heights: 1000 feet, 30 feet; power: 1 kilowatt; ground constants: $\sigma = 5 \times 10^{-14}$ emu, $\epsilon = 1.5$ esu; polarization: horizontal.



reflection as 180 degrees, for nearly all types of ground and angles of incidence. For vertically polarized waves, the reflection coefficient and phase shift vary appreciably with the ground constants and angle of incidence. (See Fig. 24 of Chapter 32.)

Measured field strengths usually show large deviations from point to point because of reflections from ground irregularities, buildings, trees, etc.

For transmission paths of the order of 30 miles and for frequencies up to about 6000 megahertz, good engineering practice should allow for possible increases of signal strength of +10 decibels with respect to free-space propagation and should allow a fading margin depending on the degree of reliability desired in accordance with the following:

- 10 decibels—90 percent
- 20 decibels—99 percent
- 30 decibels—99.9 percent
- 40 decibels—99.99 percent

Fading and Diversity*

Line-of-sight propagation at ultrahigh frequencies is affected both by signal-strength variation due to multipath transmission and by bending of the beam due to abnormal variation of refractive index with height in the lower atmosphere.

At frequencies below about 8000 megahertz, and on

paths having adequate clearance, the fading on line-of-sight paths is due to multipath transmission. Multipath fading may be divided into two main types; the first is relatively rapid and is caused by interference between two or more rays arriving by slightly different paths; this is known as *atmospheric-multipath*. The second type of fading is less rapid and is due to interference between direct and reflected rays; this is referred to as *reflection-multipath*. In general, the number of fades per unit time due to atmospheric-multipath increases with path length; however, the duration of a fade of a given depth tends to decrease with increasing path length. Fig. 20 shows the typical fading characteristics of a terrestrial line-of-sight path. See Chapter 27 for earth-to-space paths.

Either frequency or space diversity may be used to reduce the amplitude of multipath fading. In the case of atmospheric-multipath fading on line-of-sight paths, it

* Bullington, K., "Radio Propagation Fundamentals," *Bell System Technical Journal*, Vol. 36, No. 3, 1957; pp. 593-626. Pearson, K. W., "Method for the Prediction of the Fading Performance of a Multisection Microwave Link," *Proceedings of the IEE*, Vol. 112, No. 7, July 1965; pp. 1291-1300. CCIR XVth Plenary Assembly, Geneva, 1982, Vol. V, Report 338-4; pp. 279-314.

Transmission Loss Predictions for Tropospheric Communication Circuits. National Bureau of Standards Technical Note No. 101.

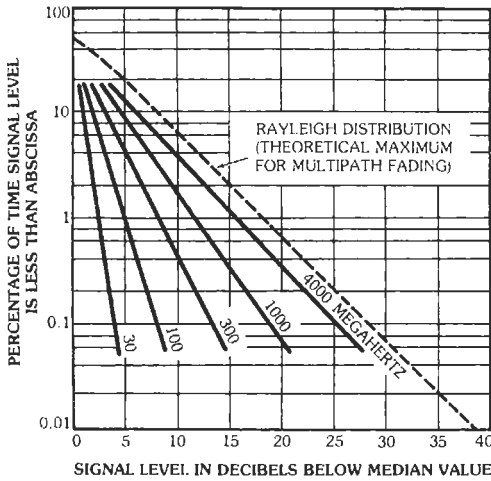


Fig. 20. Typical fading characteristics in the worst month on line-of-sight paths of 30 to 40 miles with clearance of 50 to 100 feet. (From Bullington, K. "Radio Propagation Fundamentals." Bell System Technical Journal, Vol. 36, No. 3, Fig. 4, ©1957 American Telephone and Telegraph Co.).

has been found that considerable diversity improvement can usually be obtained with a frequency difference of 100 to 200 megahertz or with a vertical antenna spacing of between 100 and 200 wavelengths.

Atmospheric Absorption

Gaseous Attenuation*—Oxygen and water vapor may absorb energy from a radio wave by virtue of the permanent electric dipole moment of the water molecule and the permanent magnetic dipole moment of the oxygen molecule. Fig. 21 shows water-vapor absorption γ_w and oxygen absorption γ_o as functions of frequency.

For a path through the atmosphere, the total attenuation depends on the scale heights of the oxygen and water-vapor components, which are approximately 8 km and 2 km, respectively, and on the elevation angle, θ . Then the total attenuation, A , in decibels is given by

$$A = (8\gamma_o + 2\gamma_w)/\sin \theta \quad \theta > 10^\circ$$

or by

$$A = 16\gamma_o(\sin^2 \theta + 16/R + \sin \theta)^{-1/2} + 4\gamma_w(\sin^2 \theta + 4/R + \sin \theta)^{-1/2} \quad \theta \leq 10^\circ$$

where R is the effective earth radius, including refraction, which is approximately 8500 km.

Rain Attenuation†—In addition to gaseous absorption, which is relatively small except in the narrow

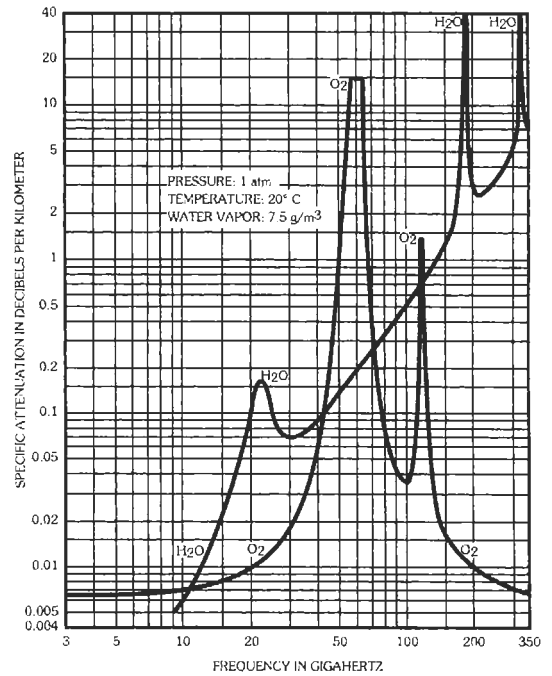


Fig. 21. Specific attenuation due to atmospheric gases (reference 4).

spectral regions shown in Fig. 21, rain attenuation must also be considered. Specific attenuation, γ_R (decibels per kilometer of path length), depends on the frequency f (GHz) and rain rate R_p (mm/hr), and is given by

$$\gamma_R = a(f) \cdot R_p^{b(f)}$$

where $a(f)$ and $b(f)$ are given in Table 5.* Fig. 22 shows this dependence for terrestrial paths.

Rain rate depends on the location of the path and on the season. For computational purposes, the world is divided into the various rain-rate regions shown in Fig. 23 and, in more detail for part of North America, in Fig. 24. The fraction of an average year during which a given rain rate is exceeded in each of these regions is shown in Table 6. In practice, annual statistics of rain rates can be misleading because of seasonal dependency and because the occurrence of rain is not random. Examples of the frequency of rain events and their duration as a function of rain rate for the United Kingdom are shown in Fig. 25. Similar data for selected locations in the USA are shown in Tables 6, 7, 8, and 9.† In temperate climates, rain rates exceeding 2.54 mm/h are not likely to occur over an area larger than about 6 km in diameter.

* Reference 4.
† Reference 3.

* Reference 2.
† Reference 1. This reference also includes other detailed rain-rate duration data.

TABLE 5. PARAMETERS FOR COMPUTING SPECIFIC ATTENUATION:
 $\gamma_R = aR^b$, 0°C, LAWS AND PARSON DISTRIBUTION*

Frequency <i>f</i> (GHz)	Multiplier <i>a</i> (<i>f</i>)	Exponent <i>b</i> (<i>f</i>)
1	0.00015	0.95
4	0.00080	1.17
5	0.00138	1.24
6	0.00250	1.28
7.5	0.00482	1.25
10	0.0125	1.18
12.5	0.0228	1.145
15	0.0357	1.12
17.5	0.0524	1.105
20	0.0699	1.10
25	0.113	1.09
30	0.170	1.075
35	0.242	1.04
40	0.325	0.99
50	0.485	0.90
60	0.650	0.84
70	0.780	0.79
80	0.875	0.753
90	0.935	0.730
100	0.965	0.715

* Reference 2.

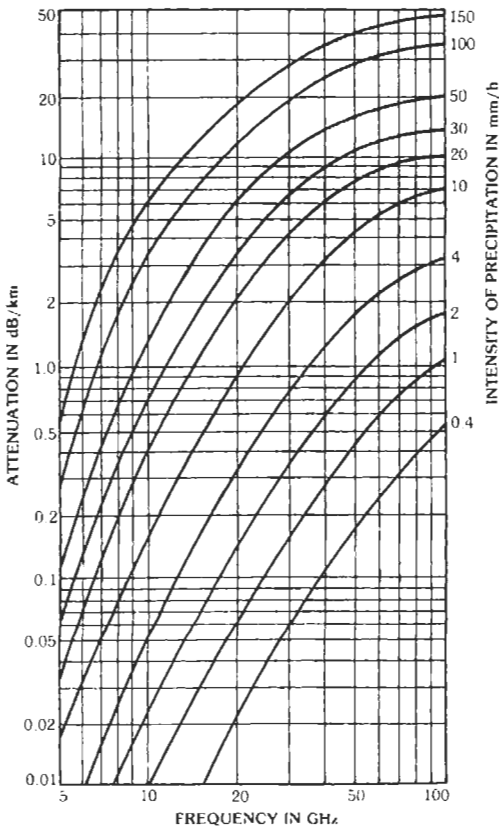


Fig. 22. Attenuation due to precipitation. (From CCIR XIIIth Plenary Assembly, Geneva, 1974, Vol. V, Report 233-3.)

Given the specific attenuation γ_R , the path attenuation, *A*, in decibels is provided by

$$A = \gamma_R L r$$

where,

L is the path length (km),

r is a path length modification factor which takes into account that the rain rate may not be uniform along the path.

The factor *r* is*

$$r = \gamma(L) \cdot R_p^{-\delta(L)}$$

where *R* is the point rain rate in millimeters per hour and the coefficients $\gamma(L)$ and $\delta(L)$ are shown in Figs. 26 and 27. The resulting dependence of *r* on the point rainfall rate and path lengths is shown in Fig. 28.

Free-Space Transmission Equations

If the incoming wave is a plane wave having a power flow per unit area equal to P_0 , the available power at the output terminals of a receiving antenna may be expressed as

$$P_r = A_r P_0$$

where A_r is the effective area of the receiving antenna.

* Reference 3.

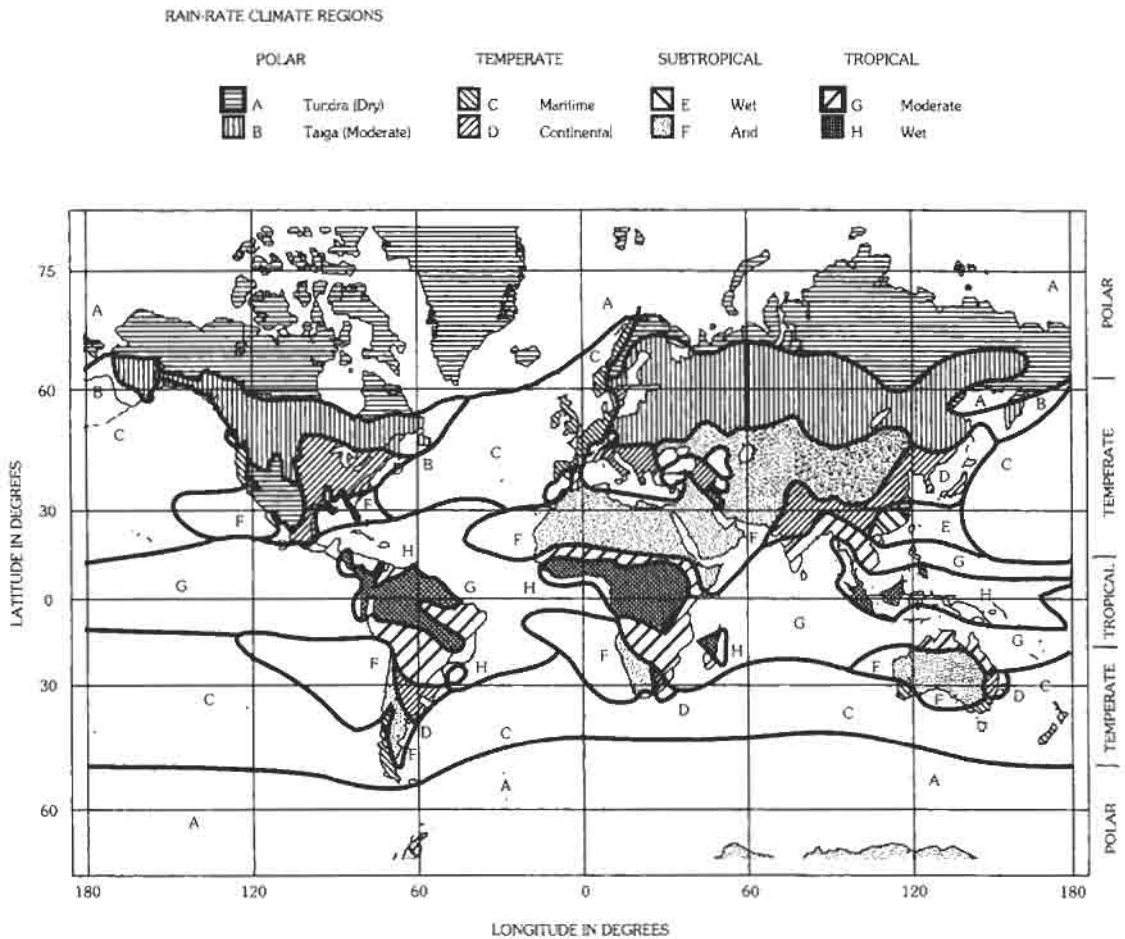


Fig. 23. Global rain-rate climate regions including the ocean areas.

The free-space path attenuation is given by

$$\text{attenuation} = 10 \log (P_i/P_r)$$

where P_i is the power radiated from the transmitting antenna (same units as for P_r). Then

$$P_r/P_i = A_r A_t / d^2 \lambda^2$$

where,

- A_r = effective area of receiving antenna,
- A_t = effective area of transmitting antenna,
- λ = wavelength,
- d = distance between antennas.

The length and surface units in the equation should be consistent. This is valid provided $d \gg 2a^2/\lambda$, where a is the largest linear dimension of either of the antennas.

Path attenuation between isotropic antennas is

$$P_i/P_r = 4.56 \times 10^3 f^2 d^2$$

where f is in megahertz and d is in miles.

Path attenuation α (in decibels) is

$$\alpha = 36.6 + 20 \log f + 20 \log d$$

A nomogram for determining the solution of α is given in Fig. 29.

Effective Areas of Typical Antennas*

Hypothetical isotropic antenna (no heat loss)

$$A = (1/4\pi)\lambda^2 \approx 0.08\lambda^2$$

* Refer to Chapter 32.

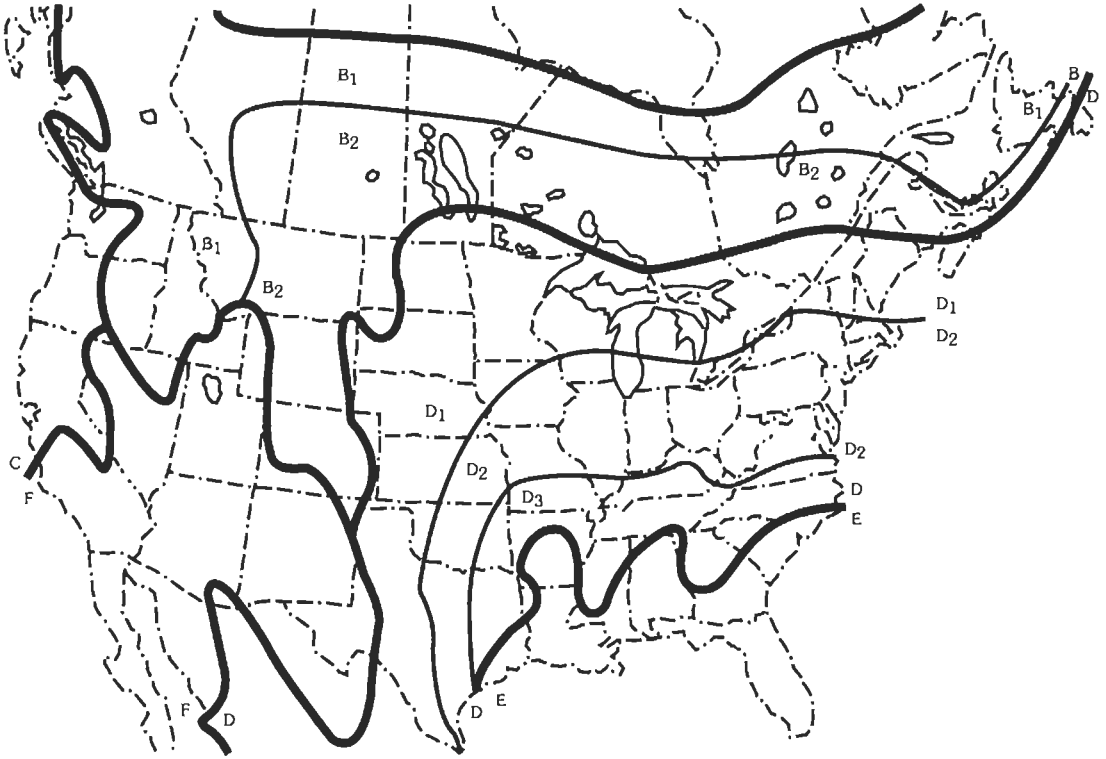


Fig. 24. Rain-rate climate regions for the continental US and southern Canada.

Small uniform-current dipole, short compared with wavelength (no heat loss)

$$A = (3/8\pi)\lambda^2 \approx 0.12\lambda^2$$

Half-wavelength dipole (no heat loss)

$$A \approx 0.13\lambda^2$$

Parabolic reflector of aperture area S (here, the factor 0.54 is due to nonuniform illumination of the reflector)

$$A \approx 0.54S$$

Very long horn with small aperture dimensions compared with length

$$A = 0.81S$$

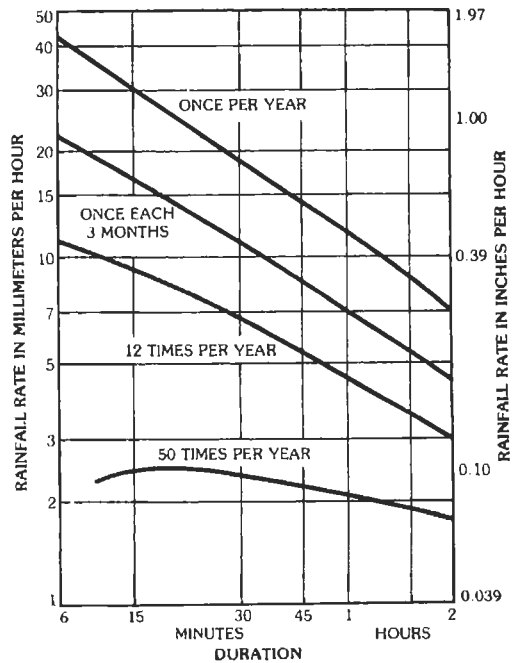


Fig. 25. Rainfall duration in England. (From Bilham, E. G. "Climate of British Isles." Toronto: Macmillan Company, 1938.)

TABLE 6. POINT RAIN-RATE DISTRIBUTION VALUES (mm/h) VERSUS PERCENT OF YEAR RAIN RATE IS EXCEEDED

Percent of Year	Rain Climate Region												Minutes per Year	Hours per Year
	A	B ₁	B	B ₂	C	D ₁	D = D ₂	D ₃	E	F	G	H		
0.001	28.5	45	57.5	70	78	90	108	126	165	66	185	253	5.26	0.09
0.002	21	34	44	54	62	72	89	106	144	51	157	220.5	10.5	0.18
0.005	13.5	22	28.5	35	41	50	64.5	80.5	118	34	120.5	178	26.3	0.44
0.01	10.0	15.5	19.5	23.5	28	35.5	49	63	98	23	94	147	52.6	0.88
0.02	7.0	11.0	13.5	16	18	24	35	48	78	15	72	119	105	1.75
0.05	4.0	6.4	8.0	9.5	11	14.5	22	32	52	8.3	47	86.5	263	4.38
0.1	2.5	4.2	5.2	6.1	7.2	9.8	14.5	22	35	5.2	32	64	526	8.77
0.2	1.5	2.8	3.4	4.0	4.8	6.4	9.5	14.5	21	3.1	21.8	43.5	1052	17.5
0.5	0.7	1.5	1.9	2.3	2.7	3.6	5.2	7.8	10.6	1.4	12.2	22.5	2630	43.8
1.0	0.4	1.0	1.3	1.5	1.8	2.2	3.0	4.7	6.0	0.7	8.0	12.0	5260	87.7
2.0	0.1	0.5	0.7	0.8	1.1	1.2	1.5	1.9	2.9	0.2	5.0	5.2	10520	175
5.0	0.0	0.2	0.3	0.3	0.5	0.0	0.0	0.0	0.5	0.0	1.8	1.2	26298	438

TABLE 7. ONE-MINUTE RAINFALL RATE (MM/MIN) VERSUS DURATION AND PROBABILITY OF AT LEAST ONE OCCURRENCE DURING THE WORST MONTH

Location	Worst Month	Duration (min)														
		5 Probability			10 Probability			15 Probability			20 Probability			30 Probability		
		0.1	0.5	0.9	0.1	0.5	0.9	0.1	0.5	0.9	0.1	0.5	0.9	0.1	0.5	0.9
Boston, MA	Aug.	1.19	0.86	0.48	0.90	0.47	0.26	0.60	0.33	0.16	0.36	0.22	0.13	0.30	0.15	0.09
Denver, CO	Aug.	1.49	0.73	0.38	1.18	0.38	0.22	0.64	0.28	0.13	0.44	0.22	0.09	0.29	0.12	0.07
Grand Junction, CO	Jul.	1.39	0.37	0.10	0.97	0.10	0.06	0.92	0.08	0.04	0.34	0.06	0.04	0.15	0.04	0.03
Key West, FL	Aug.	1.70	1.50	1.03	1.59	1.00	0.59	1.01	0.68	0.40	1.00	0.52	0.24	0.79	0.31	0.10
New Orleans, LA	Aug.	2.38	1.78	1.23	1.76	1.32	0.73	1.64	0.87	0.50	1.54	0.71	0.37	1.45	0.36	0.19
Omaha, NE	Jul.	1.60	1.09	0.68	1.00	0.69	0.39	0.90	0.45	0.24	0.65	0.39	0.15	0.65	0.20	0.09
Pittsburgh, PA	Jul.	1.80	1.02	0.59	0.99	0.52	0.28	0.98	0.29	0.16	0.74	0.20	0.12	0.35	0.14	0.09
Rapid City, SD	Jun.	1.30	0.77	0.50	0.75	0.49	0.24	0.65	0.29	0.18	0.60	0.20	0.14	0.25	0.16	0.11
San Angelo, TX	May	1.97	0.96	0.65	0.99	0.65	0.27	0.77	0.38	0.14	0.64	0.26	0.12	0.34	0.12	0.08
Santa Maria, CA	Dec.	0.57	0.38	0.26	0.44	0.27	0.18	0.39	0.19	0.14	0.30	0.18	0.12	0.20	0.13	0.09
Seattle, WA	Sep.	0.45	0.26	0.16	0.39	0.17	0.12	0.35	0.13	0.10	0.20	0.12	0.09	0.11	0.09	0.07
Urbana, IL	Jul.	1.96	1.38	1.00	1.63	0.87	0.41	0.84	0.51	0.22	0.82	0.28	0.14	0.24	0.11	0.08

TABLE 8. ONE-MINUTE RAINFALL RATE (MM/MIN) VERSUS DURATION AND PROBABILITY OF AT LEAST THREE OCCURRENCES DURING THE WORST MONTH

Location	Worst Month	Duration (min)														
		5			10			15			20			30		
		Probability			Probability			Probability			Probability			Probability		
	0.1	0.5	0.9	0.1	0.5	0.9	0.1	0.5	0.9	0.1	0.5	0.9	0.1	0.5	0.9	
Boston, MA	Aug.	0.71	0.44	0.33	0.37	0.24	0.17	0.23	0.15	0.12	0.18	0.13	0.09	0.13	0.09	0.06
Denver, CO	Aug.	0.53	0.35	0.23	0.32	0.19	0.10	0.22	0.11	0.07	0.14	0.08	0.05	0.09	0.06	0.04
Grand Junction, CO	Jul.	0.26	0.09	0.06	0.09	0.05	0.03	0.06	0.04	0.03	0.05	0.03	0.02	0.04	0.02	0.01
Key West, FL	Aug.	1.32	1.00	0.75	0.84	0.57	0.37	0.57	0.37	0.19	0.42	0.22	0.13	0.23	0.09	0.06
New Orleans, LA	Aug.	1.70	1.15	0.81	1.10	0.70	0.49	0.72	0.46	0.29	0.50	0.31	0.20	0.29	0.17	0.10
Omaha, NE	Jul.	0.96	0.65	0.44	0.56	0.36	0.20	0.38	0.22	0.11	0.26	0.13	0.08	0.14	0.08	0.05
Pittsburgh, PA	Jul.	0.86	0.56	0.39	0.46	0.27	0.16	0.26	0.15	0.11	0.14	0.11	0.09	0.11	0.08	0.06
Rapid City, SD	Jun.	0.67	0.46	0.27	0.41	0.22	0.16	0.24	0.17	0.13	0.18	0.13	0.10	0.14	0.10	0.07
San Angelo, TX	May	0.81	0.56	0.32	0.42	0.24	0.14	0.26	0.14	0.10	0.15	0.11	0.07	0.10	0.07	0.04
Santa Maria, CA	Dec.	0.33	0.25	0.18	0.23	0.17	0.13	0.17	0.13	0.10	0.15	0.12	0.09	0.12	0.09	0.06
Seattle, WA	Sep.	0.20	0.15	0.12	0.15	0.12	0.09	0.12	0.09	0.08	0.11	0.09	0.07	0.08	0.07	0.05
Urbana, IL	Jul.	1.27	0.91	0.63	0.80	0.38	0.23	0.37	0.20	0.12	0.22	0.13	0.08	0.09	0.07	0.05

TABLE 9. LONGEST DURATION (MIN) OF 1 MIN RATES AT OR ABOVE SPECIFIED THRESHOLD RATES AND THE MONTH OF OCCURRENCE

Location	Threshold Rate (mm/min)									
	0.1	0.2	0.4	0.7	1.0	1.3	1.6	2.0	2.5	
Boston, MA	275 Jan.	52 Sep.	23 Sep.	13 Sep.*	7 Jul.†	7 Jul.	3 Oct.	3 Oct.	1 Oct.†	
Denver, CO	162 Jun.	47 Aug.	20 Aug.	13 Aug.	12 Aug.	6 Aug.	4 Aug.*	4 Jul.	3 Jul.	
Grand Junction, CO	40 Jul.	24 Jul.	19 Jul.	17 Jul.	10 Jul.	8 Jul.	4 Jul.	2 Jul.	1 Jul.	
Key West, FL	156 May	74 Aug.	61 Apr.	39 Apr.*	22 Aug.	12 Aug.*	10 Aug.	7 Jul.*	4 Jul.*	
New Orleans, LA	205 Feb.	96 May*	86 Aug.	61 Aug.	53 Aug.	44 Aug.	24 Nov.	10 Jul.	6 Apr.	
Omaha, NE	142 May	67 Jul.	42 Jul.	18 Jul.	12 Jul.*	9 Jul.	5 Jul.*	3 Jul.‡	3 Aug.*	
Pittsburgh, PA	150 May	58 Apr.	27 Jul.	23 Jul.	17 Jul.	7 Jul.*	6 Jul.*	5 Oct.	4 Oct.	
Rapid City, SD	149 Jun.	55 Jun.	28 Jul.	16 Jul.	11 Jul.	8 Jul.	8 Jul.	6 Jul.	4 Jul.	
San Angelo, TX	161 Apr.	69 Apr.	35 Apr.	19 May	10 May	8 May	7 May	6 May	3 Jun.	
Santa Maria, CA	132 Jan.	46 Jan.	19 Sep.	5 Jan.*	3 Jan.	2 Nov.*	2 Nov.*	1 Nov.	0 -	
Seattle, WA	82 Feb.	59 Oct.	13 Sep.	2 Aug.*	0 -	0 -	0 -	0 -	0 -	
Urbana, IL	91 Oct.	45 May	37 May	25 Jun.	20 Jun.	14 Jun.	10 Jul.*	5 Jul.*	4 May	

* Also occurred in one other month.

† Also occurred in two other months.

‡ Also occurred in three other months.

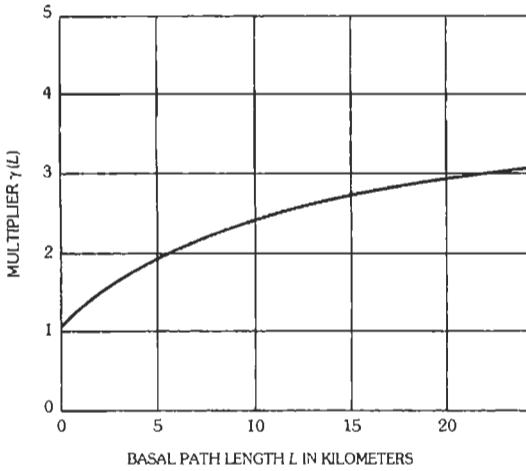


Fig. 26. Multiplier in the path averaging model.

Horn producing maximum field for given horn length

$$A = 0.45S$$

The aperture sides of the horn are assumed to be large compared with the wavelength.

Antenna Gain Relative to Hypothetical Isotropic Antennas

If directive antennas are used in place of isotropic antennas, the transmission equation becomes

$$P_r/P_t = G_t G_r [P_r/P_t]_{\text{isotropic}}$$

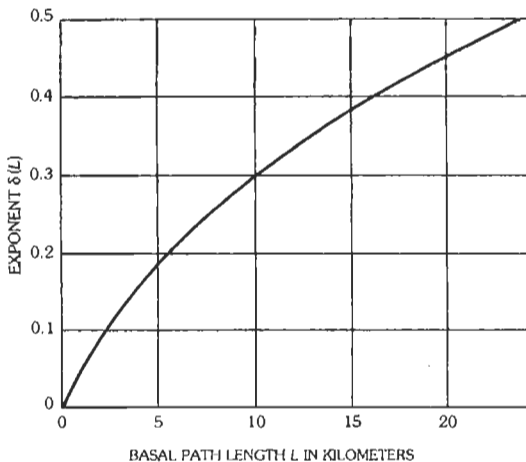


Fig. 27. Exponent in the path averaging model.

where G_t and G_r are the power gains due to the directivity of the transmitting and receiving antennas, respectively.

The apparent power gain is equal to the ratio of the effective area of the antenna to the effective area of the isotropic antenna (which is equal to $\lambda^2/4\pi \approx 0.08\lambda^2$).

The apparent power gain due to a paraboloidal reflector is thus

$$G = 0.54(\pi D/\lambda)^2$$

where D is the aperture diameter, and an illumination factor of 0.54 is assumed. In decibels, this becomes

$$G_{\text{dB}} = 20 \log f + 20 \log D - 52.6$$

where,

- f = frequency in megahertz,
- D = aperture diameter in feet.

The solution for G_{dB} may be found in Fig. 30.

Antenna Beam Width

The beam angle, θ , in degrees is related to the apparent power gain, G , of a paraboloidal reflector with respect to isotropic antennas approximately by

$$\theta^2 \approx 27\,000/G$$

Since $G = 5.5 \times 10^{-6} D^2 f^2$, the beam angle becomes

$$\theta \approx (7 \times 10^4) / fD$$

where,

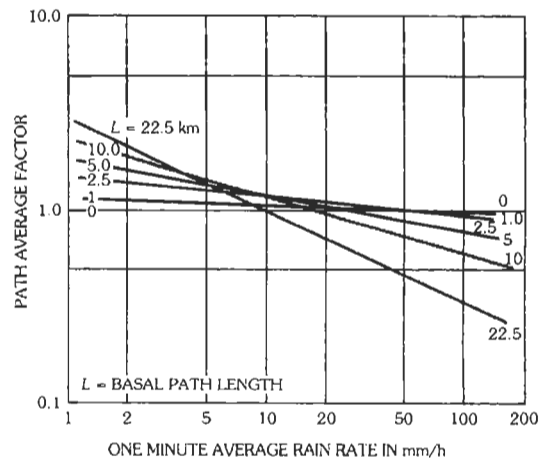


Fig. 28. Effective path average factor model for different basal path lengths.

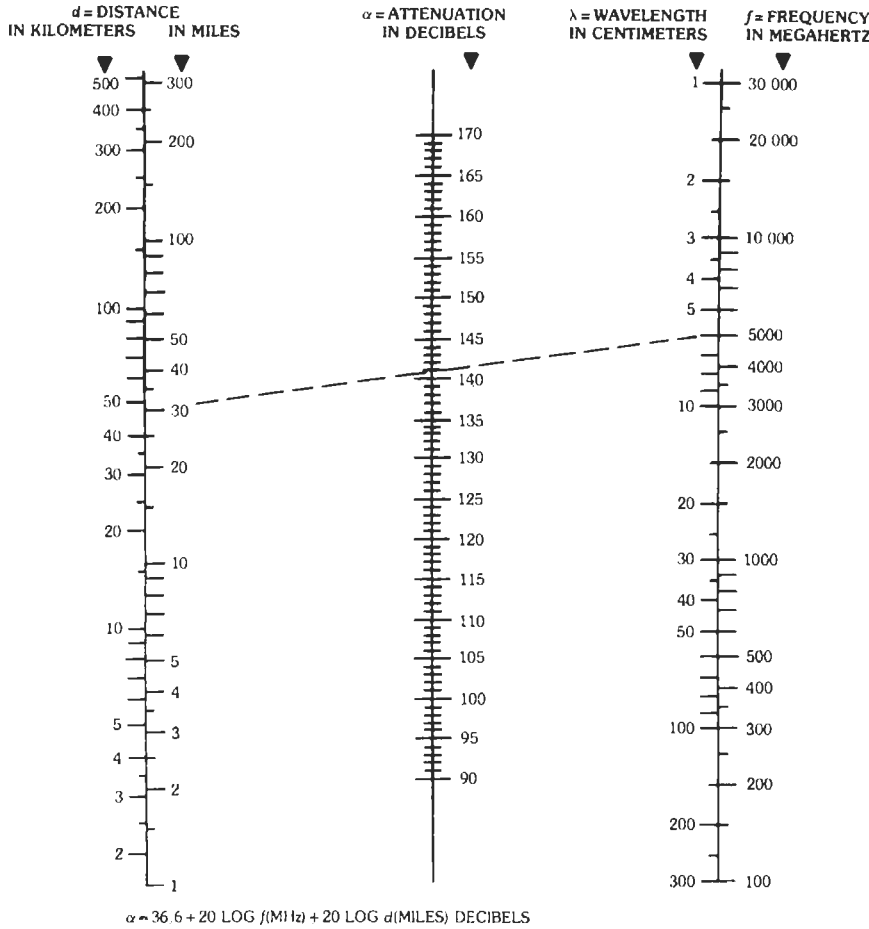


Fig. 29. Nomogram for solution of free-space path attenuation α between isotropic antennas. Example shown: distance 30 miles; frequency 5000 megahertz; attenuation = 141 decibels.

- θ = beam width between 3-decibel points in degrees,
- f = frequency in megahertz,
- D = diameter of paraboloid in feet.

Transmitter Power for a Required Output Signal/Noise Ratio

Based on the above expressions for path attenuation and reflector gain, the ratio of transmitted power to theoretical receiver noise, in decibels, is given by

$$10 \log (P_t/P_n) = A_p + (S/N) + (nf) - G_t - G_r - (nif)$$

where,

S/N = required signal/noise ratio at receiver in decibels,

- (nf) = noise figure of receiver in decibels,
- (nif) = noise improvement factor in decibels due to modulation methods where extra bandwidth is used to gain noise reduction,
- P_n = theoretical noise power in receiver,
- P_t = radiated transmitter power,
- G_t = gain of transmitting antenna in decibels,
- G_r = gain of receiving antenna in decibels,
- A_p = path attenuation in decibels.

An equivalent way to compute the transmitter power for a required output signal/noise ratio is given below directly in terms of reflector dimensions and system parameters.

(A) Normal free-space propagation

$$P_t = \frac{\beta_1 \beta_2 BL^2 F S}{40 f^2 r^4 K N}$$

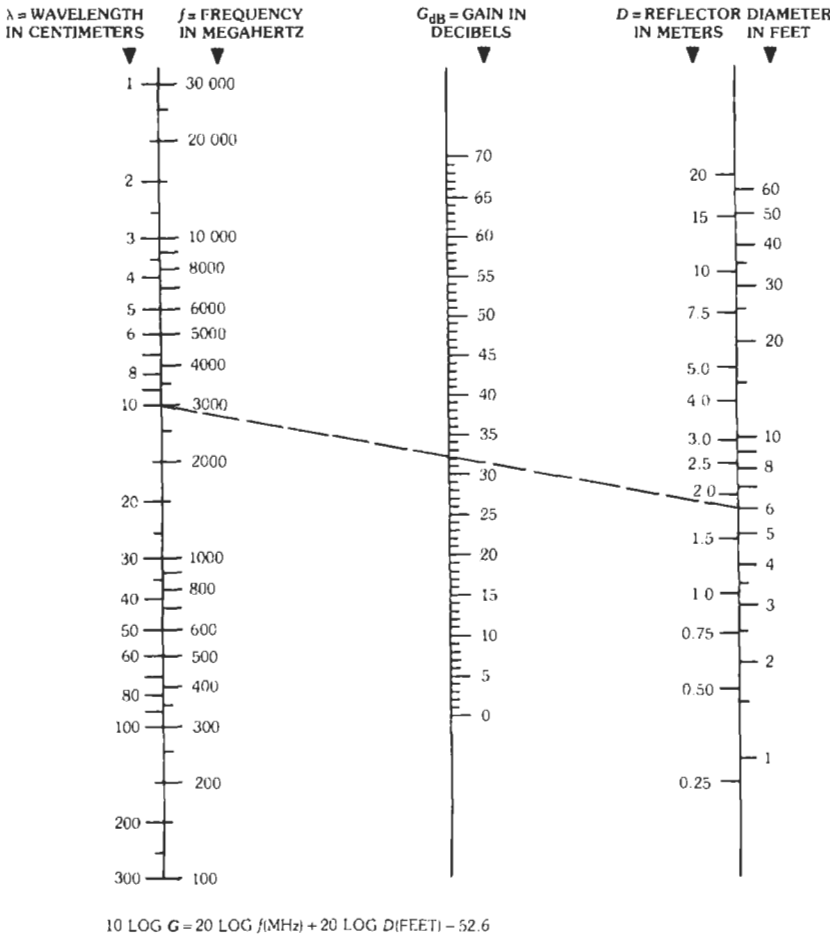


Fig. 30. Nomogram for determination of apparent power gain G_{dB} (in decibels) of a paraboloidal reflector. Example shown: frequency 3000 megahertz; diameter 6 feet; gain = 32 decibels.

(B) With allowance for fading

$$P_t = \frac{\beta_1 \beta_2 BL^2 F}{40 f^2 r^4 K} \sigma \left(\frac{S}{N} \right)_m$$

(C) For multirelay transmission in n equal hops

$$P_t = \frac{\beta_1 \beta_2 BL^2 F}{40 f^2 r^4 K} \sigma \left(\frac{S}{N} \right)_{nm}$$

(D) Signal/noise ratio for nonsimultaneous fading is

$$10 \log (S/N)_n = 10 \log \sigma (S/N)_{1m} - 10 \log \bar{n}$$

where,

P_t = power in watts available at transmitter output terminals (kept constant at each repeater point),

β_1 = loss power ratio (numerical) due to transmission line at transmitter,

β_2 = same as β_1 at receiver,

B = root-mean-square bandwidth (generally approximated to bandwidth between 3-decibel attenuation points) in megahertz,

L = total length of transmission in miles,

f = carrier frequency in megahertz,

r = radius of paraboloidal reflectors in feet,

F = power-ratio noise figure of receiver (a numerical factor),

K = improvement in signal/noise ratio due to the modulation used. (For instance, $K = 3m^2$ for frequency modulation, where m is the ratio of maximum frequency deviation to maximum modulating frequency. Note that this is the numerical power ratio.)

- σ = numerical ratio between available signal power in case of normal propagation to available signal power in case of maximum expected fading,
- S/N = required signal/noise power ratio at receiver,
- $(S/N)_m$ = minimum required signal/noise power ratio in case of maximum expected fading,
- $(S/N)_{nm}$ = same as above in case of n hops, at repeater number n ,
- $(S/N)_{1m}$ = same as above at first repeater,
- $(S/N)_n$ = same as above at end of n hops,
- n = number of equal hops,
- m = number of hops where fading occurs,
- $$\bar{n} = n - m + \sum_1^m \sigma_h$$
- σ_h = ratio of available signal power for normal conditions to available signal power in case of actual fading in hop number h (equation holds in case signal power is increased instead of decreased by abnormal propagation or reduced hop distance).

KNIFE-EDGE DIFFRACTION PROPAGATION*

Diffraction loss at an ideal knife-edge can be estimated from Fig. 31. However, the transmission loss over a practical knife-edge diffraction path depends critically on the shape of the diffracting edge. Since a natural obstacle, such as a mountain ridge, may depart considerably from an ideal knife-edge, the diffraction loss in practice is usually 10 to 20 decibels greater than that estimated for the ideal case.

A nonuniform transverse profile of the diffracting edge, or reflections on the transmission paths each side of the diffracting edge, may result in multipath transmission causing variations in the received level as a function of frequency, space, and time. The amplitude of such variations may be reduced by either space or frequency diversity and by the use of narrow-beam-width antennas.

TROPOSPHERIC SCATTER PROPAGATION†

Weak but reliable fields are propagated several hundred miles beyond the horizon in the very-high-,

ultrahigh-, and superhigh-frequency bands. An important parameter in scatter propagation is the scatter angle or angle of intersection of the transmitting and receiving antenna beams. This angle, θ , in radians is given by

$$\theta = \frac{2d - d_t - d_v}{2R} + \frac{h_t - H_t}{d_t} + \frac{h_v - H_v}{d_v}$$

where,

- d = great-circle distance between transmitting and receiving antennas,
- d_t = distance to the horizon from the transmitting antenna,
- d_v = distance to the horizon from the receiving antenna,
- h_t = height above sea level of the transmitting horizon,
- h_v = height above sea level of the receiving horizon,
- H_t = height above sea level of the transmitting antenna,
- H_v = height above sea level of the receiving antenna,
- R = effective radius of the earth.

The same units are used for distances and heights.

The effective radius of the earth is a function of the refractive index gradient and may be estimated from Fig. 32. This curve is based on the correlation found between the decrease in the refractive index in the first kilometer of altitude above the surface of the earth and the surface value of the refractive index. Fig. 33 shows typical mean values of the refractive index at sea level.

The long-term median transmission loss (in decibels) due to forward scatter is approximately

$$L(50) = 30 \log f - 20 \log d + F(\theta d) - G_p - V(d_e)$$

where $F(\theta d)$ is shown in Fig. 34 as a function of the product θd . Angular distance θ is the angle between radio horizon rays in the great-circle plane containing the antennas, and d is the distance between antennas.

A semiempirical estimate of the path antenna gain, G_p (in decibels), is provided by

$$G_p = G_t + G_r - 0.07 \exp [0.055(G_t + G_r)]$$

for values of G_t and G_r each less than 50 dB.

Fig. 35 shows $V(d_e)$, an adjustment for the indicated types of climate.

This division is, of course, rather crude, and local geographical conditions may require serious modifications. A brief description of these climates is given in Annex 1 of CCIR Report 238-2, Geneva, 1974.

Fast and slow fading is experienced on tropospheric scatter paths. Fast fading is due to multipath transmission, is in general Rayleigh distributed, and can be considerably reduced by diversity, an antenna spacing

* Bullington, K. "Radio Propagation Fundamentals." *Bell System Tech. J.* Vol. 36, No. 3, 1957; pp. 593-626.

† "Estimation of Tropospheric-Wave Transmission Loss," CCIR XVth Plenary Assembly, Geneva, 1982, Vol. V, Report 238-4. Harvey, A. F., *Microwave Engineering*, New York: Academic Press, Inc., 1963.

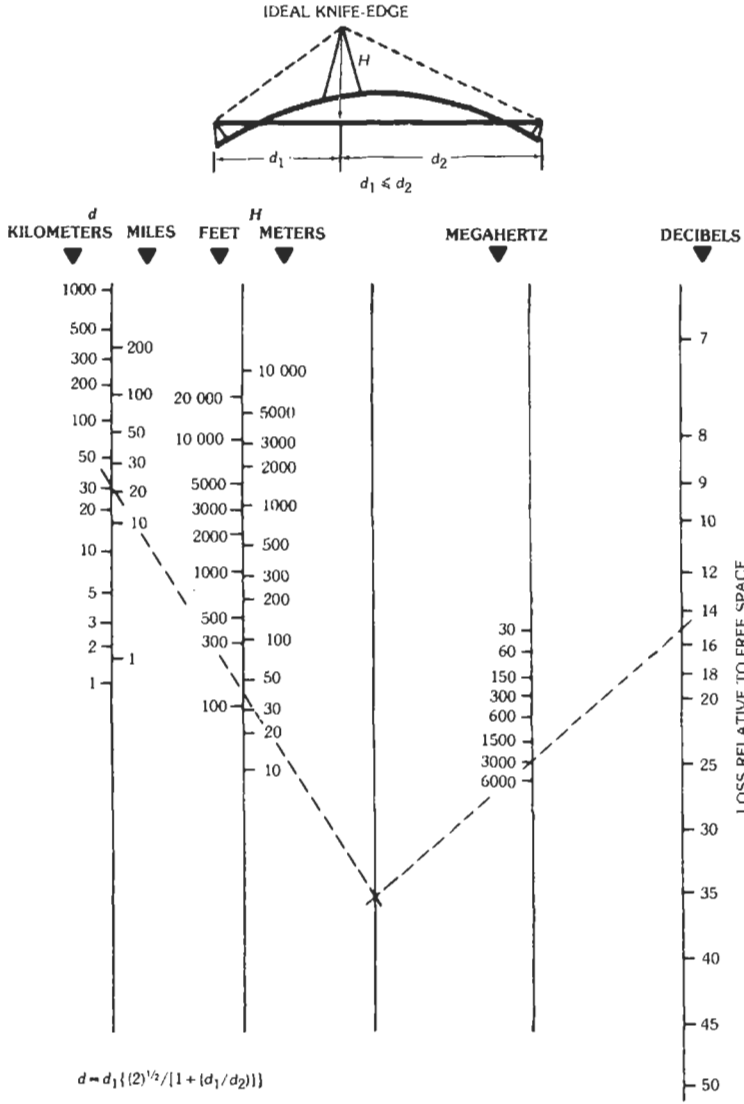


Fig. 31. Knife-edge diffraction loss relative to free space. (From Bullington, K. "Radio Propagation Fundamentals." Bell System Technical Journal, Vol. 36, No. 3, Fig. 7, ©1957 American Telephone and Telegraph Co.)

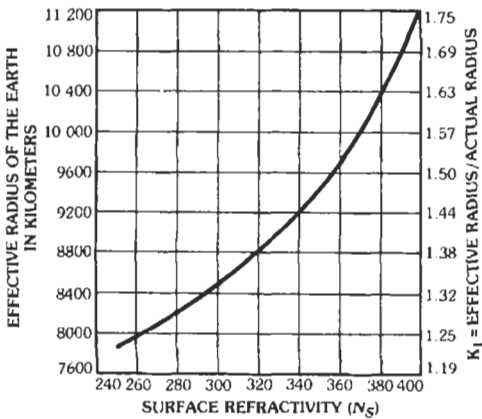


Fig. 32. Variation of the effective radius of the earth as a function of the surface refractivity, N_s . (From CCIR XVth Plenary Assembly, Geneva, 1982, Vol. V, Report 338-4.)

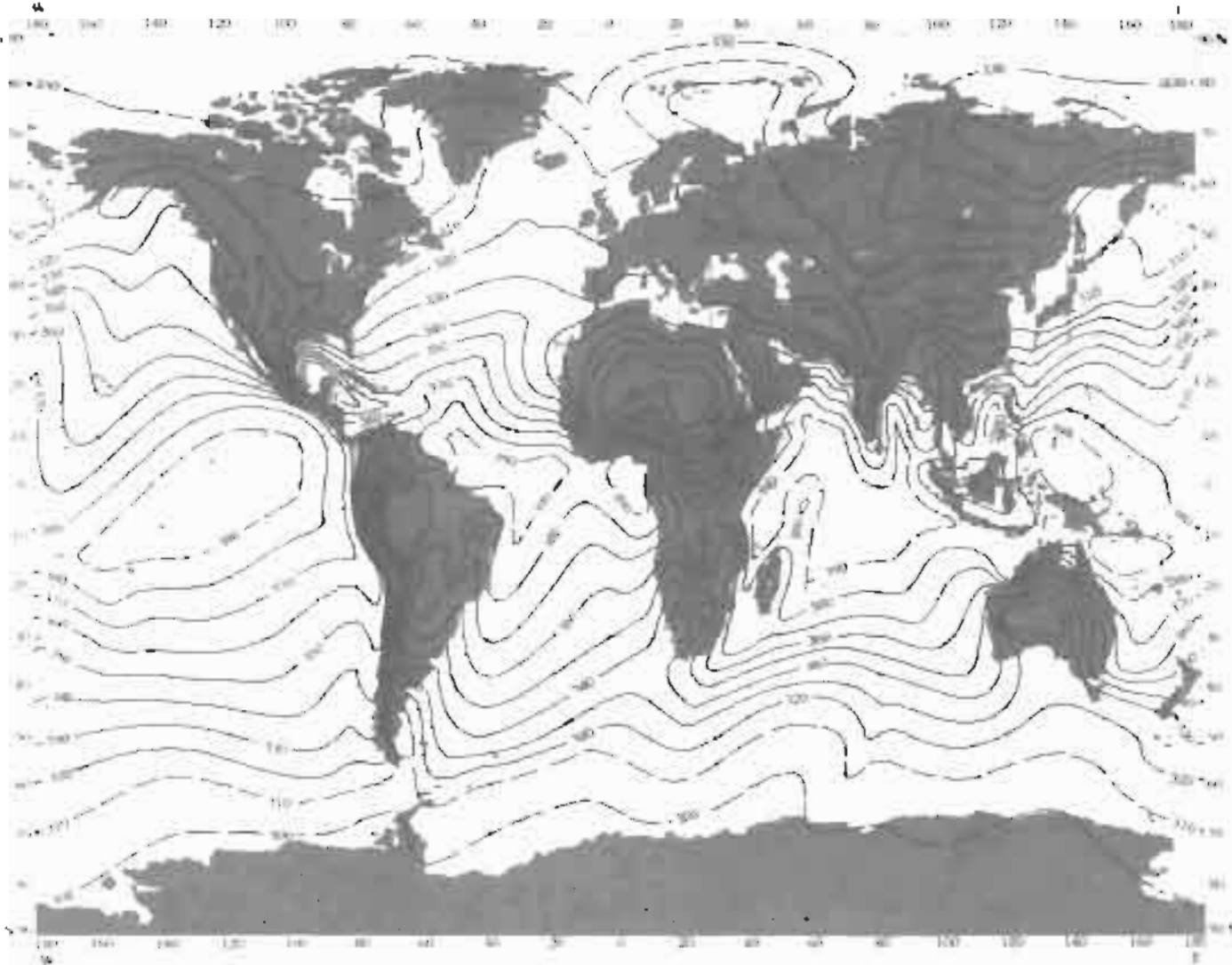


Fig. 31. Worldwide mean value of N_s for February. (From CCR XVth Plenary Assembly, Geneva, 1982, Vol. V, Report 563.2, "Influence of the Atmosphere on Wave Propagation".)

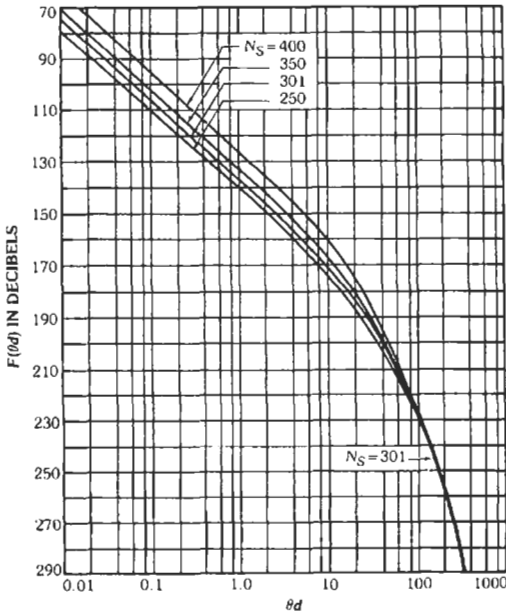
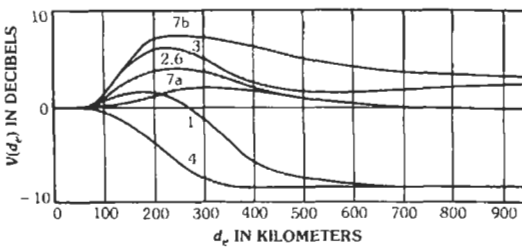


Fig. 34. Attenuation function $F(\theta d)$, where d is in kilometers and θ is in radians, for indicated values of surface refractivity N_s . (From CCIR XVth Plenary Assembly, Geneva, 1982, Vol. V, Report 238-4.)

of 60 wavelengths usually being adequate. Slow fading, with periods of hours or days, is caused by changes in the gradient of the refractive index of the atmosphere along the transmission path and is little affected by diversity.

The plane-wave gains of large antennas are not fully realized on tropospheric scatter paths. The power on

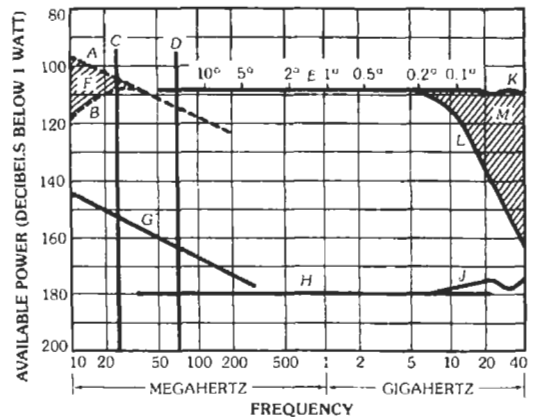


1. Equatorial (data from Congo and Ivory Coast)
2. Continental subtropical (Sudan)
3. Maritime subtropical (data from West Coast of Africa)
4. Desert (Sahara)
5. Mediterranean (no curves available)
6. Continental temperate (data from France, Federal Republic of Germany, and U.S.A.)
- 7a. Maritime temperate, over land (data from U.K.)
- 7b. Maritime temperate, over sea (data from U.K.)
8. Polar (no curves available)

Fig. 35. Function $V(d_e)$ for the types of climate indicated on the curves. (From CCIR XVth Plenary Assembly, Geneva, 1982, Vol. V, Report 238-4.)

such a path is received, not from a single point source, but from a volume in the atmosphere that subtends a solid angle at the receiving antenna. If the antenna beam angles are such as to limit the available scattering volume, then the received power will be correspondingly limited, and the antennas are said to suffer an antenna-to-medium coupling loss. The resulting median loss of received power is likely to be about 5 decibels for two 40-decibel-gain plane-wave antennas, and 17 decibels for two 50-decibel-gain antennas. The extent to which the path antenna gain is a function of the scatter angle, θ , or the height of the scatter volume has not yet been established.

Multipath transmission limits the communication bandwidth that can be used on a single carrier; however, useful bandwidths of several megahertz have been shown to be available on some 200-mile scatter paths. Narrow-beam antennas and diversity reduce the effects of multipath transmission.



- A: Signal level during ideal nighttime conditions (no absorption).
- B: Typical signal level during daytime conditions, assuming an angle of elevation of 5° .
- C: Minimum frequency to assure penetration of ionosphere: polar region, oblique path; tropical region, vertical path.
- D: Minimum frequency to assure penetration of ionosphere: tropical region, oblique path.
- E: Beam width of paraboloid between half-power points.
- F: Effect of ionospheric absorption.
- G: Minimum cosmic noise. Maximum value will be found to be higher by about 15 decibels.
- H: Noise level corresponding to a temperature of 70 kelvins.
- J: Noise due to absorption in a clear atmosphere, assuming an elevation angle of 5° .
- K: Typical signal level for a vertical path in a clear atmosphere.
- L: Typical signal level in heavy rain (16 millimeters/hour), vertical depth 1 kilometer, assuming an elevation angle of 5° .
- M: Effect of varying atmospheric conditions and elevation angles.

Fig. 36. General frequency limits in a simple earth-to-spacecraft communication system. (Spacecraft: isotropic antenna; transmitter power, 1 watt; bandwidth, 1 kilohertz; distance, 1000 kilometers. Earth station antenna: paraboloid; diameter, 20 meters; efficiency, 55%—, 15-decibel gain above isotropic antenna - - -.) (From CCIR Xth Plenary Assembly, Geneva, 1963, Vol. IV, Report 205, p. 194.)

EARTH-SPACE COMMUNICATION

Communication between earth and outer space (see Chapter 27, Satellite and Space Communications) must pass through the atmosphere of the earth, so that the optimum frequencies for this service are those that pass through the atmosphere with minimum attenuation. A range of frequency little attenuated by the atmosphere is known as a window; one such window occurs between the critical frequency of the ionosphere and the frequency absorbed by rainfall and oxygen. This frequency range extends from about 10 to 10 000 megahertz. Another window exists in the optical and infrared region of 10^6 to 10^9 megahertz. Fig. 36 shows the general frequency limits for earth-space communication.

REFERENCES

1. Tattelman, P. "New Rain Rate Analyses to Assess Rain Attenuation on Satellite EHF Communications." *International Journal of Satellite Communications*, Vol. 7, pp. 23-25, 1989.
2. Crane, R. K. "Prediction of Attenuation by Rain." *IEEE Transactions on Communications*, Vol. Com 28, pp. 1717-1723, September 1980.
3. Ippolito, L. J. *Propagation Effects Handbook for Satellite System Design*. NASA Ref Publication 1082(04), February 1989. See also Ippolito, L. J. *Radio Propagation in Satellite Communications*. New York: Van Nostrand Reinhold Co., 1986.
4. CCIR Report 719-1 *Attenuation by Atmospheric Gases*. CCIR XV Plenary Assembly, Geneva 1982, Vol. V.

34 Radio Noise and Interference

*Revised by
George W. Swenson, Jr. and A. Richard Thompson*

Natural Noise	34-2
Thermal Noise	
Atmospheric Noise	
Cosmic Noise	
Man-Made Radio Noise	34-6
Near-Zone and Far-Zone Noise Sources	
Power Line Noise	
Precipitation Static	34-8
Thermal Noise Calculations	34-8
Noise Measurements	34-8
Measurement for Broadcast Receivers	
Noise Factor of a Receiver	
Measurement of Noise Figure With a Thermal Noise Source	
Calculation of Noise Figure	
Noise Factor of Cascaded Networks	
Interference From Signals of Other Services	34-11

Radio noise and interference limit the performance of all communications systems by restricting the operating range, generating errors in messages, and in extreme cases preventing the successful operation of receivers. At locations where man-made noise is low, natural noise sources determine receiver performance. When man-made noise encroaches upon receiving sites, the performance of receiving equipment is degraded below design levels.

When noise from sources external to a receiver is involved, the gain and orientation of the antenna must be considered. For narrow-band receivers, noise is usually flat in amplitude across the bandwidth of a receiver. For such cases, noise power affecting receiver performance is proportional to the bandwidth. For wideband receivers, the noise may not be flat across the receiver bandwidth, and the determination of effective noise power requires further consideration.

Noise level can be expressed in terms of voltage or power at the terminals of a receiver, the strength of an electromagnetic field at an antenna location, or thermal

noise power at a temperature referenced to 290 kelvins. Noise that is flat in amplitude across the bandwidth of a receiver is often expressed in terms of an effective antenna noise factor, f_a , which is defined as

$$f_a = P_n/kT_0B = T_a/T_0 \quad (\text{Eq. 1})$$

where,

- P_n = noise power in watts from an equivalent lossless antenna,
- k = Boltzmann's constant,
- T_0 = reference temperature (290 kelvins),
- B = receiver noise bandwidth in hertz,
- T_a = antenna noise temperature in the presence of external noise.

NATURAL NOISE

Natural noise consists of thermal noise, atmospheric noise, and cosmic noise. These noise sources usually

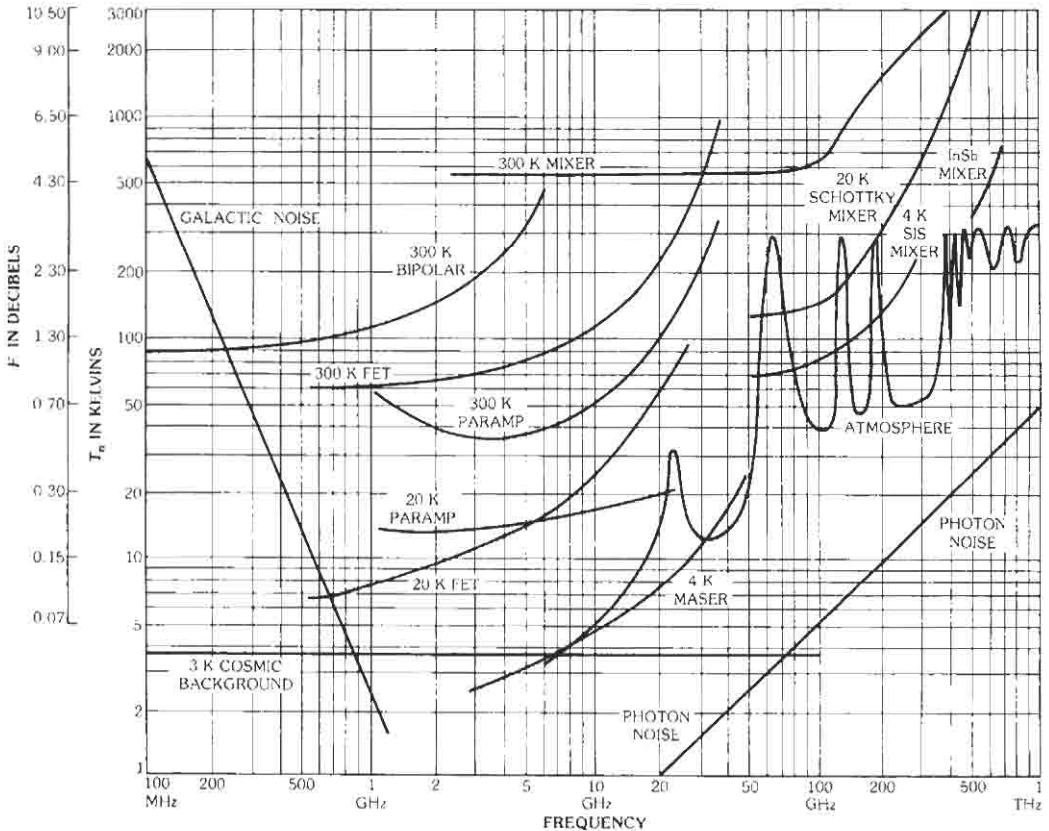


Fig. 1. Noise figure (F) and noise temperature (T_n) for various devices and natural limits—1984. (From S. Weinreb, "Low-Noise GASFET Amplifiers." IEEE Trans. on MTT, Vol. MTT-28, No. 10, October 1980, pp. 1041-1054. Material updated to 1984 by Weinreb.)

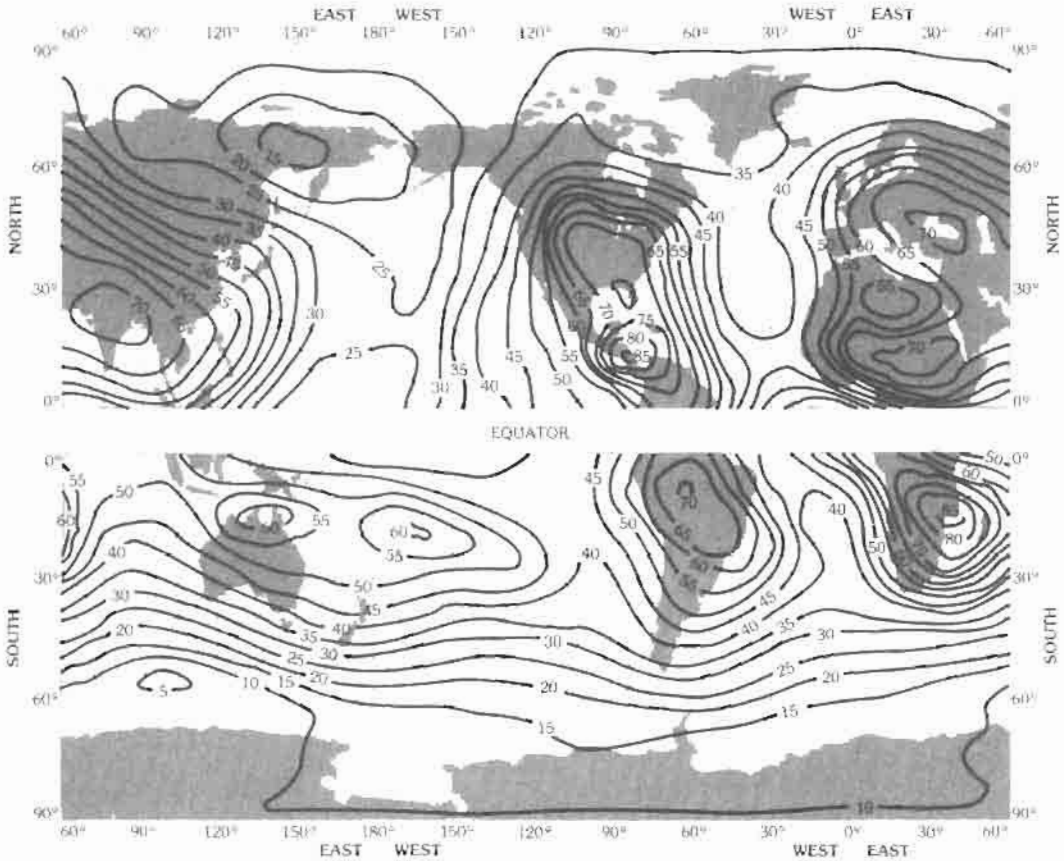


Fig. 2. Atmospheric noise levels in northern and southern hemispheres, summer, 1200–1600 hours local time. The maps show the expected values of F_{ν} at 1 MHz, in decibels above kT_0B . (From CCIR Report 322, 10th Plenary Assembly, Geneva, 1963.)

determine the minimum detectable signal level of a receiver operated in an environment free of man-made noise sources.

Thermal Noise

For many years, thermal noise in the first stage of a receiver was usually the main factor limiting the sensitivity of radar and microwave receivers. Recent advances in low-noise amplifier performance have reduced thermal noise of microwave amplifiers to very low levels, and thermal radiation from nearby objects, ground, and the sky are now major factors that must be considered in choosing sites for satellite receivers and radars. Fig. 1 shows noise temperatures for various devices (as of 1984) and natural limits at microwave frequencies. Receivers operated at frequencies below about 20 MHz usually encounter noise from other sources that is considerably above the thermal noise of conventional amplifiers; hence, low-noise amplifier performance is not usually a factor in the design of low-frequency receivers. An exception occurs for VLF receivers operating in the Arctic and Antarctic, where

atmospheric noise is extremely low and cosmic noise is screened by the ionosphere.

Atmospheric Noise

Lightning from thunderstorms produces bursts of impulsive noise. At low frequencies, these bursts are propagated to distant receivers by normal ionospheric modes. The noise is dependent on the weather, time of day, season, location of the receiver with respect to storm areas, and ionospheric propagation conditions. Atmospheric noise generally decreases with increasing latitude and increases in high-latitude equatorial areas. Atmospheric noise sources are particularly active during the rainy season in the Caribbean, the East Indies, equatorial Africa, northern India, and the Far East. An excellent summary of worldwide atmospheric noise levels is contained in CCIR Report 322.* An example of a CCIR developed map of atmospheric noise levels in the summer during daytime hours is shown in Fig. 2.

* CCIR Report 322, Vol. VI, 10th Plenary Assembly, Geneva, 1963.

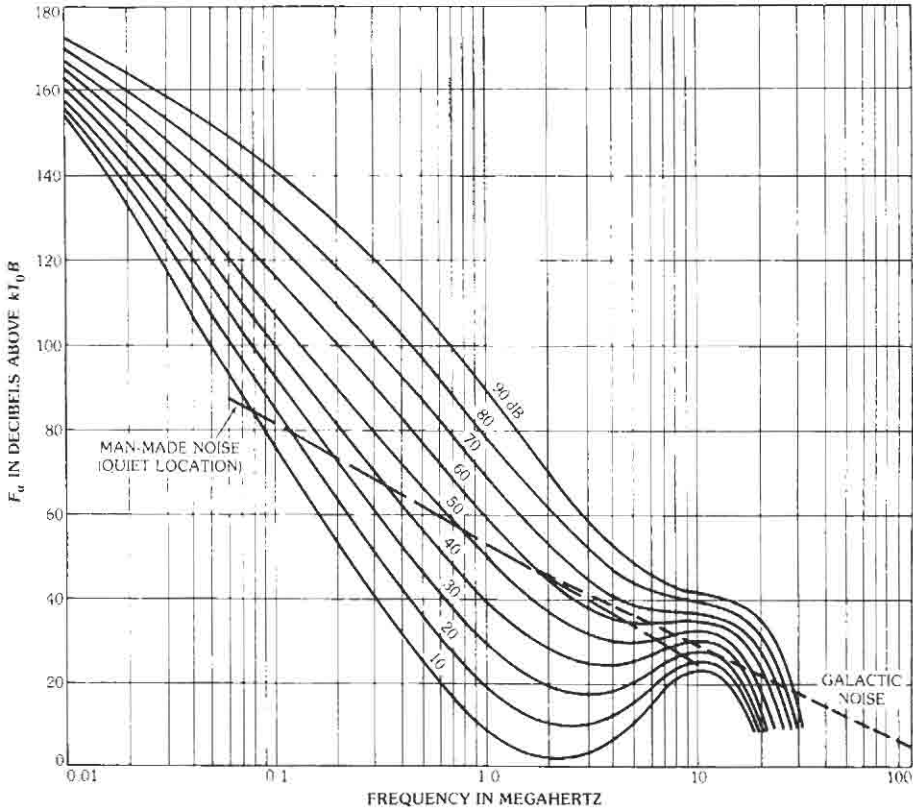


Fig. 3. Variation of radio noise with frequency, for data given in Fig. 2 legend. (From CCIR Report 322, 10th Plenary Assembly, Geneva; 1963.)

The map shows the median noise level in decibels above kTB at a frequency of 1 MHz as received on a short vertical rod antenna installed over ground (k is Boltzmann's constant, T is 290 kelvins, and B is the receiver bandwidth in hertz). This parameter is related to noise field strength by

$$E_n = F_a + 20 \log_{10} f_{\text{MHz}} - 65.5 \quad (\text{Eq. 2})$$

where,

E_n = rms noise field strength in a 1-kHz bandwidth in decibels above 1 microvolt/meter,

F_a = noise level in decibels above kTB ,

f_{MHz} = frequency in megahertz.

The level of atmospheric noise at a receiver site decreases with increasing frequency. Fig. 3 shows the frequency dependence of atmospheric noise for the data shown in Fig. 2. These levels represent the median amplitude of the noise bursts. Individual bursts can vary in amplitude from insignificant to very strong. An example of an individual atmospheric burst propagated at HF frequencies over a one-hop ionospheric mode is

shown in Fig. 4. The data were obtained with a rapidly scanning receiver and a time memory display. The noise burst originated from a lightning event that was accompanied by numerous discharges lasting about one second. The wide-bandwidth impulsive energy produced by the lightning discharge was modified by the ionospheric propagation path so that only a portion of the discharge was seen at the distant receiver (2000 km from the storm to the receiver). The maximum propagating frequency and the minimum propagating frequency of the ionosphere are shown.

Cosmic Noise

The sources of cosmic noise are external to the atmosphere of the earth; primary sources are the sun, the Milky Way galaxy, and other discrete cosmic sources. Radio astronomy activities have identified a very large number of sources of cosmic noise. Radio noise from cosmic sources must penetrate the atmosphere of the earth to reach antennas located on the surface of the earth. Ionospheric absorption limits the reception of cosmic noise at frequencies below about 20 MHz, while molecular absorption processes limit the

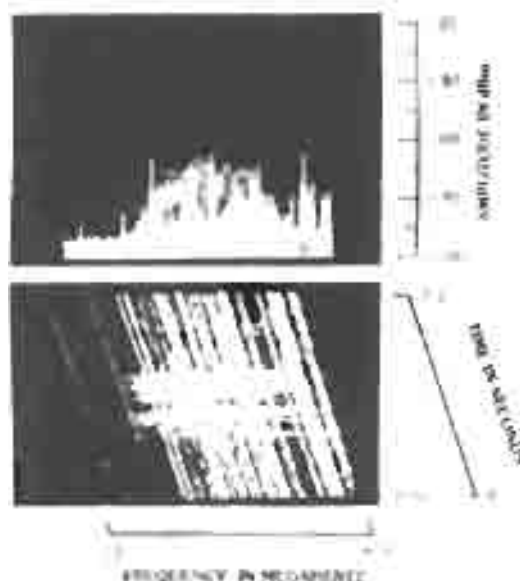


Fig. 4. Atmospheric noise band.

reception of extraterrestrial noise at frequencies above about 10 GHz. Satellite-borne receivers above about 1000 km do not encounter these limitations.

Recent advances in low-noise receiver design (see Fig. 1) and the widespread deployment of satellites and space probes have increased the importance of cosmic noise. Satellite communications systems, the broadcasting of television from satellites, and the need for data links between the space vehicles and the earth have increased the number of skyward-pointed antennas equipped with sensitive receivers that are capable of receiving cosmic noise. Cosmic noise often limits the performance of such systems.

Figs. 8 and 9 of Chapter 32 show detailed radio-sky maps of the celestial sphere for the 136-megahertz and 400-megahertz space research satellite frequency bands.

Fig. 5 shows the level of galactic noise in decibels relative to a noise temperature of 290 K when received on a half-wave dipole. The noise levels shown in this figure assume no atmospheric absorption and refer to the following sources of cosmic noise.

Galactic Plane. Galactic noise from the galactic plane in the direction of the center of the galaxy. The noise levels from other parts of the galactic plane can be as much as 12 to 15 decibels below the levels given in Fig. 5.

Quiet Sun. Noise from the "quiet" sun, that is, solar noise at times when there is little or no sunspot activity.

Disturbed Sun. Noise from the "disturbed" sun. The term "disturbed" refers to times of sunspot and solar flare activity.

Cassiopeia A. Noise from a high-intensity discrete source of cosmic noise known as Cassiopeia A. This

is one of thousands of known discrete sources. Cassiopeia A subtends a solid angle at the surface of the earth of only about 1 arc minute.

The levels of cosmic noise received by a highly-directional antenna with main lobe pointed along the galactic plane can be obtained from equations given by Kraus¹ for the antenna noise temperature (T_A) at the output terminals of an ideal, loss-free, antenna as

$$T_A = \frac{\int_{\theta_0}^{\theta_1} \sin^2 \theta \, d\theta \int_{\phi_1}^{\phi_2} T(\theta, \phi) G(\theta, \phi) \sin \theta \, d\phi}{\int_{\theta_0}^{\theta_1} \sin^2 \theta \, d\theta \int_{\phi_1}^{\phi_2} G(\theta, \phi) \sin \theta \, d\phi} \quad (5)$$

where,

$\theta = \theta'$ at zenith,

$\phi = \text{azimuth angle,}$

$T(\theta, \phi)$ = brightness-noise temperature distribution from radio-sky map, kelvins,

$G(\theta, \phi)$ = antenna radiation pattern gain distribution, assumed symmetrical,

θ_0 = minimum elevation angle between antenna main-lobe axis and the horizon, degrees.

However, for a practical antenna Taylor and Stocklin² give a simplified approximation for T_A including contributions from the main lobe, side lobes, and back lobe as

$$T_A = 0.82 T_m + 0.18 T_{sl} + T_e \quad (6)$$

angle beam, $\theta_{main} = \phi_{main} \approx 25^\circ$

where,

T_m = mean value of sky-brightness temperature within main-lobe HPBW in kelvins,

T_{sl} = mean value of sky-brightness temperature within antenna side lobes, in kelvins,

$T_e = T_0 = 290$ K, effective noise temperature of earth.

For example, a 136-megahertz, phased-array, directive antenna with main-lobe HPBW equal to 12° (pointed near Cassiopeia A, has a value of T_A equal to approximately 870 K for T_m equal to 950 K and T_{sl} equal to 400 K obtained from Fig. 8 of Chapter 32.)

¹ J. D. Kraus, *Radio Astronomy*, New York: McGraw-Hill Book Co., 1966, 2nd ed., p. 263, CH. 1, *Radio-Quantum Books*, 1966.

² R. F. Taylor and F. J. Stocklin, "VHF-UHF Signal Calibration Error Analysis," *Proceedings International Interfering Conference*, Washington, D. C., Vol. VII, pp. 553-566, September 27-29, 1971.

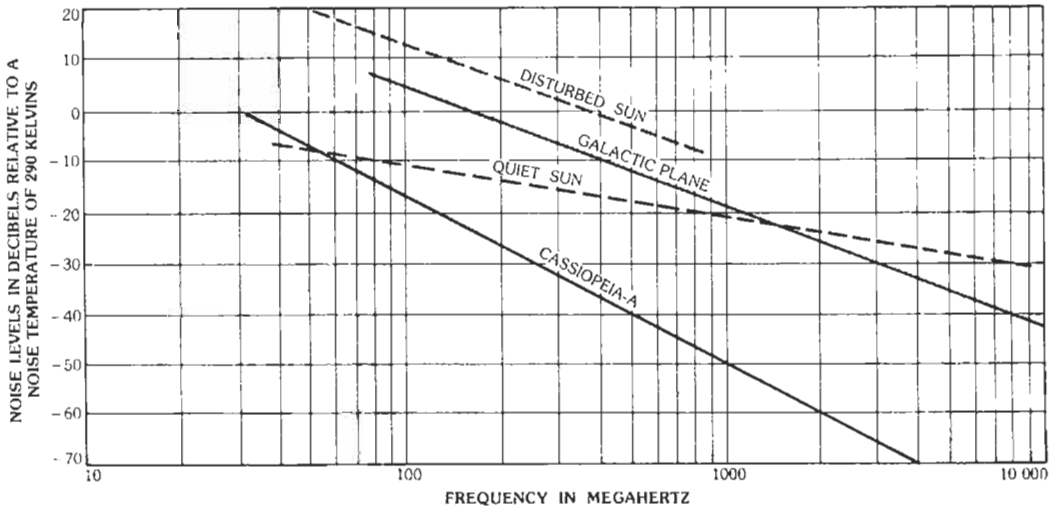


Fig. 5. Cosmic noise levels for a half-wave-dipole receiving antenna.

MAN-MADE RADIO NOISE

Man-made radio noise frequently limits the performance of receivers. This is particularly true for land-mobile communications, television reception, high-frequency radio, and other radio services below UHF. Man-made radio noise originates from a wide variety of sources; some examples are noise from ignition systems of gasoline engines, corona noise from high-voltage power lines, gap noise from utility distribution lines, noise from radio-frequency stabilized welders, and noise produced by many other electrical devices found in homes and businesses. Fig. 6 shows the frequency range that is affected by several common types of man-made noise. Sources that affect the performance of low-frequency receivers generally do not produce harmful noise at higher frequencies. The reverse is also true in that sources that produce harmful noise at VHF and UHF generally do not produce harmful noise at low frequencies. Since noise amplitude decreases with distance from the source, the magnitude of noise at a receiver site is dependent on spacial parameters, temporal variations of the noise source, directional properties of the noise-radiating elements, and the directional properties of the receiving antenna. Sometimes harmful noise is inductively coupled or conducted from its source into a receiver.

The average levels of man-made noise are higher in urban and suburban areas than in rural areas because of the larger number of sources in areas of higher population. Fig. 7 shows median values of average noise power in urban and suburban areas. In remote quiet locations, man-made noise can be extremely low, and background noise in most parts of the radio spectrum may be determined by natural noise (atmospheric noise below about 20 MHz and galactic noise above about 20 MHz).

In general, man-made noise levels decrease with increasing frequency, although a specific source may not comply with this general rule.

Near-Zone and Far-Zone Noise Sources

Most cases of harmful man-made noise affecting a receiver involve only one or two sources. A general background level of man-made noise from numerous sources is seldom encountered. Frequently, the sources of harmful noise are relatively close to a receiver. The propagation from the noise source to the receiver sometimes involves the near-zone field rather than the far-zone field of the source. For such cases, both the electric- and magnetic-field components of the noise source must be measured in order to define fully the ability of noise to affect receivers. Table 1 provides a general set of rules for the measurement of noise. When the path from the noise source to the receiver is greater than one wavelength, the far-field approximation ap-

TABLE 1. GENERAL RULES FOR MEASUREMENT OF NOISE

Distance Source to Sensor	Zone	Measure
Greater than λ	Distant	Either <i>E</i> or <i>B</i>
Less than $\lambda/6$	Inductive	Must measure both <i>E</i> and <i>B</i>
Between $\lambda/6$ and λ	Mixed	Recommend <i>E</i> and <i>B</i>

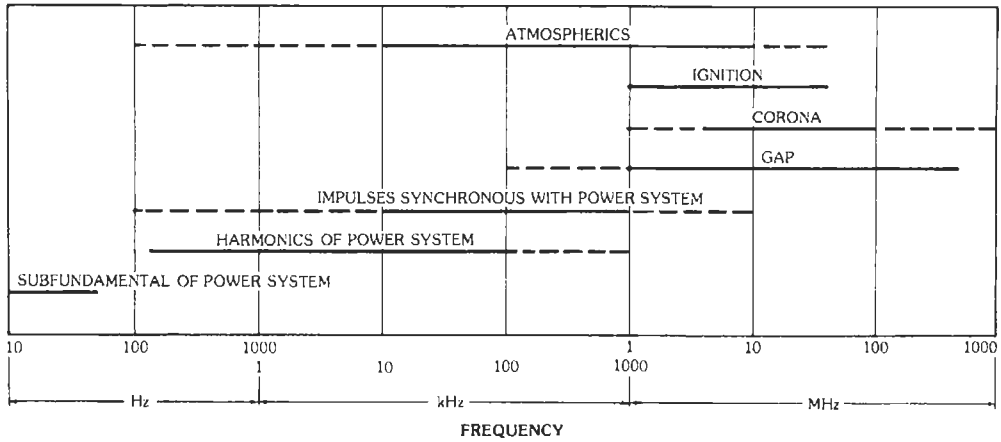


Fig. 6. Frequency range affected by man-made noise.

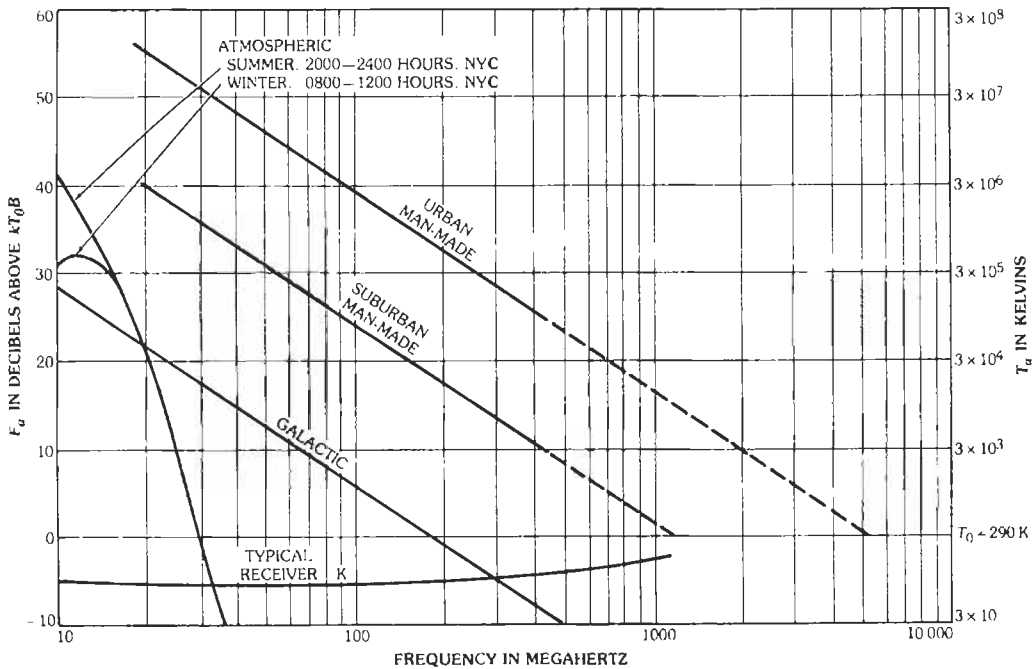


Fig. 7. Median values of average noise power expected from various sources (omnidirectional antenna near surface).

plies, and the noise can be measured with either an electric-field or a magnetic-field antenna. The electric- and magnetic-field strengths can be related by the free-space impedance of $120\pi = 377$ ohms. This applies to measurements of man-made noise at frequencies above about 30 MHz. At lower frequencies, the noise source may be electrically close to the receiving antenna. At distances from the source to the affected receiver of one sixth of a wavelength or less, i.e., $< \lambda/2\pi$, the receiver is in the near field of the source. At

these short distances, the ratio of E to H is no longer equal to 377 ohms, and both the electric and magnetic fields must be measured in order to define the impact of a noise source on a receiver.

Power Line Noise

The temporal and spectral properties of man-made noise vary considerably from one type of source to another. Noise associated with a power line usually

contains bursts of noise at intervals determined by the fundamental frequency of the power line. The temporal properties of gap noise (formed by the minute electrical breakdown between two pieces of metallic hardware exposed to a strong electric field) are shown in Fig. 8. Multiple bursts of noise occur each time the power-line voltage reaches a maximum. Individual impulses within a burst are spaced very close together (less than 1 ms apart), while the bursts are spaced 8.33 ms apart for a power line operating at 60 Hz. Individual bursts have the same amplitude. A careful inspection of the waveform in Fig. 8 shows that two gap-noise sources are present in the view.

PRECIPITATION STATIC

Precipitation static is produced by rain, hail, snow, or dust storms in the vicinity of the receiving antenna and is important chiefly at frequencies below 10 MHz. This form of interference can be reduced by eliminating sharp points from the antenna and its surroundings, and also by providing means for dissipating the charges that build up on an antenna and on its surroundings during electrical storms.

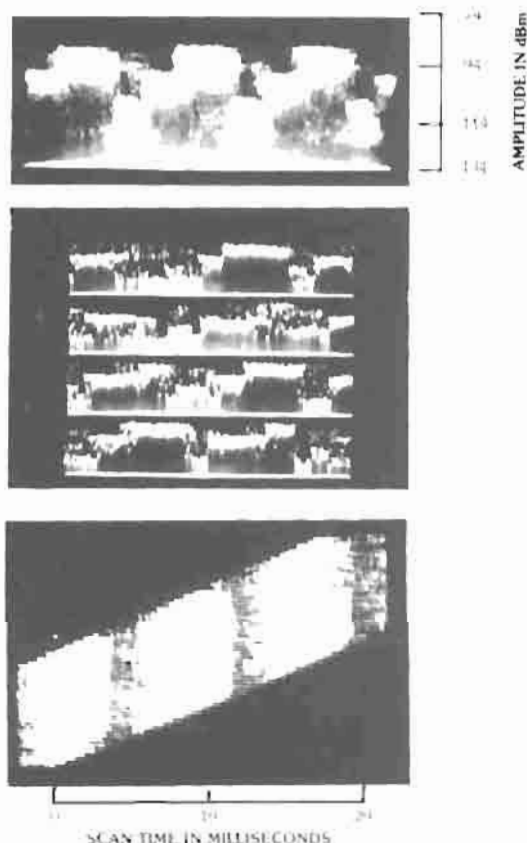


Fig. 8. Temporal structure of gap noise from 138-kV transmission line.

THERMAL NOISE CALCULATIONS

Thermal noise is caused by the thermal agitation of electrons in resistances. Let R be the resistive component in ohms of an impedance Z . The mean-square value of thermal-noise voltage is given by

$$E^2 = 4RkT \cdot \Delta f$$

where,

k is Boltzmann's constant (1.38×10^{-23} joules/kelvin),

T is the absolute temperature in kelvins,

Δf is the bandwidth in hertz,

E is the root-mean-square noise voltage.

The equation given above assumes that thermal noise has a uniform distribution of power through the bandwidth Δf .

In case two impedances Z_1 and Z_2 with resistive components R_1 and R_2 are in series at the same temperature, the square of the resulting root-mean-square voltage is the sum of the squares of the root-mean-square noise voltages generated in Z_1 and Z_2 :

$$E^2 = E_1^2 + E_2^2 = 4(R_1 + R_2)kT \cdot \Delta f$$

In case the same impedances are in parallel at the same temperature, the resulting impedance Z is calculated as is usually done for alternating-current circuits, and the resistive component R of Z is then determined. The root-mean-square noise voltage is the same as it would be for a pure resistance R .

It is customary in temperate climates to assign to T a value such that $1.38T = 400$, corresponding to about 17 degrees Celsius or 63 degrees Fahrenheit. Then $E^2 = 1.6 \times 10^{-20} R \cdot \Delta f$.

NOISE MEASUREMENTS

Measurement for Broadcast Receivers

For standard broadcast receivers the noise properties are determined by means of the equivalent noise sideband input (ensi). The receiver is connected as shown in Fig. 9. Components of the standard dummy antenna are $C_1 = 200$ picofarads, $C_2 = 400$ picofarads, $L = 20$ microhenrys, and $R = 400$ ohms.

The equivalent noise sideband input is

$$(\text{ensi}) = mE_s(P'_n/P'_c)^{1/2}$$

where,

E_s = root-mean-square unmodulated carrier-input voltage,

m = degree of modulation of signal carrier at 400 hertz,

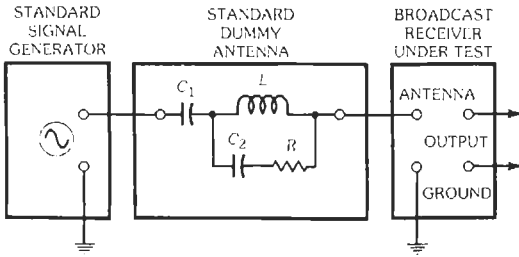


Fig. 9. Measurement of equivalent noise sideband input of a broadcast receiver.

- P'_s = root-mean-square signal-power output when signal is applied.
- P'_n = root-mean-square noise-power output when signal input is reduced to zero.

It is assumed that no appreciable noise is transferred from the signal generator to the receiver, and that m is small enough for the receiver to operate without distortion.

Noise Factor of a Receiver

A more precise evaluation of the quality of a receiver as far as noise is concerned is obtained by means of its noise factor.

It should be clearly realized that the noise factor evaluates only the linear part of the receiver, i.e., up to the demodulator.

The equipment used for measuring noise factor is shown in Fig. 10. The incoming signal (applied to the receiver) is replaced by an unmodulated signal generator with R_0 = internal resistive component, E_i = root-mean-square open-circuit carrier voltage, and E_n = root-mean-square open-circuit noise voltage produced in signal generator. Then

$$E_n^2 = 4kT_0R_0\Delta f'$$

where,

- k is Boltzmann's constant (1.38×10^{-23} J/K),
- T_0 is the temperature in kelvins,

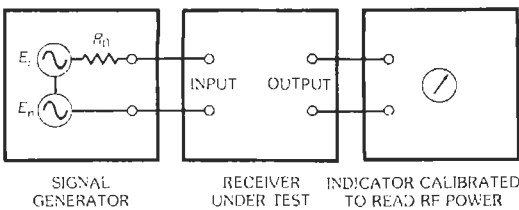


Fig. 10. Measurement of the noise factor of a receiver. The receiver is considered as a 4-terminal network. Output refers to last intermediate-frequency stage.

$\Delta f'$ is the effective bandwidth of the receiver (determined as below).

If the receiver does not include any other source of noise, the ratio E_i^2/E_n^2 is equal to the power carrier/noise ratio measured by the indicator:

$$E_i^2/E_n^2 = (E_i^2/4R_0)/kT_0\Delta f' = P_i/N_i$$

The quantities $E_i^2/4R_0$ and $kT_0\Delta f'$ are called the available carrier and noise powers, respectively.

The output carrier/noise power ratio measured in a resistance R may be considered as the ratio of an available carrier-output power P_0 to an available noise-output power N_0 .

The noise factor, F , of the receiver is defined by

$$P_0/N_0 = F^{-1}(P_i/N_i)$$

$$F = (N_0/N_i)(P_0/P_i)^{-1}$$

$$= E_{i:1}^2/4kT_0R_0\Delta f' = P_{i:1}/kT_0\Delta f'$$

where,

- P_0/P_i = available gain G of the receiver,
- $P_{i:1}$ = available power from the generator required to produce a carrier-to-noise ratio of one at the receiver output.

Noise figure is the noise factor expressed in decibels:

$$F_{dB} = 10 \log_{10} F$$

Effective bandwidth $\Delta f'$ of the receiver is

$$\Delta f' = G^{-1} \int G_f df$$

where G_f is the differential available gain. Generally, $\Delta f'$ is approximated to the bandwidth of the receiver between those points of the response showing a 3-dB attenuation with respect to the center frequency.

Measurement of Noise Figure With a Thermal Noise Source

For the case where the spurious responses of the receiver are negligible, receiver noise figure can be conveniently measured by using the noise output of a thermal noise source having an equivalent generator resistance equal to that specified for use with the receiver.

With the noise source off, but still possessing the correct output resistance, receiver gain is adjusted for a convenient amount of noise power output; then with the noise source on, and still possessing the correct output resistance, the noise power output is increased by a

convenient power ratio (N_2/N_1). The measured noise figure is then given by

$$NF = (\text{excess})_{\text{dB}} - 10 \log\{(N_2/N_1) - 1\}$$

For a thermal diode operating in the temperature limited emission mode

$$(\text{excess})_{\text{dB}} = 10 \log(20R_d I_d)$$

where,

R_d is the noise source output resistance,
 I_d is the diode current in amperes.

When the receiver has appreciable spurious responses, the correction factor that must be used with the above simple equation is a complex function of the spurious response ratios, and of the percentage of total internal receiver noise produced by the circuits preceding the mixer causing the spurious responses. For the simple case of no preselection and a diode mixer having negligible excess noise, 3 dB must be added to the measured noise figure to obtain the true noise figure.

A thermal noise source designed for a given generator impedance R_1 can be used to measure the noise figure of a receiver designed for a higher generator R_2 by adding a resistor ($R_2 - R_1$) between noise source and receiver input and using

$$NF = NF_{\text{rcad}} - 10 \log (R_2/R_1)$$

Conversion of receiver noise temperature to noise factor:

$$F = 1 + (T_R/T_0)$$

where,

T_R = receiver noise temperature in kelvins,
 $T_0 = 290 \text{ K}$,
 F = noise factor of receiver (power ratio).

Conversely,

$$T_R = (F - 1)T_0$$

Determination of effective noise temperature of receiving system (i.e., antenna, transmission line, and receiver):

$$T_E = T_A + (LF - 1)T_0$$

where,

T_E = effective noise temperature of receiving system,
 T_A = antenna noise temperature,
 L = transmission line loss (power ratio),
 F = noise factor of receiver (power ratio),
 $T_0 = 290 \text{ K}$.

Determination of the effective input noise power of the receiving system:

$$N_i = kBT_E$$

where,

N_i = effective input noise power of the receiving system,
 k = Boltzmann's constant (1.38×10^{-23} joules/kelvin),
 B = bandwidth in hertz,
 T_E = effective noise temperature in Kelvins.

$$\text{dBm}_i = -198.6 + 10 \log B + 10 \log T_E$$

Calculation of Noise Figure

The active device can be defined for noise-figure calculations as in Fig. 11.

The value of R_{eq} can be obtained experimentally by measuring, with a "zero impedance" generator, the equivalent microvolts of noise V_{sc} , in a bandwidth B , in series with the input terminals, with an almost-short-circuit on the output terminals. Then R_{eq} is given by

$$R_{\text{eq}} = |V_{\text{sc}}|^2 / 1.64 \times 10^{-20} \langle BW \rangle \quad (\text{Eq. 3})$$

The value of R_e is obtained straightforwardly by input impedance measurements with a short-circuit on the output terminals.

The value of ρ can be obtained experimentally by approximately open-circuiting the input terminals at the frequency of interest with a tuned circuit of parallel resonant resistance R_0 , and measuring the total equivalent microvolts of noise produced across the input terminals, with an almost-short-circuit on the output terminals. Then, assuming negligible correlation

$$\rho = [1 + (R_e/R_0)]^2 \{ [|V_{\text{oc}}|^2 / (1.64 \times 10^{-20} \langle BW \rangle R_e)] - (R_{\text{eq}}/R_e) \} - (R_e/R_0) \quad (\text{Eq. 4})$$

When the above characterized device is used with an input transforming circuit of parallel resonant resistance R_r , the resulting noise factor can be calculated as follows: First calculate R_1 and β from

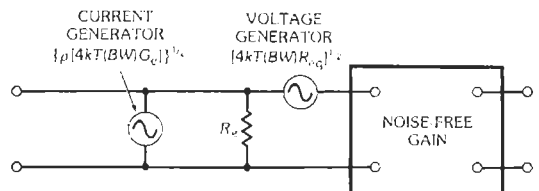


Fig. 11. Calculation of the noise figure of a receiver.

$$R_1^{-1} = R_r^{-1} + R_e^{-1} \quad (\text{Eq. 5})$$

$$\beta = [1 + \rho(R_r/R_e)]/[1 + (R_r/R_e)] \quad (\text{Eq. 6})$$

In terms of the above quantities and the transformed generator resistance R_s seen by the input terminals of the active device, the resulting noise factor is given by

$$F = 1 + 2(R_{eq}/R_1) + (R_{eq}/R_s) + (R_s/R_1) \times [\beta + (R_{eq}/R_1)] \quad (\text{Eq. 7})$$

It should be noted that to minimize noise figure the input circuit should always be tuned so as to null any part of the noise due to βR_1 , which is correlated with the noise due to R_{eq} . Equation 7 can be applied to this best noise figure tuning case if ρ is obtained, by some method, from only the uncorrelated part of the βR_1 noise.

This resulting noise figure is minimized when the transformed generator resistance has the value

$$R_{s \text{ opt}} = \{(R_1 R_{eq}/\beta)/[1 + (R_{eq}/\beta R_1)]\}^{1/2} \quad (\text{Eq. 8})$$

and with this optimum source resistance

$$F_{\text{opt}} = 1 + 2\beta(R_{eq}/R_1)^{1/2} \times \{[1 + (R_{eq}/R_1)]^{1/2} + (R_{eq}/R_1)^{1/2}\} \quad (\text{Eq. 9})$$

Noise Factor of Cascaded Networks

The overall noise factor of two networks, a and b , in cascade (Fig. 12) is

$$F_{ab} = F_a + [(F_b - 1)/G_a]$$

provided $\Delta f_b' \leq \Delta f_a'$.

The additional noise due to external sources influencing real antennas (such as galactic noise) may be accounted for by an apparent antenna temperature, bringing the available noise-power input to $kT_a\Delta f'$ instead of $N_i = kT_0\Delta f'$ (the physical antenna resistance at temperature T_0 is generally negligible in high-frequency systems). The internal noise sources contribute $(F - 1)N_i$ as before, so that the new noise factor is given by

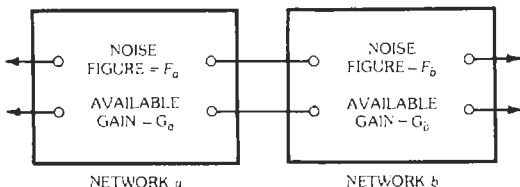


Fig. 12. Overall noise figure F_{ab} of two networks, a and b , in cascade.

$$F'N_i = (F - 1)N_i + kT_0\Delta f'$$

$$F' = (F - 1) + (T_a/T_0)$$

The average temperature of the antenna for a 6-MHz equipment is found to be 3000 kelvins, approximately. The contribution of external sources is thus of the order of 10, compared with a value of $(F - 1)$ equal to 1 or 2, and becomes the limiting factor of reception. At 3000 MHz, however, values of T_a may fall below T_0 .

INTERFERENCE FROM SIGNALS OF OTHER SERVICES

Although the usage of the radio spectrum is subject to the detailed system of regulations described in Chapter 1, there are many circumstances in which a service operating in compliance with the radio regulations can cause interference to another service. This most often results from sharing of an allocated band by more than one service. For example, the Fixed service (communications between fixed stations on the Earth) often shares a band with the Mobile service (communications between vehicles, sometimes including a fixed station). These two services also share some bands with a satellite service. Even if the systems have been designed to avoid interference most of the time, it is difficult to prevent interference when, for example, a nongeostationary satellite moving close to the horizon passes through the horizontally pointing beam of an antenna of the Fixed service, or when a mobile station is close to such an antenna. Interference can also be caused by transmitters operating in other bands, if out-of-band transmissions such as harmonics are radiated. The extensive out-of-band radiation that results from unfiltered transmitters using the direct sequence form of spread-spectrum modulation can cause harmful interference to very sensitive systems such as those used in radio astronomy, even though such side lobes are low enough to go undetected by most other services.

For any particular situation, it is possible to determine by calculation whether harmful interference will occur. First it is necessary to know the EIRP (effective isotropically radiated power) of the particular transmitter in the direction of the victim receiver. If the direction does not fall within the main beam of the transmitter, one may need to know the radiation pattern of the antenna. In the case of large parabolic antennas, the reference radiation pattern recommended by the CCIR* may be useful. This model applies to large paraboloids of diameter greater than 100 wavelengths in the frequency range 2 to 30 GHz. The envelope of the side-lobe pattern corresponds to a gain, in dB relative to an isotropic radiator, of $32 - 25 \log \phi$, where ϕ is the direction in degrees measured from the axis of the main beam. This relationship applies for values of ϕ from 1°

* Reference 1.

to 47.9° . For greater values of ϕ , the gain is -10 dBi. For other types of antennas, if specific details of the power pattern are not available, it may be necessary to use a "worst-case" estimate based on general information about the type of antenna. After the transmitted EIRP has been estimated, the next step is to find the power flux density (PFD) at the receiver location. This requires determination of the propagation loss, for which methods are described in Chapter 33. For free-space propagation the situation is simple, but propagation between two widely separated points on the surface of the Earth may involve diffraction around a spherical model earth, diffraction over one or more knife-edge ridges, or tropospheric scatter. Curves and formulae for such cases are given in Chapter 33. From the estimated PFD, the power received in the antenna of the victim receiver can be calculated if the gain of the antenna in the direction of the transmitter is known. As in the case of the transmitting antenna, this may involve an estimation of the side-lobe gain. Finally, it is necessary to compare the received power with the power level that corresponds to the threshold of harmful interference at the input of the particular receiving system involved. This harmful threshold can be calculated from the theory of operation of the receiver, or determined from laboratory measurements on the receiver. The harmful threshold is a measure of the vulnerability of a system to interference. The most vulnerable service is radio astronomy (see for example reference 2) followed by activities such as Earth exploration by passive sensing and reception of information from spacecraft at deep-space ranges. At the other end of the vulnerability scale, systems using spread spectrum techniques or digital modulation with error-correcting codes can function in the presence of interfering signals as much as 100 dB greater than those harmful to radio astronomy.

The extensive literature of the CCIR,* in the form of Reports and Recommendations of the various CCIR study groups, contains many studies of interservice interference, sharing of frequency bands, and related topics. This material provides technical input to the conferences at which the international radio regulations are revised, and is a highly authoritative source of reference.

To illustrate the consequences of overlaying an actively transmitting service on a sensitive passive service, assume a typical spread-spectrum personal communication system, radiating 0.1 watt isotropically over a 10-megahertz band centered at about 2.0 gigahertz. It would produce an interfering power density of the order of 3×10^{-19} watt/hertz in the far-out side-lobes of a receiving antenna one kilometer away. This would represent an increase in the system temperature of the service receiving interference of 20 000 kelvins. As a

result, the typical performance of a modern communication system would be degraded by a factor of 200. Such a penalty would reduce the maximum distance for communication with a spacecraft by a factor of 14. It would reduce the volume of the Universe from which an astronomical radio telescope could receive useful data by a factor of 200. Terrestrial services would suffer the same type of degradation from overlaid spread-spectrum services. Further, in addition to this interference to the overlaid services, the "direct-sequence" type of spread-spectrum system typically emits interference well outside the band it utilizes for communication unless well designed bandpass filters are employed at the output of the transmitter.

REFERENCES

1. *CCIR Recommendation 509-1, Recommendations of the CCIR*, Vol. 2. Geneva: Int. Telecommunication Union, 1990.
2. Thompson, A. R., Moran, J. M., and Swenson, G. W., Jr. *Interferometry and Synthesis in Radio Astronomy*. New York: John Wiley and Sons, 1986, and Melbourne, FL: Krieger Press, 1991, Ch. 14.
3. *Recommendations of the CCIR*, Vols. I–XII and Annexed Reports. Geneva: Int. Telecommunication Union, 1990.

General References

Morrison, Ralph. *Noise and Other Interfering Signals*. New York: John Wiley and Sons, 1992.

White, Donald R. J. *A Handbook Series on Electromagnetic Interference and Compatibility*. Germantown, MD: Don White Consultants, Inc.

Vol. I. *Electrical Noise and EMI Specifications*, 1971.

Vol. II. *Electromagnetic Interference Test Methods and Procedures*, 1973.

Vol. III. *Electromagnetic Compatibility Control Methods and Techniques*, 1973.

Vol. IV. *Electromagnetic Interference Test Instruments and Systems*, 1971.

Vol. V. *Electromagnetic Interference Prediction Techniques*, 1972.

Crane, Patrick C., and Hillenbrand, L. A. "Estimating Harmful Levels of Radio-Frequency Radiation," in Crawford, David L., ed., *Light Pollution, Radio Interference and Space Debris*, Astronomical Society of the Pacific Conference Series, Vol. 17, International Astronomical Union Colloquium No. 112, 1990, pp. 258–266.

Pakala, W. E., and Chartier, V. L. "Radio Noise

* Reference 3.

Measurements on Overhead Power Lines from 2.4 to 800 kV." *IEEE Trans on Power Apparatus and Systems*, PAS-90, 1971, pp. 1155-1165.

Skomal, E. N. *Man-Made Radio Noise*. New York: Van Nostrand-Reinhold, 1978.

Skomal, E. N., and Smith, A. A., Jr. *Measuring the*

Radio Frequency Environment. New York: Van Nostrand-Reinhold, 1985.

Cudak, M. C., Swenson, G. W., Jr., and Cochran, W. W. "Airborne Measurements of Incidental Radio Noise from Cities." *Radio Science* 26, 1991, pp. 773-781.

35 **Broadcasting, Cable Television, and Recording System Standards**

F. M. Remley, J. F. X. Browne, and S. N. Baron

Broadcast and Cable Transmission Systems	35-4
Standard Broadcasting	
Frequency-Modulation Broadcasting	
Television Broadcasting (VHF and UHF)	
Cable Television	
Other Television Services	
Network Distribution of Broadcast Program Signals	
Auxiliary Broadcast Services	
International Broadcasting Service in the United States	
Sources for Federal Communications Commission Documents	
Program Production Standards	35-26
Sound Recording Systems	
Television Recording Systems	
Selected Lists of Television Standards	
Digital Television Systems	35-34
Introduction—The Basics	
Resolution	
Sampling and Spectra	
Quantizing and Dynamic Range	
Composite and Component Signal Coding	
Bit-Parallel and Bit-Serial Interfaces	
Composite Encoded Signals: SMPTE 244 (NTSC and PAL)	

Component Coding of Signals: SMPTE 125 and CCIR 601
Synchronization and Blanking Interval Considerations
Ancillary Data Signal Format
Advanced Television Systems
References

International Standards for Broadcasting *35-46*

CCIR Documents
Selected CCIR Reports and Recommendations
IEC Publications
ISO Recommendations
Selected EBU Publications
Acknowledgements

In this age of rapid technical change, it is no longer accurate to assume that communications services directed to the public at large are "broadcasting" in the traditional sense. More and more of the general public is served by cable television and by transmissions from satellites. In addition, the video cassette recorder is an integral part of millions of home entertainment systems, and in this case the signals reaching the viewer are distributed by lengths of magnetic tape. In consideration of these facts, this chapter will not only describe conventional over-the-air broadcasting of radio and television programs, and cable and satellite television, but also will focus on the elements that enter into program recording for other media that serve the public.

The various transmission systems, and the signals that they convey to the public, must meet certain technical specifications. Much of this chapter will consist of summaries and lists of such specifications. Sources for current versions of relevant specifications are also identified.

In the United States, broadcasting is regulated by the Federal Communications Commission (FCC), which assigns frequencies and establishes technical standards. Three general classes of broadcast stations have traditionally been identified. These are standard broadcast stations (amplitude modulation in the band 535-1605

kHz), FM broadcast stations (frequency modulation in the band 88-108 MHz), and television broadcast stations (operating in the bands 54-72, 76-88, 174-216 and 470-806 MHz with vestigial-sideband amplitude modulation of the visual carrier and frequency modulation of the aural carrier). Technical specifications for these broadcast services are summarized in this chapter. Cable television systems also have channel assignment and other specifications; they are also summarized in this chapter.

Certain technical aspects of international broadcasting are described. International broadcasting ("short wave") uses frequencies between 5950 and 26 100 kHz in accordance with international agreements.

Other technical information related to broadcasting is covered under auxiliary services, intercity relay transmission, and terminal facilities.

This chapter provides a comprehensive reference for standards applicable to audio recording and video recording. North American standards are listed, and reference is made to comparable international standards. Lists of standards are provided.

The rapidly evolving field of advanced television systems using digital technologies is described in separate sections. Both national and international specifications are listed. A listing of relevant standards promulgated by CCIR, ISO, and IEC concludes the chapter.

TABLE 1. CLASSIFICATION OF STANDARD-BROADCAST STATIONS

Class of Channel	Class of Station	Permissible Power (kW)	Signal-Intensity Contour of Area Protected from Objectionable Interference (microvolts/meter)		Permissible Interfering Signal on Same Channel (microvolts/meter)	
			Day ¹	Night	Day ¹	Night ³
Clear	I-A	50	SC = 100 AC = 500	SC = 500 ² AC = 500 ¹	5	25
	I-B	10-50	SC = 100 AC = 500	SC = 500 ² AC = 500 ¹	5	25
	II-A	0.25-50 day 10-50 night	500	500 ¹	25	25
	(II-B) (II-D)	0.25-50	500	2500 ¹	25	125
	(II-B) ⁴ (II-D)	0.25-1	500	10 000	25	500
Regional	III-A	1-5	500	2500 ¹	25	125
	III-B	0.5-5 day 0.5-1 night	500	4000 ¹	25	200
Local	IV	0.25-1 day 0.25 night	500	not prescribed	25	not prescribed

Notes:

SC—same channel, AC—adjacent channel.

¹Ground wave.

²50% sky wave.

³10% sky wave.

⁴Applies to those Class II-B and II-D stations described in Part 73 of FCC Rules and Regulations.

BROADCAST AND CABLE TRANSMISSION SYSTEMS

Standard Broadcasting*

Standard-broadcast stations are licensed for operation on channels spaced by 10 kHz† and occupy the band from 535 to 1700 kHz. The major station classifications prior to 1992 are shown in Table 1.

Field-Strength Requirements—Field strength requirements have been established by the FCC to provide specified minimum signal levels in various types of locations. These vary somewhat with the class of station but more specifically with the anticipated man-made noise levels in a particular area. As of 1991 these levels generally are as follows. (As noted above, the FCC Rules were undergoing major changes in service and

interference criteria in 1992, and a current version of the FCC Rules and Regulations should be consulted.)

Primary Service:

City business, factory areas—10 to 50 millivolts/meter, ground wave.

City residential areas—2 to 10 millivolts/meter, ground wave.

Rural, all areas during winter or northern areas during summer—0.1 to 0.5 millivolt/meter, ground wave.

Rural, southern areas during summer—0.25 to 1.0 millivolt/meter, ground wave.

Secondary Service: All areas having sky-wave field strength equal to or greater than 500 microvolts/meter for 50% or more of the time.

* FCC Rules and Regulations, Part 73, Subpart A. These rules were undergoing substantial changes in early 1992, and a current version of the Rules should be consulted for latest classifications and requirements.

† 10-kHz spacing is used only in ITU Region 2. Regions 1 and 3 employ 9-kHz spacing.

For a station employing a directional antenna, all determinations of service and interference are based on the inverse field of a “standard pattern” for that station. When applied to nighttime operation, this includes the radiation pattern in the horizontal plane as well as

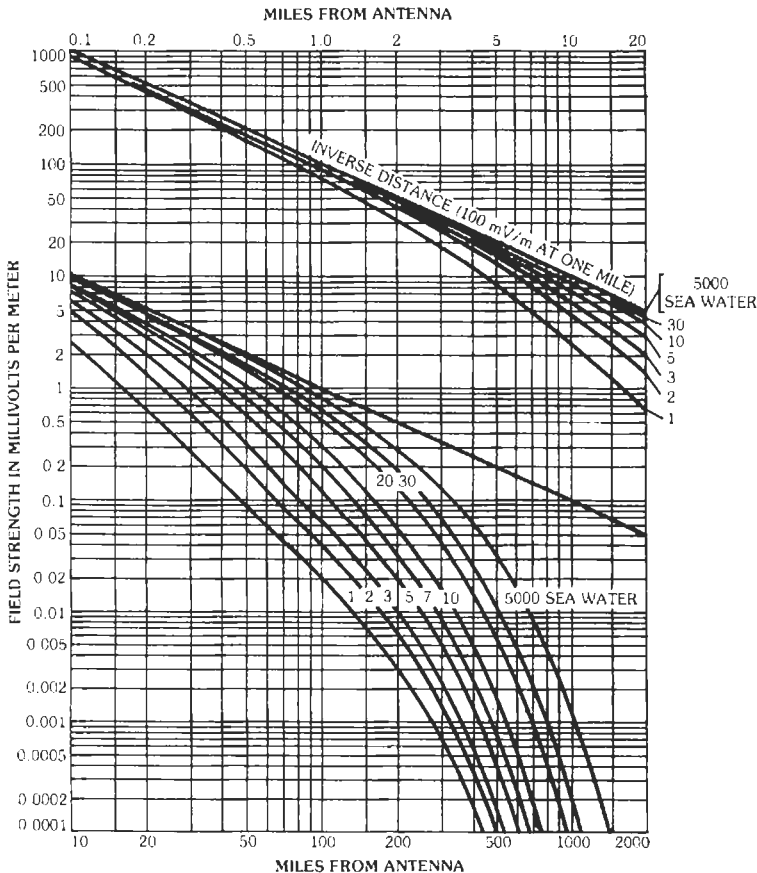


Fig. 1. Ground-wave field strength plotted against distance. Computed for 550 kilohertz. Dielectric constant = 15. Ground-conductivity values are in millimhos/meter.

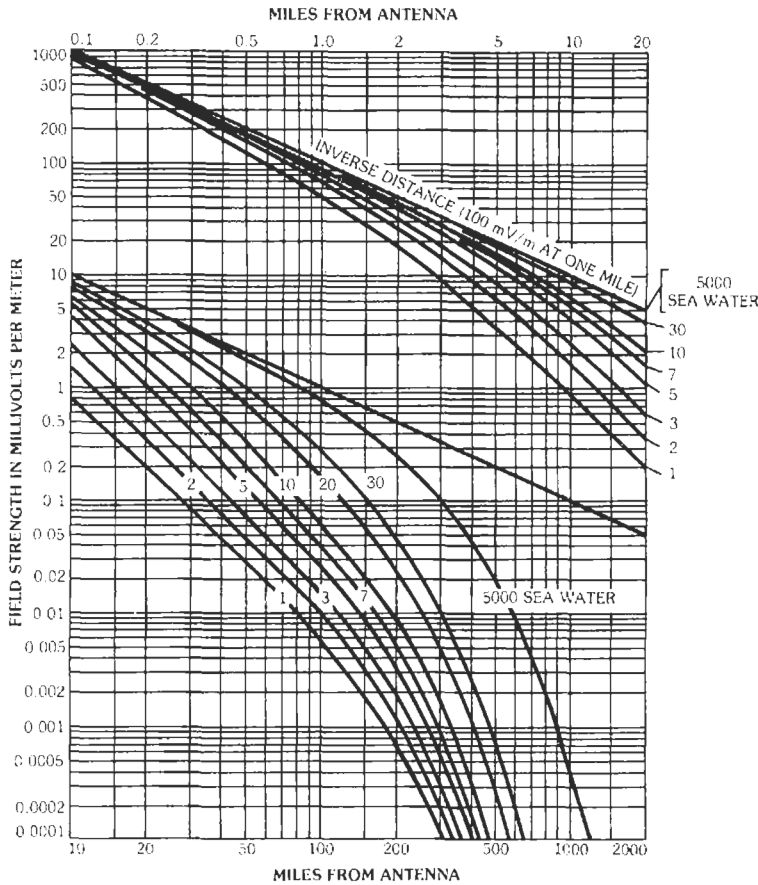


Fig. 2. Ground-wave field strength plotted against distance. Computed for 1000 kilohertz. Dielectric constant = 15. Ground-conductivity values are in millimhos/meter.

radiation at angles above the horizontal plane (vertical radiation pattern).

Table 1 outlines generally the protected contours and permissible interference for the various classes of stations. There are additional details and some exceptions in Sections 73.21–73.29 and 73.181–73.190 of Part 73 of the FCC Rules and Regulations. (See previous footnote regarding changes to these rules. The new rules will require stations to employ new interference criteria and to reduce existing interference for a site change.)

Coverage Data—Figs. 1, 2, and 3 show computed values of ground-wave field strength as a function of the distance from the transmitting antenna. These are used to determine coverage and interference. They were computed for the frequencies indicated, for a dielectric constant equal to 15 for land and 80 for sea water (referred to air as unity), and for the surface conductivities noted. The curves are for radiation from a short vertical antenna at the surface of a uniformly conductive spherical earth with an antenna efficiency and input power such that the inverse-distance field is 100 millivolts/meter at one mile. Fig. 4 shows the estimated effective field for vertical omnidirectional antennas of

various heights. Figs. 5 and 6 show the effective ground conductivity for various parts of the US and Canada, and Fig. 7 shows the sky-wave fields for 10% and 50% of the time.

Transmission System Requirements

Modulation: 85% on peaks of frequent recurrence at authorized operating power. Not greater than 100% on negative peaks or 125% on positive peaks.

Audio-Frequency Response: In accordance with NRSC recommendations.

Carrier Amplitude/Carrier Shift: Less than 5% at any level of modulation percentage.

Out-of-Band Emissions: Per FCC/NRSC requirements, see FCC Rules.

AM Stereo Transmission—Three basic transmission methods were proposed for adding stereo to the standard AM broadcast system. The three methods were:

Mixed Mode: The carrier is amplitude modulated with the (L+R) signal and phase or frequency modulated with the (L–R) signal.

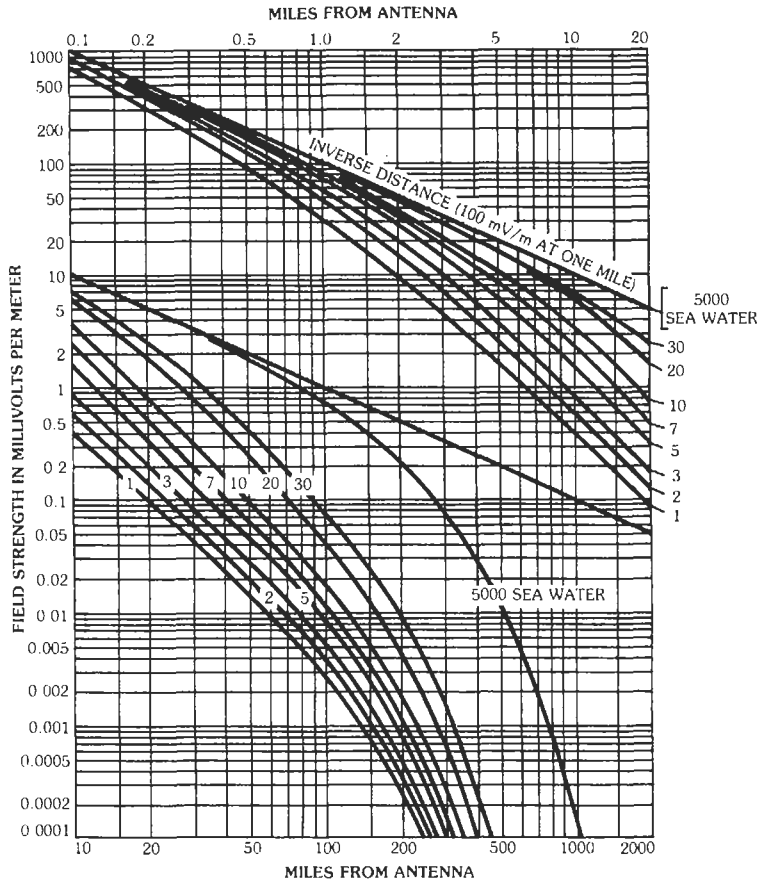


Fig. 3. Ground-wave field strength plotted against distance. Computed for 1600 kilohertz. Dielectric constant = 15. Ground-conductivity values are in millimhos/meter.

Independent Sideband: The upper and lower sidebands of the carrier are modulated with the L and R signals, respectively (Kahn system).

Quadrature: Two phase-locked carriers are amplitude modulated and combined at a fixed phase angle (Motorola system).

A "marketplace approach" has been employed in the US to select the best system, and that manufacturer's technique will become the *de facto* standard. In Canada and Japan, the Motorola system has been standardized, and most US stations also employ this system.

Frequency-Modulation Broadcasting*

Frequency-modulation (FM) broadcast stations are authorized for operation on 101 allocated channels, each 200 kHz wide, extending consecutively from channel 200 on 87.9 MHz to channel 300 on 107.9

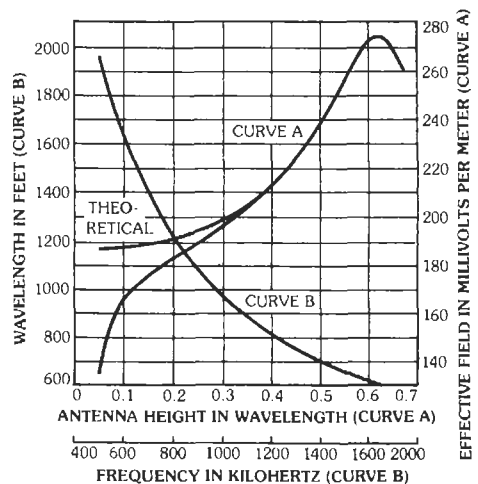


Fig. 4. Effective field at 1 mile for 1 kilowatt (curve A). Use for simple omnidirectional vertical antenna with ground system of at least 120 radials $\lambda/4$ long. From FCC Rules and Regulations.

* FCC Rules and Regulations, Part 73, Subparts B and C.

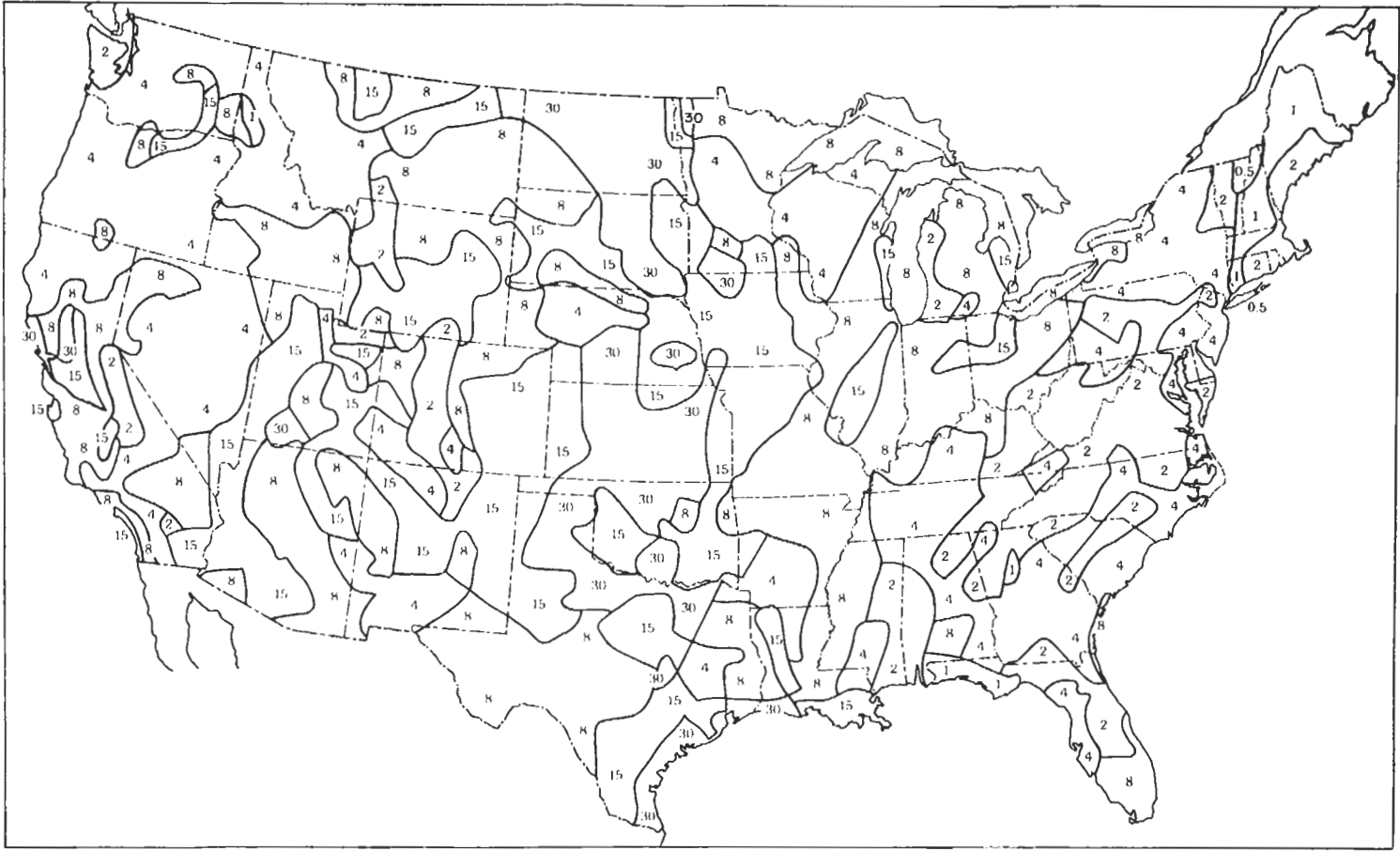


Fig. 5. Estimated effective ground conductivity in the United States. The numbers are in millimhos/meter. The conductivity of sea water (not shown) is assumed to be 5000 millimhos/meter.

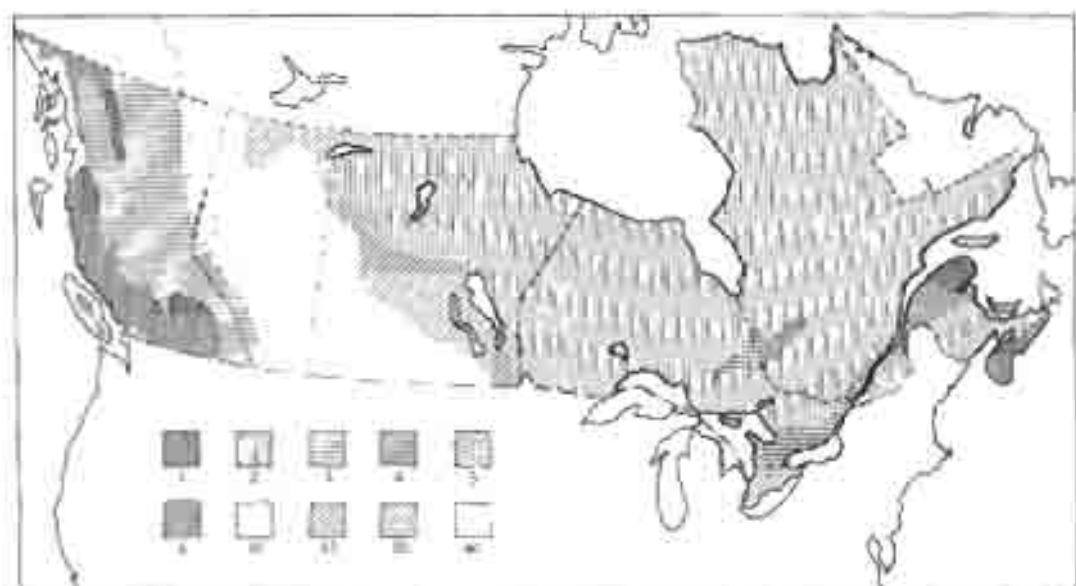


Fig. 6. Estimated effective ground conductivity in Canada.

MHz. Commercial broadcasting is authorized on channels 221 (92.1 MHz) through 300. Noncommercial educational broadcasting is permitted on any channel, but channels 200 through 220 are reserved for such use (see Part 73 of the FCC Rules and Regulations).

Station Classification—An FM station is classified by type in accordance with its authorized maximum operating power and antenna height (Table 2). A "reference distance" has been established for each class to denote the distance to the theoretical limit of its service area (the 1 mV/m contour—see "Coverage" below). A station operating with an antenna height in excess of the maximum for its class must reduce its effective radiated power (ERP) so that its 1 mV/m contour does not exceed the reference distance.

Channel Availability—Channels 221 through 300 are assigned to specific communities in accordance with the Table of Assignments (see Section 73.202 of FCC Rules). The number (and type) of channels assigned to a community is based on population, the limits imposed by cochannel and adjacent-channel interference considerations, and the area in which the community lies. The Table of Assignments may be amended to add channels upon approval of a petition for rulemaking which includes a showing of need and conformity with the separation standards shown in Table 3. Noncommercial stations operating on channels 217-220 must conform with Table 2 with respect to stations operating on channels 221-223. Stations operating on channels 200 through 216 are authorized based on mutual interference standards in lieu of a Table of Assignments. These standards are expressed in terms of the following interference ratios of desired to undesired signals.

Cochannel:	10:1, or 20 dB
1st Adjacent Channel:	2:1, or 6 dB
2nd Adjacent Channel:	1:10, or -20 dB
3rd Adjacent Channel:	1:100, or -40 dB

The ratios are the magnitude of the field strength from the undesired station, determined from the F50,10 curves* at the limit of the protected contour† of the desired station (see "Coverage" below).

Channel 200 is available only for Class D (maximum transmitter output of 10 watts) use in areas not served by a television broadcast station operating on channel 8. Stations on channels 201 through 220 are also subject to

* The protected contour is normally the 1 mV/m or 60 dBu service contour as determined by the F50,50 propagation curves. However, class B stations receive slightly greater protection. Refer to FCC Rules.

† Refer to Fig. 11 in the section on television broadcasting.

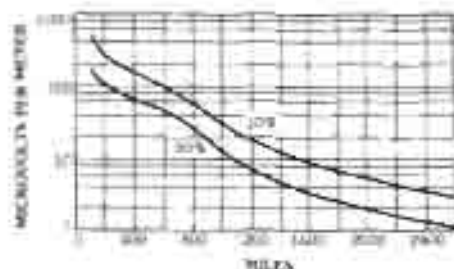


Fig. 7. Sky-wave signals for 10% and 50% of the time. The sky-wave range for frequencies from 540 to 1600 kilohertz is based on a radiated field of 100 millivolts per meter at 1 mile at the pertinent vertical angle.

TABLE 2. CLASSES OF FM STATIONS

Station Class ¹	Maximum ERP	Reference HAAT ² in Meters (ft.)	Class Contour Distance in Kilometers
A	6 kW (7.8 dBk)	100 (328)	28
B1	25 kW (14.0 dBk)	100 (328)	39
B	50 kW (17.0 dBk)	150 (492)	52
C3	25 kW (14.0 dBk)	100 (328)	39
C2	50 kW (17.0 dBk)	150 (492)	52
C1	100 kW (20.0 dBk)	299 (981)	72
C	100 kW (20.0 dBk)	600 (1968)	92

- Class B stations are available in Zone I and Class C stations in Zone II. See FCC Rules for zone boundaries.
- HAAT is height of the antenna radiation center above average terrain elevation. Average terrain is determined by a procedure set forth in the FCC Rules.

maximum power restrictions if located in or near areas served by a television station operating on channel 6.

Coverage—The estimated coverage of an FM station can be determined from the FCC F50.50 curves shown in Fig. 8.* These curves permit the prediction of

the field strength that will be exceeded for 50% of the time at 50% of the locations within the contour distance from the transmitter for various combinations of effective radiated power and antenna height above average terrain.

Station Performance—The performance of an FM broadcast station is required to be maintained in accordance with the following general specifications (see FCC

* These curves are approximate only. Official FCC data should be used for engineering designs.

TABLE 3. MINIMUM DISTANCE SEPARATION REQUIREMENTS IN KILOMETERS (MILES)

Relation	Co-channel	200 kHz	400/600 kHz	10.6/10.8 MHz
A to A	115 (71)	72 (45)	31 (19)	10 (6)
A to B1	143 (89)	96 (60)	48 (30)	12 (7)
A to B	178 (111)	113 (70)	69 (43)	15 (9)
A to C3	142 (88)	89 (55)	42 (26)	12 (7)
A to C2	166 (103)	106 (66)	55 (34)	15 (9)
A to C1	200 (124)	133 (83)	75 (47)	22 (14)
A to C	226 (140)	165 (103)	95 (59)	29 (18)
B1 to B1	175 (109)	114 (71)	50 (31)	14 (9)
B1 to B	211 (131)	145 (90)	71 (44)	17 (11)
B1 to C3	175 (109)	114 (71)	50 (31)	14 (9)
B1 to C2	200 (124)	134 (83)	56 (35)	17 (11)
B1 to C1	233 (145)	161 (100)	77 (48)	24 (15)
B1 to C	259 (161)	193 (120)	105 (65)	31 (19)
B to B	241 (150)	169 (105)	74 (46)	20 (12)
B to C3	211 (131)	145 (90)	71 (44)	17 (11)
B to C2	241 (150)	169 (105)	74 (46)	20 (12)
B to C1	270 (168)	195 (121)	79 (49)	27 (17)
B to C	274 (170)	217 (135)	105 (65)	35 (22)
C3 to C3	153 (95)	99 (62)	43 (27)	14 (9)
C3 to C2	177 (110)	117 (73)	56 (35)	17 (11)
C3 to C1	211 (131)	144 (90)	76 (47)	24 (15)
C3 to C	237 (147)	176 (109)	96 (60)	31 (19)
C2 to C2	190 (118)	130 (81)	58 (36)	20 (12)
C2 to C1	224 (139)	158 (98)	79 (49)	27 (17)
C2 to C	249 (155)	188 (117)	105 (65)	35 (22)
C1 to C1	245 (152)	177 (110)	82 (51)	34 (21)
C1 to C	270 (168)	209 (130)	105 (65)	41 (25)
C to C	290 (180)	241 (150)	105 (65)	48 (30)

Note: Different criteria apply to stations within 320 km of US-Canadian and US-Mexican borders. Consult FCC Rules and Regulations, section 73.207.

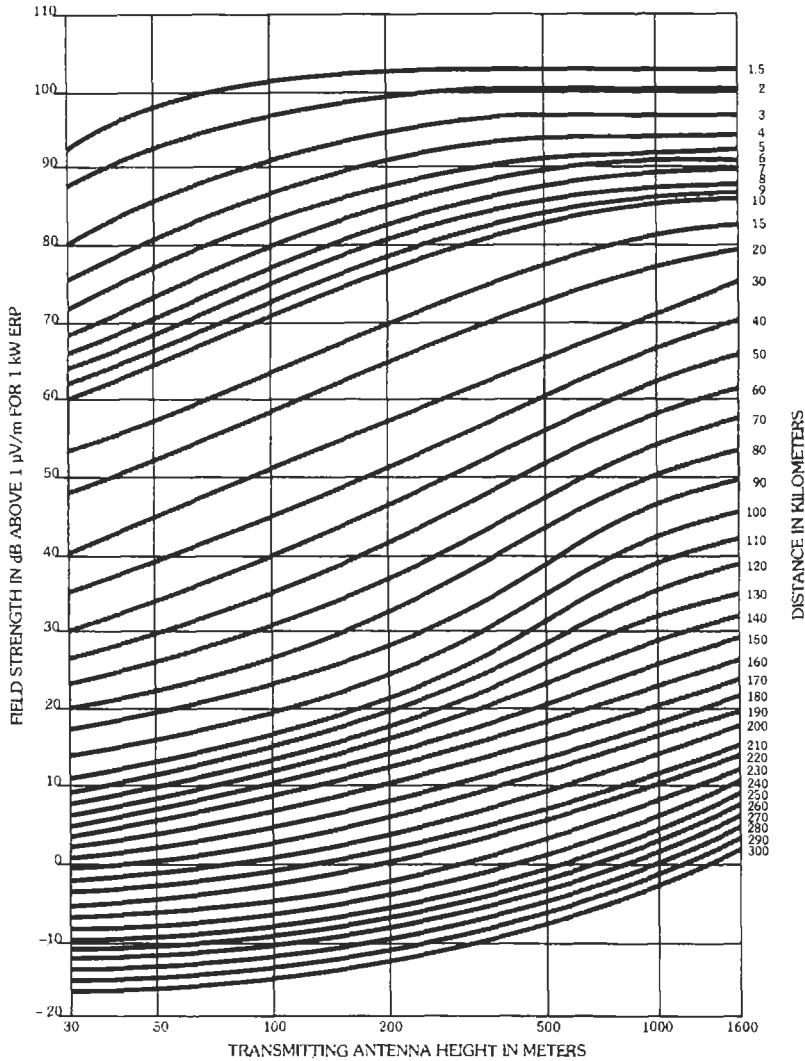


Fig. 8. Estimated field strength exceeded at 50 percent of the potential receiver locations for at least 50 percent of the time (FM channels, receiving antenna height 9 meters). After FCC Rules and Regulations.

Rules Section 73.317 for additional details). Key system criteria include:

- Modulation Capability: ± 75 kHz deviation of carrier for 100% modulation.
- Audio-Frequency Response: Within limits shown in FCC Rules for 75 microseconds pre-emphasis.
- Polarization: Typically right-hand circular; elliptical or other modes acceptable provided horizontal polarization is dominant.††

Stereophonic Transmission—Most FM stations broadcast in a stereophonic mode that is compatible

with monophonic receivers. The stereophonic signal consists of a main channel and a subchannel. The main channel is modulated by the sum of the left (L) and right (R) stereo signals. The subchannel consists of sidebands of a 38-kHz suppressed carrier that is modulated with an L minus R (difference) signal. Since the modulating frequencies can be as high as 15 kHz, the sidebands occupy the baseband spectrum from 23 kHz to 53 kHz (Fig. 9). A pilot carrier of 19 kHz is also transmitted to allow the receiver to generate a phase-locked 38-kHz carrier in order to demodulate the L-R signal. A matrix in the stereo receiver recovers the L and R signals by algebraic operation on the L-plus-R and L-minus-R signals. A monophonic receiver demodulates only the L-plus-R (main channel) signal.

†† An exception to this requirement permits noncommercial stations operating on channels 201-220 to use vertical polarization to reduce interference to television channel 6.

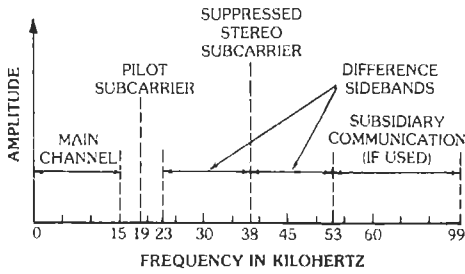


Fig. 9. Resulting stereophonic frequency spectrum.

A recent development in FM stereo modulation called FMX Stereo* employs a compressed L-R signal which modulates a 38-kHz suppressed carrier in a manner similar to the normal stereo difference signal. This signal is transmitted in quadrature to the main L-R signal and is, therefore, transparent to ordinary FM stereo decoders. Receivers equipped with an FMX decoder can demodulate and expand the compressed difference signal for an overall signal-to-noise improvement of 12–15 dB. This improvement helps to mitigate the 17–20 dB signal-to-noise penalty suffered in normal stereo transmission relative to monophonic FM transmission.

Other FM Services—A Subsidiary Communications Authorization (SCA) permits the addition of one or more subcarriers on the main channel; these subcarriers may be used for the transmission of voice or data signals related or unrelated to the broadcast-station operation. The subcarrier is amplitude modulated and is bandwidth limited. (See FCC Rules Section 73.319).

FM Translators—The FCC permits the use of very low power (1 watt east of the Mississippi river and 10 watts west) translators to rebroadcast the signals of FM stations in areas where no reception is possible. The translator station receives the main FM signal “off the air” and rebroadcasts it on a different channel; other methods of program delivery are permissible. Coverage range can be from one to five miles depending on power, antenna height, antenna gain, and terrain. Channels are assigned on a secondary basis (noninterference with operating stations). (See Subpart L, Part 74 of the FCC Rules and Regulations.)

Television Broadcasting (VHF and UHF)

Channel Designations—Television broadcast stations are authorized for commercial and educational operation on the channels shown in Table 4. Assignment of channels to specific communities is made by the FCC, and the channel assignments are designated as

commercial (unreserved) or educational (reserved). Actual carrier frequencies are assigned by a plan for carrier offsets to reduce cochannel interference (see “Transmission Standards”).

Coverage Data—The channel assignments have been made in such a manner as to facilitate maximum interference-free coverage in the available frequency bands. The radiated power of a particular station is fixed by several considerations. Minimum power is 100 watts visual effective radiated power. No minimum antenna height is specified.

Except as limited for stations using antenna heights in excess of 300 meters (600 meters for channels 14–69) in Zone I and antenna heights in excess of 600 meters in Zones II and III, the maximum visual power in decibels above 1 kilowatt (dBk) is:

Channels	Maximum Power
2–6	20 dBk = 100 kilowatts
7–13	25 dBk = 316 kilowatts
14–69	37 dBk = 5000 kilowatts

Zone I is the Northeastern United States. Zone II includes Puerto Rico, Alaska, the Hawaiian Islands, the Virgin Islands, and other parts of the US not in Zones I or III. Zone III is essentially a strip along the southeastern border of the US from Florida to Texas. Zone III cochannel separations reflect the propagation anomalies experienced in the Gulf Coastal Region. Detailed descriptions of the zones are in the FCC Rules and Regulations, Section 73.609. Stations operating with heights in excess of the class/zone maxima must reduce power in accordance with the requirements of 73.614 of the FCC Rules and Regulations.

Grades of service are designated A and B. The signal strengths in decibels above 1 microvolt/meter (dBu) specified for the various channels are:

Channel	Grade A	Grade B
2–6	68 dBu	47 dBu
7–13	71 dBu	56 dBu
14–69	74 dBu	64 dBu

The transmitter location must be so chosen that, with the effective radiated power and antenna height proposed, the minimum field strength (sometimes referred to as “city grade”) predicted over the principal community to be served will be:

Channel	Signal
2–6	74 dBu
7–13	77 dBu
14–83	80 dBu

The curves of Figs. 10 through 14* give estimated 50% and 10% time-varying field strengths for VHF and

* FMX Stereo is a registered trademark of Broadcast Technology Partners.

* The curves in Figs. 10 through 14 are approximate only. Official FCC data should be used for engineering designs.

TABLE 4. NUMERICAL DESIGNATION OF TELEVISION CHANNELS

Channel Number	Band (megahertz)	Channel Number	Band (megahertz)	Channel Number	Band (megahertz)
2	54-60	25	536-542	48	674-680
3	60-66	26	542-548	49	680-686
4	66-72	27	548-554	50	686-692
5	76-82	28	554-560	51	692-698
6	82-88	29	560-566	52	698-704
7	174-180	30	566-572	53	704-710
8	180-186	31	572-578	54	710-716
9	186-192	32	578-584	55	716-722
10	192-198	33	584-590	56	722-728
11	198-204	34	590-596	57	728-734
12	204-210	35	596-602	58	734-740
13	210-216	36	602-608	59	740-746
14	470-476	37	608-614	60	746-752
15	476-482	38	614-620	61	752-758
16	482-488	39	620-626	62	758-764
17	488-494	40	626-632	63	764-770
18	494-500	41	632-638	64	770-776
19	500-506	42	638-644	65	776-782
20	506-512	43	644-650	66	782-788
21	512-518	44	650-656	67	788-794
22	518-524	45	656-662	68	794-800
23	524-530	46	662-668	69	800-806
24	530-536	47	668-674		

The frequencies between 806 and 890 MHz, formerly allocated to television broadcasting, are now allocated to the land mobile services. Operation, on a secondary basis, of some television translators may continue on these frequencies.

UHF television channels at different heights and powers. The F(50,50) curves are used for the prediction of service, while the F(50,10) curves are used for the prediction of interference. Cochannel TV stations operating with offset carriers are afforded desired-to-undesired interference protection of 28 dB, i.e., F(50,50) - F(50,10) = 28 dB. The antenna height is the height of the radiation center of the antenna above the average terrain. "Average terrain" is determined by averaging the ground elevations between 3 and 16 km from the antenna site, taken along eight radials separated by 45° in azimuth. Effective radiated power is the product of the antenna gain and the antenna input power. Antenna input power is the peak visual output power of the transmitter less transmission-line and diplexer losses. The procedures to be followed in determining the effective radiated power used in the prediction of coverage are detailed in Section 73.684 of the FCC Rules and Regulations.

Directional antennas may be employed to improve service. The ratio of maximum to minimum radiation in the horizontal plane may not exceed 10 dB for channels 2-13, and 15 dB for channels 14-69 if the transmitter power is more than 1 kW. There is no restriction for channels 14-69 if the transmitter power is 1 kW or less.

Transmission Standards—The standards for television transmission in the US, as defined by the FCC, are:

Channel Width: 6 MHz

Picture Carrier Location: 1.25 MHz \pm 1000 Hz above lower boundary of the channel.*

Aural Center Frequency: 4.5 MHz \pm 1000 Hz above visual carrier.

Polarization of Radiation: Horizontal normal, right-hand circular optional.

Modulation: Amplitude-modulated composite picture and synchronizing signal on visual carrier, together with frequency-modulated audio signal on aural carrier. (See Figs. 15 and 16.)

Scanning Lines: 525 lines/frame interlaced two to one.

Scanning Sequence: Horizontally from left to right, vertically from top to bottom.

Horizontal Scanning Frequency: 2/455 times chrominance subcarrier frequency (15 734 Hz).

Vertical Scanning Frequency: 2/525 times the horizontal scanning frequency (59.94 Hz).

Chrominance Subcarrier Frequency: 3.579 545 MHz \pm 10 Hz.

Polarity of Transmission: Negative—a decrease in initial light intensity causes an increase in radiated power.

Transmitter Brightness Response: For luminance sig-

* Carrier frequencies may be offset \pm 10 kHz (e.g., 1.24 or 1.26 MHz above band edge) for cochannel-interference considerations.

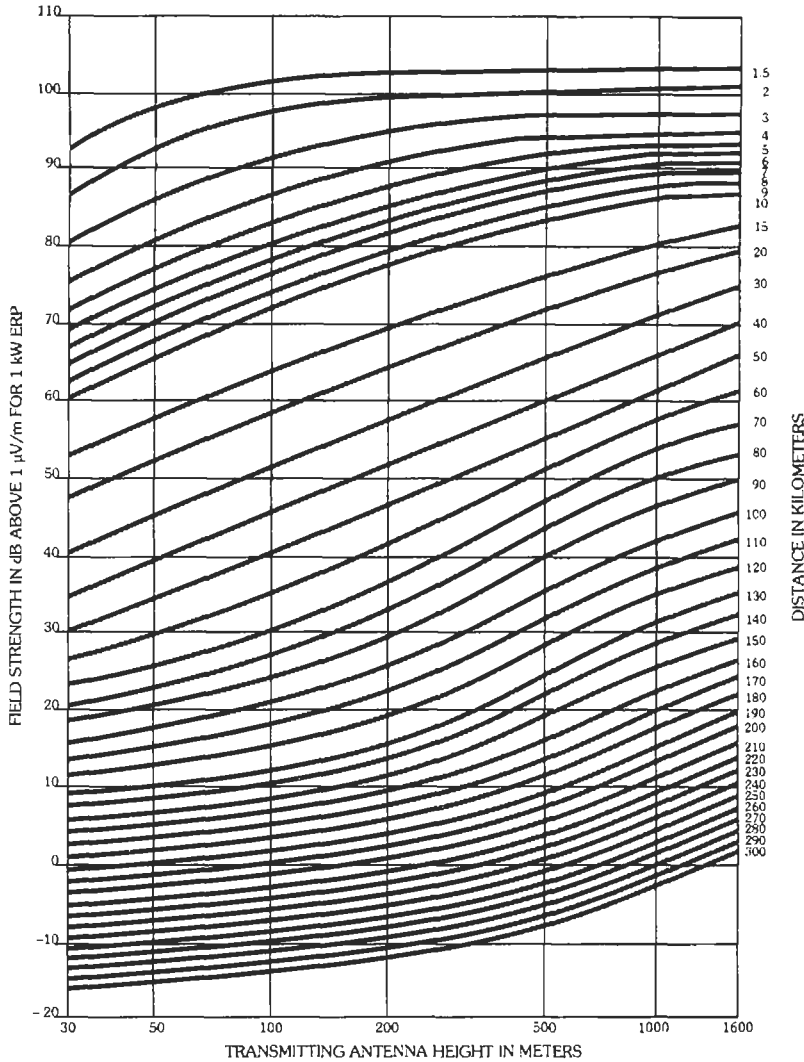


Fig. 10. Estimated field strength exceeded at 50 percent of the potential receiver locations for at least 50 percent of the time for TV channels 2-6. Receiving antenna height, 9 meters. After FCC Rules and Regulations.

nal, radio-frequency output varies in an inverse logarithmic relation to the brightness of the scene.

Aural-Transmitter Power: Maximum radiated power is 22% of peak visual effective radiated power.

For color transmission, the luminance component is transmitted as amplitude modulation of the picture carrier and the chrominance components as amplitude-modulation sidebands of a pair of suppressed subcarriers in quadrature (Fig. 17).

The interval beginning with line 17 and continuing through line 20 of the vertical-blanking interval of each field may be used for the transmission of data, test signals, and cue and control signals. Test signals may

include signals designed to check the performance of the overall transmission system or its individual components. Test signals or cue and control signals may not be transmitted during that portion of each line devoted to horizontal blanking. Line 19 in each field may be used only for transmission of the standard vertical interval reference (VIR) signal. Line 21 may be used for transmission data such as closed captioning for the deaf. Lines 10-20 may be used for data and text transmission (Teletext) except as noted above.

Aural Transmitter—The TV audio signal is frequency modulated on a carrier 4.5 MHz above the visual carrier. The power level of the carrier may not exceed 22% of the peak visual power. The carrier is

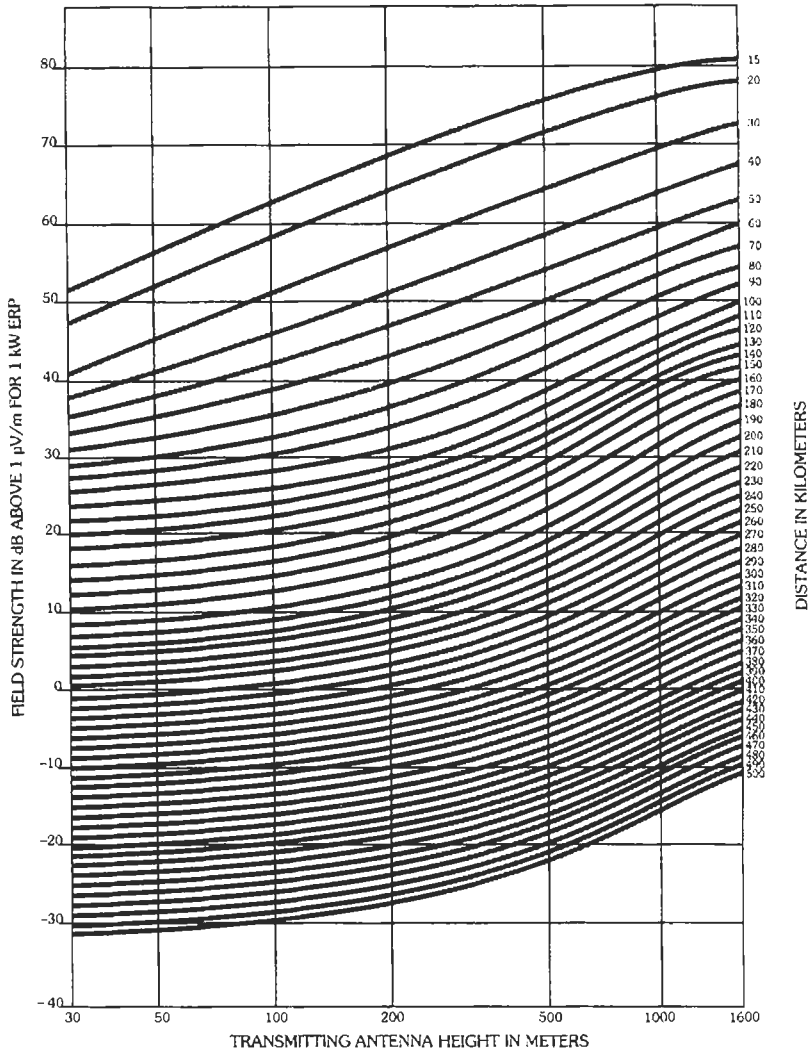


Fig. 11. Estimated field strength exceeded at 50 percent of the potential receiver locations for at least 10 percent of the time for TV channels 2-6 and FM. Receiving antenna height, 9 meters. *After FCC Rules and Regulations.*

deviated 25 kHz for 100% modulation, and the same 75-microsecond pre-emphasis specified for FM broadcasting is used. A multichannel television sound system (MTS) may be employed to provide stereophonic and other audio services. The common system in use is the BTSC (Broadcast Television Sound Committee) format which provides a compatible (L+R) signal, a companded stereo difference signal, a separate audio program (SAP) channel, and a utility audio-data channel (PRO channel). Total aural baseband information must be maintained within a 120-kHz bandwidth.

Cable Television

Cable television (CATV) systems were originally called "community antenna television systems." They

had their origin in mountainous regions of Pennsylvania where "off-air" reception of broadcast television stations was difficult, if not impossible. Beginning with crude antennas and amplifiers that delivered a few channels to homes in sparsely populated valleys, CATV systems have evolved into a broadband technology utilizing high-capacity coaxial cables to carry signals in the spectrum from 5 MHz to over 500 MHz and fiber-optic cables with seemingly limitless capacity.

Some modern systems employ multiple cables for separate subscriber and institutional networks, each with two-way capability. Cable TV systems have developed into carriers of more than "off-air" programming from broadcast stations. Modern urban cable systems carry more television signals emanating from satellites and local sources than from broadcast stations.

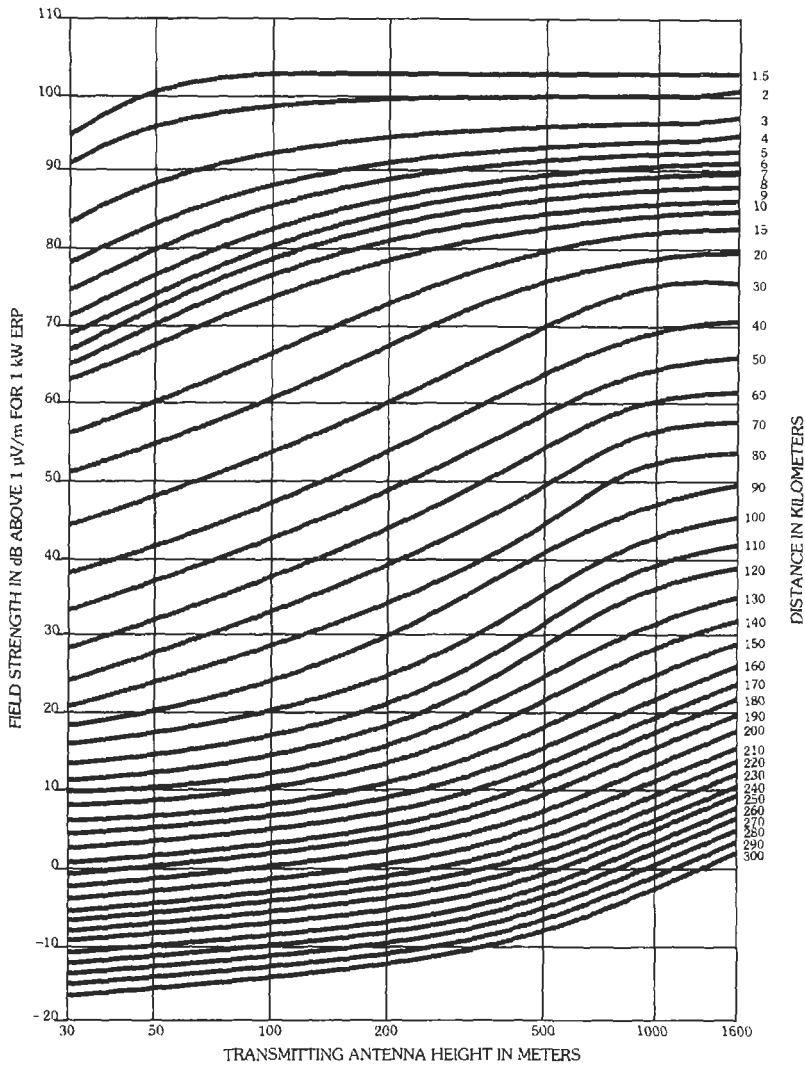


Fig. 12. Estimated field strength exceeded at 50 percent of the potential receiver locations for at least 50 percent of the time for TV channels 7-13. Receiving antenna height, 9 meters. *After FCC Rules and Regulations.*

Converters are used at subscriber locations to change nonstandard channel frequencies to a standard VHF channel (usually between 2 and 6 inclusive) so that ordinary television receivers can demodulate the signals. Most new receivers are being manufactured with "cable ready" tuners capable of tuning more than one hundred channels. The cable channel frequency assignments generally used are given in Table 5.

Bidirectional cable systems normally employ a low-band split (i.e., downstream channels above 50 MHz and upstream channels below 50 MHz) or a midband split with the division occurring near 150 MHz. As the upper limit of usable cable-system frequencies increases, the frequency where the split occurs has been adjusted.

Cable systems are not regulated to any great extent by the FCC,* and the technology has far exceeded the minimum performance requirements set by the Commission. Of particular regulatory concern, however, is the unwanted leakage from poorly constructed and maintained CATV systems of radiations such as spurious signals in radio-frequency bands used for services crucial to the safety of life. The Commission has restricted CATV operation on frequencies which coincide with aviation services (108-137 MHz) and other governmental/public-safety bands.

Some services delivered by CATV systems are tiered to varying levels of subscription fees, and access to the

* FCC Rules and Regulations, Part 76.

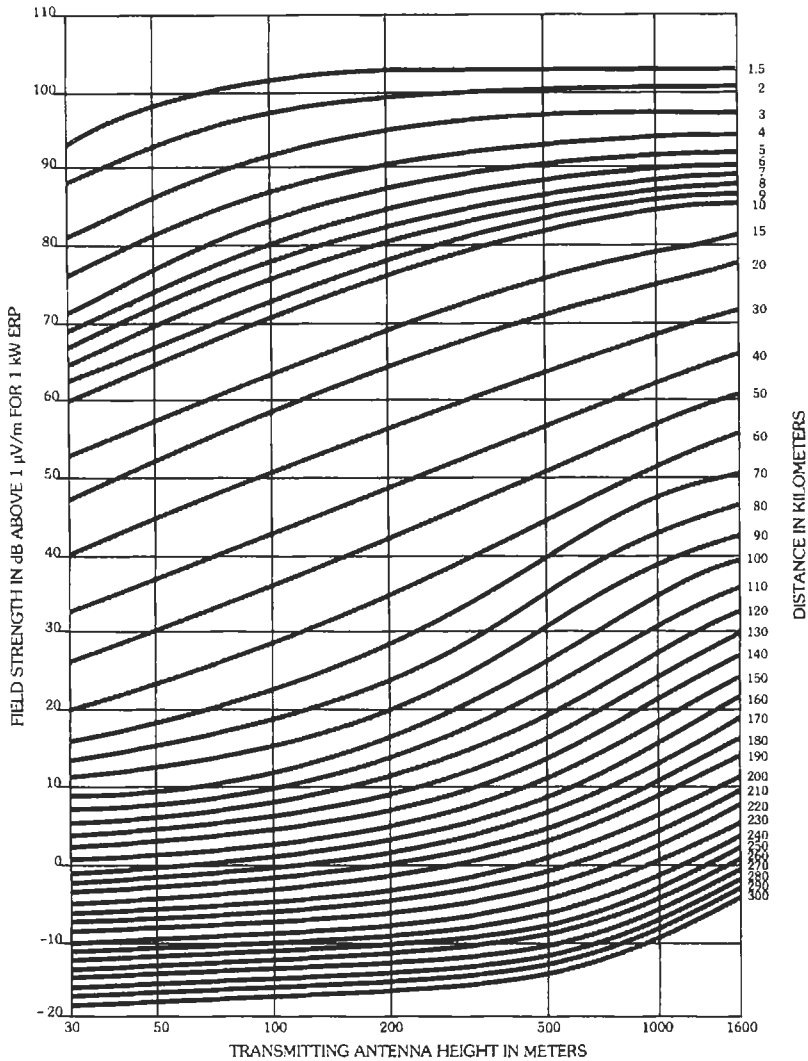


Fig. 13. Estimated field strength exceeded at 50 percent of the potential receiver locations for at least 50 percent of the time for TV channels 14-69. Receiving antenna height, 9 meters. After FCC Rules and Regulations.

higher-level tiers is limited by the use of scrambled signals and addressable converters and decoders. Newer technology incorporates provisions for "pay-per-view" access to special-events programming.

One of the limitations on the number of channels a cable system can carry is the distortion due to intermodulation and cross-modulation products resulting from system nonlinearities and the numerous carriers. *Composite triple beat* is a term applied to third-^{*} and higher-order products, which can greatly degrade the performance of a system. One method used to reduce these effects is the employment of coherent oscillators

to regenerate (and convert) carrier frequencies such that they are harmonically related (HRC). This intermodulation problem increases with amplifier output levels and builds as amplifiers are cascaded. This then becomes the limiting factor on the length of a system (number of amplifiers in cascade), and, as the number of channels is increased, the amplifier separation must be reduced to compensate for lowered output levels. Extension of systems and system interconnection into large networks are also limited by this constraint.

Fiber-optic, or light-wave cable, systems are being developed which promise to improve many of the transmission characteristics of coaxial-cable (broadband) amplifier systems. While there are few operating CATV systems delivering signals directly to subscribers

* Products of the general type $(2f_1 \pm f_2)$ or $(f_1 + f_2 \pm f_3)$.

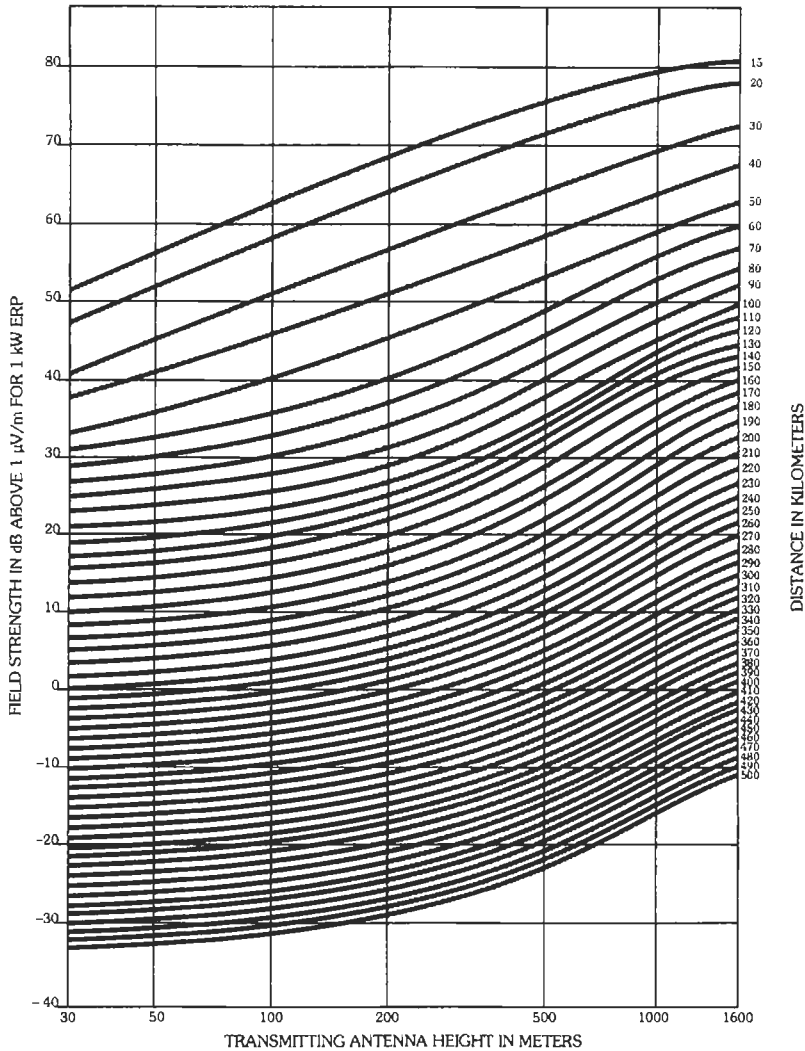


Fig. 14. Estimated field strength exceeded at 50 percent of the potential receiver locations for at least 10 percent of the time for TV channels 14-69. Receiving antenna height, 9 meters. After FCC Rules and Regulations.

via fiber optics, this technology is being employed to deliver signals to headends or as "super trunks" to deliver signals to multiple miniheadends in many systems. However, coaxial-cable hardware development has advanced to the point where it has been the most practical method of constructing economical broadband systems; fiber-optic technology has emerged from the developmental stage for this application, and greater use will be seen.

As this section was being written, the FCC had initiated regulatory action which would permit telephone companies to deliver video and wideband information services. If these companies decide to get into this business, they will employ fiber-optic systems, and this could eventually lead to the installation of light-

wave capability into most homes and businesses for all communications.

Other Television Services

Translators and LPTV—Television broadcast translators and low-power television (LPTV) stations operate under Subpart G of Part 74 of the FCC Rules and Regulations. A television broadcast translator is a station that rebroadcasts the programming of a broadcast station in an essentially unaltered form. The signals from the originating station may be delivered by direct off-the-air pickup, microwave relay, satellite transmission, or rebroadcast of another translator. An LPTV station is similar to a translator station in terms of

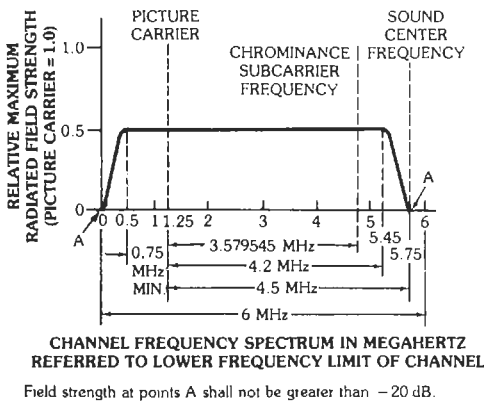


Fig. 15. Radio-frequency amplitude characteristics of television picture transmission. (Drawing not to scale.)

equipment and service area, but it is permitted to originate programming from virtually any source. The technical rules governing translators and LPTV stations are essentially identical.

For VHF stations, the transmitter output power (peak visual) may be up to 10 watts unless the station operates on an unoccupied channel assigned to its community for regular television broadcast use; in this case, 100 watts may be employed. A UHF station may be authorized on any channel with transmitter power up to 1 kW. There is no limitation on effective radiated power for either type of station. If stations employ circular polarization, authorized transmitter power may be doubled for all classes.

The technical standards for operation of these stations are generally the same as those for broadcast stations (A5/F3 modulation) with the principal exceptions relating to carrier-frequency tolerances, spurious emissions, and lower-sideband attenuation characteristics.

Channel assignments are made on the assumption that these facilities are secondary to broadcast stations and must protect existing facilities from objectionable interference. While the service area of a regular broadcast station is defined as the area within the predicted Grade B contour, the service area of an LPTV station is much smaller and generally equivalent to the area within a predicted contour approximating its Grade A service area. Existing LPTV stations are protected from interference from newly proposed LPTV stations.

High Definition Television—High definition television (HDTV) system proposals abound as of the writing of this text. The FCC is considering various systems through the Advanced Television System Committee and the Advanced Television Test Center. The latter is conducting laboratory and field tests on the systems, while the former group is dealing with the development of HDTV transmission standards. A 1993 target date has been established for the completion of the work. Major issues include compatibility of any new

system with existing NTSC systems, channel bandwidth requirements, modulation (analog/digital), signal propagation anomaly correction methodologies, and the impact of emerging international standards.

ITFS—Instructional Television Fixed Service (ITFS) stations operate under Subpart I of Part 74 of the FCC Rules in the band 2500–2686 MHz. These stations are licensed to eligible educational entities for the distribution of program material to students enrolled in instructional curricula. Public broadcast stations rendering such services are also eligible. An additional 4 MHz (2686–2690 MHz) is assigned for ITFS response stations. Response channels are intended for use as return links to the originating point for aural infor-

Notes for Fig. 16:

1. H = time from start of one line to start of next line.
2. V = time from start of one field to start of next field.
3. Leading and trailing edges of vertical blanking should be complete in less than $0.1H$.
4. Leading and trailing slopes of horizontal blanking must be steep enough to preserve minimum and maximum values of $(x + y)$ and z under all conditions of picture content.
5. Dimensions marked with an asterisk indicate that tolerances given are permitted only for long-time variations, and not for successive cycles.
6. Equalizing-pulse area shall be between 0.45 and 0.5 of the area of a horizontal-synchronizing pulse.
7. Color burst follows each horizontal pulse, but is omitted following the equalizing pulses and during the broad vertical pulses.
8. Color bursts to be omitted during monochrome transmission.
9. The burst frequency shall be 3.579545 megahertz. The tolerance on the frequency shall be ± 10 hertz with a maximum rate of change of frequency not to exceed $1/10$ hertz per second.
10. The horizontal scanning frequency shall be $2/455$ times the burst frequency.
11. The dimensions specified for the burst determine the times of starting and stopping the burst but not its phase. The color burst consists of amplitude modulation of a continuous sine wave.
12. Dimension P represents the peak excursion of the luminance signal from blanking level but does not include the chrominance signal. Dimension S is the synchronizing amplitude above blanking level. Dimension C is the peak carrier amplitude.
13. Refer to FCC standards for further explanations and tolerances.
14. Horizontal dimensions not to scale in A , B , and C .

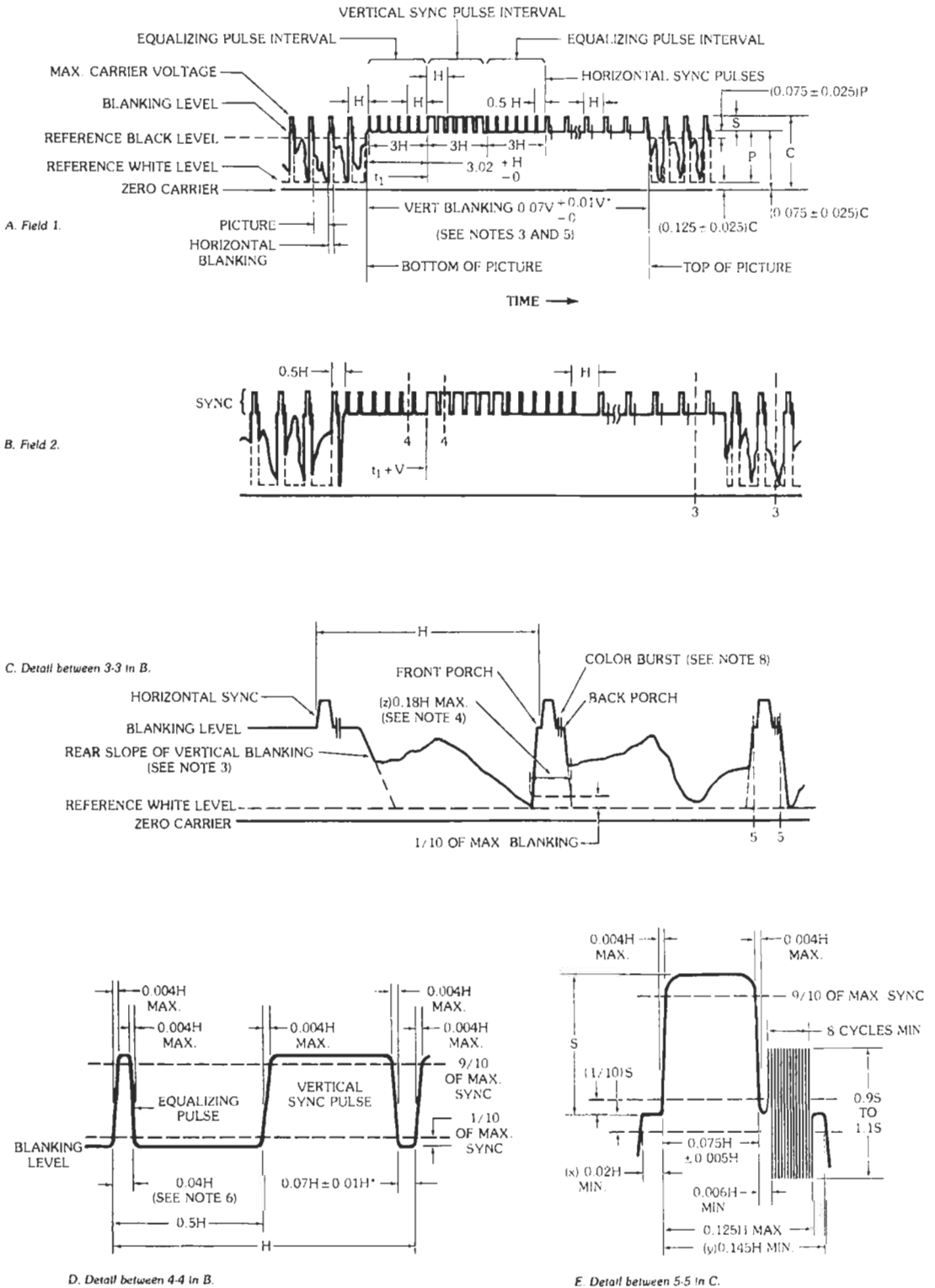


Fig. 16. Television composite-signal waveform data. (See notes on facing page.)

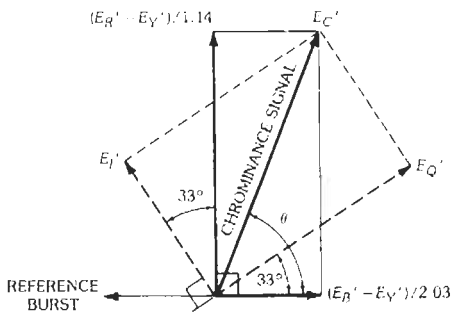


Fig. 17. Phases of color signal.

mation/data such as student questions and responses. These stations employ very low power and narrowband equipment (250 mW at 125 kHz). The ITFS channels are six megahertz wide and have a standard television broadcast signal format. The channels are presently grouped in bands of four alternately spaced channels (e.g., group A is composed of A-1, 2500–2506 MHz; A-2, 2512–2518 MHz; A-3, 2524–2530, etc.), with adjacent groups interleaved as shown in Table 6. Licensees are permitted to use multiple channels (up to four) based on need and availability. Stations are usually limited to an equivalent isotropically radiated power (EIRP) of 33 dBw unless a directional antenna is employed. The EIRP is further limited by coverage requirements and considerations of interference to other stations in the band (some “grandfathered” facilities of other services remain in the ITFS band). The band is also shared with the broadcast satellite service, and protection to that service is afforded by the Rules.

Reception of ITFS signals usually requires the use of a small parabolic antenna, which is coupled to a down-converter. The down-converter incorporates a low-noise front end ($NF < 5$ dB) and converts the signal(s) to VHF channels. A block of four channels (group) can be converted in a single device to yield four alternately spaced VHF channels (e.g., 7, 9, 11, 13) or wideband devices can down-convert all 31 channels to VHF “cable” channels in a single block. Four-foot antennas are generally used at distances up to ten miles for received carrier-to-noise ratios of approximately 50 dB.

MDS—The Multipoint Distribution Service (MDS) is a common-carrier service but is included in this chapter on broadcasting because the service is used primarily for the distribution of television program material to subscribers (“pay TV”). It is governed under Subpart K of Part 21 of the FC Rules. This service is assigned 10 MHz of spectrum in the band 2150–2160 MHz, except in the top 50 market areas where the band is extended to 2162 MHz. Channel 1 has “inverted” carrier relationships (i.e., aural carrier below visual carrier) with respect to channel 2 in order to minimize adjacent-channel interference in those locations where two channels are authorized. The band

2156–2160 MHz in smaller markets is too narrow to be used for standard television picture transmission but may be assigned for other purposes.

In all other aspects, MDS stations are very much akin to ITFS stations, including hardware. As in ITFS, reception is accomplished with a down-converter that changes the 2-GHz signal to a standard VHF channel. The higher power generally employed in the MDS service permits the use of smaller receiving antennas and higher-noise-figure down-converters, making home-receiver reception practical.

The FC has allocated groups E, F, and H to the MMDS (Multichannel MDS) service, which has evolved into an “over-the-air” or “wireless” cable service. ITFS licensees are permitted to lease excess capacity to MMDS operators, thus increasing the total channels potentially available to 31–33.

Network Distribution of Broadcast Program Signals

Terrestrial—Most network radio and television programming formerly was distributed by terrestrial microwave systems owned and operated by common carriers (see Chapter 38). Television program transmission is governed by the Network Transmission Committee Standard, NTC-7, which establishes reasonable requirements for a system that is presumed to have over 100 microwave relay stations in tandem. While most programming formerly was carried over facilities belonging to AT&T and other carriers, those facilities are being supplanted by satellite-based systems. The principal advantage of the terrestrial system was its infrastructure which provided many routing alternatives; its primary disadvantage was the poor performance of long repeater cascades (coast-to-coast) as compared to satellite transmission circuits.

Satellite Program Distribution—Space stations in the Domestic Satellite Service are being used extensively for the distribution of television programming. Space stations are located in the geostationary (or geosynchronous) orbit located approximately 23 000 miles in space above the equator. Satellites located between approximately 70° and 145° west longitude can be “seen” from most sites in the continental US except those in the most northerly latitudes. The Public Broadcasting Service (PBS), was the first to distribute its television programming exclusively by the use of satellite channels; its sister organization, National Public Radio (NPR), also distributes its programming almost exclusively by satellite. Now virtually all network TV program distribution is accomplished via satellite.

When a satellite system is viewed as a point-to-point microwave system with only one heterodyne repeater, it is readily apparent that vastly improved performance can be obtained over the terrestrial system if carrier-to-noise limitations can be overcome. With very-low-noise amplifiers (30–50 K), the weak signals from C-band

TABLE 5. CABLE TV CHANNEL FREQUENCIES

Channel	Freq. Range (MHz)	Carriers (MHz)		
		Video	Color	Sound
T- 7	5.75-11.75	7	10.58	11.5
T- 8	11.75-17.75	13	16.58	17.5
T- 9	17.75-23.75	19	22.58	23.5
T-10	23.75-29.75	25	28.58	29.5
T-11	29.75-35.75	31	34.58	35.5
T-12	35.75-41.75	37	40.58	41.5
T-13	41.75-47.55	43	46.58	47.5
2	54-60	55.25	58.83	59.75
3	60-66	61.25	64.83	65.75
4	66-72	67.25	70.83	71.75
5	76-82	77.25	80.83	81.75
6	82-88	83.25	86.83	87.75
7	174-180	175.25	178.83	179.75
8	180-186	181.25	184.83	185.75
9	186-192	187.25	190.83	191.75
10	192-198	193.25	196.83	197.75
11	198-204	199.25	202.83	203.75
12	204-210	205.25	208.83	209.75
13	210-216	211.25	214.83	215.75
FM	88-108	—	—	—
14	120-126	121.25	124.83	125.75
15	126-132	127.25	130.83	131.75
16	132-138	133.25	136.83	137.75
17	138-144	139.25	142.83	143.75
18	144-150	145.25	148.83	149.75
19	150-156	151.25	154.83	155.75
20	156-162	157.25	160.83	161.75
21	162-168	163.25	166.83	167.75
22	168-174	169.25	172.83	173.75
23	216-222	217.25	220.83	221.75
24	222-228	223.25	226.83	227.75
25	228-234	229.25	232.83	233.75
26	234-240	235.25	238.83	239.75
27	240-246	241.25	244.83	245.75
28	246-252	247.25	250.83	251.75
29	252-258	253.25	256.83	257.75

Continued on next page.

TABLE 5 (CONT). CABLE TV CHANNEL FREQUENCIES

Channel	Freq. Range (MHz)	Carriers (MHz)		
		Video	Color	Sound
30	258-264	259.25	262.83	263.75
31	264-270	265.25	268.83	269.75
32	270-276	271.25	274.83	275.75
33	276-282	277.25	280.83	281.75
34	282-288	283.25	286.83	287.75
35	288-294	289.25	292.83	293.75
36	294-300	295.25	298.83	299.75
37	300-306	301.25	304.83	305.75
38	306-312	307.25	310.83	311.75
39	312-318	313.25	316.83	317.75
40	318-324	319.25	322.83	323.75
41	324-330	325.25	328.83	329.75
42	330-336	331.25	334.83	335.75
43	336-342	337.25	340.83	341.75
44	342-348	343.25	346.83	347.75
45	348-354	349.25	352.83	353.75
46	354-360	355.25	358.83	359.75
47	360-366	361.25	364.83	365.75
48	366-372	367.25	370.83	371.75
49	372-378	373.25	376.83	377.75
50	378-384	379.25	382.83	383.75
51	384-390	385.25	388.83	389.75
52	390-396	391.25	394.83	395.75
53	396-402	397.25	400.83	401.75
54	72-78	73.25	76.83	77.75
55	78-84	79.25	82.83	83.75
56	84-90	85.25	88.83	89.75
57	90-96	91.25	94.83	95.75
58	96-102	97.25	100.83	101.75
59	102-108	103.25	106.83	107.75
60	108-114	109.25	112.83	113.75
61	114-120	115.25	118.83	119.75

satellites can be processed to produce signal-to-noise ratios well in excess of 55 dB with antennas 7 to 10 meters in diameter (depending on satellite EIRP at the particular receiving location).

The present domestic satellite system uses C-band transponders (3.7-4.2 GHz down-link and 5.925-

6.425 GHz up-link) and K_u-band (11/14 GHz). The C-band is shared with terrestrial common carriers on a coequal basis. Thus, satellite terminals must be coordinated with existing facilities to assure interference-free operation. For this reason, it is frequently difficult to coordinate C-band satellite earth stations in urban areas

TABLE 6. FREQUENCY ASSIGNMENTS OF INSTRUCTIONAL TELEVISION FIXED STATIONS

Channel	Band Limits (megahertz)
Group A	
A-1	2500-2506
A-2	2512-2518
A-3	2524-2530
A-4	2536-2542
Group B	
B-1	2506-2512
B-2	2518-2524
B-3	2530-2536
B-4	2542-2548
Group C	
C-1	2548-2554
C-2	2560-2566
C-3	2572-2578
C-4	2584-2590
Group D	
D-1	2554-2560
D-2	2566-2572
D-3	2578-2584
D-4	2590-2596
Group E*	
E-1	2596-2602
E-2	2608-2614
E-3	2620-2626
E-4	2632-2638
Group F*	
F-1	2602-2608
F-2	2614-2620
F-3	2626-2632
F-4	2638-2644
Group G	
G-1	2644-2650
G-2	2656-2662
G-3	2668-2674
G-4	2680-2686
Group H*	
H-1	2650-2656
H-2	2663-2668
H-3	2674-2680

* These frequencies shared with other services.

where terrestrial frequency congestion abounds. Techniques such as antenna shielding and IF filtering can sometimes be employed to resolve interference problems.

Satellites with K_u -band (11/14 GHz) transponders eliminate the interference problems while permitting the use of smaller antennas. Interference is not a problem since these bands are not shared with terrestrial

users; smaller antennas are usable because of the higher effective gain of both satellite transmitting antennas and earth-terminal receiving antennas. The disadvantages of the band include the much higher potential for signal degradation due to precipitation propagation losses. Distributed, interconnected earth stations can solve this problem. For nonbroadcast applications where less than 99.99% availability is acceptable, K_u -band has become the dominant service.

DBS—Direct Broadcast Satellite transmission has been approved by the FCC. The bands allotted are 12.2–12.7 GHz for down-link and 17.3–17.8 GHz for up-link. Various schemes have been proposed for channelization and modulation formats. Some proponents of high-definition television (HDTV) have proposed the use of DBS as the only available means of broadcasting such programming directly to the public. Direct Broadcast Satellite transmission will further improve on the K_u -band advantage (use of small antennas) by employing high-power transponders that will permit a reduction in antenna size to much less than one meter for practical home reception. It is predicted that as many as 20 DBS satellites could be in operation by the year 2000, offering over 200 HDTV channels for home reception. Digital video compression systems will play an important role in this application, and the development of a practical system could permit a substantial increase in the number of channels while decreasing transponder/bandwidth requirements.

Auxiliary Broadcast Services

In Part 74 of its Rules and Regulations, the FCC has made provision for various auxiliary broadcast services. These are generally channels allocated for delivering programming and associated communications from remote locations to the station. The bands and services include the following.

Remote Pickup Broadcast Stations—Stations in this service are used for the transmission of aural programming material and associated cues and data. Assigned frequencies are in bands at 153 MHz, 161 MHz, 166 MHz, 170 MHz, 450 MHz, and 455 MHz. (See Subpart D of Part 74, FCC Rules and Regulations for limitations on use.)

Aural Intercity and STL Stations—Stations in this service are to be used for relay of aural program material from studio to transmitter and between fixed facilities in other locations. Assigned frequencies are in the band 947–951 MHz. (See Subpart E of Part 74, FCC Rules and Regulations.)

Television Auxiliary Stations—Stations in this service are used for transmission of television programming (aural and visual) between studio and transmitter, for intercity relay, and for remote pickups. Various bands are available at 2 GHz, 7 GHz, 13 GHz, 18 GHz and 23 GHz. (See Subpart F of Part 74, FCC Rules and Regulations for specific channel frequencies and limitations.)

Low Power Auxiliary Stations—Stations in this service are intended for use over very short distances for audio, cues, control, etc., associated with broadcast programs (for example, wireless microphones). Frequency bands are at 26 MHz, 161 MHz, 174–216 MHz, 450 MHz, and 950 MHz. (See Subpart H of Part 74, FCC Rules and Regulations for specific frequency assignments and limitations.)

Frequency Sharing—All users of auxiliary broadcast services must share frequencies. Users must coordinate use of specific frequencies with other local broadcasters to eliminate interference problems.

International Broadcasting Service in the United States

Transmissions from international broadcasting stations located within the United States are intended to be received directly by the general public in foreign countries. International broadcasting is conducted from both government (Voice of America) and privately owned stations. Public Law 80-402, the *United States Information and Educational Exchange Act of 1948*, encourages the participation of the private sector in international broadcasting. Nongovernment international broadcasting stations are licensed by the Federal Communications Commission in accordance with Part 73, Subpart F of the FCC Rules and Regulations. A license for an international broadcasting station will be issued only after a satisfactory showing has been made in regard to the following items, among others:

1. That there is a need for the international broadcasting service proposed to be rendered.
2. That the necessary program sources are available to the applicant to render the international service proposed.
3. That the production of the program service and the technical operation of the proposed station will be conducted by qualified persons.
4. That the applicant is legally, technically, and financially qualified and possesses adequate technical facilities to carry forward the service proposed.
5. That the public interest, convenience, and necessity will be served through the operation of the proposed station.

International broadcasting stations employ frequencies in bands between 5950 and 26 100 kHz. Frequencies authorized by the FCC fall within the following bands, which are allocated exclusively for broadcasting:

Meter Band	Frequency in Kiloherzt
49	5950–6200
31	9500–9775
25	11 700–11 975
19	15 100–15 450
16	17 700–17 900
13	21 450–21 750
11	25 600–26 100

The band 7100–7300 kHz is also allocated for broadcasting, except in the western hemisphere.

The carrier frequencies assigned begin 5 kHz above the frequency specified for the beginning of each band, and are in successive steps of 5 kHz ending 5 kHz below the frequency specified as the end of each band.

Since international broadcasts cross frontiers, bridge oceans, and span continents, frequency assignments must be regulated internationally. Article 17 of the *Radio Regulations*, International Telecommunication Union, Geneva, 1979, specifies a frequency coordination procedure that member administrations must follow. The FCC has incorporated this procedure into its Rules and Regulations.

Schedules for international broadcasting stations are prepared seasonally and are implemented at 0100 UTC on the first Sunday of March, May, September, and November. The FCC requires licensees to submit their tentative schedules to the Commission six months prior to the start of each season, indicating for the season the frequency or frequencies desired for transmission to each zone or area of reception specified in the license; the specific hours of transmission to such zones or areas on each frequency; and the power, antenna gain, and antenna bearing to be used. The geographical areas to which broadcasts are directed should be designated in accordance with the ITU Geographical Zones for Broadcasting shown in Fig. 18.

Frequencies proposed should be as close as possible to the optimum working frequency, which is defined as that frequency which is returned to the surface of the earth for a specific transmission path and time of day on 90% of the days of the month, and should be chosen so that a given frequency will provide the largest period of reliable transmission to the selected zone or area of reception. The minimum transmitter power permitted for an international broadcasting station licensed by the FCC is 50 kW; the antenna power gain toward the intended reception zone must be at least 10; and the field strength incident in the reception zone, either measured or calculated, should exceed 150 microvolts/meter for 50% for the time. For purposes of calculating interference under practical operating conditions, it is assumed that the field strength in directions other than the main lobe of the antenna is not less than 222 microvolts/meter at 1 kilometer for 1 kilowatt supplied to the antenna.

Frequencies for short-wave broadcasting are very much in demand throughout the world. It has been estimated that between two and three stations often compete for an available channel during prime listening hours in many areas of the world. For this reason, international broadcasting is often subjected to a high level of cochannel and adjacent-channel interference.

International broadcasts may include commercial or sponsored programs, provided that commercial program continuities give no more than the name of the sponsor of the program and the name and general character of the commodity, utility or service, or attraction advertised and the commodity is regularly sold or is being promoted for sale on the open market in

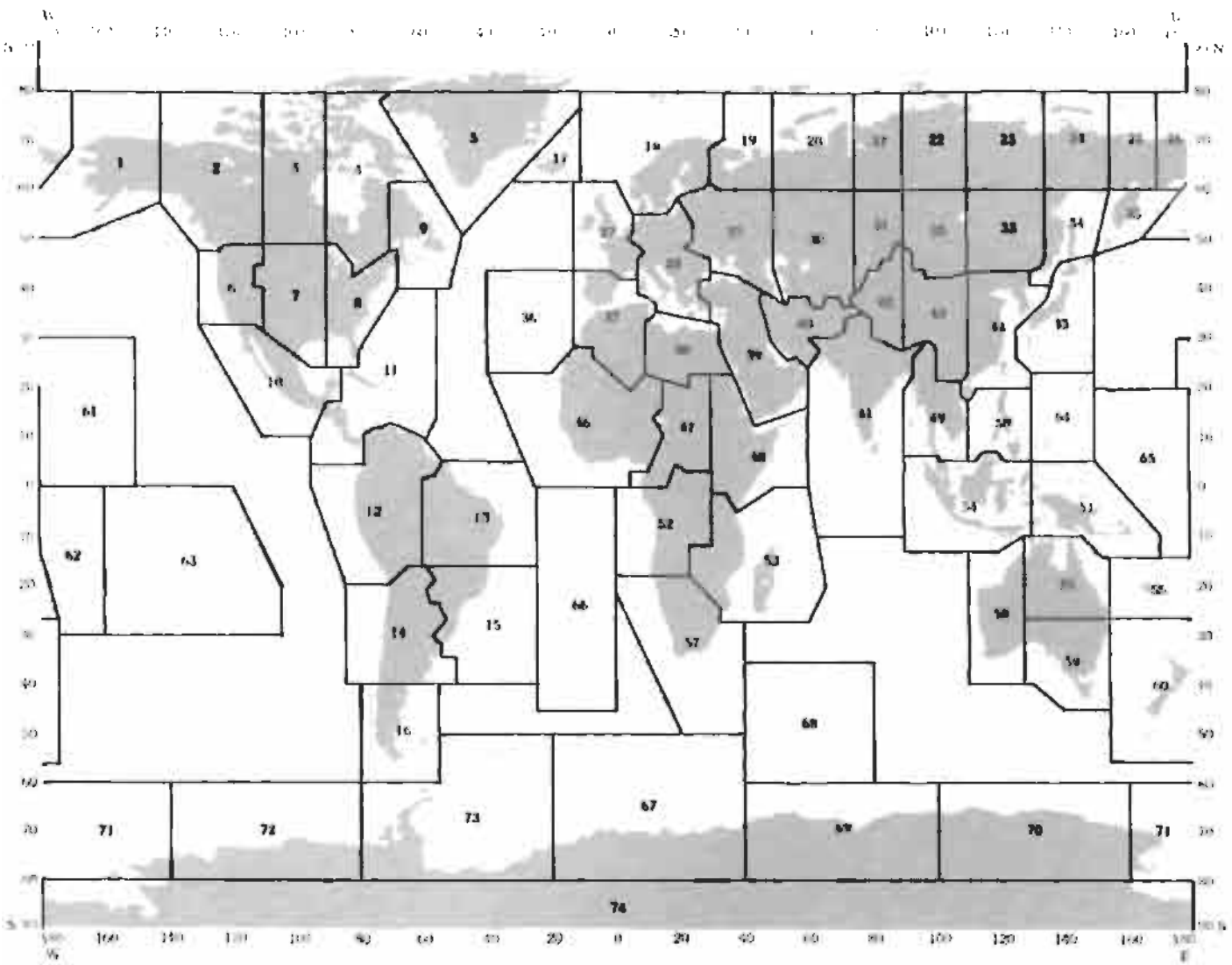


Fig. 18. Geographical zones for broadcasting.

the foreign country or countries to which the program is directed. An international broadcast station may transmit the program of a standard broadcast station or network system provided that the above commercial restrictions are met.

The 1979 World Administrative Radio Conference of the ITU allocated the following additional short-wave spectrum to broadcasting:

kHz
9775– 9900
11 650–11 700
11 975–12 050
13 600–13 800
15 450–15 600
17 550–17 700
21 750–21 850

The band 25 600–26 100 was reduced to 25 670–26 100 when the new allocations became effective.

Although the above WARC-79 extension bands may not be officially implemented until well into the 1990s, they are being used by an increasing number of international broadcast stations on a de facto basis. The FCC assigns frequencies in these bands to US licensed international broadcast stations on the basis of noninterference to other bands. The FCC also assigns frequencies on the same basis between 7305 and 7530 kHz.

At this writing, additional HF spectrum is expected to be allocated to broadcasting at a major international telecommunication conference to be held in Spain during early 1992. This conference, called a *World Administrative Radio Conference of the International Telecommunications Union (WARC-92)* has on its agenda the further expansion of the bands used by international broadcast stations.

The popularity of shortwave broadcasting is on the increase, with at least 600 million radios estimated to be in homes throughout the world. Despite advances expected in the art of broadcasting, short wave is expected to continue to enjoy its unique and important role in the field of human communications well into the 21st Century.

Sources for Federal Communications Commission Documents

The Superintendent of Documents, US Government Printing Office, Washington, D.C. 20402-9371, is the primary source for obtaining the Rules of the FCC. The GPO publishes the *Code of Federal Regulations—Title 47—Telecommunication* that contains, in five volumes, the FCC Rules and Regulations. Individual volumes or the complete set of five volumes may be purchased, and a subscription service for obtaining updates in several forms (microfiche, magnetic tape, paper) is available. In addition, the GPO publishes and sells subscriptions to the biweekly publication *FCC Record* that contains all FCC decisions and actions, specified public notices, etc.

Because of the complexities inherent in maintaining a complete up-to-date set of the often-revised FCC documents, private sources for these publications are also available. Subscriptions are obtainable for FCC and other government rules and regulations. A complete list of such nongovernment sources can be obtained from:

*Office of the Secretary
Federal Communications Commission
1919 M Street, N.W., Room 222
Washington, DC 20554
(202) 632-6410*

PROGRAM PRODUCTION STANDARDS

If sound or television programs are broadcast in "real time," such broadcasts or telecasts are described as "live." However, live broadcasts are rare occurrences at present. Most broadcasting, except news and special events, is transmitted from recordings, either audio-only (sound recordings) or audio and video (television recordings). The origination equipment, which may include microphones, television cameras, film cameras, synchronizing apparatus, transmitters, and recording systems, is designed to meet industry and government specifications. For recording systems, the primary standards specifications must ensure interchangeability of the recordings between machines. Interchangeability is mandatory if recording and replay in a variety of locations is to be achieved.

There is no inherent reason that standards designed for broadcasting should be especially suited to uses other than broadcasting. Currently, however, professional production facilities for cable television, for most educational and industrial instructional purposes, as well as some systems used for video teleconferencing and computer graphics creation, utilize essentially the same types of equipment as those used for broadcasting and so adhere to similar or identical operating standards. This situation may not prevail in the future because active development of high definition television systems and various kinds of special-purpose digital video systems is underway in many parts of the world and in many disciplines. It is likely that specialized systems, especially digitized video systems for uses such as video teleconferencing and multimedia computer use, will evolve rapidly and will require new and improved standards.

Several scientific and industrial societies and trade associations are involved in preparing documents that define agreed-to specifications for sound and television origination and recording equipment. Listed below are the names, areas of responsibilities, and addresses for many of these groups. For information on the current status of standards, it is important to ask the group concerned, since in these fields the technology is advancing at a rapid pace. Names and addresses of the most important standards organizations are listed later in this chapter.

Sound Recording Systems

The past decade has produced radically new approaches to the recording of sound material for broadcasting and for other uses, including home and automobile listening and portable personal use, perhaps while jogging. For the home market, the long-playing, mechanically reproduced phonograph record has virtually disappeared from stores. In its place has appeared the digitally recorded, laser-read compact audio disk (CD) and the digital audio tape (DAT). Analog magnetic tape recording equipment, mostly using the compact audio cassette first introduced in the mid 1960s, is also to be found in most stores and in most homes. Audio cassettes are dominant in the field of highly portable tape players intended for personal use. The audience for recorded music of all kinds has become very large and very sophisticated.

Broadcasting studios must be equipped to replay the consumer recordings listed above—analog compact audio cassette, CD, and DAT. Thus, given these factors, in recent years there has been a proliferation of standards for the audio recording and broadcasting industries.

Some advanced professional sound recording studios make digital audio recordings on magnetic or optical computer disks, but most professional sound recording facilities continue to depend on magnetic tape as the recording medium. Despite advances in digital technology represented by the CD and DAT, analog sound recording on magnetic tape remains very common. Many professional systems in service today are based on analog technology using gamma ferric iron oxide coated on a plastic tape backing material. They use recording heads of iron or ferritic materials that use high-frequency ac recording bias. Associated playback heads have narrow (short) gaps that drive the input stages of low-noise analog amplifiers. Increasingly, however, professional sound studios and professional recording equipment are being designed to make use of modern digital technology.

Digital Audio Recording—The use of digital techniques to record audio signals has the following advantages when compared to traditional analog techniques:

- A. The dynamic range and frequency response of the recording can be greater. Dynamic range is primarily determined by the number of sampling bits per digital word used in the analog-to-digital conversion process. In professional audio applications, 16- to 20-bit words are used, thus providing a theoretical dynamic range of from 96 to 120 dB. The amplitude/frequency response of digital systems is primarily determined by the digital sampling rate and the associated filters used in the digital-to-analog conversion process. For professional use, the sampling rate is usually 48 000 samples per second (48 kHz), and thus the net bit rate of the signal is between 768 000 and 960 000 bits/s. To this bit rate must be added

error-detection signals which can account for an increase of 20% or more to arrive at a figure for the gross bit rate finally to be recorded.

Other sampling rates may be used. For example, for compact disk (CD) recordings, the sampling rate is 44.1 kHz, derived originally from television-system frequency relationships, while some digital audio transmission channels, for example certain long-haul European microwave circuits, use 32-kHz sampling. Digital telephone systems typically sample speech signals at 8 kHz.

- B. Deterioration of quality when the original recordings are copied is virtually absent. This situation is quite different in traditional analog recording, where each "generation" of tape copying results in increased noise and distortion compared to the original recording. In typical digital systems, there is no incremental deterioration of the signal after the initial analog-to-digital conversion.
- C. The effects of magnetic-tape defects and signal "dropouts" due to dirt can be virtually eliminated by use of powerful digital error detection, correction, and concealment techniques. These techniques almost guarantee that the playback is a duplicate of the original recording. When the error-correction systems are overloaded by exceptionally high bit-error rates, error concealment allows a graceful cover for defective data and is designed to produce a minimal disturbance to the reproduced sound recording.

In addition to more or less conventional tape-transport systems using fixed heads and capstan-based tape-pulling mechanisms, digital audio recorders are available, such as those of the DAT design, that use tape transports based on the rotating-head scanning-assembly concept used in video recording systems. Other designs use magnetic disk mechanisms similar to the hard disks used in computer applications. Rewritable optical disk recording systems have appeared and are expected to compete strongly with magnetic tape.

Standards and Specifications for Sound-Recording Systems—Several organizations are active in drafting standards and specifications for radio-broadcasting and sound-recording studios. These include the National Association of Broadcasters (NAB), the Electronic Industries Association (EIA), and the Audio Engineering Society (AES). The fundamental specifications for radio broadcasting in the USA are promulgated by the Federal Communications Commission (FCC). There is a description of the FCC Rules and Regulations earlier in this chapter. However, many other aspects of the program production system are dealt with by these other organizations. A survey of the standards available from the various organizations is listed below.

*National Association of Broadcasters
Dept. of Science and Technology
1771 N Street, N.W.
Washington, DC 20036-2891*

TELEPHONE (202) 429-5346 or 1-800-368-5644
TELEFAX (202) 775-3515

Equipment-related parameters for audio recording and reproducing systems used in broadcasting, including monophonic and stereophonic reel-to-reel magnetic sound recorders, broadcast cassette and cartridge magnetic recording systems, and disk recordings. A complete set of the standards is contained in the publication *NAB's Recommended Engineering Technical Standards*, available from the address listed above. Their titles are:

NAB Standard: Cartridge Tape Recording and Reproduction

NAB Standard: Audio Cassette Recording and Reproduction

NAB Standard: Magnetic Tape (Reel to Reel) Recording and Reproduction

NAB Standard: Disc Recording and Reproduction

Note: NAB also publishes specifications based on reports of the National Radio Systems Committee. These include:

NSRC Interim Voluntary National Standard: Preemphasis/Deemphasis and Broadcast Audio Transmission Bandwidth Specifications. (NRSC-1)

NSRC Interim Voluntary National Standard: Emission Limitation for AM Broadcast Transmission. (NRSC-2)

NSRC Voluntary National Standard: Audio Bandwidth and Distortion Recommendations for AM Broadcast Receivers.

Electronics Industries Assoc.
2001 Pennsylvania Ave. N.W.
Washington, DC 20006
TELEPHONE (202) 457-8716
TELEFAX (202) 457-4901

Equipment-related specifications for a variety of sound broadcasting and sound recording subjects. Consumer Electronics standards are also issued. A selected partial listing of current document titles includes:

CPEB6-A, Preferred Voltage and Impedance Values for the Interconnection of Audio Products

CPEB7, Audio Rectification

SE-101-A, Amplifiers for Sound Equipment

SE-104, Engineering Specifications for Amplifiers for Sound Equipment

EIA-160, Sound Systems

EIA-243, Color Coding for Stereo Pick-Up Leads

EIA/TIA-250-C, Electrical Performance for Television Transmission Systems

EIA-288, Audio Magnetic Playback Characteristic at 7 1/2 IPS

EIA-298, Audio Transmitter Input Impedances for Single Input Transmitters

EIA-517, Car Audio Standard

EIA-157, Method for Determining Air Gap Flux Density and Energy

EIA-549, NSRC AM Preemphasis/Deemphasis and Broadcast Audio Transmission Bandwidth Specifications

EIA-490, Standard Test Methods for Audio Amplifiers

EIA-219, Audio Facilities for Radio Broadcasting Systems

EIA-297-A, Cable Connectors for Audio Facilities for Radio Broadcasting

EIA-560, Standard Method of Measurement for Compact Disc Players

Note: In addition to this selected list of EIA publications, the association publishes an extensive range of standards and other engineering publications concerned with electronics and electronic engineering. Television standards are listed later in this chapter. A catalog may be obtained from the address noted above.

Audio Engineering Society
60 E. 42nd Street
New York, NY 10017
TELEPHONE (212) 661-8528
TELEFAX 212 682-0477

Equipment-related specifications for a variety of sound broadcasting and sound recording subjects. The AES supports the work of the ANSI S4 standards committee. A selected partial listing of current document titles includes:

AES-IID-1991, Plane Wave Tubes: Design and Practice

AES-2-1984 r1991 (ANSI S4.26-1984), Specification of Loudspeaker Components used in Professional Audio and Sound Reinforcement.

AES-3-1992 (ANSI S4.40-1992), Serial Transmission Format of Two-channel Linearly Represented Digital Audio Data (Revision of AES-3-1985, ANSI S4-40-1985)

AES-5-1984 r1991 (ANSI S4.28-1984), Preferred Sampling Frequencies for Professional Digital Audio Applications Employing Pulse-Code Modulation.

AES-6-1982 r1991 (ANSI S4.3-1982), Method for Measurement of Weighted Peak Flutter of Sound Recording and Reproducing Equipment.

AES-7-1982 r1991 (ANSI S4.6-1982), Method of Measuring Recorded Flux of Magnetic Sound Records at Medium Wavelengths.

AES-10-1991 (ANSI S4.43-1991), Serial Multichannel Audio Digital Interface (MADI)

- AES-11-1991 (ANSI S4.44-1991), Synchronization of Digital Audio Equipment in Studio Operations.
- AES-14-1992 (ANSI 4.48-1992), Application of Connectors, Part 1: XLR Polarity and Gender.
- AES-15-1991 (ANSI S4.49-1991), Sound Reinforcement Systems—Communications Interface (PA 422)
- AES-17-1991 (ANSI S4.50-1991), Measurement of Digital Audio Equipment.
- AES-18-1992 (ANSI S4.52-1992), Format for the User Data Channel of the AES Digital Audio Interface.
- AES-19-1992 (ANSI S4.30-1992), Measurement of the Lowest Resonance Frequency of Loudspeaker Cones.

Note: The Audio Engineering Society provides a computer-based data exchange that may be accessed by modem. It provides fully updated information on AES standards. Details for accessing this facility may be obtained from the AES headquarters listed above.

International Sound-Recording Standards—

The contents of most of the standards and specifications listed above are duplicated in international standards. For sound recording, international standardization is performed by the International Electrotechnical Commission (IEC), a subdivision of the International Standards Organization (ISO). Responsibility for international sound-recording standards rests with the IEC Technical Committee 60—Recording (TC-60) and its subcommittee SC-60A. Loudspeakers, microphones, and audio-visual systems are standardized by TC-84.

Information on IEC and ISO standards may be obtained from:

*American National Standards Institute
1430 Broadway
New York, NY 10018*

International exchange of recorded broadcast sound programs is dealt with by the International Radio Consultative Committee (CCIR). Copies of CCIR documents may be ordered from:

*National Technical Information Service
US Department of Commerce
5285 Port Royal Road, Box 1553
Springfield, VA 22161*

A listing of international standards affecting sound broadcasting is given in the section "International Standards for Broadcasting" at the end of this chapter.

Television Recording Systems

The recording of television signals on magnetic tape, variously described as video tape recording (VTR) or television recording, is a technology that has been used professionally for about 35 years. The invention of the video recorder ranks high among the important technical and cultural developments of the 20th century.

Various recording formats are now in professional service. They include the original broadcast format using 2-inch-wide tape (now essentially obsolete) and formats using 1-inch tape, ¾-inch tape, 19-mm tape, and ½-inch tape. The last three tape widths are used for both analog and digital video recording systems, with the tape contained in a box called a cassette. Digital video recording will surely become the dominant professional technology in the immediate future. In addition to these professional formats, there are hundreds of millions of video cassette recorders (VCRs) in homes and institutions around the world. In general, VCR systems use analog recording technology and ½-inch or 8-mm tape. The tape itself is contained in a cassette that is inserted into the front of the recording machine.

Among the formats listed above, all except the obsolete 2-inch format use helical recording technology. In this system, the tape is constrained by mechanical guides to take the shape of a partial helix and is wrapped around a cylinder (the scanner) that rotates rapidly. The scanner also serves as a mounting for the video recording head(s). The tracks that result from the recording process are placed very closely together and at a shallow angle of a degree or two to the reference tape edge. The tracks are quite long, their length amounting to ten or more times the width of the tape. Because the cylinder around which the tape is wrapped rotates rapidly, usually in the 3500–3600 r/min range, the wavelength of the highest recorded modulation frequency (or of a recorded digital data symbol) ranges from 0.7 to 1.1 micron and is well accommodated by the magnetic recording process.

Older recording machine designs, like the 1-inch Type C analog composite helical recorder, use a reel-to-reel tape transport system, where two reels or spools serve to hold the tape in position for passage through the transport mechanism. The tape moves from the supply reel through the recording heads and onto the takeup reel. In all the other listed formats, the tape is contained in a plastic box called a cassette. For recording or playback use, the cassette is inserted into a loading slot in a recording or playing system. The tape is automatically extracted through a slot in the front of the cassette and guided around the helical recording assembly. When this loading process is completed, the machine is ready to receive the commands needed to record or replay a tape.

Professional video recording systems are standardized in North America by the SMPTE (Society of Motion Picture and Television Engineers) and internationally by Technical Committee 60—Recording of the IEC (International Electrotechnical Commission). The 1-inch SMPTE Type C composite recorder has been very heavily used for broadcasting and program production. This format uses full bandwidth FM recording, where the incoming video signal, in composite form, is frequency modulated onto an RF carrier which is then directly recorded on tape. For news gathering and other portable applications, the BetaCam™ and M-II™ analog component systems are frequently used. As

noted earlier, the 2-inch transverse track (quadruplex) recorder is obsolete. The 3/4-inch U-format is gradually being replaced by the analog component systems and, soon, by new digital systems still under development.

Digital formats, now in the process of supplanting analog formats, include the 19-mm Type D-1 (component) and Type D-2 (composite) formats and the 1/2-inch Type D-3 (composite) system. (The "Type" identifiers are assigned by the SMPTE).

For home recording purposes, and for institutional use, the most common video recording formats are two well known VCR systems, the 1/2-inch VHS format and the 8 mm Video format. Improved, mostly compatible variants of these systems, called S-VHS and Hi-8, are used for some semiprofessional applications. These domestic-use formats incorporate a special process to permit recording of the color part of the video signal in an inexpensive way. The process is called *color under* and is not capable of recording the full bandwidth of color signals produced in the broadcast studio. Color under is an engineering compromise that works very well for a final, distribution tape, but is not well suited for full professional applications. It is anticipated that VHS and Video 8 will encounter competition from low-cost, consumer-market digital recording formats within the next few years. However, now, most video cassettes prepared for rental and sales to the public are VHS or (to a much smaller degree) Video 8.

As highly advanced as the magnetic recording of television signals may be, very many television programs are produced by photographic processes. Motion picture technology is also standardized in North America by the SMPTE. The list of SMPTE standards printed below includes many documents defining the parameters for motion-picture films intended for television and other professional uses.

Standards for Video Magnetic Tape and for Motion-Picture Film Recording Systems for Television use—The standards documents listed below are available from

Society of Motion Picture and Television Engineers
595 West Hartsdale Ave.
White Plains, NY 10607

STANDARDS

- SMPTE 1-1990, Video Recording—2-in Magnetic Recording Tape
- V98.2-1982, Specifications of Monochrome Video Magnetic Tape Leader
- SMPTE 3-1986, Video Recording—Frequency Response and Operating Level of Recorders and Reproducers—Audio 1 Record on 2-in Tape Operating at 15 and 7.5 in/s
- SMPTE 4-1989, Television Analog Recording—2-in Magnetic Tape for Quadruplex Recording—Speed
- SMPTE 5-1989, Television Analog Recording—2-in Reels
- SMPTE 6-1988, Video Recording—Video, Audio and Tracking-Control Records—2-in Quadruplex Tape
- SMPTE 7-1988, Motion-Picture Film (16-mm)—Camera Aperture Image and Usage
- SMPTE 8-1989, Video Recording—Quadruplex Recorders Operating at 15 in/s—Audio Level and Multifrequency Test Tape
- V98.9-1983, Video Recording—Magnetic Tape—Color Leader
- SMPTE 11-1989, Video Recording—Quadruplex Recorders Operating at 7.5 in/s—Audio Level and Multifrequency Test Tape
- SMPTE 12M-1986, Television—Time and Control Code—Video and Audio Tape for 525-Line/60-Field Systems
- SMPTE 13-1988, Video Recording—Cartridge Spools—2-in Quadruplex Tape
- V98.18M-1983, Video Recording—1-in Type C Helical-Scan—Basic System and Transport Geometry Parameters
- V98.19M-1983, Video Recording—1-in Type C—Records
- SMPTE 20M-1985, Video Recording—1-in Type C Recorders and Reproducers—Frequency Response and Reference Level
- SMPTE 21M-1986, Video Recording—3/4-in Type E Helical Scan—Records
- SMPTE 22M-1986, Video Recording—3/4-in Type E Helical Scan—Cassette
- SMPTE 23M-1986, Video Recording—1/2-in Type F Helical-Scan—Records
- SMPTE 24M-1985, Video Recording—1-in Reel Dimensions
- SMPTE 25M-1989, Video Recording—1-in Magnetic Recording Tape
- SMPTE 26M-1989, Video Recording—1-in Helical-Scan Recorders—Raw Stock for Reference Tapes
- V98.27M-1983, Video Recording—1-in Type C Reference Recorders—Basic System and Transport Geometry Parameters
- V98.28M-1983, Video Recording—1-in Type C Reference Tapes—Records
- SMPTE 31M-1989, Television Analog Recording—3/4-in Type E—Small Video Cassette
- V98.32M-1983, Video Recording—1/2-in Type H Cassette—Records
- V98.33M-1983, Video Recording—1/2-in Type H Cassette—Tape and Cassette
- V98.34M-1984, Video Recording—1/2-in Type G Cassette—Records
- V98.35M-1984, Video Recording—1/2-in Type G—Cassette and Tape
- PH22.40-1984, Motion-Picture Film (35-mm)—Photographic Audio Records—Release Prints
- SMPTE 41-1989, Motion-Picture Film (16-mm)—Photographic Audio Records—Prints
- SMPTE 48-1989, Motion-Picture Film (16-mm)—Printed Area—Picture and Sound Contact Printing
- PH22.55-1983, Motion-Picture Film—Leaders and Cue Marks—35- and 16-mm Audio Release Prints
- SMPTE 56-1984, Motion-Picture Film—Nomenclature for Studios and Processing Laboratories

- SMPTE 59-1989, Motion-Picture Film (35-mm)—
Camera Aperture Images
- SMPTE 75M-1988, Motion-Picture Film—Designation of A and B Windings for Raw Stock
- SMPTE 83-1990, Motion-Picture Film (16mm)—
Edge Numbers—Location and Spacing
- PH22.86-1981, Position, Dimensions and Reproducing
Speed of Three 200-Mil Magnetic Sound Records on
35-mm and One Record on 17.5-mm Motion-Picture
Film
- SMPTE 87M-1985, Motion-Picture Film (16-mm)—
100-Mil Magnetic Striping
- SMPTE 94-1985, Television—Image Areas and
Mounts for Slides and Opaques
- PH22.95-1984, Motion-Picture Film (35-mm)—
Television Image Area
- PH22.96-1982, Dimensions for Television Image Area
on 16-mm Motion-Picture Film
- SMPTE 97-1989, Motion-Picture Film (16-mm)—
200-Mil Edge Position Magnetic Audio Record
- SMPTE 101-1988, Motion-Picture Film (16-mm)—
Perforated 2R-3000 Film—Magnetic Striping
- SMPTE 102-1986, Motion-Picture Film (35-mm)—
Perforated CS-1870
- SMPTE 183M-1985, Motion-Picture Film—Measure-
ment of Photoelectric Output Factor—Photographic
Audio Level Test Films
- SMPTE 196M-1986, Motion-Picture Film—Screen
Luminance and Viewing Conditions—Indoor Theater
Projection
- PH22.202M-1984, Motion-Pictures—B Chain Elec-
tro-Acoustic Response—Control Rooms and Indoor
Theaters
- SMPTE 203-1987, Motion-Picture Film (35-mm)—
Two-Track Photographic Audio Records—Prints
- SMPTE 204-1987, Motion-Picture Film (16-mm)—
Two-Track Photographic Audio Records—Prints
- SMPTE 207M-1984, Television—Digital Control
Interface—Electrical and Mechanical Characteristics
- PH22.208M-1984, Motion-Picture Film (35-mm)—
Recorded Characteristic—Magnetic Audio Records
- PH22.209M-1984, Motion-Picture Film (8-mm Type
S)—Recorded Characteristic—Magnetic Audio Re-
cords
- SMPTE 210M-1990, Motion-Picture Film (16-mm)—
Magnetic Audio Records—Two Records on 16-mm
Magnetic Film
- PH22.211M-1984, Motion-Picture Film—16- and 35-
mm Variable-Area Photographic Audio Records—
Signal-to-Noise Ratio
- SMPTE 213M-1984, Motion-Picture Film (16-mm)—
Recorded Characteristic—Magnetic Audio Records
- SMPTE 214M-1984, Motion-Picture Film (35-mm)—
Photographic Audio Reproduction Characteristic
- SMPTE 216-1985, Motion-Picture Film (35-mm)—
Recorded Characteristic of Magnetic Audio Re-
cords—Four-Track Striped Release Prints
- SMPTE 217-1985, Motion-Picture Film (70-mm)—
Recorded Characteristic of Magnetic Audio Records
—Striped Release Prints
- SMPTE 218M-1985, Motion-Picture Film (16-
mm)—200-Mil Center-Position Magnetic Audio Re-
cord
- SMPTE 220-1985, Motion-Picture and Television
Equipment—Camera Mounting Connections—1/4-
inch-20 Thread and 3/8-Inch-16 Thread Tripod
Screws
- SMPTE 222M-1987, Television—Monitor System
Electro-Acoustic Response—Control and Review
Rooms
- SMPTE 223M-1985, Motion-Picture Film—Safety
Film (Includes PH1.25)
- SMPTE 224M, Television Digital Component Record-
ing—19-mm Type D-1 Format—Tape Record
- SMPTE 225M, Television Digital Component Record-
ing—19-mm Type D-1 Format—Magnetic Tape
- SMPTE 226M, Television Digital Component and
Composite Recording—19-mm Type D-1 Format—
Tape Cassette
- SMPTE 227M, Television Digital Component Record-
ing—19-mm Type D-1 Format—Helical Data and
Control Records
- SMPTE 228M, Television Digital Component Record-
ing—19-mm Type D-1 Format—Time and Control
Code and Cue Records
- SMPTE 229M-1987, Video Recording—1/2-in. Type
L Mode 1—Records
- SMPTE 230M-1987, Video Recording—1/2-in Type L
Mode 1—Electrical Parameters—Video, Audio,
Time and Control Code and Tracking Control
- SMPTE 233-1987, Motion-Picture Film (16-mm)—
Projectable Image Area
- SMPTE 238M, Television Analog Recording—1/2-in
Type L—Tapes and Cassettes
- SMPTE 240M-1988, Television—Signal Parameters—
1125/60 High-Definition Production System
- SMPTE 244M, Television Digital Recording—Repre-
sentation of NTSC Encoded (System M) Video
Signal—Active Video Portion
- SMPTE 245M, Television Digital Recording—19-mm
Type D-2 Composite Format—Tape Record
- SMPTE 246M, Television Digital Recording—19-mm
Type D-2 Composite Format—Magnetic Tape
- SMPTE 247M, Television Digital Recording—19-mm
Type D-2 Composite Format—Helical Data and
Control Records
- SMPTE 248M, Television Digital Recording—19-mm
Type D-2 Composite Format—Cue Record and
Time and Control Code Record
- SMPTE 249M, Television Analog Recording—1/2-in
Type M-2—Records
- SMPTE 250M, Television Analog Recording—1/2-in
Type M-2—Tapes and Cassettes
- SMPTE 251M, Television Analog Recording—1/2-in
Type M-2—Electrical Parameters of Video, Audio,
Time and Control Code and Tracking Control
- SMPTE 252M, Television Analog Recording—1/2-in
Type M-2—Pulse Code Modulation Audio
- SMPTE 253, Television—Three-Channel Parallel
Component Analog Video Interface

SMPTE 255M, Television Analog Recording—25.4-mm Type J—High-Definition Signals
 SMPTE 256M Television—Specifications for Video Tape Leader

RECOMMENDED PRACTICES

- RP 9-1986, Dimensions of Double-Frame 35-mm 2×2 Slides for Precise Applications in Television
 RP 14-1988, Plotting Data from Sensitometric Strips Exposed on Type Ib (Intensity Scale) Sensitometers
 RP 15-1988, Calibration of Densitometers Used for Black-and-White Photographic Density Measurement
 RP 27.1-1989, Specifications for Operational Alignment Test Pattern for Television
 RP 27.2-1989, Specifications for Operational Registration Test Pattern for Multiple-Channel Television Cameras
 RP 27.3-1989, Specifications for Safe Action and Safe Title Areas Test Pattern for Television Systems
 RP 27.4-1985, Specifications for Operational Test Pattern for Checking Jitter, Weave and Travel Ghost in Television Projectors
 RP 27.5-1989, Specifications for Mid-Frequency Response Test Pattern for Television
 RP 27.7-1987, Specifications for Gray-Scale Operational Alignment Test Pattern for Telecine Cameras
 RP 37-1969, Color Temperature for Color Television Studio Monitors (R1982)
 RP 38.1-1989, Specifications for Deflection Linearity Test Pattern for Television
 RP 41-1983, Evaluation of Color Films Intended for Television
 RP 46-1990, Density of Color and Black-and-White Films, Prints and Slides for Television
 RP 47-1985, Electronic Method of Dropout Detection and Counting (R1989)
 RP 61-1989, Specifications for Azimuth Test Film for 8-mm Type S Audio Reproducers, Magnetic Type
 RP 62-1989, Specifications for Flutter Test Film for 8-mm Type S Audio Reproducers, Magnetic Type
 RP 63-1989, Specifications for Sound-Focusing Test Film for 16-mm Audio Reproducers, Photographic Type
 RP 64-1987, Specifications for Audio-Focusing Test Film for 35-mm Audio Reproducers, Photographic Type
 RP 65-1987, Step Optical Reduction Printing of 35-mm Images to 16-mm Prints and Duplicate Negatives
 RP 66-1987, Step Optical Enlargement Printing of 35-mm Images from 16-mm Images
 RP 67-1989, Specifications for Buzz-Track Test Film for 16-mm Motion-Picture Audio Reproducers, Photographic Type
 RP 68-1984, Specifications for Buzz-Track Test Film for 35-mm Motion-Picture Photographic Audio Reproducers
 RP 69-1989, Specifications for Scanning-Beam Uniformity Test Film for 35-mm Motion-Picture Audio Reproducers
 RP 70-1989, Specifications for Flutter Test Film for 16-mm Audio Reproducers, Photographic Type
 RP 71-1977, Setting Chromaticity and Luminance of White for Color Television Monitors Using Shadow-Mask Picture Tubes
 RP 72-1977, Specifications for an Illuminator of Test Pattern Transparencies for Television Studio Cameras (R1988)
 RP 85-1985, Tracking-Control Record for 1-in Type C Helical-Scan Video Tape Recording
 RP 86-1985, Video Record Parameters for 1-in Type C Helical-Scan Video Tape Recording
 RP 87-1986, Reference Carrier Frequencies, Pre-Emphasis Characteristic and Audio and Control Signals for 3/4-in Type E Helical-Scan Video Tape Cassette Recording
 RP 88-1986, Reference Carrier Frequencies and Pre-Emphasis Characteristic for 1/2-in Type F Helical-Scan Video Tape Recording
 RP 90-1979, Specifications for Magnetic-Type Audio Level and Multifrequency Test Film for 16-mm Sound Reproducers
 RP 96-1988, Specifications for Subjective Reference Tapes for Helical-Scan Video Tape Reproducers for Checking Receiver/Monitor Set-Up
 RP 101-1986, Requirements for Recording American National Standard Time and Control Code on Quadruplex Video Tape Recorders
 RP 103-1982, Care and Handling of Video Magnetic Recording Tape (R1987)
 RP 112-1983, Reference Carrier Frequencies, Pre-Emphasis Characteristic and Audio and Control Signals for 1/2-in Type H Helical-Scan Video Tape Cassette Recording (R1988)
 RP 113-1983, Supervisory Protocol for Digital Control Interface
 RP 114-1983, Dimensions of Photographic Control and Data Record on 16-mm Motion-Picture Film (R1989)
 RP 115-1983, Dimensions of Photographic Control and Data Record on 35-mm Motion-Picture Release Prints (R1989)
 RP 116-1990, Dimensions of Photographic Control and Data Record on 35-mm Motion-Picture Camera Negatives
 RP 117-1989, Dimensions of Magnetic Control and Data Record on 8-mm Type S Motion-Picture Film
 RP 118-1983, Dimensions of Photographic Control and Data Record on 8-mm Type S Motion-Picture Prints (R1989)
 RP 119-1984, Reference Carrier Frequencies, Pre-Emphasis Characteristics and Audio and Control Signals for 1/2-in Type G Helical-Scan Video Tape Recording Cassette Systems
 RP 120-1983, Measurement of Intermodulation Distortion in Motion-Picture Audio Systems (R1989)
 RP 121-1988, Tape Dropout Specifications for 1-in Types B and C Video Tape Recorder/Reproducer
 RP 125-1984, Bit-Parallel Digital Interface for Component Video Signals
 RP 127-1985, Specifications for Type U Audio Level

- and Multifrequency Test Film for 35-mm Studio Audio Reproducers, Magnetic Full-Coat Type
- RP 132-1985, Storage of Edit Decision Lists on 8-in Flexible Diskette Media (R1989)
- RP 133-1986, Specifications for Medical Diagnostic Imaging Test Pattern for Television Monitors and Hard-Copy Recording Cameras
- RP 134-1986, Polarity for Analog Audio Magnetic Recording and Reproduction
- RP 135-1986, Use of Binary User Groups in Motion-Picture Time and Control Codes
- RP 136-1986, Time and Control Codes for 24, 25 or 30 Frame-Per-Second Motion-Picture Systems
- RP 137-1986, Data Tracks on Low-Dispersion Magnetic Coatings on 35-mm Motion-Picture Film
- RP 138-1986, Control Message Architecture
- RP 139-1986, Tributary Interconnection
- RP 140-1986, Position of Photographic Audio Record for Routine Test Signals
- RP 141-1986, Background Acoustic Noise Levels in Theaters and Review Rooms
- RP 142-1986, Stereo Audio Track Allocations and Identification of Noise Reduction for Video Tape Recording
- RP 143-1986, Specifications for Type U Audio Level and Multifrequency Test Film for 35-mm Striped Four-Track Release Print Audio Reproducers
- RP 144-1987, Basic System and Transport Geometry Parameters for 1/2-in Type L Cassette
- RP 145-1987, Color Monitor Colorimetry
- RP 146-1987, Transfer of Edit Decision Lists
- RP 147-1987, Audio Channel Assignments of Multi-Channel Sub-Masters Used in Preparation for Two-Track Masters for Transfer to Video
- RP 148-1987, Relative Polarity of Stereo Audio Signals
- RP 150-1988, Channel Assignments and Test Leader for Magnetic Film Masters Intended for Transfer to Video Media Having Stereo Audio
- RP 154, Reference Signals for the Synchronization of 525-Line Video Equipment
- RP 155-1990, Audio Levels and Indicators for Digital Audio Records on Digital Television Tape Recorders
- RP 156-1990, Bar Code Labeling for Type D-1 Component and Type D-2 Composite Cassette Identification
- RP 157-1990, Key Signals
- RP 158, Basic System and Transport Geometry Parameters for 1/2-in Type M-2 Format
- RP 159, Vertical Interval Time Code and Longitudinal Time Code Relationship
- RP 160, Three-Channel Parallel Analog Component High-Definition Video Interface

ENGINEERING GUIDELINES

- EG 5-1989, Projected Image Quality of 70-mm, 35-mm and 16-mm Motion-Picture Projection Systems
- EG 6-1982, Use of 2-in Tape on Cartridge/Cassette Spools for Quadruplex Video Tape Recorders (R1987)
- EG 7-1989, Audio Sync-Pulse for 8-mm Type S

- Cameras, Magnetic Audio Recorders and Rerecording Projectors
- EG 8-1984, Specifications for Motion-Picture Camera Equipment Used in Space Environment (R1989)
- EG 9-1985, Audio Recording Reference Level for Post-Production of Motion-Picture Related Materials
- EG 10, Tape Transport Geometry Parameters for 19-mm Type D-1 Cassette for Component Digital Video Recording
- EG 13-1986, Use of Audio Magnetic Test Films
- EG 20, Tape Transport and Geometry Parameters for 19-mm Type D-2 Composite Format for Television Digital Recording
- EG 21, Nomenclature for Television Digital Recording, 19-mm Type D-1 Component and Type D-2 Composite Format
- EG 22, Description and Index of Documents for 19-mm Type D-2 Composite Television Digital Recording
- EG 23-1990, Transfer of Two-Channel Stereo Audio from Audio Magnetic Film or Tape to Video Tape
- EG 24, Video and Audio Alignment Tapes and Procedures for 1-in Type C Helical-Scan Television Analog Recorders
- EG 25, Telecine Scanning for Film Transfer to Television
- EG 26, Audio Channel Assignments for Digital Television Tape Recorders with AES/EBU Digital Audio Inputs

Video Disk Recordings—The use of optical video disk recordings is increasing. The basic design of the optical video disk system dates to the early 1980s, but various factors, including the availability of low-cost VCR machines, have slowed acceptance of the video disk. The basic process for manufacturing the optical video disk is to etch, in a continuous spiral on a master disk, millions of “pits” that are width modulated by the incoming video signal. This etching is accomplished by a modulated laser beam. Aluminized plastic replicas of the master are then made in large volume and are read (played back) by a low-power laser in the disk player. The recorded and reproduced signal is analog. The fundamental manufacturing process is similar in many ways to that used to produce compact disk (CD) recordings.

Both advantages and disadvantages are evident in the video disk system. For example, a VCR provides the user with the option of recording television programs for later replay; the video disk is, up to now, entirely a playback system. The recording situation may change as the evolution of writeable and erasable optical disks evolves, but it is not clear that costs can be sufficiently reduced to assure success in the consumer market. Even so, the disk offers a series of advantages to the professional user and to demanding consumers. The system allows a high degree of random access; the laser beam may be directed to reproduce any part of the disk, to an accuracy of a single television frame, in less than a second. In addition, the recording system is a wide-bandwidth system and is capable of excellent picture quality. Both of these useful attributes have made the

video disk popular for special applications such as the distribution of recorded instructional and advertising material by industrial users, as well as for educational purposes in schools. It is also true that critical consumers are turning to video disk recordings of motion pictures for television display in their homes. Often the picture quality obtained from the video disk is superior to that obtainable from VCR equipment. The standards for video disk systems are written and published by the IEC, through Technical Committee 60. The address of the IEC is given in the section "International Standards for Broadcasting" at the end of this chapter.

Selected Lists of Television Standards

Listed below are two groups of industrial standards that are useful in television engineering. The EIA list is abbreviated; readers are encouraged to contact EIA headquarters, at the address listed earlier in this chapter in the subsection "Standards and Specifications for Recording Systems." A comprehensive, up-to-date listing is available.

EIA Television Standards (Partial List)

- TEPI06-B, Worldwide Type Designation System for TV Picture Tubes and Monitor Tubes
- TVSB5, Multichannel TV Sound System—BTSC System Recommended Practices
- EIA/IS-6, Recommended Able Television Channel Identification Plan
- IEB1, Closed Circuit Television Definitions
- IETNTS1, Color Television Studio Picture Line Amplifier Output Drawing
- EIA-170, Electrical Performance Standards—Monochrome Television Studio Facilities
- EIA-312-A, Engineering Specification Outline for Monochrome CCTV Camera Equipment
- EIA-330, Electrical Performance Standards for Closed Circuit Television Camera 525/60 Interlaced 2:1
- EIA-343-A, Electrical Performance Standards for High Resolution Monochrome Closed Circuit Television Camera
- EIA-375-A, Electrical Performance Standards for Direct View Monochrome Closed Circuit Television Monitors 525/60 Interlaced 2:1
- EIA-403-A, Precision Coaxial Connectors for CASTV Application (75 Ohm)
- EIA-412-A, Electrical Performance Standards for Direct View High Resolution Monochrome Closed Circuit Television Monitors
- EIA-420, Electrical Performance Standards for Monochrome Closed Circuit Television Cameras 525/60 Random Interlace
- EIA-439, Engineering Specification for Color CCTV Camera Equipment
- EIA-508, Electrical Performance Standards for Television Broadcast Transmitters
- EIA-563, Standard Baseband (Audio/Video) Interface Between NTSC Television Receiving Devices and Peripheral Devices

Other EIA Standards:

- EIA-232-D, Interface Between Data Terminal Equipment and Data Circuit Terminating Equipment Employing Serial Binary Data Interchange
- EIA-422-A, Electrical Characteristics of Balanced Voltage Digital Interface Circuits

Note: In addition to this selected list of EIA publications, the association publishes an extensive range of standards and other engineering publications concerned with electronics and electronic engineering. A catalog may be obtained from the address given above in the subsection "Standards and Specifications for Recording Systems."

International Standards for Television Recording—Most of the SMPTE and EIA documents listed above have equivalent specifications in the publications of international standardizing bodies. Video recording standards are produced by the IEC Subcommittee 60B and broadcast program exchange specifications by the CCIR (published in CCIR Volume 11—Television Broadcasting). A listing of relevant IEC and CCIR standards is given later in this chapter. In addition, a set of technical publications concerning video recording of 625-line PAL and SECAM color television signals is available from the Technical Centre of the European Broadcasting Union (EBU), whose address is given later in this chapter.

International Standards for Television Film Systems—Standards defining international specifications for motion picture film are drafted by ISO Technical Committee 36 (TC 36). An abbreviated listing of these documents is contained in the last part of this chapter. In the USA, ISO standards may be obtained from ANSI. In addition, the EBU publishes specifications for motion picture film intended for television use by members of the EBU, and the CCIR publishes specifications, in CCIR Volume 11—Television, for motion picture films to be used for the international exchange of television programs.

DIGITAL TELEVISION SYSTEMS

Introduction—The Basics

Television is a medium for recording and transferring images and sound from one point to another. The image portion of the system is a "light-to-light" system, gathering variations in light at the source via a pickup device (such as a camera) and recreating those variations in light as a visual image via a display device such as a cathode ray tube (CRT). The television system transports the source information by scanning the image and generating a continuous electronic signal that is a representation of the light. Since the variations in light at both the source and the display are continuous, the television signal is inherently an analog signal, a signal whose amplitude as a function of time represents the

visual content of the television image. This analog signal can be transformed into a linear pulse-code-modulation (PCM) digital representation of the image.

Digital forms of television have certain advantages over the analog signal formats. An analog signal can be degraded by the various transport media and processes. These degradations include both linear and nonlinear distortions and the introduction of noise. The degradations are both cumulative and difficult to separate from the analog signal. Digital signals consist of a sequence of two levels of information (logical 1s and logical 0s) that can be sampled and restored as long as the distortion and noise levels do not exceed certain thresholds. Further, digital systems can be encoded with redundant information so that even in the presence of distortion and noise that exceeds the threshold of separation between a 1 and a 0, the original information can be restored. Once a television system has been digitized, the image information is in the form of a series of numbers which can be processed in a uniform and consistent manner.

A digital television system (Fig. 19) consists of an image source, a means such as a camera of converting the image into an analog electronic signal (E_A), a means of transforming the analog signal into a digital signal (analog-to-digital convertor, or A/D), some process (P_D) which is applied to the digital signal, a means of retransforming the resulting digital signal back to the analog domain (digital-to-analog convertor, or D/A), and a means such as a display of reconverting the analog signal back into "light" for the benefit of the observer.

The digital process (P_D) varies according to the application. It can be a variable time delay used for frame synchronization and time-base correction, a video processor that corrects for luminance and chrominance error, an image manipulator that permits digitally generated effects and graphics to be added to the image, a data-compression scheme to reduce the amount of data, or simply a transmission path.

The A/D converter (Fig. 20) consists of four basic components, a band-limiting, antialiasing filter (F_A), a sample-and-hold circuit (S&H) that samples the analog signal, the quantizing unit (Q) that divides the range of each analog signal sample into Q distinct levels, and finally an encoder (ENC) that places a specified code on the output data lines (D_0-D_n) for each of the Q levels. For binary systems, Q is a power of 2 ($Q = 2^B$), where B is the number of bits and the number of data lines is related to the number of bits ($n = B - 1$).

The D/A converter (Fig. 21) consists of four basic components, a digital input register (REG) in which the

bits of the word to be converted are stored for one time period, a decoder (DEC) that converts the data lines into Q distinct analog levels, a resampling (R_c) circuit that corrects for distortions due to the digital sample-and-hold process, and a band-limiting, restoration or reconstruction filter (F_b).

Resolution

The resolution of a television system is a measure of its ability to delineate picture detail. For a number of lines N , normally alternate black and white lines, the width of each line is $1/N$ times the picture height. From this definition, television resolution is expressed as R television lines per picture height (tvl/ph).

Based on early empirically derived data taken from Engstrom* and others and later by Wheeler and Loughren,† it was determined that vertical resolution (R_v) could be expressed as kN_v , where k represented a utilization factor that for 2:1 interlaced CRT devices was in the range 0.64 to 0.71. The factor k is influenced by many variables including the shape of both the gathering beam and the display beam. Work done on progressive scan displays places an upper bound on k of 0.95. Assuming $k = 0.66$, the vertical resolution of the NTSC-M system is given as $R_v = 320$ tvl/ph.

Horizontal resolution (R_h) is a function of the active line period, the aspect ratio (A) of the image, and the bandwidth of the signal (f), and is given as:

$$R_h = 2 (1/A)P(f)/(F \times L)$$

where,

- P = percent active period of the total line period,
- F = frame rate,
- L = number of scan lines.

For NTSC:

$$\begin{aligned} 1/A &= 0.75 \\ P &\approx 52.655/63.555 = 0.828 \\ F &= 29.97003 \\ L &= 525 \end{aligned}$$

and $R_h = 78.9 \times 10^{-4} (f) = 331.5$ tvl/ph for $f = 4.2$ MHz.

The image resolution and dynamic range of the digital system is limited by seven characteristics:

* Reference 1.
† Reference 2.

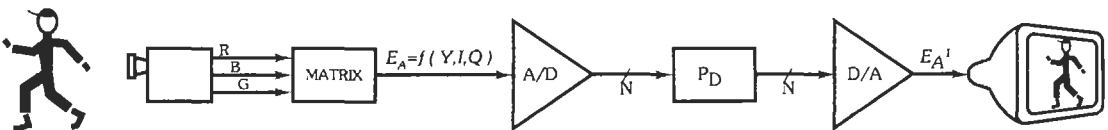


Fig. 19. Digital television system.

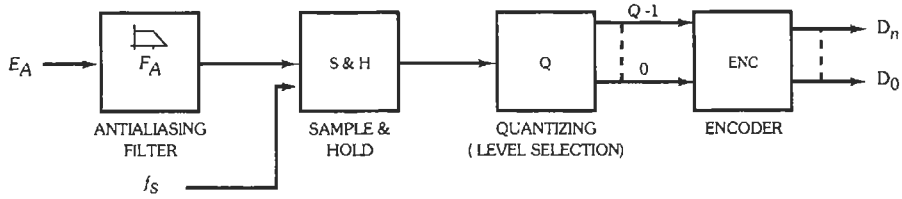


Fig. 20. A/D converter.

1. The resolution limitation (R_A) of the analog signal being converted, as noted above
2. The characteristics of the band-limiting, anti-aliasing filter (F_A) in the A/D converter
3. The frequency (f_S) of the A/D sampling pulse
4. The shape of the sampling pulse at the A/D converter
5. The number of bits (B) in the data word (D_0 - D_n) which determines the number of levels (Q) that can be defined
6. The shape of the resampling pulse (f_R) at the D/A converter
7. The characteristics of the band-limiting, restoration filter (F_D) at the D/A converter.*

Sampling and Spectra

An analysis of the television signal spectrum shows that the spectrum is composed of discrete frequency components about harmonics of the line frequency (f_H) and field rate (f_F) (Fig. 22). The analog signal is sampled by the periodic nature of the scanning structure, in that the scanning lines and the creation of image frames have the effect of sampling the analog signal in the temporal (T) and vertical (V) dimensions, respectively. Whenever an analog signal of spectrum f_A is sampled at intervals, such as with a rectangular pulse train with a period of T , (where $1/T = f_S$ and f_S is the sampling frequency) as described in Fig. 23, the sampled signal has a spectrum equal to the original spectrum plus a double sideband of width $2 \times f_A$ repeated every f_S Hz, as shown in Fig. 24. In the NTSC system, most of the energy of the luminance signal is concentrated in a set of sidebands of the line frequency. Most of the energy of the chrominance signal is concentrated in a set of sidebands of the color subcarrier frequency. Since the color subcarrier is an integer

multiple ($455/2$) of the line frequency, these latter sidebands are also related to the line frequency.† In the NTSC system, the temporal sampling frequency (f_F) is 59.94006 Hz, and the vertical sampling frequency is the line rate (f_H), 15.734 + kHz.

As shown in Fig. 20, the first step in the analog-to-digital conversion process is to sample the analog signal with the digital sampling clock of frequency f_S . This produces a third shift in the frequency domain as shown in Fig. 24 by creating a sampling structure in the horizontal (H) dimension.

A close inspection of the process in both the temporal and frequency domains indicates that the sampling process must conform to certain criteria. If the maximum frequency contained within the analog signal is f_A and the frequency of the sampling pulse is f_S , then f_S must be greater than $2 \times f_A$. Fig. 24 shows the result of sampling the spectra with a sampling frequency greater than $2 \times f_A$. When an analog signal is sampled at a rate at least twice its highest frequency component, the sampling is said to meet the Nyquist criterion. Fig. 25 shows the result of sampling the spectra with a sampling frequency less than $2 \times f_A$. It is clear that if the frequency of the sampling pulse is less than $2 \times f_A$, the spectra of the original baseband signal and the spectra of the sidebands at the multiples of the sampling frequency of the resultant Fourier transform overlap as shown by the shaded area in Fig. 25. This makes an accurate restoration of the original baseband signal extremely difficult, as the sidebands would appear to have component information that does not really exist.** The erroneous signal content is termed *aliasing*. Fig. 26 shows an example of a sub-Nyquist-sampled, analog baseband signal of a single frequency for illustrative purposes. In this instance, the reconstructed signal

* Reference 3.

† Reference 4.
** Reference 5.

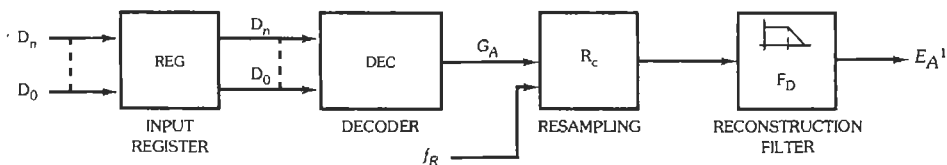


Fig. 21. D/A converter.

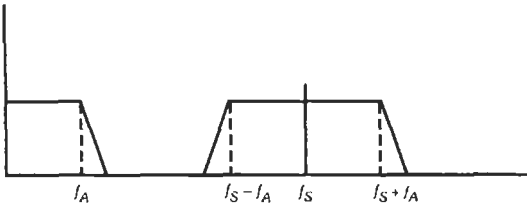


Fig. 24. Sampling—Nyquist criterion ($f_S \geq 2f_A$).

4.5 MHz as shown in Fig. 28, and therefore, the line, field, frame, and sampling frequencies have integer submultiple relationships. This also helps when samples based on one sampling structure must be translated (standards converted) into an alternative sampling structure. The sampling structures for the proposed advanced television (ATV) systems are also orthogonal.

Quantizing and Dynamic Range

The analog signal having been separated into discrete samples, the second step is to quantize the samples. Quantizing a set of values (or samples) is simply the act of describing each member of the set of values as a multiple of a basic unit, dividing the range of values into Q quantizing levels. The range of values which can be represented is 0 to $Q-1$. In binary systems, $Q = 2^B$, where B equals the number of bits. A segment of an analog signal waveform is shown in Fig. 23 with the periodic sampling signal and an eight-level quantizer ($Q = 2^3$) superimposed. If we assume that each level represents 0.125 V, then the range of the input signal that can be accurately represented is 0 to 0.875 V (0 to 7 times 0.125 V). The levels need not be uniformly quantized (equally spaced) to provide a linear sampling of the analog signal, but in most systems, the levels are uniformly quantized.

The sampled analog signal is discrete in time but remains continuous in amplitude. However, the quantizing process assigns all analog values within a specific quantizing interval the same quantizing value, the nominal value for that interval. This produces a quantizing error in the restored, reconverted analog signal; this error is known as *quantizing noise*. As shown in Fig. 23, the quantizing error can not exceed $\pm Q/2$.

The theoretical maximum peak-to-peak signal-to-quantization-noise ratio (SNR), or dynamic range, of the quantized signal is determined by the number of

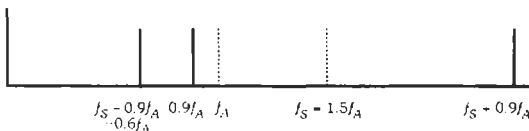


Fig. 26. Sub-Nyquist sampling—single frequency.

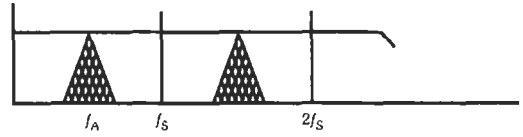


Fig. 25. Sub-Nyquist sampling ($f_S < 2f_A$).

levels (Q) selected, but is also influenced by the shape of the sampling pulse and the relationship between the maximum bandwidth of the signal, f_A , and the sampling frequency, f_S . The analog signal contains noise, and, as a very over-simplified explanation, the sampling process tends to “average” the noise in a sample. The equation for SNR of a quantized signal is:

$$SNR = K + 20 \log Q$$

and since Q can be replaced in binary systems by 2^B ,

$$SNR = K + 20 \log 2^B = K + B \times 20 \log 2 = K + B \times 6.02$$

In uniformly quantized, pulse-coded modulation (PCM) systems, K has been derived as:

$$K = K_p + 10 \log [f_S / (2 \times f_A)]$$

where K_p is the component due to the application of uniformly quantized sampling and is given as $10 \log 12$. Therefore,††

$$K = 10.79 + 10 \log [f_S / (2 \times f_A)]$$

For the example of a 4.2-MHz NTSC signal, sampled at four times the color subcarrier in a ten-bit system ($Q = 1024$),

$$SNR = 10.79 + 10 \log (14.318/8.4) + 10 \times 6.02 = 73.30$$

Composite and Component Signal Coding

The transmitted television signal format in the United States is the NTSC-M format, a composite signal consisting of the luminance (Y) information, the chrominance information as two signals (I and Q) quadrature modulated onto a color subcarrier, and the synchro-

†† Reference 10.

Line 41, Field 1	X X X X X X
Line 41, Field 2	X X X X X X
Line 42, Field 1	X X X X X X
Line 42, Field 2	X X X X X X

Fig. 27. Section of NTSC sampling array showing orthogonal arrangement of samples.

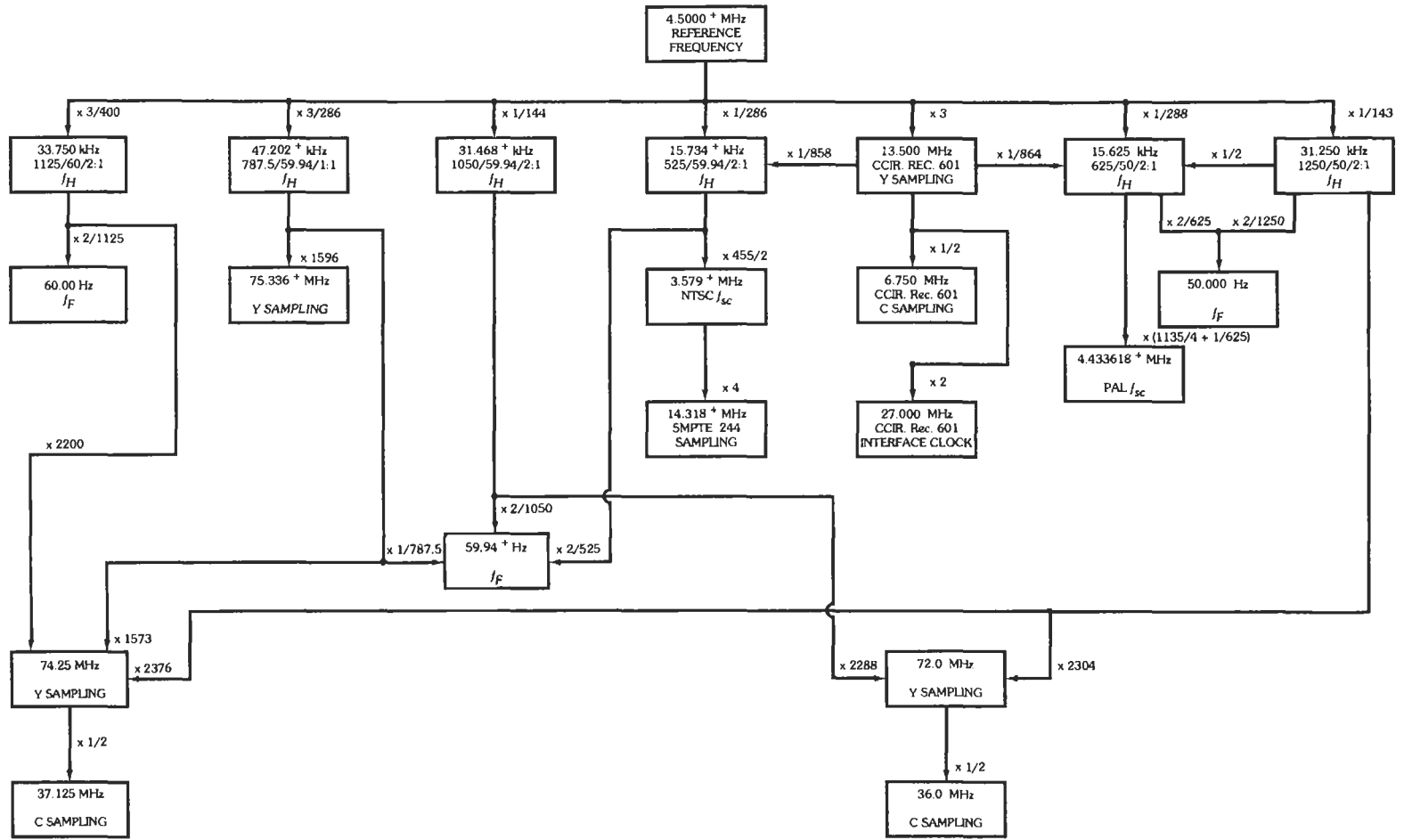


Fig. 28. Frequency relationships in NTSC, PAL, CCIR-601, and ATV systems (references 6, 7, 8, and 9).

nization pulses. This signal is also the format used in many instances in television production facilities for distribution and recording of the television signal. The analog television signal originates in the camera as a tricolor signal representing the red (R), green (G), and blue (B) components of the signal. The R , G , and B components are matrixed to form the Y , I , and Q signals as described in Fig. 19 and Table 7. The analog television signal also appears in its component form in some portions of the production facility, such as in the production of graphics or special effects in order to avoid the limitations in both luminance and chrominance signal bandwidth and cross-color and cross-luminance distortions inherent in the NTSC system and the signal degradation caused by multiple decoding of the composite signal back into its components and the re-encoding to the composite form.

The digital television signal may be implemented in either the composite or component representation of the television signal. The composite representation is based on the sampling of the NTSC encoded signal and is described in the section "Composite Encoded Signals: SMPTE 244 (NTSC and PAL)." The component representation accepts separate baseband tricolor stimulus signals and individually samples and quantizes them. The component representation of the video signal can be in one of many tristimulus formats. The most common formats are (1) a luminance signal, Y , and two color-difference signals, $k_B(B - Y) = C_B$ and $k_R(R - Y) = C_R$; and (2) the baseband R , G , and B signals. The component representations are discussed in the section "Component Coding of Signals: SMPTE 125 and CCIR 601."

Bit-Parallel and Bit-Serial Interfaces

The composite and component signal format interfaces are both electrically and mechanically similar in both the parallel data format and serial data format. The parallel interface consists of a unidirectional, eleven-pair interconnection between a transmitting equipment and a receiving equipment. Video data, timing reference information, and ancillary signals are time multiplexed and transferred on ten data pairs in NRZ form. An eleventh pair provides a synchronous clock. The

serial interface consists of a unidirectional, bit-serial single interconnection between a transmitting equipment and a receiving equipment.

The data signals are in the form of binary information coded in 8- or 10-bit words. These signals are:

- Video data
- Timing reference codes
- Ancillary data, including audio services
- Identification codes such as video index information

The signaling sense of the voltage appearing across the interconnection cable is positive binary and defined as follows (refer to Fig. 29):

1. The A terminal of the transmitter shall be negative with respect to the B terminal for a binary 0 (LOW or L or OFF) state.
2. The A terminal of the transmitter shall be positive with respect to the B terminal for a binary 1 (HIGH or H or ON) state.
3. In these interfaces, all digital-signal time intervals are specified at the half-amplitude points. All transitions are specified between the 20% and 80% amplitude points.
4. In the parallel signal interface, the data lines are designated DATA 0 through DATA 9. The group of ten signals is identified by placing parentheses around the range of subscripts included, as DATA (0-9). When 8-bit signals are conveyed by the interface, DATA (2-9) are used and DATA (0-1) are set to zero. DATA 9 is always the most significant bit.*

Composite Encoded Signals: SMPTE 244 (NTSC and PAL)

The Society of Motion Picture and Television Engineers (SMPTE) developed a common composite interface that serves the needs of both 525-line/59.94-Hz NTSC systems and 625-line/50-Hz PAL systems. The characteristics of the interface are summarized for the 525-line case in Table 8.†

* Reference 12.

† Reference 13.

TABLE 7. NTSC-M COLOR TELEVISION SYSTEM*

Parameter	Equation
Luminance signal	$Y' = 0.299 R' + 0.587 G' + 0.114 B'$
Color-difference signals	$I' = -0.27 (R' - Y') + 0.74 (B' - Y')$ $Q' = 0.41 (R' - Y') + 0.48 (B' - Y')$
Composite color signal	$M' = Y' + Q \sin (2\pi f_{sc} + 33^\circ)$ $+ I \cos (2\pi f_{sc} + 33^\circ)$ where R' , G' , B' are the gamma-precorrected signals.

* Reference 11.

TABLE 8. ENCODING PARAMETERS FOR NTSC COMPOSITE TELEVISION*

Parameter	Composite NTSC-M
Coded signal	Composite NTSC
Samples per total line	910
Sampling structure	Orthogonal (line, field, and picture repetitive)
Sampling frequency	4 times $f_{sc} = 14.318$ MHz
Serial interface data rate	143.18 Mbits/s
Form of coding	Uniformly quantized PCM, 8 or 10 bits per sample
Number of samples per active line	768
Correspondence between video signal levels and quantization levels	<p>8-bit samples: The tip of the synchronization pulse (-40 IRE) corresponding to level 4 (04_H), the blanking level corresponding to level 60 ($3C_H$), and the peak white level ($+100$ IRE) corresponding to level 232 ($C8_H$)</p> <p>10-bit systems: The tip of the synchronization pulse corresponding to level 16 (010_H), blanking level corresponding to level 240 ($0F0_H$) and the peak white level corresponding to level 800 (320_H)</p>
Alignment	The sample at interval 0, line 10, field 1, color frame A is an I axis ($+123^\circ$) sample.

* Reference 14.

Component Coding of Signals: SMPTE 125 and CCIR 601

The Society of Motion Picture and Television Engineers (SMPTE) and the European Broadcasting Union (EBU), working jointly, developed a common component interface that serves the needs of both 525-line/59.94-Hz systems and 625-line/50-Hz systems. The interface is documented in CCIR (International Radio Consultative Committee) Recommendation 601 and in SMPTE 125M. The characteristics of the interface are summarized in Tables 9 and 10.

The Y , C_B , and C_R signals are transmitted at the 4:2:2 level equivalent of the CCIR document 601 digital hierarchy for television signals, with a nominal sampling frequency of 13.5 MHz for the luminance signal and 6.75 MHz for each of the color-difference signals.

Because of the existence of both 8-bit and 10-bit equipment, all synchronizing signals (EAV, SAV, ANC) must be detected by reference to the 8 most significant bits only.

In the parallel data format, data is transmitted across the interface on ten data pairs: DATA 0 through DATA 9. DATA 9 is the most significant bit (MSB). For the 10-bit system, 1016 of the 1024 levels [digital levels 4 (004_H) through 1019 ($3FB_H$)] of the ten-bit word are

used to express quantized values. For the 8-bit system, data values 1 through 254 (01_H – FE_H) may be used to express quantized signal values.

The digital code words that describe the video signal are conveyed as a 27-megaword/second multiplex in the following order:

$$C_B, Y, C_R [Y], C_B \dots$$

where the three words C_B , Y , and C_R refer to the co-sited samples; the following $[Y]$ is an isolated luminance-only sample. The C_R , C_B samples are co-sited with the odd (1st, 3rd, 5th, . . .) Y samples on each line, as shown in Fig. 30. The first video data word in the active line period is C_B .

Synchronization and Blanking Interval Considerations

The digital interfaces provide for the transmission of synchronization and ancillary data that may be multiplexed into the data stream during video blanking intervals. In each of the component systems, four specific sample words or interface clock intervals are reserved for the end-of-active-video (EAV) timing reference, and four specific sample words are reserved for

TABLE 9. ENCODING PARAMETERS FOR COMPONENT TELEVISION AT 4:2:2 LEVEL

Parameter	525-Line, 59.94-Fields/s	625-Line, 50-Fields/s
Coded signals	Y, C_B, C_R	
No. of samples, total line		
Luminance signal (Y)	858	864
Each color-difference signal (C_B, C_R)	429	432
Sampling structure	Orthogonal (line, field, and picture repetitive). C_B and C_R samples cosited with odd Y samples in line.	
Sampling frequency		
Luminance signal	13.5 MHz	
Each color-difference signal	6.75 MHz	
Multiplexed parallel data interface clock rate	27.0 MHz	
Serial interface clock rate	270 MHz	
Form of coding	Uniformly quantized PCM, 8 or 10 bits per sample, for the $Y, C_B,$ and C_R signals.	
Number of samples per digital active line		
Luminance signal	720	
Each color-difference signal	360	

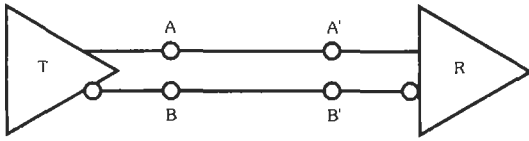
the start-of-active-video (SAV) timing reference. The timing reference signals are contiguous with the video data, when present, and continue through the vertical-blanking interval. The EAV and SAV timing references are the digital horizontal-synchronization signals and

occur on every line. The interval starting at EAV and ending with SAV is the digital horizontal-blanking period.

Any interface clock intervals not required to carry video information may be used to carry ancillary

TABLE 10. QUANTIZATION LEVELS FOR COMPONENT TELEVISION AT THE 4:2:2 LEVEL

Parameter	Quantization
Correspondence between video signal levels and quantization levels	
Luminance signal	8 bit systems: 220 quantization levels with the black level corresponding to level 16 (10_H) and the peak white level corresponding to level 235 (CB_H)
Each color difference signal	10 bit systems: 877 quantization levels with the black level corresponding to level 64 (040_H) and the peak white level corresponding to level 940 ($3AC_H$) 8 bit systems: 225 quantization levels symmetrically distributed about level 128 (80_H) corresponding to the zero level 10 bit systems: 897 quantization levels symmetrically distributed about level 512 (200_H) corresponding to the zero level
Quantization level assignments	8 bit systems: 254 of 256 levels (digital levels 1 through 254) of the eight-bit word are used to express quantized values. Data levels 0 (00_H) and 255 (FF_H) are reserved to indicate timing references. 10 bit systems: 1016 of 1024 levels (digital levels 4 through 1019) of the ten-bit word are used to express quantized values. Data levels 0 (00_H) through 3 (03_H) and 1020 ($3FC_H$) through 1023 ($3FF_H$) are reserved to indicate timing references.



T = Transmitter
 R = Receiver
 A, A' = Data line
 B, B' = Return line

Fig. 29. Positive binary signal convention.

information. In each system, a specific interface clock interval is designated line-word 0 for the purpose of reference only. All the other sample words per total line are therefore numbered 1 to N , where N is one less than the ratio of the interface clock frequency to the line frequency. For example, in SMPTE 125M, the ratio of the 27-MHz interface clock frequency to the 15.734-kHz line frequency is 858. The sample words are, therefore, numbered 0 to 857.

The EAV and SAV references consist of a four-word sequence in the following format:

3FF 000 000 XYZ

The first three words are a fixed preamble. The fourth word contains information providing even field (field 2) identification, the state of vertical blanking, and the state of horizontal blanking.

Because of the existence of both 8-bit and 10-bit equipment, for detection purposes all values in ranges 000-003 and 3FC-3FF must be considered equivalent to 000 and 3FF, respectively.

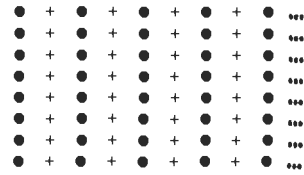
In each system, the half-amplitude point of the leading or falling edge of the analog horizontal-sync signal (OH) is described as coincident with a specified sampling point.

Ancillary Data Signal Format

Ancillary data may be inserted in any portion of the data stream not occupied by timing-reference signals or video data. Two categories of the ancillary data, horizontal ancillary data (HANC) and vertical ancillary data (VANC), are defined for different portions of the data stream. The same three-word header is used to identify both HANC and VANC ancillary data, although 8-bit representation of the header is permitted for VANC only.

Horizontal ancillary data is permitted in all horizontal intervals. It is not permitted in the active portion of the lines. HANC data is of 10-bit format, and each block of HANC data is preceded by the three-word ancillary-data header, 000 3FF 3FF. The ancillary header may occur multiple times during each horizontal-blanking period if different blocks of data are transmitted.

All permitted data identification words and data formats will protect the values 000 to 003 and 3FC to 3FF.



● = cosited C_B, Y, C_R sample
 + = isolated [Y] sample

Fig. 30. Siting of video samples.

Vertical ancillary data is permitted only in the active portion of specified lines of the vertical interval. In some systems, certain lines are reserved for digital vertical-interval time code (DVITC) and video index. VANC data is of the 8-bit format, and each block of VANC data is preceded by the three-word ancillary-data header, 000 3FF 3FF. The ancillary header may occur multiple times during each line period if different blocks of data are transmitted.

Advanced Television Systems

ATV Systems—The abbreviation ATV refers to television systems that provide improved audio and video quality or enhance the current television systems. The term ATV embraces both enhanced definition television (EDTV) and high definition television (HDTV). The term HDTV is generally accepted as referring to systems which provide approximately twice the resolution of the current NTSC services in both the horizontal and vertical dimensions. The term EDTV is generally accepted as referring to systems which provide improved resolution when compared to the current NTSC services but which do not qualify as HDTV systems. The audio performance of an ATV system should approximate that of current CD technology.

An interface signal format has been defined for a 525-line, 59.94-fields/s, wide-screen (wide screen aspect ratio = 16:9 = 1.7777:1) digital color television system. The system is an upscaled version of SMPTE 125M, modified to accommodate the requirements for the wide aspect-ratio images. The standard NTSC aspect-ratio image is 4:3 = 1.3333:1. The ratio of 1.7777 to 1.3333 is 1.3333 = 4/3. Therefore, the values in Table 9 and those in Table 11 have the relationship of 4/3.

HDTV—At least three different HDTV production systems have been developed for possible use in NTSC countries. The three systems are the 1125/60/2:1, 1050/59.94/2:1, and 787.5/59.94/1:1 systems, where the first value is the number of lines per frame, the second value is the number of fields per second, and the third quantity (X:Y) is the interlace factor. All three systems feature a wide-screen aspect ratio (16:9 = 1.7777:1).

An interface signal format has been defined for the 1125-line, 60-fields/s, 2:1-interlaced digital color

TABLE 11. ENCODING PARAMETERS FOR WIDE-SCREEN COMPONENT TELEVISION AT 4:2:2 PRODUCTION LEVEL*

Parameter	525-Line, 59.94-Fields/s	625-Line, 50-Fields/s
Coded signals	Y, C_B, C_R	
No. of samples, total line		
Luminance signal (Y)	1144	1152
Each color-difference signal (C_B, C_R)	572	576
Sampling structure	Orthogonal (line, field, and picture repetitive). C_B and C_R samples cosited with odd Y samples in line.	
Sampling frequency		
Luminance signal	18.0 MHz	
Each color-difference signal	9.00 MHz	
Multiplexed parallel data interface clock rate	36.0 MHz	
Serial interface clock rate	360 MHz	
Form of coding	Uniformly quantized PCM, 8 or 10 bits per sample, for the $Y, C_B,$ and C_R signals.	
Number of samples per digital active line		
Luminance signal	960	
Each color-difference signal	480	

* Reference 15.

TABLE 12. ENCODING PARAMETERS FOR WIDE-SCREEN HDTV 1125/60 COMPONENT TELEVISION AT 4:2:2 PRODUCTION LEVEL*

Parameter	1125/60/2:1 System
Coded signals	Y, C_B, C_R $Y = 0.701 G + 0.087 B + 0.212 R$ $C_B = (B - Y)/(1.826)$ $C_R = (R - Y)/(1.576)$
No. of samples, total line	
Luminance signal (Y)	2200
Each color-difference signal (C_B, C_R)	1100
Sampling structure	Orthogonal (line, field, and picture repetitive). C_B and C_R samples cosited with odd Y samples in line.
Sampling frequency	
Luminance signal	74.250 MHz
Each color-difference signal	37.125 MHz
Multiplexed parallel data interface clock rate	148.50 MHz
Form of coding	Uniformly quantized PCM, 8 or 10 bits per sample, for the $Y, C_B,$ and C_R signals
Number of samples per digital active line	
Luminance signal	1920
Each color-difference signal	960

* Reference 16.

TABLE 13. ENCODING PARAMETERS FOR COMPONENT 1050- AND 1250-LINE TELEVISION AT 4:2:2 LEVEL*

Parameter	1050-Line, 59.94-Fields/s	1250-Line, 50-Fields/s
Coded signals	Y, C_B, C_R	
No. of samples, total line Luminance signal (Y) Each color-difference signal (C_B, C_R)	2288 1144	2304 1152
Sampling structure	Orthogonal (line, field, and picture repetitive). C_B and C_R samples cosited with odd Y samples in line.	
Sampling frequency Luminance signal Each color-difference signal Multiplexed parallel data interface clock rate Serial interface clock rate	72.0 MHz 36.0 MHz 144 MHz 1440 MHz	
Form of coding	Uniformly quantized PCM, 8 or 10 bits per sample, for the $Y, C_B,$ and C_R signals.	
Number of samples per digital active line Luminance signal Each color-difference signal	1920 960	

* Reference 17.

HDTV electronic production system. The suggested parameter values are given in Table 12.

An interface signal format has been defined for the 1050-line, 59.94-fields/s, 2:1-interlaced digital color HDTV television system. Since the 1050-line system is

a doubled, upscaled version of the 525-line system described in SMPTE 125M, the digital-system parameter values that apply to the 1050-line system could also apply to a 1250-line, 50-fields/s, 2:1-interlaced system. The suggested parameter values are given in Table 13.

TABLE 14. ENCODING PARAMETERS FOR COMPONENT 787.5-LINE TELEVISION AT 4:2:2 LEVEL*

Parameter	787.5-line, 59.94-fields/s
Coded Signals	Y, C_B, C_R
No. of samples, total line Luminance signal (Y) Each color-difference signal (C_B, C_R)	1596 798
Sampling structure	Orthogonal (line, field, and picture repetitive). C_B and C_R samples cosited with odd Y samples in line.
Sampling frequency Luminance signal Each color-difference signal Multiplexed parallel data interface clock rate Serial interface clock rate	75.336 MHz 37.668 MHz 150.672 MHz 1506.72 MHz
Form of coding	Uniformly quantized PCM, 8 or 10 bits per sample, for the $Y, C_B,$ and C_R signals.
Number of samples per digital active line (square pixels system) Luminance signal Each color-difference signal	1280 640

* Reference 18.

An interface signal format has also been defined for a 787.5-line, 59.94-fields/s, 1:1 progressive scanned digital color HDTV television system. The suggested parameter values are given in Table 14.

References

References Cited in Text

1. Engstrom. "Study of Television Image Characteristics." *Proceedings of the IRE*, Vol. 21, 1933, pp. 1631-1651.
2. Wheeler and Loughren. "Fine Structure of Television Images." *Proceedings of the IRE*, Vol. 26, 1934, pp. 1246-1265.
3. Baron, S., and Pursell, S. "Digital Video Processing Systems." CBS Laboratories Report, Stamford, 3 April 1972, pp. 17-18.
4. *Ibid.*, pp. 37-39.
5. Pratt, William. *Digital Image Processing*. John Wiley and Sons, 1978, pp. 96-98.
6. SMPTE T14.391.
7. SMPTE T14.392.
8. SMPTE 125.
9. SMPTE 244.
10. Baron, S. "Dynamic Range of Digital Video Signals." SMPTE Working Group on Digital Video, 12 August 1977.
11. CCIR Report 624.
12. SMPTE 125M.
13. SMPTE 244M.
14. SMPTE 244.
15. Baron, S. "Proposed Studio Component Digital Representation and Bit-Parallel Digital Interface for 525-line Wide Screen Standard." August 1991.
16. SMPTE 240M.
17. SMPTE S17.391.
18. SMPTE S17.392.

Normative References

- SMPTE S17.394 Standard for Digital Component Video Signal Wide Screen.
- SMPTE 170M Studio Standard for Composite Video Signal, System M/NTSC.
- SMPTE 125M Bit-Parallel Digital Interface for Component Video Signals.
- SMPTE 240M Signal Parameters, 1125/60 High-Definition Production System.
- SMPTE 244M Bit-Parallel Digital Interface for Composite Video Signals.
- SMPTE RP-XXWM Recommend Practice for Video Index
- CCIR Recommendation 601: Encoding Parameters of Digital Television for Studios.
- CCIR Report 624: Characteristics of Television Systems.
- CCIR Recommendation 656: Interfaces for Digital Component Video Signals in 525-line and 625-line Television Systems.
- CCIR Report 962: The Filtering, Sampling, and Multi-

plexing for Digital Encoding of Colour Television Signals.

MIL-C-24308C, Connectors, Electrical, Rectangular, Nonenvironmental, Miniature, Polarized Shell, Rack and Panel, General Specifications for

INTERNATIONAL STANDARDS FOR BROADCASTING

CCIR Documents

The CCIR (International Radio Consultative Committee) is a branch of the International Telecommunications Union (ITU), a specialized agency of the United Nations. Thus the CCIR is a treaty organization related directly to the United Nations. The CCIR is concerned with generating documents dealing with the preparation, transmission, and reception of all kinds of information using radio signals, with the term "radio" being taken in the broadest sense. The CCIR has established several Study Groups, each dealing with specific aspects of radio communications. Each Study Group is further subdivided into smaller groups, each of which is assigned responsibility for the documents in a specialized area of interest. The scopes of the various groups are described below.

The Study Groups of interest in radio and television broadcasting are Study Group 10—Sound Broadcasting Service, and Study Group 11—Television Broadcasting Service. Also related and of interest to broadcasters is the work of the CMTT (Joint Committee on Television Transmission), a group set up jointly by the CCIR and the CCITT (International Telephone and Telegraph Consultative Committee). The latter organization is a branch of the ITU that deals with telegraph and telephone transmission systems. The CMTT is concerned with the transmission of signals over long distances through facilities that include a mixture of the responsibilities of the CCIR and CCITT.

The CCIR and CMTT produce documents of various kinds. Categories include Recommendations and Reports (the two most important), as well as Questions, Study Programs, Decisions, Resolutions, and Opinions. The documents are reviewed every four years. During this interval, there are two or three preliminary meetings of the Study Groups and their subgroups, to review proposals submitted by members. These proposals may take the form of modifications to existing documents, or the submission of entirely new documents. The agreements of the Study Groups after the last preliminary meeting are reviewed and approved at the final, or Plenary, Assembly of the CCIR, convened each four years.

The CCIR Plenary Assembly documents are published within a year of the Plenary in thirteen volumes, one for each Study Group. The volumes are bound in paper, with a green cover traditionally used for the English text. Hence, they are colloquially known as the "Green Books." They are available from the International Telecommunication Union in Geneva, Switzerland.

land and, in the United States, from the National Technical Information Service (see the earlier subsection "International Sound-Recording Standards" for the address).

The work of the CCIR is organized into the following groups and committees:

Study Group 10—Broadcasting Service (Sound)

Study Group 10A-1, Amplitude-modulation sound broadcasting

Study Group 10A-2, Sound broadcasting in the tropical zone

Study Group 10B, Frequency-modulation sound broadcasting

Study Group 10C, Audio-frequency characteristics of sound broadcasting signals

Study Group 11—Broadcasting Service (Television)

Study Group 11A, Characteristics of systems for monochrome and colour television

Study Group 11B, Ancillary television services

Study Group 11C, Control, measurement and international exchange of television programmes

Study Group 11D, Picture quality and parameters affecting it

Study Group 11E, Planning of television networks, protection ratios, television receivers and antennas

Study Group 11F, Digital methods of transmitting television information.

Joint Working Group 10/11F, Exchange of recorded sound programmes

Joint Working Group 10/11G, Exchange of recorded television programmes

Joint Working Group 10/11H, Use of film in television

Joint Working Group 10/11I, Utilization and synchronization of different programme supports

Selected CCIR Recommendations

Listed below are some of the more important Recommendations ("Rec") issued by the CCIR in the broadcasting, recording, and transmission fields.

Rec 265-6, Standards for the international exchange of programmes on film for television use

Rec 407-4, International exchange of sound programmes recorded in analogue form

Rec 408-5, Standards of sound recording on magnetic tape for the international exchange of programmes

Rec 450-1, Transmission standards for FM sound broadcasting at VHF

Rec 469-5, *Analogue television tape recording*. Standards for the international exchange of television programmes on magnetic tape

Rec 470-2, Television systems

Rec 500-4, Method for the subjective assessment of the quality of television pictures

Rec 501-2, Appraisal of programmes on colour film intended for television use

Rec 562-3, Subjective assessment of sound quality
Rec 601-2, Encoding parameters of digital television for studios

Rec 602-1, Exchange of television recordings for programme evaluation

Rec 642-1, Limiters for high-quality sound-programme signals

Rec 646, Source encoding for digital sound signals in broadcasting studios

Rec 647-1, A digital audio interface for broadcasting studios

Rec 648, Digital recording of audio signals

Rec 649, Measuring methods for analog audio disk and tape recordings

Rec 653-3, Teletext systems

Rec 656, Interfaces for digital component video signals in 525-line and 625-line television systems

Rec 657-1, *Digital television tape recording*. Standards for the international exchange of television programmes on magnetic tape

Rec 709, Basic parameter values for the HDTV standard for the studio and for international programme exchange

Rec 710, Subjective assessment methods for image quality in high definition television

Rec 711, Synchronizing reference signals for the component digital studio

Rec 713, Recording of HDTV images on film.

Rec 714, International exchange of programmes electronically produced by means of high definition television

Rec 715, International exchange of ENG recordings

IEC Publications

The International Electrotechnical Commission (IEC) publishes standards on a wide range of electrical and electronic subjects. IEC publications may be purchased through ANSI in New York. The headquarters address of the IEC is:

International Electrotechnical Commission
3, rue de Varembe
PO Box 131 - 1211 Geneva 20
Switzerland
TELEFAX (+41 22) 7333843

Important to this chapter are the IEC publications produced by Technical Committee 60, Recording. Some of these are listed below:

Pub. 94, Magnetic tape sound recording and reproducing systems

Pub. 347, Transverse track video recorders

Pub. 386, Method of measurement of speed fluctuations in sound recording and reproducing equipment

Pub. 461, Time and Control code for video tape recorders

- Pub. 503, Spools for 1 in (25.4 mm) video magnetic tape
- Pub. 558, Type C helical video tape recorders
- Pub. 574-10, Audio cassette systems
- Pub. 602, Type B helical recorders
- Pub. 698, Measuring methods for television tape machines
- Pub. 712, Helical-scan video-tape cassette system using 19 mm (3/4 in) magnetic tape, known as U-format
- Pub. 735, Measuring methods for video tape properties
- Pub. 756, Time base stability of non-broadcast video recorders
- Pub. 767, Helical-scan video tape cassette system using 12.65 mm (0.5 in) magnetic tape on type beta format
- Pub. 774, Helical-scan video tape cassette system using 12.65 mm (0.5 in) magnetic tape on type VHS
- Pub. 841, Audio recording—PCM encoder/decoder system
- Pub. 843, Helical-scan video tape cassette system using 8 mm magnetic tape—Video 8
- Pub. 883, Measuring method for chrominance signal-to-random noise ratio for video tape recorders
- Pub. 1016, helical-scan digital component video cassette recording system using 19 mm magnetic tape (format D-1)

In addition to the recording specifications shown above, IEC also publishes the work of related technical groups, including TC 12, *Radiocommunications (including SC-12A, Receivers)*, TC-84, *Equipment and systems in the field of audio, video and audiovisual engineering*, and of JTC-1, *ISO/IEC Joint technical committee for information technology*.

ISO Recommendations

The International Organization for Standardization (ISO) has primary responsibility for the standardization of mechanical systems. ISO publications are available through ANSI in New York. The headquarters address for ISO is:

International Organization for Standardization
1, rue de Varembe
Case postale 56
CH-1211 Geneva 20
Switzerland

The standardization of electrical and electronic systems is primarily the province of the IEC. Standards for motion picture films are drafted by ISO Technical Committee 36—Cinematography. A partial list of ISO TC 36 documents for television broadcasting appears below:

- ISO 23:1976, Camera usage of 35 mm motion-picture film

- ISO 25:1976, Camera usage of 16 mm motion-picture film
- ISO 26:1985, Projector usage of 16 mm motion-picture films for direct front projection
- ISO 490:1978, Magnetic stripes and magnetic recording head gaps for sound records on 16 mm motion-picture film perforated along one edge
- ISO 1188:1984, Recorded characteristic for magnetic sound records on full-coat 16 mm motion-picture film
- ISO 1189:1986, Recorded characteristic for magnetic sound records on 35 mm motion-picture film excluding striped release prints
- ISO 1223:1985, Picture areas for motion-picture films and slides for television—position and dimensions
- ISO 3640:1982, Motion-picture prints and sound records for international exchange of television programmes
- ISO 4242:1980, Recording head gaps for two sound records on 16 mm magnetic film—positions and width dimensions
- ISO 6035:1983, Viewing conditions for the evaluation of films and slides for television—colours, luminances and dimensions

Selected EBU Publications

The European Broadcasting Union publishes a well regarded technical journal (*The EBU Review—Technical*) and a wide variety of technical documents. The EBU supplies a comprehensive catalog of these publications that can be obtained from their Geneva headquarters:

European Broadcasting Union
Ancienne Route, 17A
Case postale 67
CH-1218 Grand Saconnex/Geneva
Switzerland
TELEFAX (+41 22) 798 58 97

Included in the scope of EBU publications are documents in the following categories:

- Lists of broadcasting stations in the European area, updated frequently*
- “Tech 3000” documents that reproduce committee reports and study results*
- EBU Technical Monographs, documents of a thorough tutorial nature*
- EBU Official Technical Texts, including Technical Standards, Technical Recommendations, Technical Statements, and Technical Information Sheets*

More than 200 of these technical documents are available by mail order from EBU. Examples of a few of the EBU documents include:

Tech 3094, Specifications for the basic signals recommended by the EBU for the synchronization of television sources

Tech 3097, EBU time-and-control codes for television tape-recordings (625-line television systems)

Tech 3238, Methods for measuring the main characteristics of television cameras

Tech 3252, Standard for recording digital television signals on magnetic tape in cassettes

N6-1989, Helical-scan television recording on 25.4 mm tape

N12-1990, Time-and-control codes for television tape recording

N13 and supplements, Digital television tape recording

R20-1987, Exchange of television programmes

recorded on transverse-track, format B or format C tapes

R48-1991, Allocation of audio channels in the D-1 and D-2 digital television tape recording formats

Acknowledgements

The following individuals offered invaluable assistance in compiling the lists of standards:

Mark Hyman, SMPTE
Derek Mullen, IEC
Ralph Justus, EIA
Jane Pessall Frock, NAB
Daniel Queen, AES
George Waters, EBU

36

Radar

Merrill I. Skolnik

Introduction to Radar	36-3
Prediction of Radar Range	36-4
Fluctuating Target Models	36-5
Radar Cross Section of Targets	36-7
Other Forms of the Radar Equation	36-7
Surveillance Radar Equation	
Tracking Radar Equation	
Surface Clutter Range Equation	
Volume Clutter Radar Equation	
Noise Jamming Radar Equation (Surveillance)	
Noise Jamming Radar Equation (Tracking)	
Self-Screening Range Equation	
Weather Radar Equation	
Synthetic Aperture Radar Equation	
HF Over-the-Horizon Radar Equation	
Symbol Definitions	
Radar Letter Bands	36-11
Radar Antennas	36-11
Antenna Gain	
Cosecant-Squared Antenna Loss	
Antenna Errors	

Coverage	36-13
Doppler Frequency Shift	36-14
Pulse Radars That Employ Doppler	36-14
MTI Radar	36-16
MTI Performance Limitations	
Pulse Compression	36-18
Clutter	36-18
Radar Measurement Accuracy	36-21
RMS Range Measurement Error	
RMS Doppler Frequency Measurement Error	
RMS Angle Measurement Error	
Tracking with Radar	36-21

INTRODUCTION TO RADAR

Radar is an electromagnetic sensor that greatly extends one's ability to detect reflecting objects (or targets) at long or short range and to accurately locate targets in fair, as well as poor weather. Since its introduction during the second world war, radar has been a vital part of air defense in its many forms and for other military missions such as battlefield surveillance, fighter/attack aircraft, ballistic missile defense, and antisubmarine warfare. It has also been important for many nonmilitary applications including weather observation (precipitation, severe storms, winds, and wind shear), observing beneath the ever-present clouds of a planet like Venus, probing below the surface of the Earth, high-resolution imaging of the Earth's surface in three dimensions, and for mapping of sea ice for the more efficient routing of shipping in northern regions. Targets of interest to radar have been aircraft, ships, missiles, spacecraft, vehicles, people, birds, insects, as well as the natural environment.

Radar operates by radiating from an antenna, a known waveform, usually a series of short-duration pulses. After a portion of the radiated energy is reflected by a target and returned to the radar, it is received by the antenna and processed in the receiver to detect the presence of a target and to determine something about its nature. Two of the basic measurements made by a radar are *range* (distance) and *angu-*

lar location of a target. By observing the location of a target over time, the radar can establish its trajectory, or *track*, and predict the target's future location. Many modern radars use to advantage the shift in frequency (relative to the frequency that was transmitted) of the echo signal from a moving target. The shift of the echo signal frequency from a moving target is caused by the doppler effect, something familiar from high school or college physics. The doppler frequency shift is proportional to the *radial velocity* of the target, so it can be used to separate the frequency-shifted signals from moving targets (such as aircraft) from large undesired stationary *clutter* echo signals from the land, sea, and weather. It is an important part of MTI (moving target indication), pulse doppler, and CW (continuous wave) radars that have to detect small moving targets in the midst of very large clutter echoes. The doppler shift is also the basis for meteorological radars that detect and recognize hazardous weather effects to provide information about the environment not readily available by other means. In addition to the usual measurements of range, angular location, relative velocity, and target track, radar sometimes can obtain information about the size, shape, symmetry, and surface properties of a target.

A simple block diagram illustrating the major subsystems of a radar that might be used for air surveillance is shown in Fig. 1. The transmitter is usually a *power amplifier* such as a klystron, traveling wave tube,

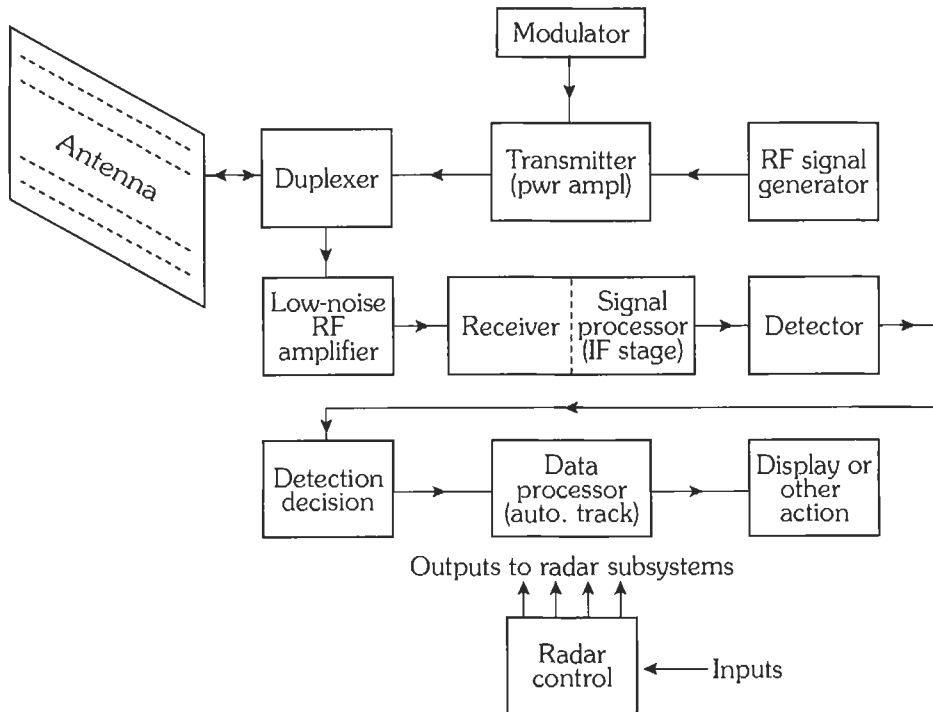


Fig. 1. Simple block diagram of a generic radar system.

or transistor. Although the magnetron oscillator was widely used in the early days of radar, its limited average-power, poor stability, and inability to generate sophisticated modulated waveforms restrict its application to radars with only modest capability. The first stage of the receiver is often a *low-noise transistor amplifier*. In a superheterodyne receiver, the echo is converted by a *mixer* and *local oscillator* (not shown) to an intermediate frequency (IF), where the signal is amplified and subject to signal processing to extract the desired signal and reject or attenuate undesired signals and noise. An important example of a *signal processor* is the *matched filter* that maximizes the ratio of the peak-signal-to-mean-noise output of the receiver, which in turn maximizes the detectability of the desired signal. In a receiver with a matched filter, the peak-output-signal-to-mean-noise-power ratio is $2E/N_0$, where E = signal energy and N_0 = noise power per unit bandwidth. Thus detectability of a radar signal when a matched filter is used does not depend on the shape of the signal or its bandwidth, but only on its total energy. The *detector* stage following the IF stage extracts the signal modulation from the carrier frequency. In a radar where there are no undesired clutter echoes to compete with the detection of the desired target echoes, the detector stage is an *envelope detector* (also called the second detector). In a radar that employs the doppler frequency shift to separate (by the use of filters) desired moving targets from undesired fixed clutter echoes, the detector stage is a *phase detector*. It requires a reference signal (not shown) that is a faithful representation of the transmitted signal so as to recognize that the echo signal has experienced a doppler shift. A *video amplifier* (not shown) following the detector amplifies the signal and a decision is made whether the receiver output is due to a target (signal plus noise) or is due to noise alone. The *detection decision* is based on observing when the receiver output exceeds a predetermined threshold whose level depends on achieving an acceptable probability of false alarm. In early radars the detection decision was made by an operator viewing a radar display, but in modern radars the decision whether or not a target has been detected is made automatically without direct operator intervention.

The received signal is digitized for processing, either after the detector stage (in the video) or before the detector (in the IF), especially in radars that depend on the doppler frequency shift for detection of moving targets. Digital processing makes it possible to automatically detect and accurately track many hundreds or thousands of targets so as to present fully processed tracks rather than individual detections or "raw" (unprocessed) radar data. The automatic tracker is an example of a *data processor*. The processed output of the radar or the established tracks of targets might be displayed to an operator or used to perform some automated operation. The *antenna* can be one of several different forms of mechanically steered parabolic reflectors, a mechanically steered planar array, or one of several types of electronically steered phased

arrays. The *duplexer* is the device that allows a single antenna to be time-shared between the transmitter and the receiver.

A typical long-range air-surveillance radar might have a resolution in the range dimension of about one or two hundred meters. When required, a radar can have a range resolution of a small fraction of a meter. The beamwidth of a radar antenna might typically be one or two degrees, but some operational radars have had beamwidths as small as 0.3 degree. Thus the resolution in the cross-range dimension (determined by the beamwidth and the target range) is usually much worse than the range resolution. It is possible, however, to achieve high resolution in the cross-range dimension comparable to the resolution achieved in the range dimension by employing synthetic aperture radar (SAR). Here the resolution of a large antenna is obtained by utilizing a small antenna on a moving platform, such as an aircraft, to store the received echoes over a relatively long time so as to synthesize (virtually) in a digital processor the equivalent of a large antenna. The output of a SAR is usually a high resolution map or image of a target scene.

Radar is generally found within what is known as the microwave region of the electromagnetic spectrum from about 400 MHz to 40 GHz; but there have been many operational radars in the VHF region (30 to 300 MHz) as well as in the HF region (3 to 30 MHz). An HF over-the-horizon radar can reach out to ranges of about 2000 nmi by utilizing refraction from the ionosphere. Radar has also been considered for use at frequencies higher than the microwave region, at millimeter wavelengths. Laser radars are found in the IR and optical regions of the spectrum, where they can provide precision range and radial-velocity measurement.

PREDICTION OF RADAR RANGE*

The radar range equation is important not only for predicting the range performance of a radar, but to act as a focus for radar design and for better understanding the factors that affect radar performance. The simple form of the radar range equation is

$$P_r = P_t G A_e \sigma / [(4\pi)^2 R^4] \quad (\text{Eq. 1a})$$

where

P_r = received echo signal power in watts,

P_t = transmitted signal power in watts,

G = antenna gain,

A_e = antenna effective area in square meters,

σ = radar cross section of the target in square meters,

R = range to the target in meters.

* Reference 1.

If a single antenna is used for both transmitting and receiving, as is usually the case, $G = 4\pi A_c/\lambda^2$, where λ is the radar wavelength in meters. Then

$$P_r = P_t G^2 \lambda^2 \sigma / [(4\pi)^3 R^4] \quad (\text{Eq. 1b})$$

$$= P_t A_c^2 \sigma / [4\pi \lambda^2 R^4]$$

The maximum range R_{\max} of a radar occurs when the received signal $P_r = S_{\min}$, the minimum detectable signal. The minimum detectable signal is a statistical quantity limited by receiver noise. It can be written as

$$S_{\min} = kT_0 B F_n (S/N)_1 \quad (\text{Eq. 2})$$

where

- k = Boltzmann's constant,
- T_0 = standard temperature (290 K),
- $kT_0 = 4 \times 10^{-21}$ W/Hz
- B = receiver bandwidth in hertz,
- F_n = receiver noise figure,
- $(S/N)_1$ = minimum signal-to-noise ratio required for reliable detection.

The received echo signal power can be increased by integrating (adding) a number of echo signal pulses n . This can be incorporated into the radar equation by dividing S_{\min} by $nE_i(n)$, where $E_i(n)$ is the efficiency with which the n pulses can be integrated. Since the average power P_{av} is more indicative of radar capability than is the peak power, it is introduced via the relation

$$P_{\text{av}} = P_t \tau f_p \quad (\text{Eq. 3})$$

where

- τ = pulse width in seconds,
- f_p = pulse repetition frequency in hertz.

With the above, the form of the radar equation suitable for calculating the range is

$$R_{\max} = \left[\frac{P_{\text{av}} G^2 \lambda^2 \sigma n E_i(n)}{(4\pi)^3 kT_0 F_n (B\tau) f_p (S/N)_1 L_s} \right]^{1/4} \quad (\text{Eq. 4})$$

The radar system losses L_s (number greater than one) have been included. For most radars designed with a matched filter receiver (a filter that maximizes the output signal-to-noise ratio), the product $B\tau \approx 1$. [In Eq. 4, $(S/N)_1/nE_i(n)$ is the required signal-to-noise ratio per pulse $(S/N)_1$.]

Fig. 2 shows the relationship of the required signal-to-noise ratio $(S/N)_1$ to the probability of detection and the probability of false alarm. The probability of detection is usually taken as 0.90, but sometimes it is quoted as 0.5 or 0.8. Its choice is usually the prerogative of the customer. The probability of a false alarm is given here as

$$P_{\text{fa}} = 1/BT_{\text{fa}}$$

where

- B = receiver bandwidth in hertz,
- T_{fa} = average time between false alarms.

The reciprocal of P_{fa} is n_f , the false-alarm number. The false-alarm time T_{fa} is usually specified for radar performance rather than the probability of false alarm or the false-alarm number.

Fig. 3 is a plot of the integration-improvement factor $nE_i(n)$ as a function of n . The number of pulses returned from a target when an antenna of beamwidth θ_B degrees rotates at a rate of ω_m revolutions per minute, with a pulse repetition rate of f_p Hz is

$$n = \theta_B f_p / 6\omega_m \quad (\text{Eq. 5})$$

Failure to include the many factors that contribute to the system losses L_s can result in considerable difference between the calculated range and the actual range. Losses include:

- Loss in the transmission line connecting the antenna to the transmitter and receiver.
- Loss in the duplexer, rotary joint, and other microwave components.
- Beam-shape loss, to account for the fact that the radar equation employs the maximum gain rather than a gain that changes pulse to pulse as the antenna is scanned past the target.
- Signal processing losses, which can sometimes be surprisingly large.
- Loss due to degradation of transmitter power and receiver noise figure.

The system losses from all factors might be from 10 to 20 dB, or even greater. (A loss of 16 dB reduces the radar range by a factor of two.)

FLUCTUATING TARGET MODELS

The complex nature of most radar targets causes the radar cross section to vary with changing aspect. Fluctuations in the radar cross section are difficult to specify precisely; however, four simple statistical models first described by P. Swerling (Reference 2) are often used for computing the radar range. The statistical nature of the radar cross section σ for cases 1 and 2 is described by the probability density function

$$p(\sigma) = (1/\sigma_{\text{av}}) \exp[-\sigma/\sigma_{\text{av}}] \quad (\text{Eq. 6a})$$

where σ_{av} is the average cross section. Cases 3 and 4 are described by

$$p(\sigma) = (4\sigma/\sigma_{\text{av}}) \exp[-2\sigma/\sigma_{\text{av}}] \quad (\text{Eq. 6b})$$

In cases 1 and 3, it is assumed that the echo pulses received from a target on any one scan are of constant amplitude throughout the entire scan but are independent (uncorrelated) from scan to scan. In cases 2 and 4, the cross-section is assumed to be independent from

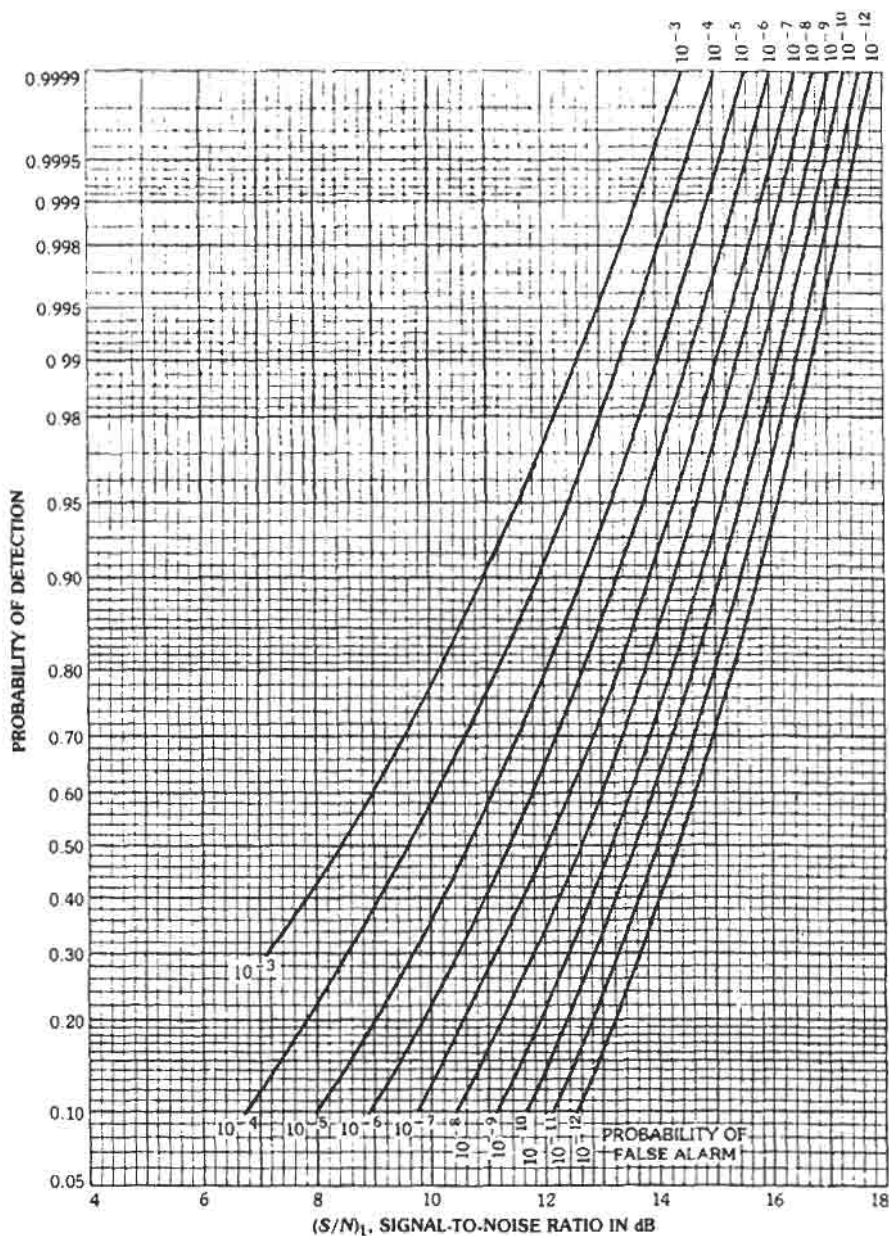


Fig. 2. Probability of detection for a sine wave in noise as a function of the signal-to-noise (power) ratio and the probability of false alarm. (Courtesy McGraw-Hill Book Co.)

pulse to pulse. When these models are used, the signal-to-noise ratios and the integration improvement factors inserted in the radar equation are different from those used for a constant cross section. The required values can be found in several references (References 3-7); however, a simple approximate method (Reference 1) for most purposes can be obtained using Figs. 2, 4, and 5 with the following procedures:

- Find the signal-to-noise ratio from Fig. 2 corresponding to the desired values of detection probability P_d and false-alarm probability $P_{fa} = 1/BT_{fa}$.
- From Fig. 4, determine for the desired Swerling case the correction factor to be applied to the signal-to-noise ratio found from Fig. 2.
- Find the integration improvement factor $nE_c(n)$ from Fig. 5.

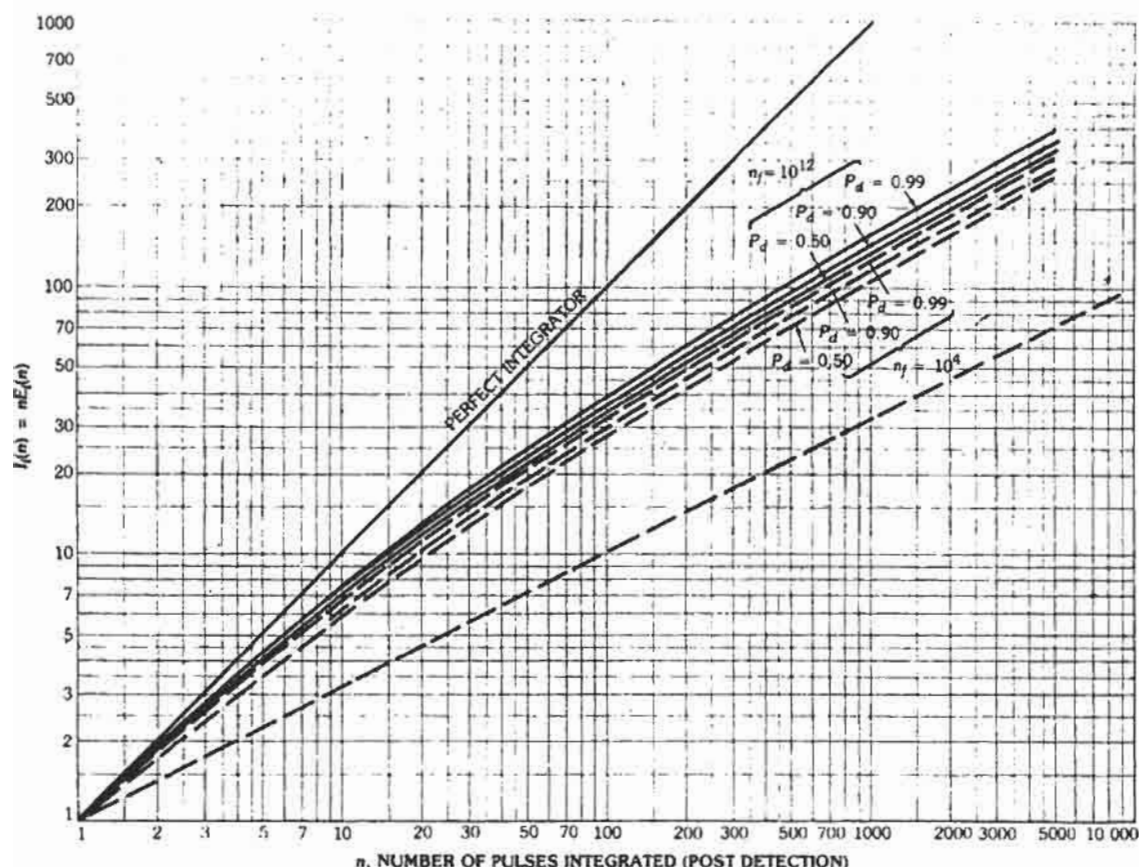


Fig. 3. Integration-improvement factor, assuming square-law detector. P_d = probability of detection, $n_f = T_{fa}B$ = false alarm number, T_{fa} = average time between false alarms, B = bandwidth. (Courtesy McGraw-Hill Book Co.)

Note that in cases 2 and 4 when $P_d \leq 0.9$ and the number of pulses are large ($n > 20$ for case 2 and $n > 10$ for case 4), cases 2 and 4 can be replaced by the nonfluctuating target model (indicated in Fig. 5 as case 5).

A more general expression for describing the statistical properties of the radar target cross section is the chi-square distribution of degree $2m$, whose probability density function is

$$p(\sigma) = [m!(m-1)! \sigma_w] (m\sigma/\sigma_w)^{m-1} \exp(-m\sigma/\sigma_w) \quad (\text{Eq. 7})$$

The Swerling cases 1 and 2 are given by the chi-square distribution with $m = 1$; the other two cases are given with $m = 2$.

RADAR CROSS SECTION OF TARGETS*

The radar cross section of a complex target, such as an aircraft or ship, varies considerably with change in aspect or change in frequency so that a single number cannot adequately describe the radar cross section of a

target. Nevertheless, Table 1 lists "example" values for various targets at microwave frequencies. These are for illustrative purposes to show the relative "sizes" of common targets as "seen" by radar.

OTHER FORMS OF THE RADAR EQUATION

The radar equation is used for the calculation of range; but it is also used as a basis for assessing the tradeoffs in radar design. The simple forms of the radar equation given above (Eqs. 1 and 2) are seldom sufficiently complete, however, and they must be extended. Each specific radar application has some particular requirements or constraints that can result in a slightly different form of the radar equation. Examples are presented below. All symbols used in the equations are defined at the end of this section.

* References 1 and 8.

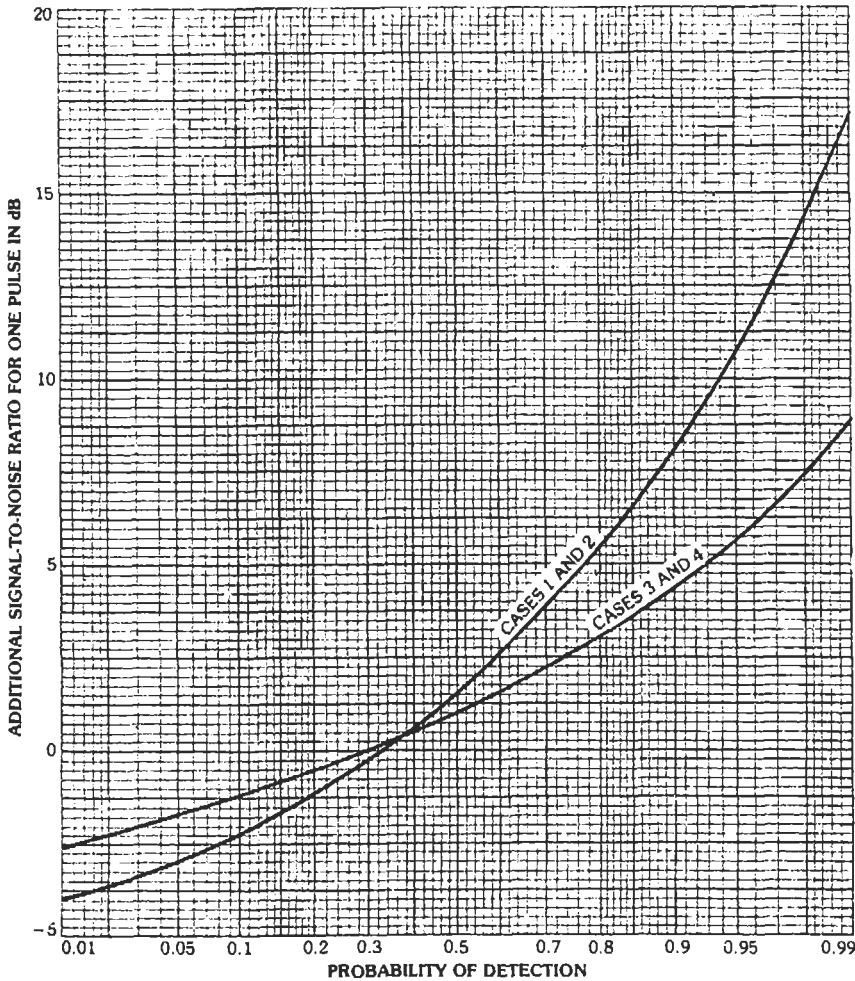


Fig. 4. Additional signal-to-noise ratio required to achieve a particular probability of detection, when the target cross section fluctuates according to the Swerling models, as compared with a nonfluctuating target; single hit, $n = 1$. To be used with Fig. 2 to find $(S/N)_1$. (Courtesy McGraw-Hill Book Co.)

Surveillance Radar Equation

$$R_{\max}^4 = \frac{P_{\text{av}} A_e \sigma E i(n)}{4\pi k T_0 F (S/N)_1 L_s} \cdot \frac{t_s}{\Omega} \quad (\text{Eq. 8})$$

This equation applies to a radar that must observe all targets within an angular region of solid angle Ω steradians once every t_s seconds. When the surveillance radar utilizes a conventional rotating fan beam whose elevation beamwidth is θ_e , the solid angle Ω equals $2\pi \sin \theta_e$, and t_s is the azimuth rotation period (or revisit time).

Tracking Radar Equation

Equation 4 is basically the tracking radar equation, where $n/f_p = t_0$ is the signal integration time. (It has also been called the searchlight equation.)

Surface Clutter Range Equation

$$R_{\max} = \sigma n_c J [(S/C)_0 \sigma^o \theta_a (c\tau/2) \sec \psi] \quad (\text{Eq. 9})$$

This equation describes the detection of a target when viewed at a low grazing angle ψ in the presence of surface clutter of radar cross section per unit area σ^o . It assumes that the received clutter echo power is much greater than receiver noise. The effective number of pulses integrated, n_c , will depend on the decorrelation time of the clutter. (With completely stationary clutter, there is no integration improvement so that $n_c = 1$.)

Volume Clutter Radar Equation

$$R_{\max}^2 = \sigma G n_c J [(S/C)_0 \eta (\pi^3/4) (c\tau/2)] \quad (\text{Eq. 10})$$

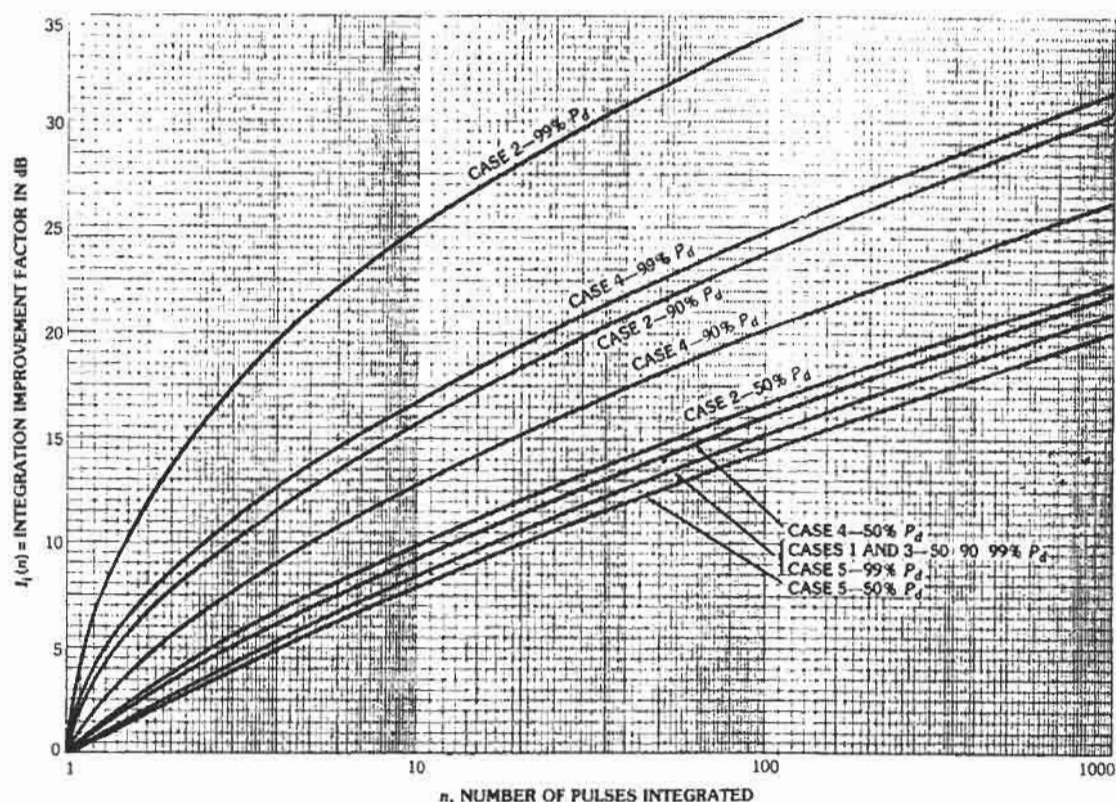


Fig. 5. Integration-improvement factor as a function of the number of pulses integrated for the four Swerling cases (cases 1 to 4) and the constant cross section (case 5). (Courtesy McGraw-Hill Book Co.)

The reflectivity η is the radar cross section of the clutter per unit volume.

Noise Jamming Radar Equation (Surveillance)

$$R_{\max}^2 = \frac{P_{av} G E_i(n)}{G_{SL} L_s} \cdot \frac{\sigma}{(S/N)_1} \cdot \frac{t_r}{\Omega} \cdot \frac{B_j}{P_j G_j} \quad (\text{Eq. 11})$$

This equation assumes that the jamming noise enters the antenna sidelobes whose gain is G_{SL} . When the jamming enters the main beam, $G_{SL} = G$. The jammer power P_j is spread over a bandwidth B_j and is radiated by an antenna whose gain is G_j .

Noise Jamming Radar Equation (Tracking)

$$R_{\max}^2 = \frac{P_{av} G^2 E_i(n) t_0}{4\pi G_{SL} L_s} \cdot \frac{\sigma}{(S/N)_1} \cdot \frac{B_j}{P_j G_j} \quad (\text{Eq. 12})$$

When the jamming noise enters the radar via the main beam, $G_{SL} = G$.

Self-Screening Range Equation

This is the range at which the radar echo signal S received from a target exceeds the received jamming noise power J by the amount S/J . It is also called the cross-over range. The self-screening range is found from either Eq. 11 or 12 (depending on the application) by setting $G_{SL} = G$, setting $(S/N)_1 = S/J$, and calling R_{\max} the self-screening range R_{ss} . The value of required S/J is often taken to be the same as $(S/N)_1$ found for receiver noise.

Weather Radar Equation

$$\bar{P}_r = 2.4 P_t G \pi r^{1.6} / [R^2 \lambda^2 L_s] \quad (\text{Eq. 13})$$

This equation is employed by radar meteorologists to relate the average echo signal power \bar{P}_r , to the rainfall rate r (mm/hr). It assumes that rain uniformly fills the radar resolution cell.

TABLE 1. "EXAMPLE" VALUES OF RADAR CROSS SECTION

Target	σ (square meters)
Conventional unmanned winged missile	0.1
Small single-engine aircraft	1
Small fighter or 4-passenger jet	2
Large fighter	6
Medium bomber or medium jet airliner	20
Large bomber or large jet airliner	40
Jumbo jet	100
Helicopter	3
Small open boat	0.02
Small pleasure boat	2
Cabin cruiser	10
Ship, grazing angle greater than zero	Displacement tonnage expressed in m^2
Pickup truck	200
Automobile	100
Bicycle	2
Man	1
Large bird	10^{-2}
Medium bird	10^{-3}
Large insect (locust)	10^{-4}
Small insect (fly)	10^{-5}

Synthetic Aperture Radar Equation

$$\frac{S}{N} = \frac{2P_{av} \rho_a^2 \sigma^\circ \delta_{cr} \delta_r}{\pi f k T_0 F_n R S_w L_s \sin^2 \psi} \quad (\text{Eq. 14})$$

This equation relates the signal-to-noise ratio of a resolution cell (sometimes called a pixel) with range resolution δ_r and cross-range resolution δ_{cr} located within a swath S centered at a range R . The above takes account of the combined restriction on cross-range resolution and swath necessary to avoid ambiguities in either range or cross range.

HF Over-the-Horizon Radar Equation

$$R_{max}^2 = \frac{P_{av} G_t G_r \lambda^2 \sigma F_p^2 T_c}{(4\pi)^3 N_0 (S/N)_1 L_s} \quad (\text{Eq. 15})$$

The transmitting antenna gain G_t and the receiving antenna gain G_r are shown separate since two different antennas are often used for transmit and receive. The

propagation loss is accounted for by F_p (number less than unity), and T_c is the coherent processing time. The noise power per unit bandwidth N_0 (W/Hz) at the receiver is determined by external noise.

Symbol Definitions

The symbols used in the above radar equations are defined as follows.

- A_e = antenna effective aperture in square meters
- B = receiver bandwidth in hertz
- B_j = jammer bandwidth in hertz
- c = velocity of propagation in meters/second
- δ_{cr} = cross-range resolution in meters
- δ_r = range resolution in meters
- $E_i(n)$ = efficiency in integrating n pulses
- η = volume clutter of reflectivity, or radar cross section of clutter per unit volume, in meters⁻¹
- f = radar frequency in hertz
- f_p = pulse repetition frequency in hertz
- F_n = receiver noise figure
- F_p = propagation factor
- G = antenna gain
- G_j = jammer antenna gain
- G_r = radar receiving antenna gain
- G_{SL} = antenna sidelobe gain
- G_t = transmitting antenna gain
- k = Boltzmann's constant = 1.38×10^{-23} joules/kelvin
- L_s = system losses
- λ = wavelength in meters
- n = number of echo pulses received per target
- n_e = effective number of pulses integrated
- N_0 = noise power per unit bandwidth
- Ω = solid angular region (steradians) of radar coverage
- P_{av} = average power in watts
- P_j = jammer power in watts
- P_r = average received signal power in watts
- P_r = received signal power in watts
- P_t = peak power in watts
- ψ = grazing angle
- r = rainfall rate in millimeters/hour
- R = range in meters
- R_{max} = maximum radar range in meters
- ρ_a = antenna efficiency
- $(S/C)_0$ = minimum signal-to-clutter ratio necessary to detect a target with a specified probability of detection and probability of false alarm, for a single pulse
- S/N = signal-to-noise ratio in a SAR resolution cell
- $(S/N)_1$ = minimum signal-to-noise ratio necessary to detect a target with a specified probability of detection and probability of false alarm, for a single pulse
- S_w = swath width in meters
- σ = radar cross section of target in square meters
- σ° = radar cross section of surface clutter per unit area
- t_s = scan time, or revisit time, in seconds

- $t_0 = n/f_p$ = signal integration time in seconds
- T_c = coherent processing time in seconds
- T_0 = standard temperature = 290 K
- τ = pulse width in seconds
- θ_a = azimuth beamwidth in radians
- θ_e = elevation beamwidth in radians
- v = velocity in meters/second

RADAR LETTER BANDS

The frequency bands in which radar operates have traditionally been designated by letters of the alphabet. This nomenclature is given in Table 2 along with the International Telecommunication Union frequency allocations for "radiolocation." There are other letter-band nomenclatures that have been employed in the past for various purposes, but those of Table 2 have been accepted as a Standard by the Institute of Electrical and Electronic Engineers (IEEE) and are the only ones that should be employed for *radar*.

RADAR ANTENNAS

The half-power beamwidth of a radar antenna of dimension D is

$$\theta_B = \beta\lambda/D \quad (\text{Eq. 16})$$

where beamwidth θ_B is in degrees, and wavelength λ and dimension D are in the same units. The constant β depends on the shape of the illumination (current distribution) across the aperture. It is sometimes called the normalized beamwidth. The choice of aperture illumination determines the normalized beamwidth and the peak side-lobe level. Equation 16 is plotted in Fig. 6 for $\beta = 65$ which is a typical value for horn-fed parabolic reflector antennas. The Taylor illumination is often used in radar-antenna design. It is characterized as having its first $\bar{n} - 1$ sidelobes equal, beyond which the sidelobe amplitude decreases with increasing angle. The variation of the beamwidth as a function of the sidelobe level of a planar antenna designed with a Taylor circular aperture illumination is shown in Fig. 7 for various \bar{n} (References 9 and 10). For a Taylor line source, the beamwidths of Fig. 7 are reduced by the (approximate) factor $\bar{n}/(\bar{n} + 0.25)$.

Table 3 illustrates how the antenna beamwidth, null width, peak sidelobe level, and the relative antenna gain vary as the aperture illumination becomes more tapered at the ends of the aperture compared to the illumination at the center of the antenna. The example in this table is that of an aperture illumination defined as $\cos^N(\pi x/2)$, where N is an integer, x is the normalized distance from the center of the aperture, and $|x| \leq 1$. The $\cos^N(\pi x/2)$ aperture illumination is not commonly

TABLE 2. RADAR LETTER-BAND NOMENCLATURE

Band Designation	Nominal Frequency Range	Specific Radiolocation (Radar) Bands Based on ITU Assignments for Region 2
HF	3 – 30 MHz	
VHF	30 – 300 MHz	138 – 144 MHz 216 – 225
UHF	300 – 1000 MHz	420 – 450 MHz 850 – 942
L	1 – 2 GHz	1215 – 1400 MHz
S	2 – 4 GHz	2300 – 2500 MHz 2700 – 3700
C	4 – 8 GHz	5250 – 5925 MHz
X	8 – 12 GHz	8500 – 10 680 MHz
K _u	12 – 18 GHz	13.4 – 14.0 GHz 15.7 – 17.7
K	18 – 27 GHz	24.05 – 24.25 GHz
K _s	27 – 40 GHz	33.4 – 36.0 GHz
V	40 – 75 GHz	59 – 64 GHz
W	75 – 110 GHz	76 – 81 GHz 92 – 100
mm	110 – 300 GHz	126 – 142 GHz 144 – 149 231 – 235 238 – 248

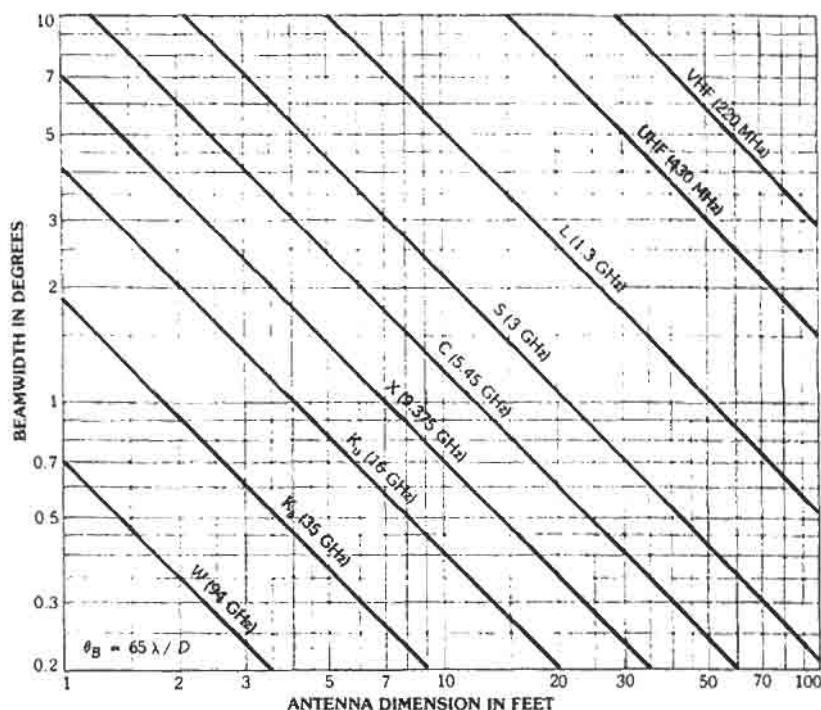


Fig. 6. Antenna beamwidth as a function of the antenna dimension for the various radar letter bands.

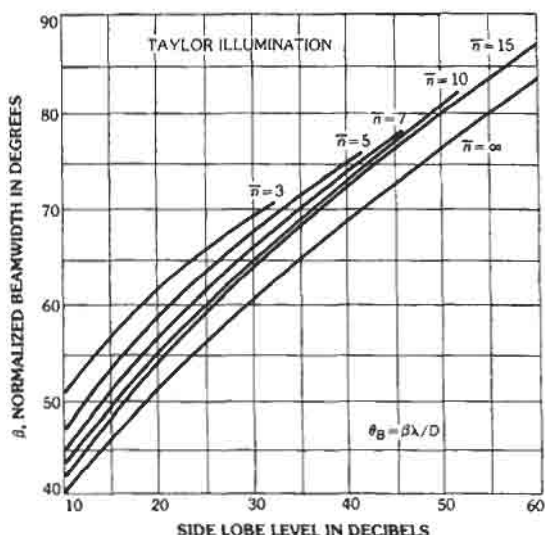


Fig. 7. Normalized beamwidth, β , as a function of the peak sidelobe level for the Taylor circular aperture illumination in which the first $\bar{n} - 1$ sidelobes are all equal. (After Reference 10.)

used for radar antennas, but it is a convenient analytical representation that illustrates the effect that the shape of the aperture illumination has on the character-

istics of the antenna pattern. When $N = 0$, the aperture illumination is constant (uniform) across the aperture. It has the greatest gain and the narrowest beamwidth for a given size aperture, but its peak sidelobe is high (-13.2 dB below the peak of the main beam). As N increases, the peak sidelobe decreases, but the antenna gain decreases and the beamwidth increases. Thus low antenna sidelobes are obtained by tolerating a reduction in gain and a wider beamwidth (or by using a bigger antenna).

Antenna Gain

A useful expression for the gain of "typical" reflector antennas, when no better information is available about the antenna other than its beamwidths, is (Reference 11)

$$G_D = 26,000 / \theta_B \phi_B \quad (\text{Eq. 17})$$

where θ_B and ϕ_B (in degrees) are the half-power beamwidths measured in the two principal (orthogonal) planes.

Cosecant-Squared Antenna Loss

The cosecant-squared antenna has a conventional beam shape in elevation from $\phi = 0$ to $\phi = \phi_0$ ($\phi =$ elevation angle), but its gain is proportional to $\csc^2 \phi / \csc^2 \phi_0$ from $\phi = \phi_0$ to $\phi = \phi_m$. It is used mainly to shape the

TABLE 3. ANTENNA CHARACTERISTICS FOR COSINE APERTURE ILLUMINATION

Aperture illumination	3-dB beamwidth (radians)	Null width (radians)	Peak sidelobe (dB)	Relative gain (relative loss)
N = 0	0.88 λ/D	2 λ/D	13.2	1.00 (0 dB)
N = 1	1.20 λ/D	3 λ/D	23	0.81 (0.92 dB)
N = 2	1.45 λ/D	4 λ/D	32	0.667 (1.18 dB)
N = 3	1.66 λ/D	5 λ/D	40	0.575 (2.40 dB)
N = 4	1.93 λ/D	6 λ/D	48	0.515 (2.88 dB)

elevation pattern of an air-surveillance radar antenna so as to more efficiently cover the region where aircraft are expected to fly. For small ϕ_0 and large ϕ_m the loss in gain of a cosecant-squared antenna relative to the gain of a conventional antenna of beamwidth ϕ_0 (in radians) is approximately

$$L_c \approx 2 - \phi_0 \cot \phi_m \quad (\text{Eq. 18})$$

Antenna Errors

The gain of an antenna when the rms phase error (radians) of the aperture illumination is $(\delta^2)^{1/2}$ is approximately (Reference 12)

$$G = G_0 \exp(-\delta^2) \quad (\text{Eq. 19})$$

where G_0 = no-error antenna gain. In a reflector antenna, the phase error might be due to deformation of the reflecting surface from its true value. For a circular paraboloidal reflector antenna of diameter D with antenna efficiency ρ_a , and with an rms mechanical tolerance ϵ measured in the same units as the wavelength λ , the maximum gain that can be achieved due to the resulting phase errors is

$$G_{\max} = (\rho_a/43)(D/\epsilon)^2, \quad (\text{Eq. 20})$$

which occurs at a wavelength

$$\lambda_{\max} = 4\pi\epsilon \quad (\text{Eq. 21})$$

The gain of a phased array antenna with δ^2 = mean-square phase error, Δ^2 = mean-square (relative) amplitude error, and P_c = fraction of the elements operative, is approximately (Reference 12)

$$G \approx G_0 P_c / (1 + \delta^2 + \Delta^2) \quad (\text{Eq. 22})$$

where G_0 = no-error gain.

When discrete phase shift is employed, as in digital phase shifters, the sidelobe level due to the quantization is approximately (Reference 13)

$$\text{rms sidelobe level} \approx 5/(2^{2B}N) \quad (\text{Eq. 23})$$

where

- B = number of bits in the phase shifter,
- N = total number of elements in the array.

The above assumes a random distribution of the quantization phase error across the aperture. If the quantized phase error is periodic, the peak lobe relative to the main beam is

$$\text{Peak quantization lobe} = 1/2^{2B} \quad (\text{Eq. 24})$$

Its position is

$$\sin \theta_q = (1 - 2^B) \theta_0 \quad (\text{Eq. 25})$$

where θ_0 = angle to which the beam is steered. The maximum pointing error due to quantization is

$$\Delta \theta_0 = (\pi/4)(1/2^B) \theta_B \quad (\text{Eq. 26})$$

where θ_B = beamwidth.

COVERAGE

The curvature of the earth limits the coverage of an earth-based radar. From simple geometrical considerations, the distance d (in nautical miles) to the horizon from a radar at height h (in feet) is

$$d = 1.23\sqrt{h} \quad (\text{Eq. 27})$$

This assumes that the refraction of the radar energy by the atmosphere can be represented by an effective earth radius 4/3 times the actual radius. If the target is at a height h_t the distance between the radar and target when the line of sight just grazes the surface of the earth is given by

$$d_0 = 1.23(\sqrt{h} + \sqrt{h_t}) \quad (\text{Eq. 28})$$

where d_0 is in nautical miles and h and h_t are in feet.

The presence of the earth causes part of the energy to be reflected from the surface. This reflected wave can interfere either constructively or destructively with the direct wave to produce a series of maxima (lobes) and minima (nulls). The angle of the first (lowest) lobe in radians is given by

$$\theta_1 = \lambda/4h \quad (\text{Eq. 29})$$

where

- λ = wavelength
- h = radar antenna height.

Both λ and h are measured in the same units. This equation assumes a flat earth.

An example of the elevation lobing pattern resulting from constructive and destructive interference for a 900 MHz radar is shown in Fig. 8.

DOPPLER FREQUENCY SHIFT

The *doppler frequency shift* of an echo signal reflected from a moving target is

$$f_d = (2v \cos\theta)/\lambda \quad (\text{Eq. 30})$$

where

- f_d is the frequency shift in hertz,
- v is the velocity of the target in meters/second,
- λ is the wavelength in meters,
- θ is the angle defined by the direction of target travel and the radar line of sight to the target.

The relative velocity is $v_r = v \cos \theta$. The doppler frequency shift per knot of relative velocity (f_d/v_r) is plotted in Fig. 9. When v_r is in knots and λ in meters, an approximate expression (to an accuracy of about 3 percent) for the doppler frequency shift is the following:

$$f_d (\text{Hz}) \approx v_r (\text{kt})/\lambda (\text{m}) \quad (\text{Eq. 31})$$

In a pulse radar, the measurement of the doppler frequency shift is ambiguous if the doppler frequency is

greater than the Nyquist rate, which in this case is twice the pulse repetition frequency (prf). When the prf is equal to the doppler frequency f_d or some multiple, the target velocity cannot be distinguished from stationary clutter; *i.e.*, it will appear to have no doppler shift. The relative velocity that produces a doppler frequency equal to the prf or some multiple is called a *blind speed*. The blind speeds are given by

$$v_b = n\lambda f_p/2 \quad (\text{Eq. 32})$$

where

- v_b is the blind speed in meters/second,
- λ is the wavelength in meters,
- f_p is the pulse repetition frequency in hertz,
- n is an integer.

The first blind speed (Fig. 10) is approximately (to an accuracy of about 3 percent)

$$v_1 (\text{kt}) \approx \lambda (\text{m})f_p (\text{Hz}) \quad (\text{Eq. 33})$$

PULSE RADARS THAT EMPLOY DOPPLER

There are three types of pulse radars that extract the doppler frequency shift, or relative velocity, in addition to the range information. These are:

1. A *moving target indicator* (MTI) radar in which the pulse repetition frequency (prf), denoted by f_p ,

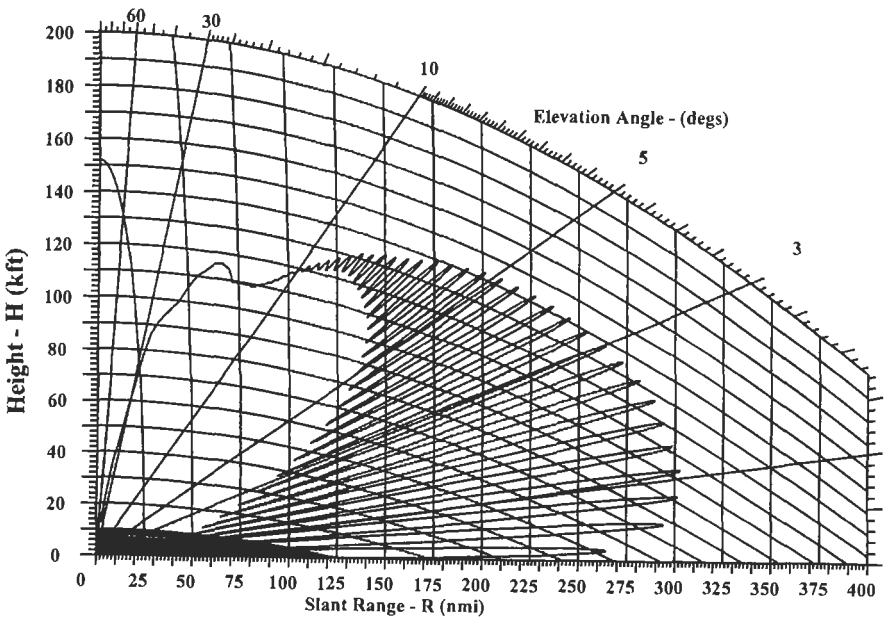


Fig. 8. Calculated elevation coverage of an air-surveillance radar at a frequency of 900 MHz, showing the uneven coverage due to the lobing that results from constructive and destructive interference between the direct and the surface-reflected waves. (Courtesy of Vilhelm Gregers-Hansen, Naval Research Laboratory.)

is low enough to have no range ambiguities; *i.e.*, its maximum unambiguous range $R_{unamb} = c/2f_p$. On the other hand, the MTI radar has ambiguities in doppler and many blind speeds (due to doppler frequency ambiguities) where moving targets are not detected.

2. A high pulse repetition frequency pulse doppler radar in which the prf is large enough to have no blind speeds [as given by Eq. (32)] within the expected values of target doppler frequencies. There will be, however, many range ambiguities because of the high prf. Range ambiguities are

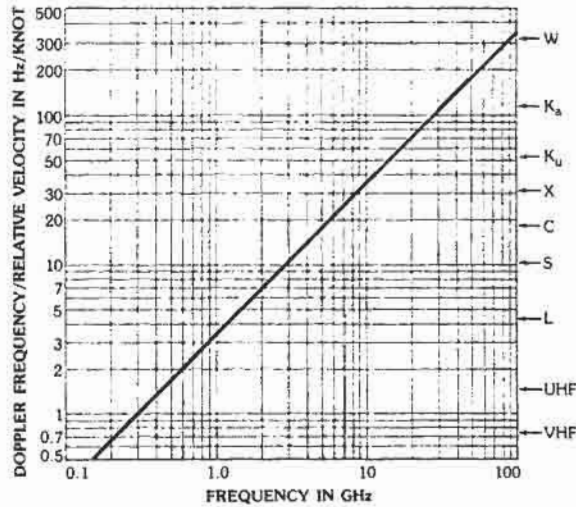


Fig. 9. Doppler frequency shift per unit relative velocity (hertz per knot) as a function of radar frequency.

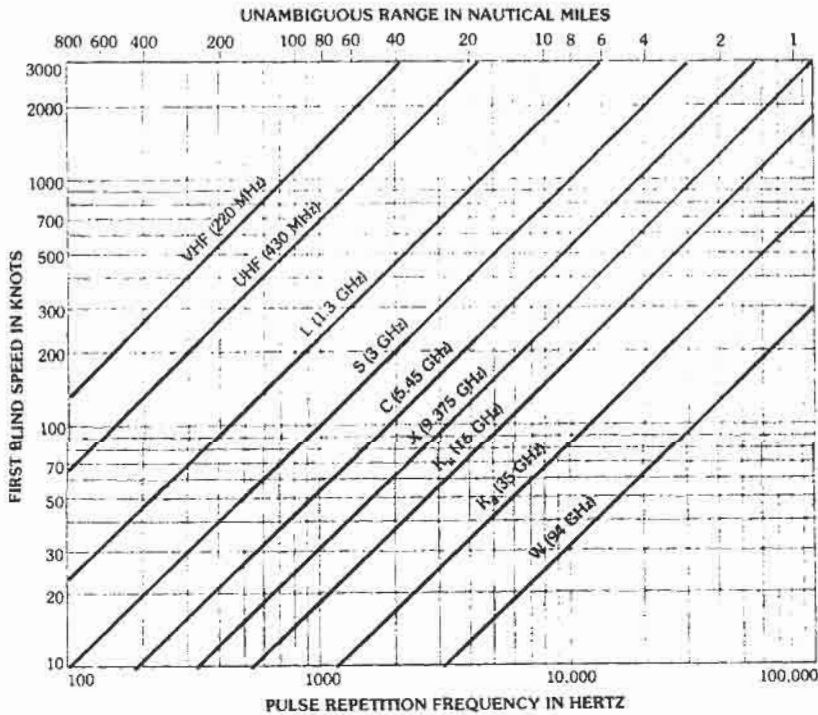


Fig. 10. First blind speed in an MTI or pulse doppler radar as a function of the pulse repetition frequency and the radar frequency.

resolved and the true range is found by transmitting three separate waveforms, each at a different prf. (In theory, only two different prfs are required to resolve range ambiguities, but in practice at least three are needed.) An advantage of a high-prf pulse doppler radar is that it can readily detect targets with a high relative-velocity since such echoes do not compete with the echoes from clutter that are at lower doppler frequencies. It has much poorer performance, however, against low relative-velocity targets which compete with the large clutter echoes seen by this type of radar in its antenna sidelobes because of the high prf.

3. A *medium prf pulse doppler radar* has both range and doppler ambiguities. It will not detect high relative-velocity targets as well as can a high-prf pulse doppler radar, but it will detect low relative-velocity targets better because its low prf sees less clutter than does the high prf.

Most ground-based air-surveillance radars that must see aircraft in the midst of clutter are generally MTI radars. When flown in an aircraft for purposes of air-surveillance, these are called AMTI, or *airborne moving target indicator radar*. AMTI radars are very good for airborne air-surveillance radars at UHF, but result in too many blind speeds when employed at the higher microwave frequencies. The high-prf and the medium-prf pulse doppler radars can both be used for airborne air-surveillance purposes at the higher frequencies. (The S-band AWACS, or AN/APY-1, airborne air-surveillance radar utilizes a high-prf pulse doppler waveform.) An X-band fighter/attack radar in a modern military aircraft might employ on a time-shared basis three widely differing prfs, depending on the operational situation. When no clutter is present (as when the antenna beam is looking up above the surface of the earth), a low prf waveform without any doppler processing might be employed. (The low prf waveform will provide the greatest range, if no clutter is present.) When searching for a target in clutter, the high prf and the medium-prf pulse doppler waveforms might be interleaved. The high prf would be looking for high-speed approaching targets at long range and the medium prf would be looking for slower-speed targets at shorter ranges.

The doppler frequency shift is important for many purposes in radar. In addition to the MTI, AMTI, and pulse doppler radars, the doppler frequency is the key to such radars as

- *Doppler weather radars*
Nexrad, whose output is regularly shown on TV weather broadcasts.
Terminal Doppler Weather Radar (TDWR). These are located in the vicinity of major airports to warn aircraft, that are landing or taking-off, of the presence of dangerous wind shear.
Wind profiler, that measures as a function of altitude the wind speed and direction for both weather prediction and the efficient routing of aircraft.

Airborne weather-avoidance doppler radar, that detects dangerous wind shear to warn the pilot of danger during take-off or landing.

- *Synthetic aperture radar (SAR)*, (for mapping a scene on the surface of the earth, and the *inverse synthetic aperture radar (ISAR)*, for imaging a target well enough to recognize it from other similar targets. (SAR is more usually thought of as an antenna synthesized in a digital processor, but it was originally invented as a doppler radar and can be described as depending on the doppler effect.)
- *Doppler navigator*, a multi-beam radar which can provide the vector velocity of an aircraft which carries it.
- *IIF over-the-horizon radar*, which depends on the doppler frequency to detect aircraft and ships in the presence of large clutter echoes. The doppler spectrum of the sea echo obtained by an HF radar can also be used to extract the strength and direction of the winds over the ocean.

The continuous wave (CW) radar also depends on the doppler frequency shift to detect targets in clutter and/or to measure relative velocity.

MTI RADAR

The MTI radar uses a digital implementation of a time-domain filter composed of "delay lines" to pass desired moving targets, but to reject stationary clutter. The simplest type of MTI time-domain filter uses a single delay line to subtract the received radar echoes from the echoes received one pulse repetition period previously. A single delay-line canceler has a frequency response given by

$$H_1(f) = 2 \sin \pi f_d T_p \quad (\text{Eq. 34})$$

where

$$\begin{aligned} f_d &= \text{doppler frequency,} \\ T_p &= \text{pulse repetition period.} \end{aligned}$$

When n single delay-line cancelers are used in cascade the frequency response is

$$H_n(f) = 2^n \sin^n \pi f_d T_p \quad (\text{Eq. 35})$$

An arrangement of n delay lines in cascade with $N = n + 1$ taps (one at the input of each delay line and one at the output of the last delay line), whose outputs are weighted by the coefficients of the binomial series with alternating sign before summing, produces a response equivalent to that of Eq. 35 for the n cascaded delay-line cancelers. This is sometimes called an N -pulse canceler and is an example of a *transversal filter* (also called a *nonrecursive filter*). The transversal filter may be used with weightings other than the binomial coefficients with alternating sign to give a tailored frequency response. Feedback can also be employed in a cascade of delay-line cancelers to shape the frequency response further. This is called a *recursive filter*. Its "infinite" duration transient response, however, can

limit its usefulness, especially for military radars subject to hostile electronic countermeasures.

Many air-surveillance radars used for air-traffic control employ a staggered prf rather than a constant prf. A staggered or multiple prf MTI radar is one that uses more than one pulse repetition period. Four or five stagger periods might be used. For example, the periods of an *L*-band air surveillance radar might be in the ratio of 28:30:27:31 (Reference 15). A staggered prf is used to increase the first blind speed beyond that given by Eq. 32. If the individual periods of the multiple PRF waveform are designated by T_1, T_2, \dots, T_N , and if they are related by the ratios $n_1/T_1 = n_2/T_2 = \dots = n_N/T_N$, where n_1, n_2, \dots, n_N are integers, then the first blind speed when N different periods are used is

$$v_1 = [(n_1 + n_2 + \dots + n_N)/N]v_B \quad (\text{Eq. 36})$$

where v_B is the first blind speed as given by Eq. 32 for a nonstaggered waveform with a constant period equal to the average of the N periods, or $(T_1 + T_2 + \dots + T_N)/N$. In the above example with four staggered periods, $n_1 = 25, n_2 = 30, n_3 = 27, \text{ and } n_4 = 31$, the first blind speed of the staggered waveform is 28.25 times that of a waveform with constant pulse repetition period. The staggered prf can also be used with a transversal filter of unequal spacings and with amplitude weights to shape the frequency response of the doppler filter.

Still another method for providing doppler filtering in an MTI or pulse doppler radar is to generate a bank of narrow-band doppler filters. These can be obtained using a transversal filter with phase weights rather than amplitude weights. N pulses processed as a filter bank will produce N contiguous narrow-band filters. The filter bank has the advantage over other methods in (1) separating desired moving targets from undesired moving clutter (such as rain), (2) providing a measure of the target's radial velocity, and (3) improving the signal-to-noise ratio compared to other filters. Of all these advantages, the first (separating aircraft echoes from rain) has been the most important.

In the above, several methods for providing doppler filtering have been briefly mentioned. Each has its advantages and limitations, and no one method is superior in all cases to the rest. Thus "engineering judgement" is needed when selecting the type of filtering to be used in any doppler radar application.

MTI Performance Limitations

The performance measure used to describe the ability of an MTI to see targets in clutter is the *MTI Improvement Factor*, defined as the signal-to-clutter ratio at the output of the MTI processor divided by the signal-to-clutter ratio at the input, averaged uniformly over all target relative velocities of interest. *Clutter Attenuation* is the ratio of the clutter power at the canceler input to the clutter residue at the output, normalized to the attenuation of a single pulse passing through the unprocessed channel of the canceler. *Sub-clutter Visibility* is the ratio by which the target echo

power may be weaker than the clutter echo power and still be detected with specified detection and false alarm probabilities. All target relative velocities are assumed equally likely.

The performance of an MTI radar is limited by clutter fluctuations, antenna scanning, and equipment instabilities.

Clutter Fluctuations—The improvement factor of an N -pulse canceler with $n = N - 1$ delay lines, or n cascaded delay-line cancelers, is

$$I_N = (2^n/n!) (f_p/2\pi\sigma_c)^{2n} \quad (\text{Eq. 37})$$

where

- f_p = pulse repetition frequency,
- σ_c = standard deviation of the clutter spectrum which is assumed to be of Gaussian shape.

This is plotted in Fig. 11. The standard deviation of the clutter spectrum in hertz is related to the standard deviation of the velocity spread in meters/second by the expression $\sigma_c = 2\sigma_v/\lambda$. "Typical" values of σ_v , are (Reference 14):

- Heavily wooded hills, 20 mph wind: 0.2 m/s
- Sparsely wooded hills, calm day: 0.02 m/s
- Sea echo, windy day: 0.9 m/s
- Rain clouds: 2 m/s
- Chaff: 1 m/s

The Gaussian model (and Eq. 37) does not apply for heavy clutter. The exponential clutter model needs to be used when improvement factors are greater than 20 or 30 dB, as has been described by Billingsley (Reference 16).

Antenna Scanning—The limitation to the improvement factor for n delay-line cancelers in cascade is

$$I_n = (2^n/n!) \times (0.6n_b)^{2n} \quad (\text{Eq. 38})$$

where n_b = number of pulses received within the half-power beamwidth as the antenna scans by the target.

Equipment Instabilities—The improvement factor in a doppler radar can be limited by the instabilities in such things as the transmitter frequency, stalo and coho oscillator phase or frequency, phase shift in the transmitter, and changes in the pulse width and pulse amplitude. Simple expressions can be found in the literature for the limitation on the improvement factor due to these equipment instability effects when a single or a double delay-line canceler is used; but most modern radars employ more sophisticated doppler filters and thus require more sophisticated analysis of the effect of equipment instabilities.

Digital Quantization—The limit on the improvement factor in a digital MTI due to the finite size of the quantization interval of the A/D converter is approximately

$$I_q \approx 6 \times \text{number of bits} \quad (\text{Eq. 39})$$

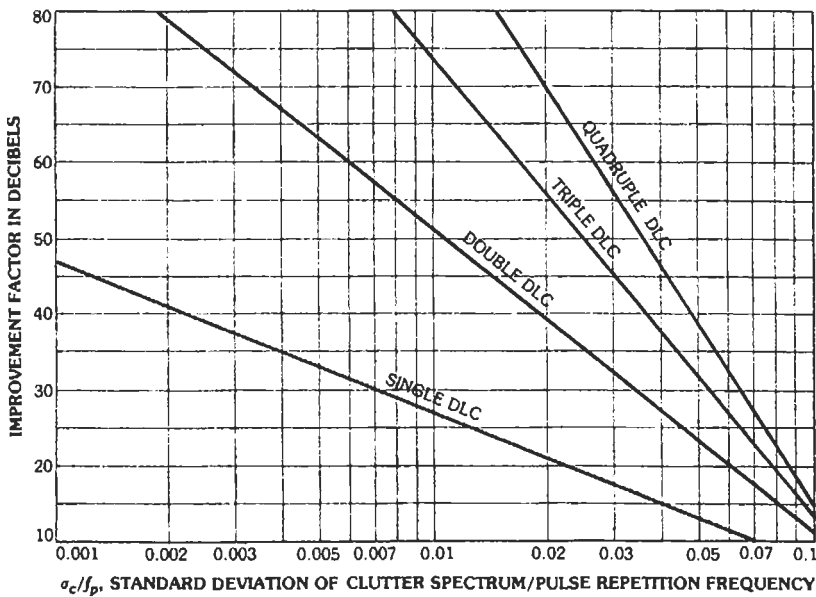


Fig. 11. The improvement factor as a function of σ_c/f_p (ratio of the standard deviation of the clutter and the pulse repetition frequency) for single and cascaded delay-line cancelers (DLC) as given by Eq. 37.

where I_q is in decibels. For example, an eight-bit A/D converter limits the improvement factor to 48 dB. In practice one or two additional bits must be added to account for receiver noise in the lower quantization levels and the limitations introduced by a limited dynamic range. In the previous example, a 10 bit A/D converter would likely be used if the limitation on the improvement factor had to be 48 dB.

PULSE COMPRESSION

Pulse compression allows a radar to utilize a long pulse to achieve a large radiated energy, but with the range resolution of a short pulse of wide bandwidth. It achieves this by modulating the long pulse of width T to obtain a bandwidth $B \gg 1/T$. The received signal is passed through a matched filter to produce a compressed pulse of width $1/B$. The pulse compression ratio, which is the duration of the long pulse divided by the duration of the short (compressed) pulse, is equal to BT . Frequency and phase modulations have both been used for pulse compression. Amplitude modulation could also be employed, in principle, but it is seldom found in practical pulse compression systems.

Table 4 gives the major types of pulse compression waveforms that have been employed in radar.

Doppler tolerant waveforms, such as the linear-FM modulation and the polyphase code, can be processed using a single matched filter even though the received signal experiences a wide range of doppler frequency shifts. *Complementary codes* consist of pairs of equal-

length codes which have the property that the time sidelobes of one code are the negative of the other so that they can be canceled to result in no sidelobes, theoretically. (This is generally difficult to obtain in practice.) The *pulse burst* is a series of pulses transmitted as a group before any echo signals are received (as when a long minimum range can be tolerated). Further information about pulse compression may be found in Chapter 6 of Reference 1.

CLUTTER

Clutter echoes are those received from undesired scatterers such as land, sea, rain, snow, chaff, clouds, birds, insects, aurora, and meteors. Good radar practice requires that these clutter echoes be eliminated or minimized so as to prevent degraded detection of desired target echoes. It is especially important to limit the amount of clutter echoes presented to automatic detection and tracking (ADT) systems so that overloading of the tracking computer does not occur by trying to establish tracks based on the echoes from clutter rather than from targets.

Surface clutter is described by the dimensionless parameter σ^0 , defined as the clutter cross section per unit area illuminated. Fig. 12 is a plot of the mean value of σ^0 for sea clutter as a function of the grazing angle. It applies for a medium sea (10–20 knot winds). At the lower grazing angles, sea clutter increases slowly with increasing winds above 15 to 20 knots (a few tenths of a decibel per knot). Below about 5 knots

TABLE 4. PULSE COMPRESSION WAVEFORMS

Type	Description	Pulse Compression Ratio	Peak (Time) Sidelobe	Comments
Linear FM (chirp)	Linear modulation of the frequency over the range Δf in a time T .	$\Delta f \times T$	-13.2 dB; but with weighting in the frequency domain on receive, sidelobes can be -30 dB with loss in SNR of about 1 dB.	Widely used. Especially applicable to very high resolution. Doppler tolerant.
Binary phase code; also known as linear recursive sequences, pseudorandom sequences, P-N sequences, and maximal length sequences.	A pulse of width T is divided into N subpulses of width τ , where $N = T/\tau$. The phase of each subpulse is either 0 or π radians, chosen at random or pseudorandom.	N , the number of subpulses.	Approximately $0.5/N$ (in the best of cases)	Sidelobes are not low, especially with a doppler shift.
Barker code	A form of binary phase code that has all its time sidelobes equal to $1/N^2$ that of the compressed pulse peak power.	N	$1/N^2$	Limited to $N \leq 13$. Longer pulse compression ratios can be had with compound Barker codes in which each segment of a Barker code is modulated by a Barker code, but no decrease in sidelobes is obtained.
Frank polyphase code with variants due to Lewis and Kretschmer known as P codes. Phase quantization less than π radians.	A pulse of width T is divided into M^2 equal subpulses. The phase of each subpulse is chosen so that a linear FM waveform is approximated with a phase quantization equal to $2\pi/M$ radians.	M^2	$1/\pi^2 M^2$	Doppler tolerant for aircraft velocities. Lower sidelobes than other pulse compression methods.
Costas codes; also known as discrete frequency shift codes and time-frequency codes	A pulse of width T is divided into M subpulses, and the carrier frequency is changed subpulse to subpulse, with frequency steps separated by the reciprocal of the subpulse width. Frequencies are selected in a definite order from the M frequency steps over a band $\Delta f = M/\tau$, where $\tau = T/M$.	$\Delta f \times T = M^2$	$1/M^2$ in regions away from the central peak response; $(2/M)^2$ near the central peak response.	Produces a thumbtack ambiguity diagram. For a given value of M , there can be a large number of different Costas codes.
Nonlinear FM	Frequency is varied nonlinearly over a band Δf in a time T .	$\Delta f \times T$	Depends on the nonlinear waveform used.	Low sidelobes can be obtained with matched filter processing without a loss in SNR. Can obtain a thumbtack ambiguity diagram.
Stretch	Transmits a linear FM waveform of bandwidth Δf over a time T , but processes the received signal with circuitry of bandwidth $\Delta f/a$ (where $a > 1$). The region where pulse compression occurs is over a restricted range interval $1/a$ that of conventional FM pulse compression.	$\Delta f \times T$	As with linear FM pulse compression.	The time waveform is "stretched" by a factor of a so as to achieve a reduction of processing bandwidth by $1/a$. Used for target recognition applications.

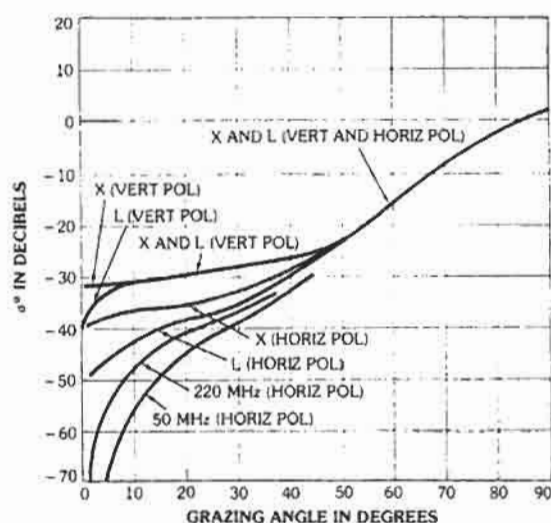


Fig. 12. Composite of σ^0 data for a "medium" sea, as a function of the grazing angle.

wind velocity, sea clutter decreases rapidly with decreasing wind.

Land clutter is difficult to describe because of the many factors that affect it. Fig. 13 is a plot of the mean value and variance of rural clutter measured at low grazing angles (Reference 16). It is seen that the mean is relatively independent of frequency, but the variance decreases with increasing frequency. The following are a few of the many findings from the extensive measurements reported by Billingsley in Reference 16 for land clutter when viewed at low grazing angle:

- Most of the significant clutter echoes come from spatially localized or discrete vertical features;

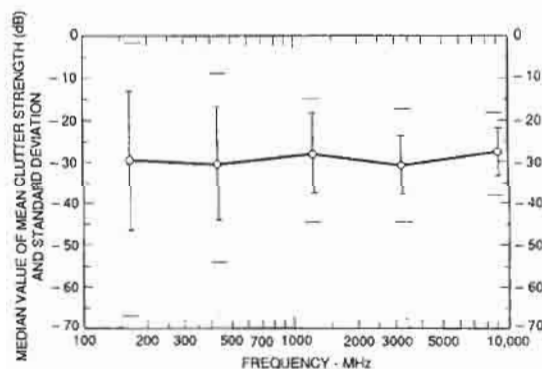


Fig. 13. Ground clutter of rural terrain at low gazing angles as a function of frequency, showing mean values of clutter strength (open circles), standard deviation (vertical bars), and the extreme measured values (horizontal bars). (After J. Billingsley, Reference 16.)

such as trees, fences, buildings, and high regions of terrain.

- The effect of vertical discrete objects on the overall clutter strength is large even when these objects are relatively sparse.
- The variation of the average clutter echo with weather, season, and from day-to-day is small.
- There is little difference in the clutter echo between vertical and horizontal polarization.
- Measured amplitude distributions "almost never pass rigorous statistical hypothesis tests for belonging to Weibull, lognormal, or other theoretical distributions that have been tried."
- The chief characteristic of clutter is its variability.

The radar echo from volume clutter, such as rain, is described by the volume reflectivity, η , defined as the radar cross section per unit volume. The volume reflectivity of rain in the microwave region is

$$\eta = 7f^4 r^{1.6} \times 10^{-12} \quad (\text{Eq. 40})$$

where

- η is the volume reflectivity in meters²/meter³,
- f is the radar frequency in gigahertz,
- r is the rainfall rate in millimeters/hour.

This expression assumes that the backscatter from rain is described by Rayleigh scattering, which is applicable when the radar wavelength is large compared to the circumference of the raindrop.

The approximate number of hours per year that various rainfall rates are exceeded in Washington, DC are as follows:

Rainfall Rate, mm/hr	Number of hrs/yr
0.25 (drizzle)	500
1 (light rain)	230
4 (moderate rain)	65
16 (heavy rain)	10
40 (excessive rain)	2

An alternate description of rainfall rate, instead of millimeters/hour, is the reflectivity factor $Z = 200r^{1.6}$, where r is the rainfall rate in mm/hr and Z is in mm⁶/m³. This measure is often used by radar meteorologists. It is usually expressed in decibels and abbreviated dBZ. For example, a rainfall rate of 4 mm/hr is equivalent to 32.6 dBZ.

The following brief statements review the general nature of radar clutter:

- Clutter can be more harmful than receiver noise in limiting the ability of a radar to detect targets.
- It is far more difficult to characterize clutter than to characterize receiver noise.
- Land clutter is a more serious limitation than sea clutter.
- The available information on land clutter at low grazing angles is quite good. With the possible exception of vertical incidence, the description of land clutter at higher grazing angles is not as

good. Theoretical models of clutter echo are poor and remote sensing results over land might be misleading.

- The theory of radar sea clutter at HF and VHF is readily described by Bragg scatter. On the other hand, Bragg scatter does not apply at the higher microwave frequencies. “Sea spikes” are the dominant clutter mechanism at these frequencies.
- Rain clutter can seriously affect the performance of radars at L band and higher. Unless doppler processing is used, S band and higher frequency radars are seriously degraded (clobbered) in rain.
- Snow generally is not as serious as rain in reducing radar performance.

RADAR MEASUREMENT ACCURACY

The rms error, δM , in making a measurement of a radar parameter M in the presence of Gaussian noise and in the absence of any bias can be written conceptually in the form

$$\delta M = kM/(2E/N_0)^{1/2} \tag{Eq. 41}$$

where

- E is the received signal energy,
- N_0 is the noise power per unit bandwidth,
- k is a constant of the order of unity.

(Note that $2E/N_0$ is also the peak-signal-to-mean-noise ratio at the output of a matched filter.) For a range measurement (time delay), k depends on the shape of the signal spectrum and M can be taken as the rise time of the pulse. For a doppler frequency measurement, k depends on the shape of the time waveform and M is the spectral resolution (or the reciprocal of the observation time). For an angle measurement, k depends on the shape of the antenna aperture illumination and M is the beamwidth.

RMS Range Measurement Error

Rectangular pulses are seldom truly rectangular in that they have a finite rise and a finite fall time. A quasi-rectangular pulse (Reference 1) is more realistic than a perfectly rectangular pulse since it is a pulse which has passed through a filter so that it has a rounded shape rather than have unrealistic rise and fall times with infinite slopes. The theoretical rms range error for such a pulse when the noise is Gaussian is

$$\delta R = (c\tau/2)/[1.3(2E/N_0)^{1/2}] \tag{Eq. 42}$$

where

- c = velocity of propagation
- τ = half-power time duration.

RMS Doppler Frequency Measurement Error

Rectangular pulse of width τ .

$$\delta f = \sqrt{3}/[\pi\tau(2E/N_0)^{1/2}] \tag{Eq. 43}$$

Quasi-rectangular pulse of half-power width τ .

$$\delta f = 1/[2.6\tau(2E/N_0)^{1/2}] \tag{Eq. 44}$$

RMS Angle Measurement Error

Uniform aperture illumination:

$$\delta\theta = 0.628\theta_B/(2E/N_0)^{1/2} \tag{Eq. 45}$$

where θ_B = half-power beamwidth.

Cosine aperture illumination:

$$\delta\theta = 0.73\theta_B/(2E/N_0)^{1/2} \tag{Eq. 46}$$

TRACKING WITH RADAR

There are at least four different radar methods for obtaining the track of a target, as described below.

Single-Target Tracker (STT): The single-target tracker continually engages one target with a high data rate. Ten target observations per second is a typical data rate for military tracking radars. Tracking in range and angle is performed with a closed-loop servo system so as to keep the radar positioned on a moving target. *Monopulse* trackers are the most accurate tracker and have a high degree of resistance to deliberate countermeasures, which is why it has been the preferred tracking method for military air defense systems based on the STT. In a monopulse radar, the angle measurement in one coordinate is determined by using two beams slightly displaced (squinted) in angle. Transmission occurs on the sum of the two squinted beams and reception on both the sum and the difference of the two beams. The angle measurement is obtained with the difference pattern. The sum pattern provides detection and the measurement of range. The sum pattern signal also acts as a reference signal to determine in which direction is the angle error. To provide angle measurements in two angle coordinates, four squinted antenna beams are required. A monopulse tracker has three receiving channels: a range channel and two angle channels.

The *conical scan* STT is much simpler than a monopulse tracker in that it needs only one receiving channel to extract the angle measurement in two angle coordinates. A single offset (squinted) antenna beam is rotated typically at a rate of about 30 Hz. The conical scan tracker time-shares a single beam to sequentially obtain the angle measurements in two coordinates; whereas, a monopulse radar obtains its angle measurements simultaneously. Conical scan radars are used where accuracy and vulnerability to electronic counter-

measures are less important than a simpler, lower cost system.

Automatic Detection and Track (ADT): A conventional rotating air-surveillance radar such as employed for air-traffic control obtains an azimuth angle measurement on each rotation, or scan, of the antenna. Tracks can be established based on this lower scan rate (revisit times from 4 to 12 s). Automatic tracking is combined with automatic detection and is called automatic detection and track. The data rate is considerably slower than that of a single-target tracker, but a single ADT system might be able to hold in track many hundreds or thousands of aircraft targets. Tracking is done open loop in that the position of the antenna is not controlled by the processed tracking data as it is in a closed-loop STT.

Phased Array Radar Tracking: Multiple targets can be tracked on a time-shared basis with an agile electronically steered phased array radar. Tracking is not closed loop, but it combines the high data rate of a STT with the ability of an ADT to hold a large number of targets in track.

Track While Scan (TWS): This term has been applied to both the ADT and to a radar which rapidly scans a limited angular sector to maintain tracks with a moderate data rate on a modest number of targets. TWS is now almost always applied to the latter meaning rather than to an ADT system. Angle information can be obtained in a TWS radar by using two scanning fan beams perpendicular to each other, one for the azimuth angle measurement and the other for the elevation angle measurement. A TWS radar can also be obtained with a single cluster of four monopulse beams scanning in a raster fashion to provide the required volume coverage.

REFERENCES

1. Skolnik, M. I. *Introduction to Radar Systems*, Chapter 2, 3rd ed. New York: McGraw-Hill Book Co., 2001.
2. Swerling, P. "Probability of Detection for Fluctuating Targets." *IRE Trans. IT-6* (April 1960), pp. 269–308.
3. Marcum, J. I. "A Statistical Theory of Target Detection by Pulsed Radar, Mathematical Appendix." *IRE Trans. IT-6* (April 1960), pp. 145–267.
4. Meyer, D. P., and Mayer, H. A. *Radar Target Detection*. New York: Academic Press, Inc., 1973.
5. Blake, L. V. "Prediction of Radar Range." Chapter 2 of *Radar Handbook*, 2nd ed. Edited by M. I. Skolnik. New York: McGraw-Hill Book Co., 1990.
6. Blake, L. V. *Radar Range-Performance Analyses*. Lexington, MA: Lexington Books, 1980.
7. Barton, D. K., and Barton, W. F. *Modern Radar System Analysis Software and User's Manual: Version 2.0*. Norwood, MA: Artech House, 1993.
8. Knott, E. F. "Radar Cross Section." Chapter 11 of *Radar Handbook*, 2nd ed. Edited by M. I. Skolnik. New York: McGraw-Hill Book Co., 1990.
9. Hansen, R. C. "Aperture Theory." Chapter 1 of *Microwave Scanning Antennas*, Vol. 1. Edited by R. Hansen. New York: Academic Press, Inc. 1964.
10. White, W. D. *Desirable Illuminations for Circular Aperture Arrays*. Research Paper P-351, December 1967. Arlington, VA: Institute for Defense Analyses. (Approved for public release, but not generally available.)
11. Stutzman, W. L. "Estimating Directivity and Gain of Antennas," *IEEE Antennas and Propagation Magazine* 40 (August 1998), pp. 7–11.
12. Ruze, J. "Antenna Tolerance Theory—A Review." *Proc. IEEE* 54 (April 1966). pp. 633–640.
13. Cheston, T. C., and Frank, J. "Array Antennas." Chapter 7 of *Radar Handbook*, 2nd ed. Edited by M. I. Skolnik. New York: McGraw-Hill Co., 1990.
14. Barlow, E. J. "Doppler Radar." *Proc. IRE*. 37 (April 1949), pp. 340–355.
15. Shrader, W. W., and Gregers-Hansen, V. "MTI Radar." Chapter 15 of *Radar Handbook*, 2nd ed. Edited by M. I. Skolnik. New York: McGraw-Hill Book Co., 1990.
16. Billingsley, J. B. *Low Angle Radar Land Clutter*. Norwich, NY: William Andrew Publishing, 2001.
17. Nathanson, F. E. *Radar Design Principles*. 2nd ed. New York: McGraw-Hill Book Co., 1991.

37

Radio Navigation Systems

*Revised by
D. G. Jablonski*

Introduction	37-2
Major Navigation Agencies	
Propagation	37-2
Issues	
Pulse Transmission	
Space Diversity	
Frequency Diversity	
Radio Navigation Systems	37-3
Omega	
Loran-C	
Radio Beacons	
VOR (VHF Omnidirectional Range)	
ILS (Instrument Landing System)	
DME (Distance Measuring Equipment)	
TACAN (Tactical Air Navigation)	
MLS (Microwave Landing System)	
ATCRBS (Air Traffic Control Radar Beacon System)	
Transit	
TCAS (Traffic Alert and Collision Avoidance Systems)	
NAVSTAR/Global Positioning System	
Differential GPS	
Future Trends	

INTRODUCTION

Electromagnetic waves at radio and microwave frequencies can be used for the precise measurement of time, distance, and direction. Several factors make this possible. First, one can generate extremely accurate timing signals using a variety of techniques, including atomic clocks. In addition, radio and microwave signals can be generated coherently, so that their frequency and phase are well controlled and easily measured. Because of the low ambient thermal noise at these frequencies, radio signals can be detected at extraordinarily low power levels (~ -160 dBW using time integration techniques). In addition, the diverse propagation characteristics of radio waves, as a function of frequency, permit engineers to exploit over-the-horizon transmission at low frequencies, and sight-limited propagation at high frequencies. The high resolution afforded by the short signal wavelengths at higher microwave frequencies, the effects of constructive interference at all frequencies, and the always useful Doppler effect are also utilized to good advantage in modern radio navigation systems.

However, in recent years, there has been considerable change in the radio navigation systems deployed throughout the world. In particular, several time-honored systems, including Omega and the Transit satellite system, are no longer operational. Efforts to develop the Microwave Landing System (MLS), which was intended to replace the existing Instrument Landing System (ILS) used by aircraft, have been reduced greatly, at least in the United States. In their stead, the Global Positioning System (GPS) has been deployed and declared operational, and is now the system of choice for worldwide radio navigation. The Russian counterpart to GPS, the GLONASS system, is also operational, and an additional worldwide satellite-based radio navigation system, called Galileo, is under development in Europe. In addition, various military communications systems, such as the Joint Tactical Information Distribution System, or JTIDS, have an incipient radio navigation capability for their users.

In addition to the new systems, the venerable VOR, TACAN, DME, and Air Traffic Control Radar Beacon System (ATCRBS) remain stalwarts of the commercial and military aircraft navigation infrastructure.

Major Navigation Agencies

Airlines Electronic Engineering Committee (AEEC), Annapolis Science Center, Annapolis, Maryland: A division of Aeronautical Radio, Inc. (ARINC) and owned by the scheduled US airlines. Publishes technical standards for avionics purchased by the scheduled airlines.

Airline Owners and Pilots Association (AOPA), Washington, DC: Defends the needs of the airline industry and pilots with regard to safety-of-life issues,

frequency allocation that affects radio navigation systems, etc.

Department of Transportation, United States Coast Guard (USCG), Washington, DC: Operates the Loran-C navigation system for marine and aeronautical navigation; operates and maintains a beacon system for differential GPS in the coastal regions of the United States and along the Mississippi River.

Federal Aviation Administration (FAA), Washington, DC: Operates navigation aids and air traffic control systems for both civil and military aircraft in the US and its possessions.

Federal Communications Commission (FCC), Washington, DC: The agency that licenses transmitters and operators in the United States and aboard US registered ships and aircraft.

International Air Transport Association (IATA), Montreal, Canada: The international association representing scheduled airlines.

International Civil Aviation Organization (ICAO), Montreal, Canada: A United Nations agency that formulates standards and recommended practices, including navigation aids, for all civil aviation.

International Telecommunication Union (ITU), Geneva, Switzerland: An agency of the United Nations that allocates frequencies for best use of the radio spectrum.

RTCA, formerly Radio Technical Commission for Aeronautics, Washington, DC: Supported by contributions from industry and government agencies. Participation by manufacturers, users, and others in the recommended standards for aviation electronics. The ICAO and the FAA adopt many of these standards, at least in part.

Radio Technical Commission for Maritime Services (RTCM), Washington, DC: Functions similar to those of RTCA; however, addresses primarily marine issues.

PROPAGATION

Issues

A detailed understanding of the propagation characteristics of radio waves is essential to the implementation of radio navigation systems. In low-frequency navigation systems, complex propagation models may be used to correct for nonhomogeneous propagation paths. In addition, in most radio navigation systems, errors may result from the contamination of the signals with multipath signals that have traveled by a non-direct path that is often variable. This is of increasing importance as the operational bandwidth of systems increases. This is due to a phenomenon known as frequency-selective fade, in which a portion of the signal spectrum suffers from frequency-specific multipath effects. To reduce multipath effects, the following techniques are commonly used in navigation systems.

Pulse Transmission

Through the use of an appropriate pulse length and repetition rate combined with a means in the receiver to recognize the leading edges of the pulses, the desired direct signal may be separated from a signal that has traveled a longer path (sky wave for low-frequency systems and multipath for higher-frequency systems). Pulse transmission techniques are effectively used in Loran-C, DME, with ATRCBS transponders, and radar systems.

Space Diversity

The larger the aperture of the antenna system, the greater is the statistical probability that the desired signals will add linearly, while the multipath signals add randomly. This approach is effectively used in Doppler VOR systems. Antenna directivity and polarization are frequently used to reduce interference from undesired multipath signals. Horizontal polarization is used in the ILS system, and vertical polarization is used in VOR and DME systems.

Frequency Diversity

While the line-of-sight path remains the same at all radio frequencies, indirect paths may vary with frequency. In such cases, spectrum-spreading techniques, such as Code Division Multiple Access (CDMA), may achieve the same result as space diversity.

RADIO NAVIGATION SYSTEMS

The following subsections provide information on radio navigation aids ranging from the now obsolete low-frequency Omega system through the current satellite system, the Navstar Global Positioning System (GPS).

Omega

The Omega system was a worldwide VLF navigation system used for marine and enroute air navigation. In Omega, the system CW transmitting stations sequentially transmit long, but precisely timed, pulses at four frequencies: 10.2 kHz, 11.3 kHz, 13.6 kHz, and 11.05 kHz. Position information is obtained by measuring the relative phase differences of the received signals. Since lines of constant phase for the signal from a particular transmitter are hyperbolic in shape, OMEGA is known as a hyperbolic system. The inherent accuracy of the Omega system is limited by propagation corrections that must be applied to the receiver. These corrections vary depending on location and the time of day. In many cases, accuracies of 2 nautical miles (rms) day and 4 nautical miles night could be

achieved in most of the coverage areas. Note that worldwide coverage is achieved by the choice of VLF for the operating band. Only three transmitter stations are needed to achieve worldwide coverage. The choice of VLF results in the complex propagation issues reflected in the dependence of system accuracy on time-of-day.

Loran-C

Loran-C is a long-range (as opposed to worldwide) hyperbolic radio navigation system that possesses an inherent high degree of accuracy at ranges of 800 to 1000 nautical miles. The Loran-C system transmits synchronized, phase-coded pulses from a master station and two or more secondary stations at 100 kHz. Note the difference in operating frequency compared to the VLF OMEGA system. The Loran-C transmitting stations form a chain characterized by the group repetition interval (GRI) in which the pulses are repeated. A GRI starts with the master station transmitting eight pulses, each spaced one millisecond apart, followed by a ninth pulse two milliseconds later. The master station transmission is followed after a prescribed coding delay by transmissions from each of the secondary stations in the chain, each transmitting eight pulses at one-millisecond intervals. Phase coding is used to differentiate the master pulses from those of the secondaries. The pulse spacing and phase code allow the ground wave (direct propagation path) to be differentiated from the varying sky wave.

A typical Loran-C receiver makes use of a microprocessor for signal processing, navigation computation, and control. Through the use of microprocessors and other signal-processing and timing integrated circuits, the cost of high-quality navigation has continued to decrease. Although there have been plans to deactivate Loran-C, it is still in use in the United States and elsewhere. The accuracy of the system is one quarter of a mile absolute and less than 100 feet for measurement of relative position.

Radio Beacons

Radio beacons are nondirectional transmitters that operate in the low-frequency and medium-frequency bands. A radio direction finder is used to measure the relative bearing to the transmitter with respect to the heading of an aircraft or marine vessel. The aeronautical nondirectional beacons (NDBs) operate in the 190–415 kHz and 510–535 kHz bands. These beacons transmit either a coded or modulated CW signal for station identification. The coded signal is generated by modulating the carrier, and the upper carrier is keyed for the Morse-code identification. The aeronautical-beacon accuracy is in the range ± 3 degrees to ± 10 degrees, and the marine-systems accuracy is maintained to within ± 3 degrees. Aeronautical NDBs are used to

supplement the combined VOR/DME system for transition from enroute to precision approach facilities and as nonprecision approach aids. These aeronautical systems are considered one of the ICAO standard radio navigation aids.

VOR (VHF Omnidirectional Range)

The VOR transmits continuous-wave signals on one of the 20 assigned channels in the 108 to 118 MHz bands with 100-kHz channel separation. A nondirectional 30-Hz reference signal with a ± 480 -Hz frequency modulation on a 9960-Hz subcarrier is transmitted along with a carrier radiating from a rotating antenna with a horizontal cardioid pattern. The cardioid antenna pattern rotates at a 30-Hz rate, allowing the airborne receiver to determine its bearing from the station as a function of phase between the reference and the rotating signal. The VOR system has line-of-sight limitations in that at altitudes above 5000 feet the range is approximately 100 nautical miles and above 20 000 feet the range is approximately 200 nautical miles. The enroute VOR stations are rated at 50 watts. The accuracy of the VOR ground station is better than ± 1.4 . (However, the station magnetic declination is usually allowed to increase to 2 before being reset.) The total system error (rss of the ground and airborne system plus flight technical error) is less than ± 4.5 degrees.

Distance measuring equipment (DME) installations are often colocated with the VOR stations to provide ranging information. In the United States and other countries, TACAN (Tactical Air Navigation) installations are colocated with the VORs to provide a navigation system utilized primarily by the military. A combined VOR/DME/TACAN installation is known as a "VORTAC" station. They are typically located in cleared fields with good line-of-sight in all directions and are readily visible to passersby.

ILS (Instrument Landing System)

At present, the instrument landing system operating in the 108–112 MHz band is the primary worldwide, ICAO approved precision landing system. An alternate system, the microwave landing system (MLS) was scheduled to replace the ILS by the year 2000. However, in the United States, development of the Microwave Landing System was deferred in most cases in favor of the long-term adoption of GPS as a precision landing system. However, approval of GPS as the sole means of navigation for precision landing has not yet occurred, and ILS remains in widespread use as the system of choice.

An ILS normally consists of two or three marker beacons, a localizer, and a glide slope to provide verti-

cal and horizontal guidance information. The localizer, operating in the 108–112 MHz band, is normally located 1000 feet beyond the stop end of the runway. The glide slope is normally positioned 1000 feet after the approach end of the runway and operates in the 328.6–335.4 MHz band. Marker beacons operating along the extension of the runway centerline at 75 MHz are used to indicate decision height points for the approach or distance to the threshold of the runway.

Azimuth guidance provided by the localizer is accomplished by use of a 90-Hz-modulated left-hand antenna pattern and a 150-Hz-modulated right-hand pattern as viewed from the aircraft on an approach. A 90-Hz signal detected by the aircraft receiver will cause the course deviation indicator (CDI) to deviate to the right. A 150-Hz signal will drive the CDI vertical needle to the left when the aircraft is right of the centerline course. When the aircraft is on the centerline, the CDI vertical needle will be centered. A total of 40 channels are provided by the ILS localizer system, each being paired with a possible glide-slope channel.

Vertical guidance is provided by the glide-slope facility that is normally located to the side of the approach end of the runway. A total of 40 channels are provided in the 328.6–335.4 MHz band; each is paired with one of 40 ILS localizer channels. The carrier radiated in the antenna pattern below the glide slope is amplitude modulated with a 150-Hz signal. The pattern above the glide slope produces a signal with 90-Hz amplitude modulation. When the approaching aircraft is on the glide slope, the CDI horizontal (glide slope) needle will be centered.

The marker-beacon facilities along the course provide vertical fan markers to mark the key productions along the approach. The inner marker is normally at the runway threshold; the middle marker is about 3500 feet from the threshold; and the outer marker is usually 5 miles from the runway. A DME on one of the 20 paired channels with the localizer channels may also be used for indicating position during the approach.

The quality of the ILS installations varies depending on the equipment, terrain, and calibration. The ICAO has established the categories in Table 1 for minimum approach ceiling and forward visibility.

TABLE 1. CATEGORIES FOR MINIMUM APPROACH CEILING AND FORWARD VISIBILITY

Category	Minimum Ceiling (feet)	Forward Visibility (feet)
I	200	2600
II	100	1200
III-A	50	700
III-B	35	150
III-C	0	0

DME (Distance Measuring Equipment)

The airborne equipment (interrogator) generates a pulsed signal that is recognized by the ground equipment (transponder); the transponder then transmits a reply that is identified by the tracking circuit in the interrogator. The distance is computed by measuring the total round-trip time of interrogation, reply, and fixed delay introduced by the ground transponder. The airborne interrogator transmits about 30 pulse pairs per second on one of the 126 allocated channels between 1025 and 1150 MHz. The ground transponder replies on one of the paired channels in the 962–1024 MHz band or 1151–1213 MHz band. A DME and a collocated VOR constitute the ICAO standard ρ - θ system.

TACAN (Tactical Air Navigation)

The TACAN system provides both omnibearing and distance-measuring capability. The rotating directional horizontal-plane radiation pattern produces the azimuth signal, which contains a coarse (15 Hz) and a fine (135 Hz) azimuth element. The rotation of the pattern at 15 Hz results in a modulation of the carrier with a composite 15-Hz sine wave. Reference signals are transmitted by coded pulse trains to provide the phase reference. The airborne receiver obtains bearing by comparing the 15-Hz and 135-Hz sine waves with the reference pulse groups. The TACAN system operates in the 960–1215 MHz band with 1-MHz channel separations.

MLS (Microwave Landing System)

The microwave landing system is an ICAO-approved replacement for the current ILS system. However, the advent of GPS has placed many of the plans for deployment of MLS on hold. In the United States, there are no plans for widespread deployment of MLS. The MLS system is based on time-referenced scanning beams (TRSB), referenced to the runway, which enable the airborne unit to determine precise azimuth angle and elevation angle. The angular position of the aircraft is determined by measuring the time intervals between the TO and FROM azimuth antenna beam scan and the UP and DOWN scan of the elevation antenna pattern. The time interval represents a unique position within the range of the scanning beams. The azimuth scan is typically 60 degrees either side of the runway centerline, and the elevation scan is from 0 to 30 degrees. The signal format provides for 360-degree azimuth coverage for future implementation. 200 channels in the 5000–5250 MHz band provide the azimuth and elevation angle functions. Range information is provided by DMEs operating in the

960–1215 MHz band. An option is included in the signal format to permit a special-purpose system operating in the 15 400–15 700 MHz band.

ATCRBS (Air Traffic Control Radar Beacon System)

The ATCRBS ground interrogator transmits at 1030 MHz with a rotating (5 Hz enroute and 2.5 Hz terminal areas) directional antenna pattern. The interrogator transmits approximately 400 pulse pairs per second and receives replies from aircraft transponders that are within the beam of the antenna pattern. The airborne transponder replies at 1090 MHz with one of the 4096 pulse codes available. The decoded replies are displayed on the surveillance radar ppi along with primary radar returns. An omnidirectional pulsed pattern is also radiated from the ground to suppress unwanted sidelobe replies. This system is often referred to as secondary surveillance radar.

Transit

The Transit satellite navigation system consisted of a constellation of satellites in approximately 600-nautical-mile polar orbits. The system, now deactivated, was operated by the US Navy. The satellites broadcast ephemeris information continuously at 150 MHz and 400 MHz. A receiver measures successive Doppler shifts of the signal as the satellite approaches or passes the user. The geographic position of the receiver is then calculated from the satellite position information (transmitted from the satellite every two minutes) and the measurement of the Doppler shift. Normally only one frequency (400 MHz) is used for navigation, and the other is used when ionosphere correction (for improved accuracy) is performed. Coverage was worldwide but not continuous due to the relatively low altitude and the low number of satellites. The update of navigation (time between satellite fixes) was as short as 1 hour and up to 8 hours, depending on latitude. The Transit system has been replaced by the Global Positioning System (GPS).

TCAS (Traffic Alert and Collision Avoidance Systems)

The TCAS collision avoidance system permits aircraft to interrogate each other's transponder beacons. The response to an interrogation is processed to provide relative range information. Relative altitude information is derived from a comparison of the altimeter information coded into the transponder's reply and the interrogating aircraft's knowledge of its own altitude. When automatic computer computations indicate that a collision is probable, either a traffic advisory (TA) or a resolution advisory (RA) is provided to the pilot of the interrogating aircraft. The more sophisticated RA

is available only in the TCAS-II implementation. The RA directs the pilot to perform a specific maneuver to avoid a collision. TCAS-I is installed on small commercial aircraft (10–30 passengers) and TCAS-II on large transport aircraft.

NAVSTAR/Global Positioning System

The Global Positioning System (GPS) is a worldwide satellite navigation system developed by the US Department of Defense. Shown in notional form in Fig. 1, it is currently operational with 27 satellites. Operation of the system is based on accurate knowledge of the position of each satellite with respect to the time and distance from a satellite to the user. Each satellite transmits a unique ephemeris data table, which is periodically updated by the master control station. The user's position is determined relative to the satellites by processing signals received from at least four satellites. The time-of-arrival difference measurements are used to obtain the distance to each of the satellites. A time correction then relates the satellite system to earth coordinates.

The satellite signals are transmitted at two L-band frequencies (L_1 at 1575.42 MHz and L_2 at 1227.6 MHz) to permit corrections for ionospheric delays in propagation. The L_1 signal is modulated with a P and a C/A pseudorandom noise code in phase quadrature, in a manner similar to quadrature phase shift keying (QPSK). The L_2 signal broadcasts the P code, but will be upgraded to include the C/A code.

The P code is a long-period ($T = 38$ weeks) code operating at 10.23 Mb/s. The C/A code is a short code ($T = 1$ ms) readily acquired but operating at 1.023 Mb/s, which provides a less accurate measurement. The C/A

codes are different for each of the deployed satellites. Gold codes are used to minimize the cross-correlation between the signals from different satellites. This prevents a GPS receiver from inadvertently interpreting the signal from one satellite as being from another.

The same P-code sequence is used for each satellite. However, the 38 week-long code is broken into separate week-long offsets to enable the receiver to differentiate the signals from the various satellites. The P/Y code is reset weekly at midnight, Greenwich mean-time (GMT).

The P code designation stands for "precise," whereas C/A is an abbreviation for "Coarse/Acquisition." The latter designation is because of the common practice in P-code receivers of first acquiring the C/A code, then transitioning to P code operation. Some military receivers, however, have a "direct-P" acquisition capability.

Simultaneous tracking of the pseudo-random noise (PRN) codes from several satellites allows both identification of the satellites and calculation of the signal transit time from each satellite. The latter is accomplished by measurement of the phase shift required to match the codes from each satellite with versions of the same codes generated within the GPS receiver. This is accomplished using digital delay-lock loop technology. Note that this is a classic example of a direct-sequence spread spectrum system. Since each satellite transmits its pseudorandom noise sequence at the same carrier frequency, GPS is also an archetypal Code Division Multiple Access (CDMA) system.

In addition to permitting re-use of the same frequency by each satellite, the use of PRNs yields significant "processing gain." This permits GPS receivers to operate effectively with extremely small (~ -160 dBW) satellite signals. The despreading process used in the receiver also yields significant immunity to jamming and interference.

Navigation is accomplished by using the measured time-delays of the received satellite signals to infer the range from the GPS receiver to each satellite. With precise knowledge of the orbital location of the satellites, the receiver is able to compute its position.

In order to obtain the most accurate possible information about the location of each satellite, the L_1 and L_2 signals are also continuously modulated with a 50 b/s data stream. Known as the navigation message, this signal contains the "almanac" and "ephemeris" which together comprise a tabulation of the orbital parameters of the various satellites, as well as timing corrections for the atomic clocks onboard each satellite.

The operation of a typical GPS receiver is illustrated in Fig. 2. The received signal is detected using a series of nested tracking loops (which are located in the IFU portion of the system shown in the figure). Separate loops are used for each satellite. The pseudorandom noise sequences are first tracked in the code loop. Once this occurs, a second loop locks on to the carrier signal at the L_1 or L_2 frequency. In addition to providing improved range calculation, the carrier loop mea-

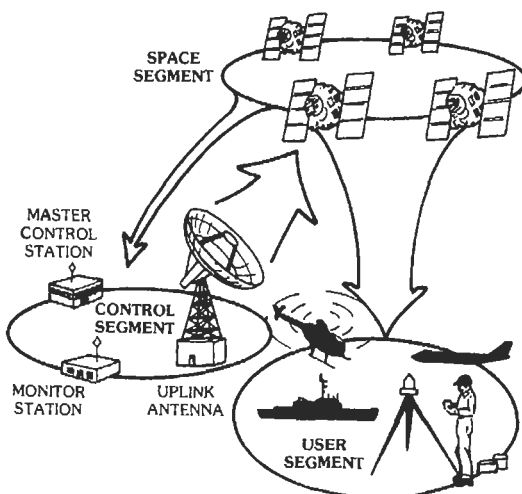
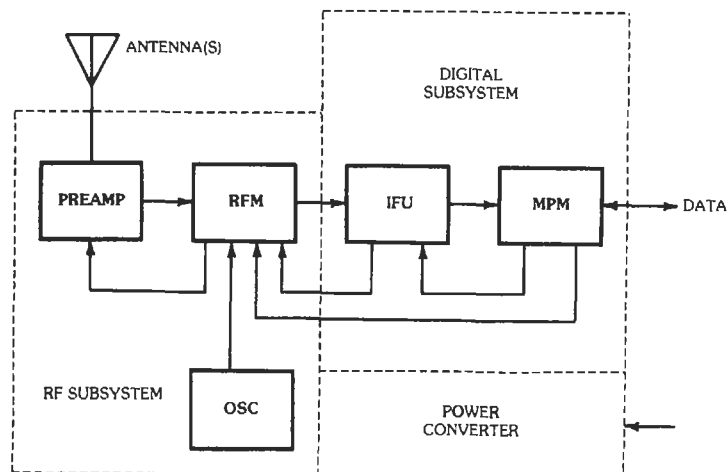


Fig. 1. Global Positioning System.



NOTE:
 RFM—RF Module (Dual-conversion from L-band to 1st IF)
 IFU—IF Module (IF signal processing, code correlation, and generation of I/Q signals along with carrier/frequency track)
 MPM—Microprocessor Module (Microprocessor for signal processing, navigation solution, control, and display functions)

Fig. 2. Global Positioning System receiver.

asures the Doppler shift of the satellite signals. Thus, carrier loop tracking provides direct information concerning the relative velocity between the satellites and the GPS receiver.

Once code and carrier loop tracking is accomplished for each of at least four satellites, range information is fed into a Kalman filter. The Kalman filter algorithm is implemented in software on a microprocessor contained in the GPS receiver. In sophisticated navigation systems, the Kalman filter may accept additional inputs from an inertial navigation system.

The P code is typically encrypted before being broadcast by the satellites. This prevents its utilization by the non-military community. Encryption also prevents an unauthorized party from "spoofing" the system by generating an incorrect signal that is perceived by a GPS receiver to be legitimate. When encrypted, the P code signal is known as the Y code, and the encryption process is known as "anti-spoof," or AS.

In addition to anti-spoof, the operators of the GPS constellation can introduce deliberate timing errors into the C/A code signal to reduce its accuracy. Known as selective availability, or "SA," this deliberate reduction in system accuracy was deactivated by Presidential order on May 1, 2000. Consequently, the accuracy of C/A code receivers is now much better than it used to be. The accuracy is now limited by the uncertainty associated with the frequency dependent propagation delays introduced as the satellite signals pass through the ionosphere and troposphere.

The P/Y code compares measurements made at both the L_1 and L_2 frequencies to estimate these additional delays. Since the C/A code is not currently transmitted

as part of the L_2 signal, there is not currently a convenient means for correcting for the ionosphere delays as part of C/A code operation. However, planned future improvements to the GPS signal include transmitting the C/A code on the L_2 frequency. In deference to military needs, a third signal, the "M" code, added to the L_1 and L_2 channels, is planned for exclusive military use. A new civilian signal, L_5 , will be added at a frequency of 1176.45 MHz.

The navigation message is formatted in five subframes, each of six seconds duration, which makes up a data frame that is 30 seconds long and contains 1500 bits. The data are nonreturn to zero (NRZ) at 50 b/s and are common to the P and C/A signals on both the L_1 and L_2 frequencies. Each data subframe starts with a telemetry word and the C/A to P code handover word. The initial eight bits of the telemetry word contain a preamble that facilitates the acquisition of the data message.

The accuracy of a position fix varies with the capability of the user's equipment and with the satellite geometry with respect to the user. Now that the Selective Availability feature has been disabled, low cost C/A code receivers have accuracies that are limited to the uncertainties associated with ionospheric delays, which can be as large as 30 feet. The P-code, when used with signals at both the L_1 and L_2 frequencies, eliminates this source of error, and is fundamentally more accurate than the C/A receivers due to its higher chipping rate (10.23 MBPS as opposed to 1.023 MBPS). This yields real-time accuracies that can be better than 10 meters. However, the actual accuracy achieved at any point in time depends on the relative

geometry between the user and the observed satellites. For example, an ideal situation would be to have one or more satellites be directly overhead, and for a minimum of three additional satellites to be slightly above the horizon at 120-degree intervals.

Accuracy can be further improved by utilizing additional satellites. Many GPS receivers now have the capability to track up to 12 satellites simultaneously. Because the GPS constellation, as currently deployed, has 27 satellites, it is conceivable that more than 12 might be visible simultaneously. However, the accuracy improvement due to tracking more than 12 satellites is likely negligible.

Differential GPS

Another approach to improving accuracy is to use differential GPS, in which one measures the difference in position between two receivers. It is quite common for benchmark locations to be surveyed with considerable accuracy. A GPS receiver located at the benchmark then provides real-time calibration of the GPS-derived position versus its "ground-truth" fix. Correction signals, which are useful for receivers as far as a few hundred miles away, are then broadcast using a variety of beacon systems. For example, the Coast Guard operates a beacon system at frequencies of approximately 300 kHz in the coastal areas of the United States and along the Mississippi River. Many low-cost receivers are already pre-equipped to process this differential correction signal. Accuracies to better than 6 feet are routine with differential GPS. With additional time-integration and/or post-processing, accuracies to the sub-centimeter level have been achieved.

Future Trends

In order for GPS to serve as the single-navigation source for commercial air transportation, its accuracy, availability, and reliability must be improved. The plan for accomplishing this is known as the Wide Area Augmentation Service, or WAAS. This will be a network of ground-based reference systems whose correction signals are broadcast via a network of geostationary satellites. WAAS-capable receivers will also utilize integrity checking algorithms to verify that the instantaneous operation of the system meets the strict requirements of a safety-of-life navigation system.

To further improve the accuracy of C/A code receivers, it is planned to include the C/A signal on the L_2 broadcast, thus allowing for elimination of ionospheric delays. To further separate the needs of military users from those of civil users, a third signal, the M-code, is planned for inclusion at the L_1 and L_2 frequencies. A new signal, L_5 , will be broadcast for civilian use at 1176.45 MHz.

In a similar fashion, a Local Area Augmentation Service (LAAS) will enable GPS to be used for Category-I and higher instrument landings without the need for conventional ILS equipment.

An additional alternative to the ILS, for the military, is the Joint Precision Approach and Landing System (JPALS). Communication networks, such as the Joint Tactical Information Distribution System (JTIDS) will also provide an incipient radio navigation capability to its users.

A system known as automatic dependent surveillance (ADS) is being developed in which aircraft broadcast their own navigation data to other, nearby aircraft.

Finally, the existing Russian GLONASS System and the Proposed European Galileo System provide similar capabilities to GPS. The GLONASS system, which is a frequency division multiple access (FDMA) system, as opposed to a CDMA system like GPS, has a questionable long-term future for economic reasons. Likewise, development and deployment of Galileo is likely to be affected by economic considerations.

REFERENCES

1. Beser, J., and Parkinson, B. W. "The Application of NAVSTAR Differential GPS in the Civilian Community." *Navigation*, Vol. 29, No. 2, Summer 1982.
2. Borty, J. E., Gupta, R. R., Scull, D. C., and Morris, P. B. "OMEGA Signal Coverage Prediction." *Navigation*, Vol. 23, No. 1, Spring 1976.
3. Bowditch, N. *American Practical Navigator*. US Navy Hydrographic Office, Publ. 9, 1958.
4. Copps, E. M., Geier, G. J., Fidler, W. C., and Grundy, P. A., "Optimal Processing of GPS Signals." *Navigation*, Vol. 27, No. 3, Fall 1980.
5. Frank, Robert L. "History of Loran C." *Navigation*, Vol. 29, No. 1, Spring 1982.
6. Fritch, V. J., and Sanders, L. J. "Instrument Landing Systems." *IEEE Communications*, May 1973, pp. 435-454.
7. Garrett, P. H. "Advances in Low-Frequency Radio Navigation Methods." *IEEE Transactions on Aerospace and Electronic Systems*, July 1975, pp. 562-574.
8. Handbook: VOR/VORTAC Siting Criteria. FAA, Department of Transportation, 1968.
9. *IEEE Transactions on Aerospace Navigation Electronics*. Special issue on the VOR/DME navigation system—Its present capabilities and future potential, Vol. ANE-12, March 1964.
10. Johnson, C., and Ward, P. "GPS Applications to the Seismic Oil Exploration." *Navigation*, Vol. 26, No. 2, Summer 1979.
11. Kayton, M., and Fried, W. R. *Avionics Navigation Systems*, 2nd ed. New York: John Wiley & Sons, Inc., 1997.

12. Kelly, R. J. "Time Reference Microwave Landing System Multipath Control Techniques." *Navigation*, Vol. 23, No. 1, Spring 1976.
13. Logsdon, Tom. *Understanding the Navstar; GPS, GLS, and IVHS, 2nd ed.* New York: Chapman and Hall, 1995.
14. *Loran C User Handbook COMDTINST M16562.3*, US Coast Guard (G-NRN-3), Washington, DC.
15. Microwave Landing System Signal Format and System Level Functional Requirements, FAA-ER-70008C. April 1979.
16. *Navigation: Journal of the Institute of Navigation*. Vol. 25, No. 2, Summer 1978; Special issue with coverage of NAVSTAR Global Position System.
17. Pierce, J. A. "OMEGA." *IEEE Transactions on Aerospace and Electronic Systems*, Vol. AES-1, December 1965, pp. 206-215.
18. Pierce, J. A., and Woodward, R. H. "The Development of Long Range Hyperbolic Navigation in the United States." *Navigation*, Vol. 18, No. 1, Spring 1971.
19. Pisacane, V. L., Holland, B. B. and Black, H. D. "Recent Improvements in the Navy Navigational Satellite System." *Navigation*, Vol. 20, 1973.
20. Poppe, Martin. "The Loran-C Receiver, A Functional Description." *Navigation*, Vol. 29, No. 1, Spring 1982.
21. Stansell, T. A., Jr. "The Navy Navigation Satellite System; Description and Status." *Navigation*, Vol. 15, No. 3, Fall 1968.
22. Swanson, E. R. "OMEGA Possibilities: Limitations, Options, and Opportunities." *Navigation*, Vol. 26, No. 3, Fall 1979.
23. Vass, E. R. "OMEGA Navigation System: Present Status and Plans 1977-1980." *Navigation*, Vol. 25, No. 1, Spring 1978.
24. Watt, A. D. *VLRF Radio Engineering*. Elmsford, N.Y.: Pergamon Press, Inc., 1967.
25. Winnick, A. B., and Brandewie, D. M. "Recent VOR/DME System Improvements." *Proceedings of the IEEE*, Vol. 58, March 1970.
26. Zimmerman, W. "Optimum Integration of Aircraft Navigation Systems." *IEEE Transactions on Aerospace and Electronic Systems*, Vol. AES-5, September 1969.

38

Common Carrier Transmission

*Revised by
Ludwell Sibley*

The Switched Telecommunications System	38-3
Conventional Voiceband Service	
Integrated Services Digital Network Service	
The Exchange Plant	38-3
Subscriber Sets	
Subscriber Loops	
Local Interoffice Trunks	
The Toll Transmission Plant	38-6
Exchange-Access Trunks	
Intertoll Trunks	
Digital Subscriber Line	38-7
IDSL	
ADSL	
ADSL “Lite”	
RADSL	
VDSL	
HDSL	
SDSL	
International Recommendations	38-8
Overall System Design	
Loss	
Echo Control	

Quality Objectives and Definitions 38-9

- Noise
- Levels and Loss
- Echo and Stability
- Crosstalk
- Digital Performance

Facilities 38-11

- Wire
- Repeaters
- Four-Wire Repeaters
- Hybrid Repeaters
- Fiber-Optic Facilities
- Carrier Systems

Network-Control Signaling 38-16

- Subscriber-Loop Signaling
- Interoffice Signaling

THE SWITCHED TELECOMMUNICATIONS SYSTEM

The telecommunications system combines many elements that work together to provide a variety of services. It is a network of networks, owned and operated by multiple parties. In this section we consider two applications of the switched network: voice service and ISDN, and a special case, the Digital Subscriber Line.

Conventional Voiceband Service

Switched voiceband connections for speech or data use (including facsimile and data approaching 56 kb/s) are set up by giving appropriate instructions from the subscriber's telephone set or other terminal. Every subscriber is connected to a local end-office switching system via a loop; consequently, there is a loop at each end of a switched connection. To connect between loops served by different end offices, it is necessary to use a trunk route.

The route may be a single trunk connected directly between end-office switching systems. Other local trunks provide access to and from "wireless" carriers for cellular-mobile, paging, and similar services; to a Public Safety Answering Point for "911" emergency service; or to an Internet Service Provider (ISP).

Alternatively, the trunk route may consist of several trunks routed through one or more tandem switching systems. Such a call may involve trunks between the switching facilities of exchange and interexchange carriers (domestically) or of different operating administrations (internationally). Direct trunks and trunks to tandem switching systems are generally limited to the exchange area, covering tens of miles or occasionally a few hundred miles. When a connection must be made over greater distances, the trunk route passes through toll-connecting and intertoll trunks and is switched by toll switching systems in accordance with a routing plan which is usually dynamically controlled according to traffic load.

As a matter of terminology, the local end office is commonly referred to as a "Class 5" office. This identification is a remainder of a one-time switching plan that used a hierarchy of toll switching offices, up through "Class 1." The hierarchy has been replaced by a plan in which all toll switching offices have essentially equal rank, but the "Class 5" office remains.

Integrated Services Digital Network Service

ISDN is a multi-use network in which wholly digital transmission is provided between customer locations, with digital telephones and data terminals being used. Access from the subscriber to the serving central office

is defined in terms of "B" channels (64 kb/s each) and "D" channels (16 or 64 kb/s). The "B" channels carry voice, data, and video traffic in digital form, while the subscriber controls service (signaling and supervision) via a two-way message link on the "D" channel. There may be two "B" channels and one "D" (in basic-rate access) or 23 "B" channels and one "D" (in primary-rate access).

Digital connections are established in a way comparable to that for conventional voice service. Pairs or multiples of "B" channels may be used ("bonded") together for 128 kb/s or higher speeds.

The following material gives information about the various transmission elements (e. g., loops and trunks) and about the signaling systems used to set up, hold, and disconnect the channels.

THE EXCHANGE PLANT

Subscriber Sets

The 500-type subscriber set is the historical basis for standards for analog subscriber instruments (ANSI/TIA/EIA-470-B-97), even though later telephones are usually designed quite differently. Fig. 1 illustrates transmission performance of a 500-type set on two hypothetical connections involving relatively short and relatively long loops. Fig. 2 shows the test configuration for this measurement. Fig. 3 gives the relative transmitting and receiving volume levels for this set as a function of loop length.

In current standards, the performance of telephone sets is defined in terms of loudness rating. Two parameters are of particular interest: the transmitting and receiving objective loudness ratings.

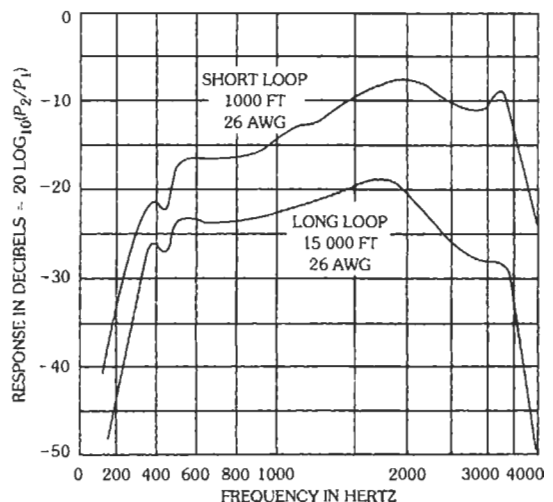


Fig. 1. Comparison of overall response. (From W. F. Tuffnell, "500-Type Telephone Set," *Bell Laboratories Record*, Vol. 29, September 1951, pp. 414-418. © 1951 Bell Telephone Laboratories, Inc.)

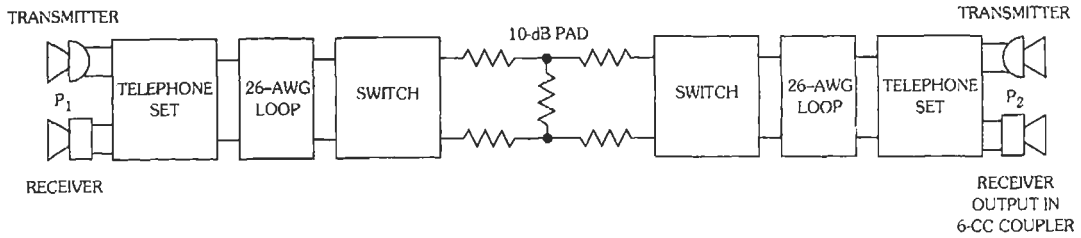


Fig. 2. Test arrangement for Fig. 1. (From W. F. Tuffnell, "500-Type Telephone Set," *Bell Laboratories Record*, Vol. 29, September 1951, pp. 414-418. © 1951 Bell Telephone Laboratories, Inc.)

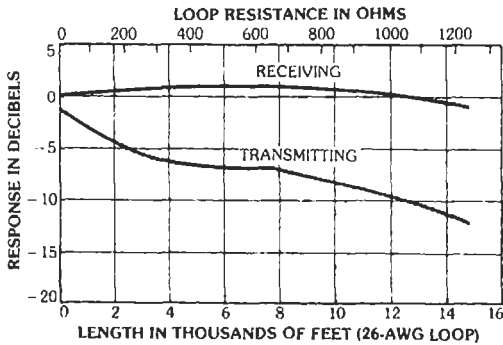


Fig. 3. Relative volume levels. (From W. F. Tuffnell, "500-Type Telephone Set," *Bell Laboratories Record*, Vol. 29, September 1951, pp. 414-418. © 1951 Bell Telephone Laboratories, Inc.)

Transmitting objective loudness rating (TOLR) is a measure, in decibels, of the efficiency with which a telephone set converts acoustic pressure at the set transmitter to an electrical output voltage. This loss-like acoustic-to-electric quantity is given by:

$$TOLR = -20 \log_{10} V_t / S_m$$

where

V_t = Loudness output voltage of the transmitting component (in millivolts across 900 ohms resistive),

S_m = Loudness sound pressure at the mouth reference point (in pascals).

Receiving objective loudness rating (ROLR) is a measure, in decibels, of the efficiency with which a telephone set converts an electric source voltage to an acoustic pressure output at the receiver. This electric-to-acoustic loss-like quantity is given by:

$$ROLR = -20 \log_{10} S_e / (\frac{1}{2} V_w)$$

where

V_w = Loudness open-circuit voltage of the electric source (in millivolts),

S_e = Loudness sound pressure at the ear (in pascals).

Fig. 4 shows the average TOLR versus loop length, measured for nine recent-vintage subscriber sets with electronic transmitters. The figure includes performance standards from standard ANSI/TIA/EIA-470-B-97 (upper and lower desirable limits, and mean performance), and includes the 500-type carbon-transmitter set as a reference model. The figure pertains to operation with a 400-ohm 48-volt battery-feed circuit as found in many central offices, and a 900-ohm ac termination. Results are similar with a regulated (constant-current) battery-feed circuit also commonly used in central offices. Fig. 5 gives the corresponding ROLR figures.

Typical subscriber-set dc internal resistance is 200 to 300 ohms, up to a maximum of 430 ohms for all equipment and wiring at the customer's premises.

Transmission standards for digital telephone sets are intended to give comparable performance. They are given in standard ANSI/TIA/EIA-579-A-98.

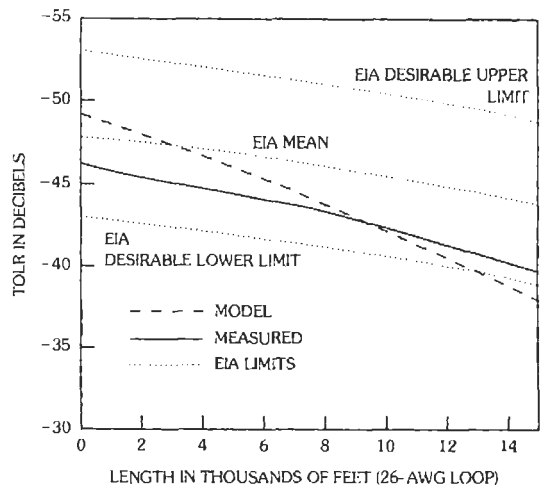


Fig. 4. Transmitting objective loss rating.

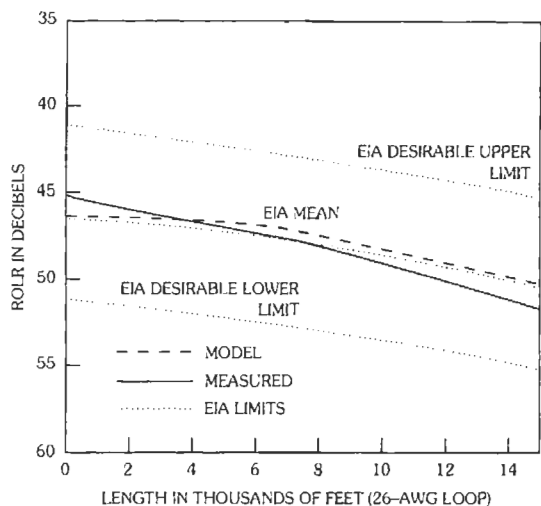


Fig. 5. Receiving objective loss rating.

Subscriber Loops

Design Objectives—The performance objectives for loops are generally as follows:

- Loss*—The loss of an analog subscriber loop, including any carrier link, should be less than 9.0 dB at 1000 Hz.
- Noise*—The noise objective for subscriber loops is typically 20 dBmC at the customer interface.
- Current*—For satisfactory transmission and signaling performance, subscriber-loop dc current is generally held between 20 and 65 milliamperes. The customer-premises equipment is assumed to have a dc resistance up to 430 ohms.

Resistance Design—The design practice for most subscriber loops using wire facilities, before the wide use of digital loop carrier (DLC) facilities, was the

“Resistance Design” plan. It used a maximum loop length of 24 kilofeet and a maximum loop resistance of 1500 ohms, with the objective to achieve a distribution of loss values clustered broadly in the vicinity of 3 to 4 dB. Limits on bridged-tap (temporarily unused extensions of the loop for flexibility in future use) and rules for loading were included to provide satisfactory transmission. On loops extending beyond 18 kft (including bridged tap), inductive loading was applied with 88-mH coils at 3, 9, 12, etc. kft from the central office. The added inductance was intended to reduce the transmission loss due to shunt capacitance up to about 3500 Hz within the voice band, at the expense of transmission at higher frequencies. Central-office switches, either digital or electronic analog, typically have signaling and supervision limits of 1600 ohms or more, leaving loop resistance limited only by the need to provide 20 mA or more of loop current. For longer loops, resistance above 1500 ohms was permitted, with the use of signaling-range extension and transmission gain. However, such long wire loops are now unusual because loops of nine kilofeet or even less are usually candidates for carrier facilities.

Digital Loop Carrier—Although it was first introduced to serve customers at great distances from an end office, the declining relative cost of electronics, and the demand for digital services like ISDN or DSL, have made DLC attractive economically. As a result, roughly half of North American loops use it. Where carrier is applied to loops, transmission and signaling are generally better than with loops on wire facilities. (A typical carrier channel has a 1004-Hz loss of 2 dB, but some later designs insert extra loss when the distribution cable is short, in order to prevent excessive loudness.) With carrier facilities, or with wire loops close to the central office, a Carrier Serving Area (CSA) plan is commonly used for new or rehabilitated facilities. Fig. 6 shows a typical carrier system used with CSA loops. With the CSA plan, the carrier remote terminal is placed no farther from the subscriber than 12 kft (9 kft if purely 26-gauge cable), and bridged

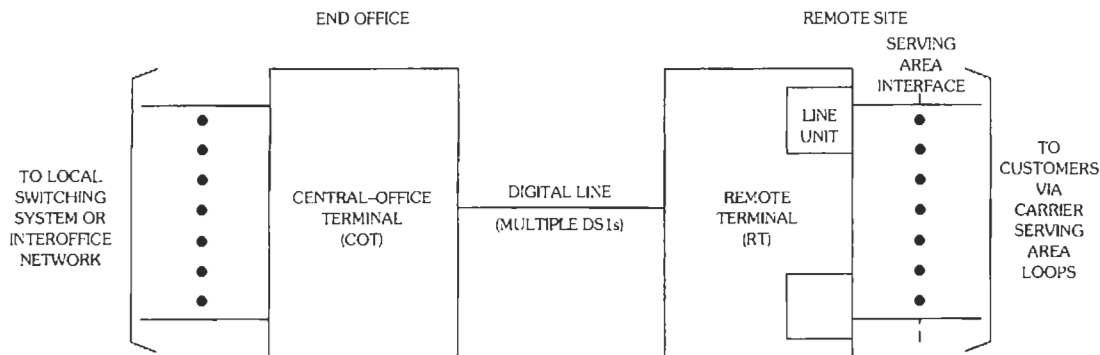


Fig. 6. Universal digital loop carrier system.

taps are restricted. Loading coils are not needed over these lengths and, in the interests of high-speed digital services, are not used. The figure shows a "universal" system in which each channel emerges individually from the CO terminal (COT) for physical cross-connection to the switch or (for a special service) an inter-office facility. With integrated carrier, the COT function is integrated into the switch, and there is no terminal as such.

The remote terminal may be placed on the premises of a customer needing many lines. The system normally uses conventional 64-kb/s PCM coding, although 32-kb/s coding is available as an option in some systems.

High-Speed Digital Services—The CSA Plan is fully compatible with the digital subscriber line for ISDN basic-rate access, while the Resistance Design plan is compatible up to 18 kft. The standard ISDN two-wire loop carries the two "B" channels, the "D" channel, and a maintenance channel at 12 kb/s, for a total of 160 kb/s. The line signal is quaternary, resulting from 2B1Q coding (two binary pulses recoded into one quaternary). A midspan repeater, and removal of loading coils, is sometimes used for range extension beyond 18 kft. Most versions of the DSL family (see below) are intended for use up to the 12–18 kft maximum zone.

Fiber in the Loop (FITL)—Where conventional telephone service is provided by means of optical fibers, a terminal device comparable to a DLC terminal is used at or near the subscriber's location. It is designed to accommodate the very low dc resistance and transmission loss of the metallic drop cable and station wiring.

Wireless Access—In deep-rural areas, there is currently an application of digital radio, operating at an 8- or 16-kb/s coding rate, for subscriber loops in the Basic Exchange Telecommunications Radio Service. This involves access from a digital radio terminal on the subscriber's premises to a "radio port" connecting with the central office. The overall design intent is to provide voice performance equivalent to that of wire facilities although the low bit rate precludes use for high-speed data or facsimile service.

Local Interoffice Trunks

Local interoffice trunks interconnect end-office switching systems and tandem switching systems. Telecommunications administrations use objectives for transmission loss that may depend on the type of trunk. Digital trunks between digital switches are inherently lossless, so the loss of a built-up connection is fixed by the decoding level at the end office (for an analog connection) or at the customer equipment (for a digital connection). Local trunks are short enough that talker echo is not a problem. Fixed-loss design objectives are usually used. In an analog network, a direct

trunk between end-office switching systems is the only trunk in a connection and may have a loss objective between 3 and 6 dB. Objectives for trunks between analog end-office switches and tandem switching systems are typically 3 dB, with 0 dB as the objective for trunks between tandems.

THE TOLL TRANSMISSION PLANT

Exchange-Access Trunks

Toll connecting trunks (in the United States, exchange-access trunks) interconnect an end-office switching system and a toll switching system. Present transmission design objectives require a loss of between 2 and 4 dB in such trunks, or use of echo cancellers, to compensate for the imperfect impedance match to an analog subscriber loop at the end office. The match is relatively poor because one or two compromise impedances cannot closely simulate the variety of impedances of loops of different constructions and lengths.

Intertoll Trunks

Overall Connection Loss—Intertoll trunks connect toll switching systems. The connection loss of a digital trunk or combination of digital trunks, as stated before, is set by the decode level. On long connections it is common practice to include digital echo cancellers to improve singing margin and block echoes. Echo-canceller technique involves formulating an estimate of the echo that will result from a signal reflected back onto the four-wire side of a hybrid that converts from four-wire trunking to a two-wire loop. That estimate is subtracted from the signal on the return path. The echo-canceller device produces the echo estimate by means of a signal processor (filter) for which parameters are set by correlation techniques. Previous practice, on analog intertoll trunks, was to use a "Via Net Loss" transmission plan with loss increasing with trunk length so as to mask echoes. On very long trunks, above 1850 miles, the VNL plan converted to zero loss and relied on an echo suppressor, a switching device which momentarily inserted loss into the idle direction of transmission in order to block the echo path.

Delay—Delay, by itself, is seldom annoying in speech communication unless it reaches a value of approximately 600 milliseconds. Delays encountered in terrestrial transmission plant are well below this value, but a round-trip delay of this magnitude occurs on circuits operating via synchronous-orbit satellites. Further delays occur in low-bitrate speech coders, echo cancellers, digital cross-connection systems, voice packetizers, cryptographic devices, etc.

Echo Objectives—In the absence of echo cancellation, echo return loss (ERL) must be held to high levels to assure satisfactory echo performance of toll

facilities. The ERL objective for analog loops at the end office is typically 11 dB. With the general use of four-wire switching, this is the only relevant echo source in a built-up connection. (Historically, in the important special case of a two-wire interface to a toll connecting trunk at an analog toll office, the ERL was held to a minimum 16 dB, with an objective of 22 dB.)

DIGITAL SUBSCRIBER LINE

An emerging application of network elements having great near-term importance is the Digital Subscriber Line. This facility does not necessarily provide switched service in the usual sense—its most visible application is for “always-on” access to an ISP. In that sense, it is a private-line channel from one user to a DSL access multiplexer (DSLAM). The DSLAM aggregates the data traffic from a group of users, usually by packet switching via Asynchronous Transfer Mode (ATM) techniques in the interoffice environment. (See Chapter 26.)

This is an emerging technology whose standards are still evolving, with a variety of pre-standard and proprietary versions being made available. Some implementations of them use equipment built into the line card of a digital switch or DLC terminal to minimize the effort of installation. The core idea is to send bi-directional digital data, usually on a two-wire loop, at previously unattainable speeds by use of combinations of signal processing, automatic equalization, error correction, and echo cancellation. The base facility is the nonloaded copper subscriber loop. The terminals compensate for signal attenuation at high frequencies, irregularities in transmission and echoes caused by bridged taps and changes in gauge at splices, crosstalk from the other direction of transmission and other DSL systems in the same cable, errors caused by impulse

noise, and even radio-frequency interference from broadcast stations.

On a given loop, the speed attainable may be controlled by either the technical conditions of the loop or by marketing considerations of the service provider.

“DSL,” also called xDSL, is actually a family of approaches having somewhat different uses. Table 1 is a snapshot of the major present contenders. The distance and speed capabilities are subject to change as the technology matures.

ISDL

The ISDN DSL uses the technology discussed above under ISDN basic-rate access (160 kb/s bi-directional, no analog telephone on the same pair, 2B1Q coding). It is usable for non-ISDN applications, at the same speed, to distances beyond the capabilities of other DSLs.

ADSL

For many important uses, the flow of data to/from the subscriber is asymmetric: a small traffic of files and controls toward the network (upstream), and a large volume of data downstream. Such applications include Internet downloads, major file transfers, and streaming compressed video. This implies the value of an Asymmetric DSL.

Fig. 7 illustrates such a line. At the central office or carrier remote terminal, a low-pass filter connects the voicegrade end-office switch to the loop, while a high-pass filter attaches an ADSL transceiver. Data to/from the transceiver may be aggregated from multiple users and sent to an ISP or other site.

At the subscriber end, a second low-pass/high-pass filter (splitter) separates voiceband and ADSL signals.

TABLE 1. DIGITAL SUBSCRIBER LINES

Name	Cable Pairs Used	Max. Loop Length kft	Speed Downstream Mb/s	Speed Upstream Mb/s	Operates Above Voice ?
ISDL	1	18*	0.16	0.16	No
ADSL	1	18	1.54	0.064-1.0	Yes
ADSL “Lite”	1	18	1.0	0.16	Yes
HDSL	2	18	1.544	1.544	No
SDSL	1	10	0.16-2.048	0.16-2.048	No
RADSL	1	21.3	0.64	0.128	Yes
	1	18	1.544	0.176	Yes
	1	12	12	1.0	Yes
VDSL	1	1	52	1.52-2.3	Yes

* Without repeater.

Speeds shown are examples / ranges only.

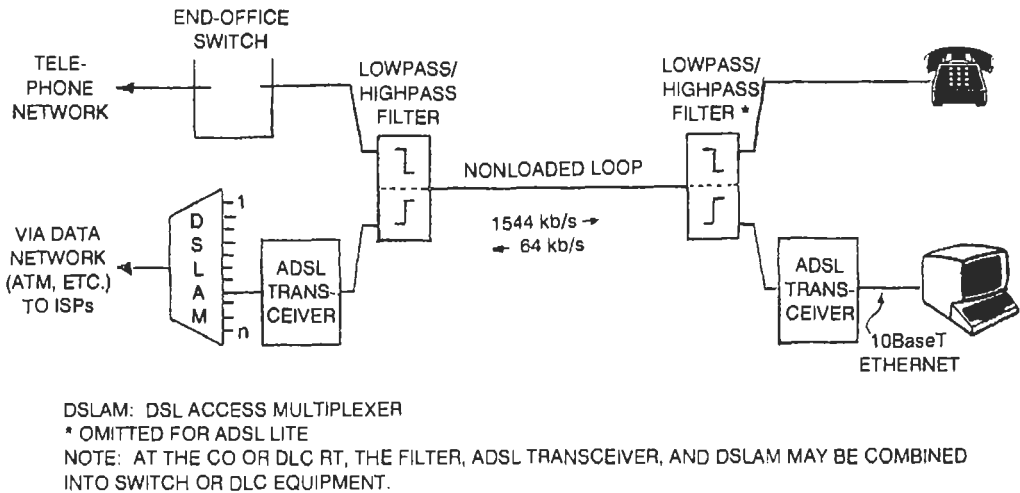


Fig. 7. Asymmetric digital subscriber line.

ADSL "Lite"

The ADSL concept assumes a common filter at the subscriber end to separate the telephones from the data signal. In mass-market applications this requires installation effort and possibly added wiring. By giving up the use of the filter, it is possible to obtain service with some modest sacrifice in speed. (The ADSL transceiver must contend with varying line impedance, depending on the usage of the telephones. Some telephones may rectify the data signal and generate interfering noise, requiring addition of a plug-in low-pass-only filter.)

RADSL

The Rate Adaptive DSL is an ADSL variant intended to supply downstream service at speeds up to 12 Mb/s (to 12 kft), 1.544 Mb/s (to 18 kft), or 640 kb/s (to 21 kft). Upstream speeds vary proportionately: 1 Mb/s, 176 kb/s, and 128 kb/s, respectively. The terminals automatically "fall back" or raise speed upon occurrence of changes in the loop.

VDSL

The Very-High-Speed DSL is intended to deliver asymmetric data at speeds up to 52 Mb/s over relatively short distances, typically in a campus situation or as the distribution link from a multiplexed fiber terminal.

HDSL

The High-Speed DSL variant is intended to provide a symmetric DS1-level facility (1.544 Mb/s) over two cable pairs, thus avoiding the time and expense of add-

ing a string of intermediate regenerators. Symmetric DSLs may be attractive as a means of multiplexing several voice lines, plus data, over a common facility.

SDSL

The Symmetric DSL is intended to supply bi-directional data over one wire pair at, depending on loop condition, up to the E1 (2.048 Mb/s) rate.

INTERNATIONAL RECOMMENDATIONS

Overall System Design

A summary of the ITU-T (formerly CCITT) recommendations for system objectives and design criteria for circuits used in international telecommunications service is given in Chapter 2.

Loss

The ITU-T recommends that international analog toll circuits be given an insertion loss of 0.5 dB for each 500 kilometers of length or fraction thereof, assuming that the international circuit is derived from carrier systems in cable or radio relay. The result approximates VNL within 0.33 dB for lengths up to 2500 km.

Echo Control

The ITU-T recommends that the overall loss of a connection be adjusted so that echo signals are sufficiently attenuated or, alternatively, that an echo canceler or suppressor be fitted if the loss adjustment results in an excessive insertion loss.

QUALITY OBJECTIVES AND DEFINITIONS

Noise

Noise—Noise, in its broadest definition, is any undesired signal in a communication channel. Noise may be classified as quantizing, or signal-dependent, noise; thermal, or white, noise; impulse noise; crosstalk; tone interference; and miscellaneous.

When noise is measured on a voiceband communication channel, a weighting network is often inserted ahead of the detector to account for the varying subjective annoyance from noise of different frequencies into telephone sets, or differing interference into a data modem (MODulator–DEMODulator). For noise measurement on channels of bandwidth greater than voice, weighting networks with correspondingly greater bandwidths are used.

Quantizing Noise—Noise occurring in the presence of signal, resulting from a mismatch between the exact value of an analog signal and the closest available quantizing step in a digital coder.

Thermal Noise—Noise arising from random electron motion, occurring on all transmission media and in all communications apparatus. It is characterized by uniform energy distribution over the frequency spectrum and by a normal or Gaussian distribution of voltage or current.

Impulse Noise—Noise consisting of irregular pulses of short duration and relatively high amplitude. Some sources of impulse noise induced in communication channels are an erroneous digital coding bit caused by an error on a transmission facility, crosstalk from dc signaling systems, lightning, or transients due to switching in adjacent power circuits.

Tone Interference—Interference due to single tones or complex periodic waveforms.

Miscellaneous Noise—Interferences that cannot readily be placed in any of the preceding categories.

Reference Noise—One picowatt (10^{-12} watt) of power. Also commonly stated as -90 dBm (90 dB below one milliwatt).

dBm—Decibels above reference noise. The dBm is the unit of measurement of noise power used in the IEEE Standard 743-1995 noise-measuring set. For measuring noise on voice communication channels, the noise-measuring set is normally equipped with a “C-message” weighting network. When this filter is used, the unique network response causes the reading to deviate from what would be obtained with some other network. Readings in dBm taken with the C-message network are designated “dBmC.” The calibration tone is 0 dBm (1 mW) at 1000 Hz, which reads 90 dBm with or without the C-message weighting network.

One milliwatt of white noise, limited to the band of 0-3 kHz, reads 88 dBm.

The noise objective in present use is that, on any digital connection, the noise should not exceed 17 dBmC, as controlled mainly by the decoder at the receiving end.

pWp—Picowatts, psophometrically weighted. The unit of noise power measured with the ITU-T-recommended psophometer (“noise meter”). The meter is frequency-weighted by a curve comparable to the C-message weighting curve. The reference tone is -90 dBm (1 picowatt) at 800 Hz. The psophometer is defined as measuring the internal (open-circuit) voltage of an equivalent noise generator having impedance of 600 ohms and delivering noise power to a 600-ohm load. For convenience in comparison, the psophometric voltage may be converted to dBm. In such a case, a tone of 0 dBm at 1000 Hz will read +1 dBm on the psophometer; one mW of white noise, limited to the 0-3 kHz band, reads -2 dBm.

Levels and Loss

Reference Level Point—An analog point in a communication circuit arbitrarily chosen as a reference location for signal-level measurements. Common equivalent terms are “0 dB transmission level point,” “zero level,” “zero level point,” “0 dB TL,” “0 dB TLP,” and “0 TLP.”

Relative Level—The relative level at any point in a circuit is the power gain or loss in decibels between the 0 TLP and the point under consideration. Reference of signal and interference powers to the 0 TLP is convenient in system design and maintenance. Relative level is expressed in such terms as “ -3 dB TLP.” Common practice is to define the input of a digital coder in a switch or transmission system as the 0 TLP. At the distant end, the corresponding decoder is conventionally at a TLP of -3 dB, -6 dB, etc. Signal and/or interference powers may be referred to the 0 TLP with a phrase such as “a signal power of -16 dBm0,” which indicates the power of the signal if it were measured at the 0 TLP. The 0 TLP may not be accessible for measurement and, in fact, need not even exist in a given system.

Net Loss—The net loss of a transmission channel is the ratio of the signal powers at the input and the output of the channel, usually expressed in decibels (dB). By custom, the net loss of a channel is understood to be measured at 1004 hertz in the American and Canadian plant and at 1020 (formerly 800) hertz in international practice (ITU-T). (The measurement frequency is displaced slightly from 1000 Hz to keep it from being a subharmonic of 8000 Hz, the sampling rate in a digital coder, to avoid generation of distortion products.)

Insertion Loss—The ratio of the power delivered from a source to a load, to the power delivered from

the same source to the same load through a transducer (network, channel). The definition may be applied to the loss or gain effect caused by the insertion of a gain element (repeater) or equalizer into a transmission channel. The concept is equally applicable to the channel as a whole.

Transducer Loss—The ratio of the maximum power available from a source to the power delivered by that source to a load through a transducer (network, channel). The transducer loss of a channel differs from the insertion loss only because the maximum power from the source is substituted for the power from that source to that load.

Echo and Stability

Return Loss—A measure of the match between the two impedances on either side of a junction point. Return loss is defined by:

$$RL \text{ (dB)} = 20 \log_{10} [(Z_1 + Z_2)/(Z_1 - Z_2)],$$

where Z_1 and Z_2 are the complex impedances of the two halves of the circuit. When the impedances are not matched, the junction becomes a reflection point. The return loss expresses the ratio of incident to reflected signal power. It is often made effectively infinite by introducing echo cancellation into a four-wire portion of the circuit.

Echo Return Loss (ERL)—The weighted power-averaged return loss at the reflection point. The ERL expresses in decibels the ratio of the power of a broadband incident signal to that of the correspondingly broadband reflected signal. Both powers are measured through a weighting network covering the band of approximately 500–2500 Hz.

Singing Return Loss—The same as echo return loss, but over a narrower band near the lower or upper edge of the voice band, e.g., 200–500 Hz or 2500–3200 Hz. The extremes of the band are used because singing usually takes place there.

Talker Echo—A signal returned to the talker after making one or more round trips between the talker and a distant reflection point. Fig. 8 depicts the echo paths associated with a connection between two end offices connected via two toll or tandem offices. The first talker echo is generally the most important. Echoes that are sufficiently loud and sufficiently delayed are annoying and interfere with the speech process.

Listener Echo—A signal first returned toward the talker at a distant reflection point and then reflected again toward the listener (Fig. 8). The listener echo mixes with the original signal, adding to its strength at some frequencies and diminishing it at others. It can result in amplitude distortion and/or reduced singing margin.

Singing Margin—The singing margin of a two-wire circuit is defined as the amount by which the combined losses of the two directions of transmission may be reduced before oscillation (singing) occurs. Inadequate singing margin results in distortion that is often described as sounding hollow or "like a rain barrel."

Crosstalk

Crosstalk—Interference from other communication channels is called crosstalk. It is classified as near-end and far-end crosstalk.

Near-end crosstalk is measured on a channel at a receiving point near the sending point of the interfering channel. Far-end crosstalk is measured on a channel at a receiving point near the receiving point of the interfering channel. In far-end crosstalk, the two channels transmit in the same direction.

Crosstalk may occur via capacitive or inductive coupling, in a cable or item of terminal equipment; via imprecise timing of sampling gates, in a multichannel digital terminal using a shared coder; via intermodulation of channels assigned to different frequency bands, in an analog transmission system; by insufficient antenna discrimination between two cochannel signals

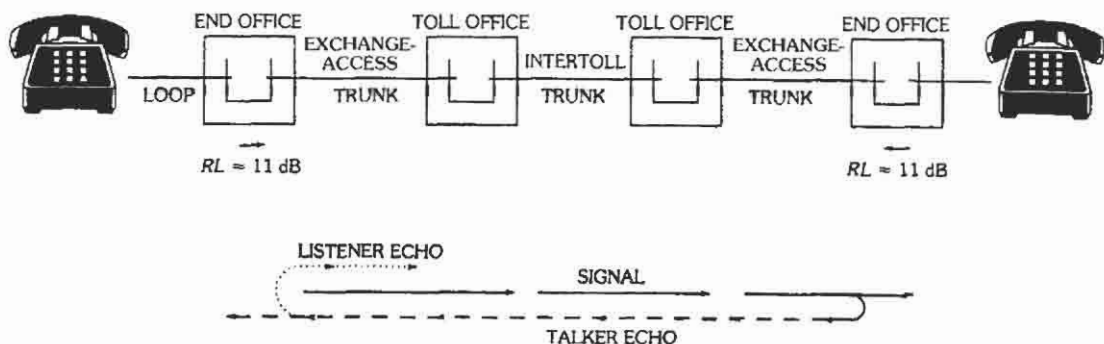


Fig. 8. Return losses and echo paths in a toll connection.

using different polarizations, in a radio system; etc. Crosstalk between digital transmission lines manifests itself by degrading the error margin of one line or both.

Crosstalk is also classed as intelligible and unintelligible. Intelligible crosstalk can be understood by the listener; because it diverts attention, it has more interfering effect than unintelligible cross talk. Crosstalk into a voice-frequency circuit from adjacent voice-frequency circuits may be intelligible. Crosstalk due to incomplete suppression of sidebands, or to intermodulation of two or more frequency-multiplexed channels, is generally unintelligible. Such crosstalk is often classed as miscellaneous noise. Where a coupling path between transmission facilities gives rise to intelligible or nearly intelligible cross talk, normal practice is to design the system so that the probability that a customer will hear a "foreign" conversation does not exceed one percent.

Digital Performance

On digital facilities and services, transmission quality is commonly defined in terms of three parameters:

Bit Error Ratio (BER)—The number of errors in a given number of transmitted bits (e. g., one error in 10^8 or a BER of 10^{-8}).

Errored Second (ES)—A one-second transmission interval that contains one or more bit errors.

Error-Free Second (EFS)—The converse of ES.

Severely Errored Second (SES)—Conventionally, a one-second transmission interval having a BER of 10^{-3} or more.

As examples, the following performance objectives and limits have been used:

Parameter	ES	EFS	SES
	BER	%	%
DSL Design Objective	10^{-7}		
Local DS1 Svc.	10^{-7}	216	99.75
Local DS3 Svc.	2×10^{-10}	432	99.5

FACILITIES

Wire

Multiconductor copper cable is widely used in pair counts up to 2700 or even higher. Cable insulated with paper pulp or polyethylene, in sizes of 22, 24, and 26 American Wire Gauge, are most common for voiceband, DSL, and carrier frequencies up to several megahertz. The transmission characteristics of typical paired cables are given in Table 2. The data cover the voiceband and higher frequencies: 40 kHz (DSL), 225 kHz (HDSL), and 772 kHz (T1 carrier spans). While

most longer loops and interoffice trunks utilize carrier transmission, it is possible to use loaded cable, by itself or with voice-frequency repeaters. Open-wire facilities have generally disappeared, even in rural areas. For carrier transmission frequencies up to 60 MHz, multiple coaxial cable facilities have been used. Such cables have occasionally been converted to digital transmission.

Repeaters

General—If cable loss exceeds limits, it is necessary to add gain to the circuit. Amplifiers for this purpose are termed "voice-frequency repeaters." They are used mainly for transmission-range extension on some local loops, special-service circuits, etc. Two types of repeater are typically used: four-wire and hybrid. Both provide independent gain and equalization for each direction of transmission.

Four-Wire Repeaters

Four-wire repeaters are associated with four-wire cable facilities, generally as voice-frequency extensions of carrier transmission systems in full-duplex data circuits and the like. The associated repeaters can be used in applications requiring more gain and/or equalization than can be obtained with two-wire repeatered facilities.

Hybrid Repeaters

Two-wire hybrid repeaters are found in special subscriber-loop applications. For two-wire facilities, the gain and equalization elements are coupled to the facility with four-wire hybrids (of transformer or electronic types), as shown in Fig. 9. To assure stability of the internal gain loop, hybrid balance is controlled by design and adjustment of the balance networks which should exactly match the impedance seen at the cable-facility interface. Any mismatch limits the gain capability

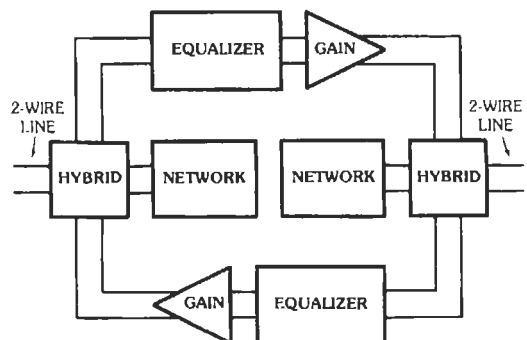


Fig. 9. A typical hybrid repeater.

TABLE 2. CHARACTERISTICS OF TYPICAL TELEPHONE CABLES*

	Frequency f (Hz)	Primary Constants				Nonloaded Secondary Constants				H88 Loaded** Secondary Constants			
		R (Ω /mi)	L (mH/mi)	G (μ mho/mi)	C (μ F/mi)	α (dB/mi)	β (deg/mi)	$ Z $ (Ω)	$\angle Z$ (deg)	α (dB/mi)	β (deg/mi)	$ Z $ (Ω)	$\angle Z$ (deg)
19 AWG Pulp	500	86.1	0.887	0.506	0.084	0.91	6.2	571	-44.0	0.42	14.9	1009	-10.3
	1000	86.2	0.886	1.219	0.084	1.27	8.9	405	-43.1	0.42	29.8	1017	-5.3
	2000	86.4	0.884	3.133	0.084	1.74	13.0	287	-41.3	0.43	62.2	1183	-2.7
	3000	86.6	0.883	5.520	0.084	2.07	16.5	236	-39.5	0.52	105.5	1914	-3.4
	5000	87.1	0.881	11.35	0.084	2.52	22.6	186	-36.1	13.78	158.7	926	-88.5
	40K	103.8	0.834	177.1	0.083	4.45	123.2	106.2	-13.4				
	225K	205.9	0.735	1577	0.083	10.10	633.3	95.2	-5.3				
	772K	370.1	0.668	7570	0.082	21.00	2062	90.4	-2.7				
22 AWG Pulp	500	173.1	0.871	0.494	0.082	1.29	8.6	819	-44.5	0.76	15.3	1098	-18.0
	1000	173.2	0.870	1.190	0.082	1.81	12.3	580	-44.0	0.79	29.7	1051	-9.9
	2000	173.3	0.869	3.058	0.082	2.52	17.6	411	-43.1	0.79	61.5	1196	-5.1
	3000	173.5	0.867	5.388	0.082	3.04	21.9	336	-42.2	0.88	103.5	1871	-5.6
	5000	173.8	0.866	11.08	0.082	3.81	29.2	262	-40.4	13.54	159.5	960	-87.4
	40K	185.4	0.872	177.1	0.084	7.40	132.5	117.9	-20.6				
	225K	298.5	0.798	1577	0.084	13.85	666	99.3	-7.2				
	772K	528.9	0.714	7570	0.084	28.02	2154	93.4	-3.8				
24 AWG Pulp	500	274.0	0.951	0.506	0.084	1.64	10.9	1019	-44.6	1.13	16.4	1201	-24.3
	1000	274.0	0.950	1.219	0.084	2.31	15.6	721	-44.3	1.21	30.7	1079	-14.6
	2000	274.1	0.949	3.133	0.084	3.24	22.2	510	-43.7	1.23	62.7	1200	-7.7
	3000	274.2	0.948	5.520	0.084	3.92	27.5	417	-43.0	1.35	105.6	1906	-8.7
	5000	274.3	0.947	11.35	0.084	4.95	36.3	324	-41.8	13.81	160.8	927	-86.2
	40K	285.4	0.892	177.1	0.084	10.52	142.6	133.2	-26.3				
	225K	391.4	0.830	1577	0.084	17.54	682.7	102.3	-9.0				
	772K	674.2	0.753	7570	0.083	34.14	2216	95.8	-4.8				
26 AWG Pulp	500	440.0	0.995	0.476	0.079	2.03	13.4	1330	-44.7	1.58	17.6	1444	-30.5
	1000	440.0	0.995	1.146	0.079	2.85	19.0	942	-44.5	1.79	31.0	1199	-20.8
	2000	440.0	0.995	2.946	0.079	4.01	27.1	666	-44.1	1.86	61.3	1247	-11.8
	3000	440.1	0.994	5.191	0.079	4.87	33.4	544	-43.7	1.97	101.0	1820	-11.8
	5000	440.2	0.993	10.67	0.079	6.20	43.7	422	-42.9	13.21	161.7	1018	-84.4
	40K	446.3	0.943	177.1	0.081	14.3	156.8	161	-31.2				
	225K	531.2	0.889	1577	0.080	22.2	697.3	109.9	-11.3				
	772K	854.2	0.826	7570	0.080	40.0	2266	102.9	-5.5				

* Information for 5000 Hz and below taken from *Transmission Data--Exchange Area Cables*, ©1962 AT&T; for 40 kHz and above, from *Provision of Metallic Pairs to Other Common Carriers for Digital Data Transmission*, ©1974 AT&T. Data may be used for PIC cables in the voiceband.

** H88 loading involves use of 88-mH loading coils at 6000-ft intervals. Data apply to cables with 3000-ft end sections on both ends.

of the repeater. Where precise automatic balancing of the hybrid is desired, a design may be used in which the four-wire path includes a digital coder and decoder in each direction, along with a digital echo canceller. The canceller provides the equivalent of a very high degree of hybrid balance.

Fiber-Optic Facilities

Glass-fiber optic transmission facilities (see Chapter 22) have generally displaced wire facilities as a more economical and higher-performance alternative for carrier transmission. Fiber-optic facilities are used for both digital trunk carrier and digital subscriber carrier systems and, in sizes up to 864 fibers, are finding direct application in subscriber loops as well. Transmission speeds up to OC-192 (10 Gb/s) on a single wave length are commonplace in the long-haul environment, and extensions of capacity are available by use of wavelength-division multiplexing (WDM), in which 40 wavelengths or more may be carried simultaneously. Optical amplifiers may be used on particularly long fiber lines, especially in submarine cables. Fiber facilities, in "self-healing" ring configurations, are often used to improve service reliability.

Carrier Systems

General—As suggested before, when many telecommunications channels are needed between two points, it is usual to use a carrier system to multiplex many channels over one medium. Media used for such systems include glass fibers, cable pairs, and microwave radio (terrestrial or satellite).

Multiplexing Techniques—Two basic techniques are used for the transmission of multiple channels over a single transmission medium:

Time-division systems, usually digital, in which each communication channel is allotted a discrete time slot within a sampling frame, occupying essentially the entire wideband frequency spectrum for the allotted time. The original practice was to multiplex directly, with a fixed destination; another technique of growing importance is to convert the bitstream into packets for transmission and routing, usually to multiple destinations.

Frequency-division systems, in which a unique band of frequencies within the wide frequency spectrum of the medium is allotted to each communication channel on a continuous-time basis. This technique is found mainly in radio systems, and in wavelength-division multiplexing in fiber-optic systems.

Modulation Techniques—Several modulation methods may be employed with either of the multiplexing techniques. Pulse-code modulation is the most common time-division modulation technique, although adaptive differential PCM, delta modulation, and other techniques are also found occasionally. Several frequency-division modulation techniques are in use,

including frequency modulation and amplitude modulation, both double- and single-sideband. See Chapters 23 and 24 for a more detailed discussion of modulation techniques.

Trunk Carrier—Trunk carrier systems operate between switching offices, using either trunk terminations built into digital switches (now usual) or discrete carrier terminals (channel banks) that deliver analog voice channels.

Circuit Multiplication—For voice connections through digital submarine cables, digital circuit-multiplication systems may transcode the 64-kb/s speech signal to 32 kb/s and perform speech interpolation (connecting a talker to the channel only when speech is present). This yields a typical circuit multiplication of 5:1.

Common-Carrier Hierarchies—There are several agreed-upon bundle sizes or multiplex levels for carrying channels in common-carrier systems. Most of the standards are covered in ITU-T documents. In digital transmission, the basic entity is a 64 kb/s channel in which the nominal 4-kHz voice signal has been filtered (to minimize aliasing) and sampled at 8 kHz, with each sample encoded into eight bits. Digital data at up to 64 kb/s can be substituted for the digitized voiceband signal. Fig. 10 gives the digital hierarchy levels corresponding to bit rates, and the numbers of usable 64-kb/s time slots, recommended by ITU-T. The figure is couched in terms of the basic DS1 transmission rate of 1544 kb/s as used in North America. Not defined by ITU-T, but also occasionally used in the United States and Canada, is a DS1C level intermediate between the first and second levels. Fig. 11 gives the corresponding hierarchy for the basic E1 rate of 2048 kb/s as used in Europe. The figures diagram the levels, multiplexes, digital line systems, and references to ITU-T documents that give more detailed specifications for the two generic hierarchies. These various hierarchy levels are not timed from a common clock; they are termed the Plesiochronous Digital Hierarchy.

In North America most local interoffice trunks, and virtually all toll connecting trunks, are digital due to the synergy between digital transmission and switching. At digital-signal levels 1 and 2, the bit stream may be carried on wire pairs with regenerators, or via channels on a higher-level multiplex system using fiber or radio as its base facility. Radio may be used for higher bit rates between 45 and 155 Mb/s, and optical-fiber media up to 40 Gb/s. Rapid advances in technology, such as low-loss/low-dispersion single-mode fiber, have led to the introduction of additional bit rates into the digital hierarchy. The Synchronous Optical Network (SONET) [or, in ITU-T terms, Synchronous Digital Hierarchy] involves ten unique transmission speeds for use on fiber or suitable digital microwave systems. These are synchronous multiples of 51.840 Mb/s, as listed in Table 3. SONET multiplexers accept and deliver DS1- and DS3-level signals, on either a

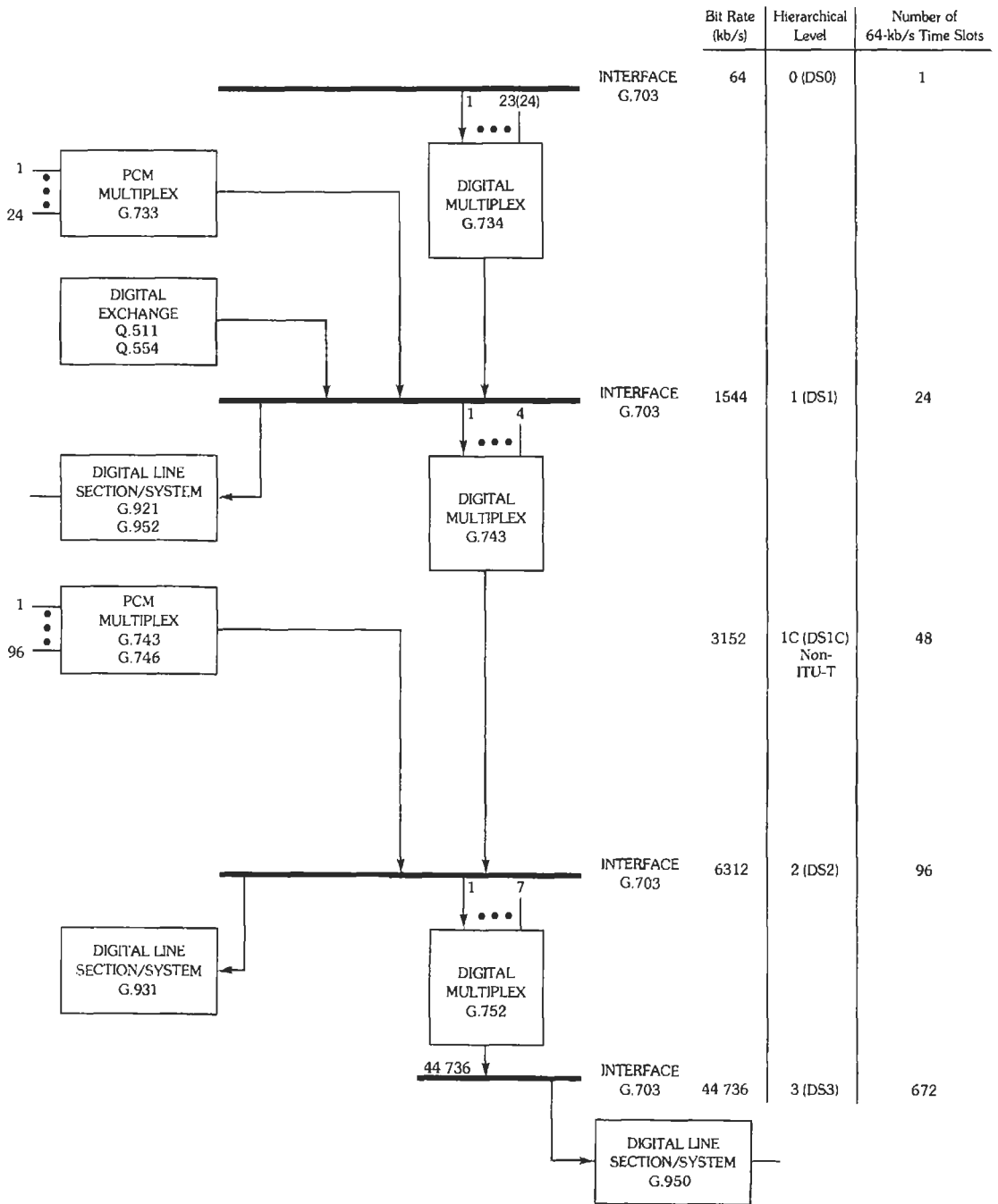


Fig. 10. Hierarchical bit rates for digital networks with the hierarchy based on a first-level bit rate of 1544 kb/s, with references to ITU-T recommendations.

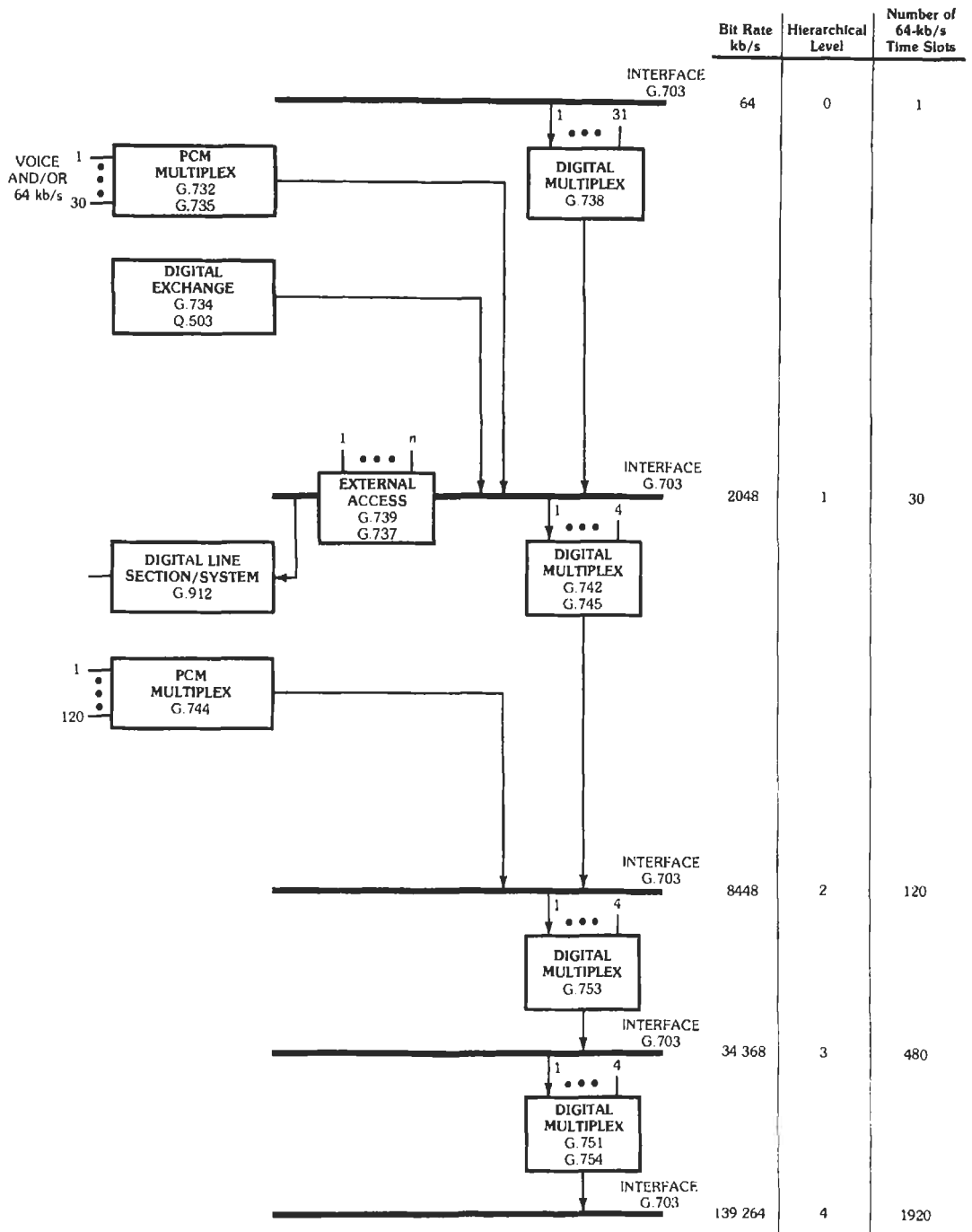


Fig. 11. Hierarchical bit rates for digital networks with the hierarchy based on a first-level bit rate of 2048 kb/s, with references to ITU-T recommendations.

TABLE 3. SONET/SDH TRANSMISSION RATES

Optical Carrier	STS Level	STM* Level	Line Rate (Mb/s)
OC-1	STS-1	N/A	51.840
OC-3	STS-3	STM-1	155.520
OC-9	STS-9	STM-3	466.560
OC-12	STS-12	STM-4	622.080
OC-18	STS-18	STM-6	933.120
OC-24	STS-24	STM-8	1244.160
OC-36	STS-36	STM-12	1866.240
OC-48	STS-48	STM-16	2488.320
OC-192	STS-192	STM-64	9953.280
OC-768	STS-768	STM-256	39813.120

* Synchronous Transfer Mode.

terminal or a drop-and-add basis, and produce optical signals at the OC-1 through OC-768 levels. Pairs of channels may be "concatenated" to be passed together through multiplexers, making a pair of OC-3s, for example, into an OC-3c.

NETWORK-CONTROL SIGNALING

Subscriber-Loop Signaling

Supervision and Ringing—In conventional (non-ISDN) use, subscriber-loop signaling uses direct current in the loop to provide supervisory signals. These (on-hook: loop open; off-hook: loop closed) are used to detect when the calling subscriber is requesting service and when the called subscriber answers. A ringing signal for summoning the called subscriber to the telephone or activating other terminal equipment is commonly transmitted as a high-voltage, low-frequency signal (at least 40 V rms, 20 Hz, at the subscriber terminal) to actuate a bell or electronic ringer.

Address Signaling—Dual-Tone Multifrequency Signaling (DTMF)—A tone signaling technique is in wide use for subscriber lines. Its advantages are speed of service, reduced holding time of registers in the switching system, and ease of end-to-end signaling for such applications as voicemail or radio paging. At the subscriber set, the keypad, when pressed by the caller, sends pairs of audio-frequency tones to the central office, each pair representing a numeric digit. Eight frequencies are provided, with the six combinations beyond the 10 required for numerics being used for special signals. The DTMF signaling code is given in Table 4.

Address Signaling—Dial Pulsing—In this system, numerics are transmitted as dial pulses obtained by opening and closing the loop at a rate of 8 to 12 pulses

per second. The number of pulses in a train represents the dialed digit (10 pulses for the digit 0). This system originated with the rotary subscriber dial, but commonly uses pulses generated electronically.

ISDN—Subscriber-loop signaling in ISDN applications employs the D-channel to handle data packets for call-control functions.

Interoffice Signaling

Common-Channel Signaling—Trunks between switching systems (whether exchange or toll) may use either signaling on a per-trunk basis or common-channel signaling (CCS). The principle of CCS, introduced in an early form in the mid-1970s, is to transmit all signaling and address information pertaining to a group of trunks over a dedicated data network. This is particularly attractive with electronic switching systems because the signaling information is a stream of data from the stored-program controller in the first place. Advantages include speed, flexibility to provide processing-intensive services, and low cost. By removing the need to signal by time-sharing any transmission bits on digital trunks (see below), CCS makes it possible to offer "clear-channel capability." This capability suits applications such as ISDN, where the user may send any digital bit pattern without constraint, including long strings of zeroes. Fig. 12 shows the signaling network used for CCS, in which individual switches have redundant pairs of data links to Signal Transfer Points (STPs). The pairs of STPs are duplicated in a quad configuration. For access to routing information for advanced (intelligent-network) services like number portability between carriers, other data links connect to Service Control Points (SCPs) giving access to data bases.

DC Signaling—Short voice-frequency trunks, consisting primarily of wire facilities, may use dc signaling for both supervisory and numeric signals.

Loop Reverse-Battery Signaling—Two-wire one-way trunks do not require full-duplex signaling. Loop reverse-battery signaling has frequently been used on them. At the terminating end of the trunk, the switching-system trunk circuit (a signaling interface) provides central-office battery, nominally 48 volts, over the facility to the originating end. The originating end applies a loop closure to seize the circuit and sends address information toward the terminating end in the form of dial pulses. The terminating trunk circuit provides supervision toward the originating end by reversing the polarity of battery applied to the trunk in an off-hook condition.

E and M Lead Signaling—E and M signaling is symmetric, and is thus suited to two-way trunks. With it, the signaling paths are separated from the trunk circuit and reside between the trunk circuit and the line. Signaling between a trunk circuit and the signal chan-

TABLE 4. DTMF AND MF SIGNALING CODES

Signal	Frequencies in Hz													
	Dual-Tone Multifrequency						Multifrequency							
	697	770	852	941	1209	1336	1477	1633*	700	900	1100	1300	1500	1700
Start Pulsing											■			
End Pulsing														
0				■		■						■	■	■
1	■				■				■	■				
2	■					■			■		■			
3	■						■			■	■			
4		■			■				■			■		
5		■				■				■		■		
6		■					■				■	■		
7			■		■				■				■	
8			■			■				■			■	
9			■				■				■		■	
*				■	■						■			
#				■			■							
a	■							■						
b		■							■					
c			■							■				
d				■							■			

* 1633 Hz is used with the other frequencies for special-category signals.

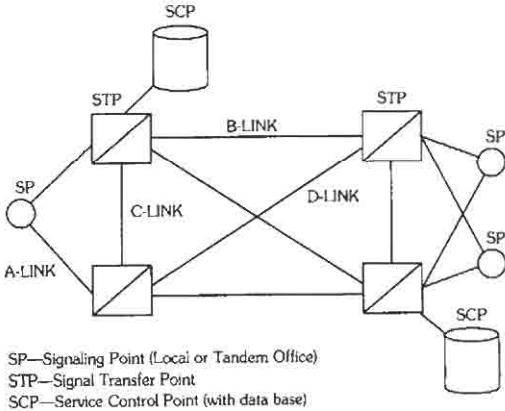


Fig. 12. Basic CCS network.

nel of a transmission facility occurs over two independent pairs of leads: the M lead and its return wire, which send signals outgoing from the trunk circuit; and the E lead and its return wire, which receive signals incoming to the trunk circuit. "E and M" historically stems from designations of the interconnecting leads on signaling-circuit drawings.

Signaling on Carrier-Derived Channels—Where CCS is not used, circuits without dc continuity use techniques whereby signaling information is transmitted either by designated bits in the pulse stream of a

digital carrier system or by audio frequencies within the voice channel (in-band signaling).

A- and B-Bit Signaling—In digital carrier systems, including DLC, supervisory and address information may be assigned to specific bits in the pulse stream. These bit positions are derived by time-sharing the least significant speech bit during one frame out of six, obtaining a 1.333-kb/s signaling bit stream. This stream is subdivided into A and B bits for multistate signaling, or may be further subdivided into A, B, C, and D bits.

Single-Frequency (SF) Signaling—Single-frequency in-band signaling is now rare; its use is effectively limited to frequency-division multiplex carrier systems and special private-line applications. It uses a tone frequency of 2600 hertz. The tone is transmitted at a level of -8 dBm0 during pulsing, and -20 dBm0 for continuous tone in the on-hook state.

Multifrequency (MF) Signaling—For speed in setting up interoffice connections in a non-CCS environment, multifrequency signaling may be applied on trunk circuits. Digit information is sent by combinations of two of five audio tones: 700, 900, 1100, 1300, and 1500 hertz. A sixth tone, 1700 Hz, is used in combination with 1100 Hz as a "start [of] pulsing" (KP) signal, and in combination with 1500 Hz as an "end of pulsing" (ST) signal. Receipt of the ST signal at the end of the digit signals starts call processing. Table 4 gives the MF signaling code. The three tone pairs not listed are used for specialized purposes in operator-services systems

and special networks for "Enhanced 911" service. Each tone is customarily transmitted at a level of -6 dBm0. Tone pairs are sent at 7 to 10 pulses per second and the KP signal is nominally 90–120 ms in duration.

Signaling Systems—ITU-T—For ITU-T Recommendations, refer to Chapter 2.

REFERENCES

1. *American National Standard for Telecommunications, Integrated Services Digital Network (ISDN)—Basic Access Interface for S and T Reference Points (Layer 1 Specification)*, ANSI T1.605-1991, New York: American National Standards Institute, 1991.
2. *Asymmetric Digital Subscriber Line (ADSL) Transceivers*, ITU-T Recommendation G.992.1, June 1999.
3. Bennett, A. F. "An Improved Circuit for the Telephone Set," *Bell System Tech. J.*, Vol. 32, May 1953, pp. 611–626.
4. Bohn, P. P., et al., "Fiber in the Loop," *AT&T Tech. J.*, Vol. 71, No. 1 (Jan.–Feb. 1991), pp. 31–43.
5. *High Bit Rate Digital Subscriber Line (HDSL) Transceivers*, ITU-T Recommendation G.991.1, Oct. 1998.
6. *IEEE Standard Equipment Requirements and Measurement Techniques for Analog Transmission Parameters for Telecommunications*, IEEE Standard 743-1995.
7. *IEEE Standard Telephone Loop Performance Characteristics. Telecommunications*, IEEE Standard 820-1984 (reaffirmed 1999).
8. *Objective Instrumentation for the Determination of Loudness Ratings*, ITU-T Recommendation P.56, 1988.
9. *Splitterless Asymmetric Digital Subscriber Line (ADSL) Transceivers*, ITU-T Recommendation G.992.2, June 1999.
10. *Telcordia Notes on the Networks*, Chapter 7 ("Transmission"), Chapter 12 ("Distribution Technology"), and Chapter 16 ("Exchange Access"), Special Report SR-2275, Issue 4, Oct. 2000.
11. *Telecommunications—Telephone Terminal Equipment—Performance and Compatibility Requirements for Telephone Sets with Loop Signaling* (ANSI/TIA/EIA-470-B-97). Washington, DC: Telecommunications Industry Association, 1997.
12. *Telecommunications—Telephone Terminal Equipment—Transmission Requirements for Digital Wireline Telephones* (ANSI/TIA/EIA-579-A-98). Washington, DC: Telecommunications Industry Association, 1998.
13. *Telecommunications Transmission Engineering*, Vol. 2—Facilities and Vol. 3—Networks and Services. Telcordia Technologies, Inc., 1990.
14. Tuffnell, W. F. "500-Type Telephone Set," *Bell Laboratories Record*, Vol. 29, Sept. 1951, pp. 414–418.

39

Switching Networks and Traffic Concepts

Revised by Amos E. Joel, Jr.

PART 1: COORDINATE SWITCHING NETWORKS

Definitions of Terms	39-2	
Space- and Time-Division Switching	39-2	
Properties of Coordinate Switching Networks		39-2
Single-Stage Coordinate Switches	39-2	
Blocking and Nonblocking Networks	39-4	
Nonblocking 3-Stage Networks		
Networks With Internal Blocking		
Time-Division Switching	39-6	

PART 2: TRAFFIC CONCEPTS

Definitions of Terms	39-10	
Traffic Measures	39-11	
Traffic Quantity		
Traffic Intensity		
Call Intensity		
Traffic Equations and Tables	39-11	
Grade of Service		
Choice of Formula		
Traffic Equations		
Delays		

PART 1; COORDINATE SWITCHING NETWORKS

DEFINITIONS OF TERMS

Concentration: The function associated with a switching network having fewer outlet than inlet terminals.

Coordinate switch: A rectangular array of crosspoints in which one side of the crosspoint is multiplied in rows and the other side in columns.

Crosspoint: A two-state switching device containing one or more elements that have a low transmission impedance in one state and a very high one in the other.

Expansion: The function associated with a switch or switching network having more outlet than inlet terminals.

Folded network: A network in which each terminal can serve as either inlet or outlet and which is capable of completing a path between any pair of inlet-outlets.

Full availability: Property of a switch or switching network capable of providing a path from every inlet terminal to every outlet terminal.

Internal blocking: The inability to interconnect an idle inlet to an idle outlet because all possible paths between them are already in use.

Nonblocking network: A network in which there is at all times at least one available path between any idle inlet and any idle outlet, regardless of the number of paths already occupied.

Nonfolded network: A network in which inlets and outlets are separate and which is capable of completing a path between any inlet and outlet.

Single-linkage array: The mesh or spread of interconnections between the stages of a switching network whereby every switch of one stage has one connection to every switch of the adjacent stage.

Space-division switching network: A switching network in which the transmission paths are physically distinct.

Switching network: That part of a switching system that establishes transmission paths between pairs of terminals.

Switching stage: Those switches in a switching network that have identical parallel functions.

Time-division switching network: A switching network in which the transmission paths are separated in time.

SPACE- AND TIME-DIVISION SWITCHING

Most of the switching systems of the world have utilized space-division analog switching. Recently,

combinations of space and time division have become almost universal in new switching systems, employing primarily pulse-code modulation on multiplexed lines. Because there is duality between time-division and space-division switching and the principles are more easily understood in terms of space division, space division will be used in describing switching-network principles. This will be followed by a description of how to map time-division into space-division networks.

PROPERTIES OF COORDINATE SWITCHING NETWORKS

The simplest coordinate switch has a number of "crosspoints" that can give an inlet to the switch a connection to an outlet. (In most applications, only one connection will be desired at a time, but this is not always a requirement.) It is convenient to give a number of inlets identical access to a number of outlets so that the network can be arranged in a single piece of apparatus (switch) as shown in Fig. 1A. Such a switch is known as a rectangular coordinate switch.

In electromechanical coordinate switches, the connecting devices may be individual contact-making relays, in which case the number of complete relays (coil and a set of contacts) required is the product of the numbers of inlets and outlets. Alternatively, the whole crosspoint array may be provided by a crossbar switch in which a single relay coil is associated with each row and column of the switch, and the concurrent energizing of a row coil and column coil closes an individual set of contacts. In electronic switches, the crosspoints may be solid-state elements, usually transistors.

The rectangular switch discussed so far is nonblocking—any idle inlet can reach an idle outlet. For small numbers of inlets and outlets, this is an efficient arrangement, but when the inlets and outlets are more than twenty to thirty, the number of crosspoints can be reduced significantly by replacing the single switch by a number of interconnected smaller coordinate switches, arranged in a multistage coordinate switching network in one of many possible ways, as will be described. If blocking is acceptable, the number of crosspoints can be reduced even more.

SINGLE-STAGE COORDINATE SWITCHES

Fig. 1A shows a rectangular-coordinate switch interconnecting inlets from N lines and outlets to M other lines. When interconnection is possible at every crosspoint, the switch provides full availability and is said to be nonblocking. The particular switch shown in Fig. 1A acts as a concentrating switch so that, although the switch is nonblocking, an idle inlet may be blocked because there is no idle outlet. Such concentration switches are often used to bring higher traffic loads per circuit to the subsequent stages of a multistage system in order to use switches more efficiently. A rectangular switch with full availability requires NM crosspoints.

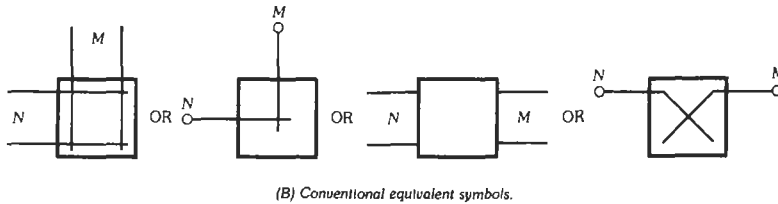
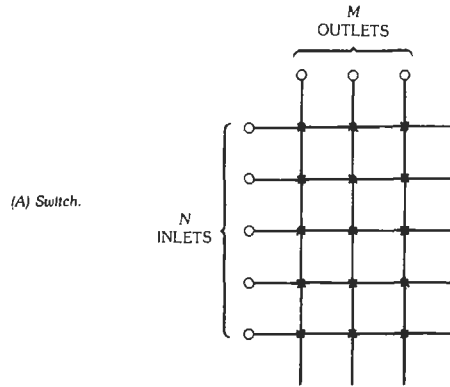


Fig. 1. Single-stage rectangular coordinate switch (full availability).

If full availability is not a requirement, an economy of crosspoints can be achieved by use of a limited-availability (restricted-access) coordinate switch. Fig. 2 shows a square switch with five inlets and five outlets in which every inlet has access to only three of the outlets. The grading is said to be homogeneous when, as shown, each set of three outlets is unique.

When the outlets from a switch are connected to the same lines or trunks as are the inlets to the switch, a triangular “folded” arrangement may be used. Fig. 3 shows a triangular switch for N inlets-outlets. The switch provides full availability and is nonblocking.

$$\text{Number of crosspoints required} = N(N - 1)/2$$

Fig. 4 shows a symmetrical two-stage network with square switches. This formed the basic frame arrangement of many electromechanical switching systems in which the limited size of the crossbar switch was compensated for by introducing two stages of switching—in the decimal planned system, expanding the access of a ten-by-ten crossbar switch to 100. This increase is achieved at the cost of introducing blocking; that is, an idle inlet can no longer reach every idle outlet regardless of the number of calls in progress. In fact, it is easy to see that every call from an inlet on primary switch 1 to an outlet on secondary switch 1 blocks all possible calls from other inlets on primary switch 1 to other outlets on secondary switch 1. Such an arrangement would usually give intolerable service, so three or more stages are generally used to provide parallel paths between inlet and outlet.

Fig. 5 shows a three-stage network where now there are several possible paths between any inlet and any outlet. The ability to search for one of several paths reduces blocking materially below the single path of the primary-secondary arrangement previously discussed.

It should be noted that while the number of crosspoints of a particular network is a first-order approach to its cost, there are other factors such as control cost, path setup time, and physical design that may make the most favorable design for a particular application quite different from the one that minimizes crosspoints.

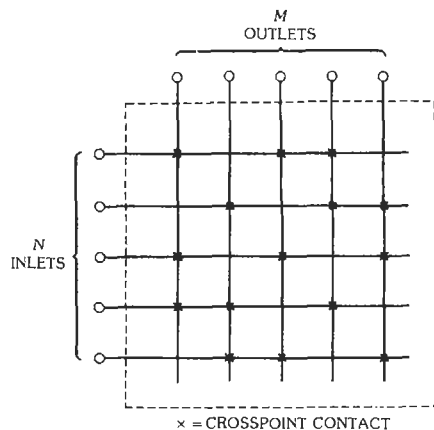


Fig. 2. Single-stage square coordinate switch (limited availability).

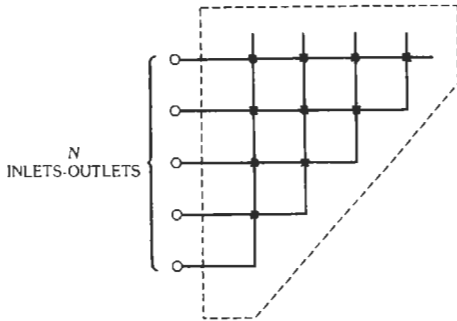


Fig. 3. Single-stage folded (triangular) switch.

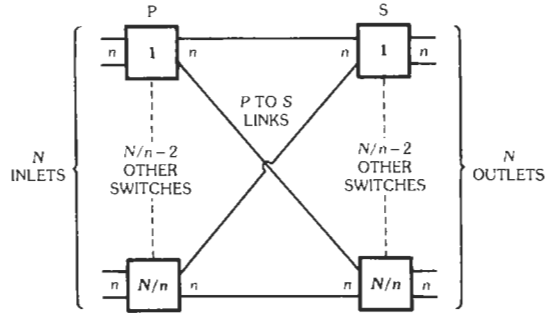


Fig. 4. Symmetrical two-stage switch.

BLOCKING AND NON-BLOCKING NETWORKS

In the following paragraphs, configurations for blocking and nonblocking networks will be described. The choice of network depends to a large extent on the service objectives, cost objectives, control capability, and, where a network covers a wide range of quantities of inlets and outlets, how additional switches and terminals can be added to the network without expensive rearrangement and changes in control methods.

Nonblocking 3-Stage Networks

A. General Case—Fig. 6 shows part of a network that presents no internal blocking (see definition). Let N = number of inlets, M = number of outlets, S = number of center-stage switches, n = number of inlets per first stage switch, and m = outlets per final-stage switch. When $N \neq M$, the condition for nonblocking is given by $S = n + m - 1$.

The number of crosspoints required is given by

$$X = (n + m - 1)[N + M + (NM/nm)]$$

A minimum number of crosspoints is obtained when $m = n$ and when n satisfies the equation

$$NM/(N + M) = n^3/(n - 1) \quad (\text{Eq. 1})$$

Columns 1 and 2 of Table 1 show corresponding values of n and $NM/(N + M)$ that satisfy equation 1, from which the optimum value of n for given values of N and M may be selected. Fig. 7 shows part of the network.

The number of crosspoints required is then given by

$$X_{\min} = (2n - 1)[N + M + (NM/n^2)]$$

B. Symmetrical Case—Fig. 8 shows part of a symmetrical nonblocking network, where $M = N$ and $m = n$. The condition for nonblocking is given by

$$S = 2n - 1$$

The number of crosspoints required is given by

$$X = N(2n - 1)[2 + (N/n^2)]$$

A minimum number of crosspoints is obtained when n satisfies the equation

$$N = 2n^3/(n - 1) \quad (\text{Eq. 2})$$

Columns 1 and 3 of Table 1 show corresponding values of n and N satisfying equation 2, from which the optimum value of n may be selected.

For large values of n , beyond the range of practical three-stage networks (see C below)

$$2n^2 \rightarrow N$$

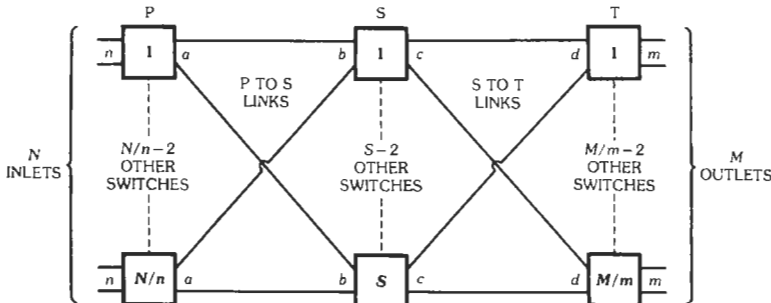


Fig. 5. Three-stage network—general case.

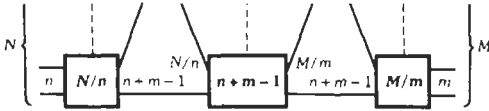


Fig. 6. Nonblocking three-stage network.

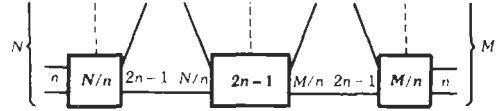


Fig. 7. Nonblocking three-stage network with minimum crosspoints.

Columns 1 and 4 of Table 1 show corresponding values of n and N satisfying the equation

$$N = 2n^2 \tag{Eq. 3}$$

The optimum value of n selected by using column 3 is either equal to or one less than the value indicated by column 4.

C. Comparison of 3-Stage Network and Single-Stage Switch—A single-stage, full-access switch is inherently nonblocking. The most favorable nonblocking three-stage network requires fewer than NM crosspoints if

$$NM/(N + M) > n^2(2n - 1)/(n - 1)^2$$

Table 2 shows, for some practical values of n , the limiting value of $NM/(N + M)$ below which a single-stage switch requires fewer crosspoints.

Table 3 compares single-stage switches and three-stage symmetrical nonblocking networks for typical values of N and n , where N/n is integral (see **D** below), to illustrate the trends of design choices.

D. Practical Nonblocking 3-Stage Network with Minimum Crosspoints—When a nonblocking network with a minimum number of crosspoints is sought, the indicated optimum value of inlets, n , per primary-stage switch may be such that N/n is not integral. The desired result may be achieved by providing some of the primary-stage and tertiary-stage switches with $(n - 1)$ inlets and outlets, respectively, by adjusting the sizes of the secondary-stage switches, and by superimposing two sets of interstage linkages.

The method is illustrated in Fig. 9, a nonblocking network for 100 inlets and 100 outlets requiring a minimum number of crosspoints. The value $n = 6$ is selected from column 3 of Table 1. The nearest multiple of 6 that exceeds 100 is 102. Thus, 17 primary-stage switches are required, 15 with 6 inlets each and 2 with 5 inlets. The larger switches require $2n - 1 = 11$ outlets each, cross-linked to 11 secondary-stage switches (links shown by continuous lines). The smaller switches require $2n - 2 = 10$ outlets each, cross-linked to 10 only of the secondary-stage switches (links shown by dashed lines). Thus, the secondary-stage switches also are of two sizes. The total number of crosspoints is 5291 (as compared with 5423 for a nonblocking network for 102 inlets and outlets, with $n = 6$). The crosspoint saving may not be worth the added control complexity.

E. Extension to 5-Stage and 7-Stage Networks—If the number of inlets and outlets on each secondary-stage switch of a three-stage network is large, it is advantageous to use a five-stage network. One possible arrangement is shown in Fig. 10, in which each secondary switch of a symmetrical three-stage network conforming to Fig. 8 is expanded into a nonblocking three-stage subnetwork.

If $N > 160$, a five-stage nonblocking network can be designed—by judicious selection of parameters n and a —which requires fewer crosspoints than the most favorable nonblocking three-stage network. The advantage increases slowly with N ; at $N = 240$, the advantage is less than 5%.

The number of crosspoints is given by

$$X = N(2n - 1)\{2 + (2a - 1)[(2/n) + (N/n^2a^2)]\}$$

TABLE 1. RELATIONSHIP OF n , N , AND M IN NONBLOCKING 3-STAGE NETWORKS

n	$NM/(N+M)$ (Eq. 1)	N (Eq. 2)	N (Eq. 3)
2	8.0	16*	8
3	13.5	27*	18
4	21.3	42.7	32
5	32.2	62.5	50
6	43.2	86.4	72
7	57.2	114.3	98
8	73.1	146.3	128
9	91.1	182.2	162
10	111.1	222.2	200
11	133.1	266.2	242
12	157.1	314.2	288

*The only two integral solutions of equation 2.

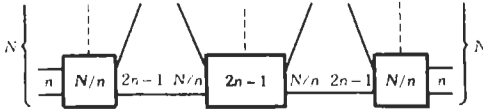


Fig. 8. Symmetrical nonblocking network with minimum crosspoints.

A minimum number of crosspoints is obtained when n and a satisfy the equations

$$N = 2na^3 / (a - 1) \\ = [na^2(2n^2 + 2a - 1)] / [(2a - 1)(n - 1)]$$

When the number of inlets and outlets, N , is such that the switches of the center stage of a five-stage network become large, each center-stage switch can with advantage be expanded in like manner into a nonblocking three-stage assembly. The result is a seven-stage network that is nonblocking overall.

F. "Folded" Nonblocking 3-Stage Networks— Fig. 11 shows the case in which the outlets from the tertiary-stage switches of a symmetrical three-stage network are connected to the same lines or trunks as are the inlets to the primary-stage switches. The number of crosspoints required to ensure no internal blocking is significantly less than in the case described in **B** above, provided the inlet-outlets are able to handle both incoming and outgoing traffic and the network control can set up paths in either direction. Such a network is said to function as a "folded" network.

The condition for nonblocking is given by

$$S = n$$

The number of crosspoints required is given by

$$X = N[2n + (N/n)]$$

A minimum number of crosspoints is obtained when n satisfies $2n^2 = N$. Corresponding values of n and N may be selected from columns 1 and 4 of Table 1.

As in the previous networks, when the network becomes large, further crosspoint savings can be achieved by substituting three-stage nonblocking networks for large rectangular switches.

G. Source Reference—Further information for the design of nonblocking multistage coordinate switching networks may be obtained from C. Clos, "A Study of Non-Blocking Switching Networks," *Bell System Tech. J.*, March 1953.

Networks With Internal Blocking

General Case—When internal blocking is allowed, it is possible to design a multistage network that requires many fewer crosspoints while providing full availability between the inlets and the outlets—but introducing a measure of blocking. The number of secondary-stage switches, S , required in a three-stage network to interconnect N inlets and M outlets, with a selected value of n , may be represented in a practical case as shown in Fig. 12 by $S = kn$, where k (constant) $< (2n - 1)/n$ determines the blocking.

The number of crosspoints required is given by

$$X = kn[N + M + NM/n^2]$$

A minimum number of crosspoints is obtained if n satisfies

$$n^2 = NM / (N + M)$$

The number of crosspoints is then given by

$$X_{min} = 2k[NM(N + M)]^{1/2}$$

When the inlets and outlets of a symmetrical network are connected to the same lines or trunks, it is advantageous to arrange the network to function as a "folded" network (see **F** above).

Traffic-Carrying Capability of Blocking Networks—Exact calculation of the probability of blocking in a multistage network is difficult. Order-of-magnitude estimates can be made by using the method described by C. Y. Lee, in which the number of parallel paths, r , possible for a connection is used. If $a_1, a_2,$ and a_{s-1} designate the occupancy of links leaving stages 1, 2, and $s - 1$ of an s -stage network, the probability, P , of blocking is approximated by:

$$P = [1 - (1 - a_1)(1 - a_2) \dots (1 - a_{s-1})]^r$$

More detailed methods of calculating blocking in switching networks may be found in the publications of the Institute of Switching and Data Technics, University of Stuttgart, Stuttgart, Germany.

TIME-DIVISION SWITCHING

Time division with pulse-code modulation has for many years been used as a method of multiplexing a number of telephone conversations on a digital carrier. With the same methods, a great reduction in active

TABLE 2. LOWER LIMITS OF $NM/(N + M)$ FOR NONBLOCKING 3-STAGE NETWORKS

n	2	3	4	5	6
$NM/(N+M) >$	12	11.25	12.45	14.06	15.84

TABLE 3. COMPARISON OF TYPICAL 1-STAGE AND 3-STAGE NONBLOCKING COORDINATE NETWORKS
($M = N, m = n, N/n$ INTEGRAL)

N	1-Stage		3-Stage		N	1-Stage		3-Stage		N	1-Stage		3-Stage	
	X	n	X	n		X	n	X	n		X	n	X	n
8*	64	2	96	50*	2500	5	1800	100	10 000	5	5 400			
15	225	3	275	54	2916	6	2079	100	10 000	10	5 700			
16	256	2	288	55	3025	5	2079	105	11 025	7	5 655			
16	256	4	336	56	3136	4	2156	108	11 664	6	5 940			
18*	324	2	351	56	3136	7	2236	108	11 664	9	6 120			
18	324	3	360	60	3600	5	2376	110	12 100	5	6 336			
20	400	2	420	60	3600	6	2420	110	12 100	10	6 479			
20	400	4	455	64	4096	4	2688	120	14 400	6	7 040			
24	576	3	560	64	4096	8	2880	120	14 400	8	6 975			
24	576	4	588	70	4900	5	3024	128*	16 384	8	7 680			
25	625	5	675	70	4900	7	3120	130	16 900	10	8 151			
27	729	3	675	72*	5184	6	3168	140	19 600	7	8 840			
30	900	3	800	75	5625	5	3375	140	19 600	10	9 044			
30	900	5	864	80	6400	5	3744	144	20 736	8	9 180			
32*	1024	4	896	80	6400	8	3900	144	20 736	9	9 248			
35	1225	5	1071	81	6561	9	4131	150	22 500	10	9 975			
36	1296	4	1071	84	7056	6	4004	160	25 600	8	10 800			
36	1296	6	1188	84	7056	7	4056	160	25 600	10	10 944			
40	1600	4	1260	90	8100	6	4455	162*	26 244	9	11 016			
40	1600	5	1296	90	8100	9	4760	170	28 900	10	11 951			
45	2025	5	1539	91	8281	7	4563	180	32 400	9	12 920			
48	2304	4	1680	96	9216	6	4928	180	32 400	10	12 996			
48	2304	6	1760	96	9216	8	5040	190	36 100	10	14 089			
49	2401	7	1911	98*	9604	7	5096	200*	40 000	10	15 200			

Legend:
 N = Inputs, outputs per network
 n = Inputs per P-switch, outputs per T-switch
 X = Crosspoints required
 * = Optimum configuration (where $N = 2n^2$, except when $n = 3$).

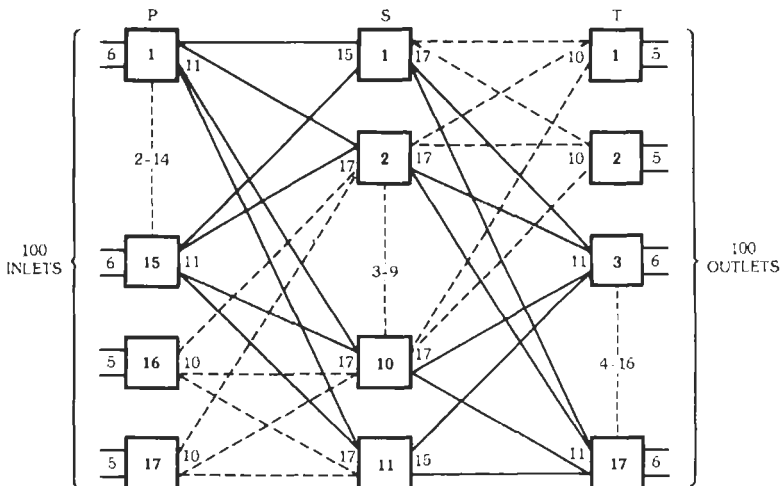


Fig. 9. Nonblocking three-stage network for 100 inlets, 100 outlets.

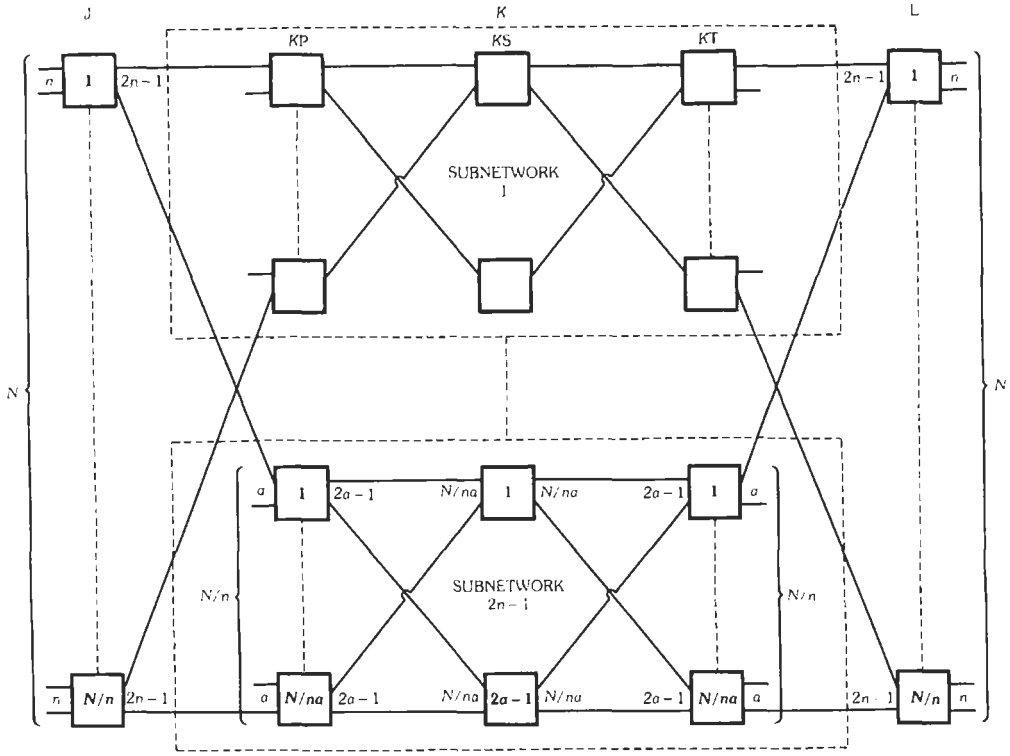


Fig. 10. Typical nonblocking five-stage network.

elements in a switching network can be achieved; in fact, it is advantageous to switch channels from digital carriers without demodulating them. The prevalence of digital transmission systems and the rapid advance of digital technology have made the present trend completely toward digital switching.

Conceptually, each time slot in a digital carrier can be thought of as a separate traffic path separated by the dimension of time rather than that of space as in the networks discussed so far. There are three usual time-division network configurations. In the first, the time

slots in one digital stream are inserted into similarly positioned time slots in an output stream. This forms the basis of the "space division" part of many time-division networks, as will be discussed.

In the second configuration, the information in one time-slot position of the input digital stream is inserted in a different position in the output digital stream (which may have a different bit rate).

Hybrid between the first two configurations is the taking of a sample of an analog inlet signal, encoding it and placing it in a selected time slot in an output digital

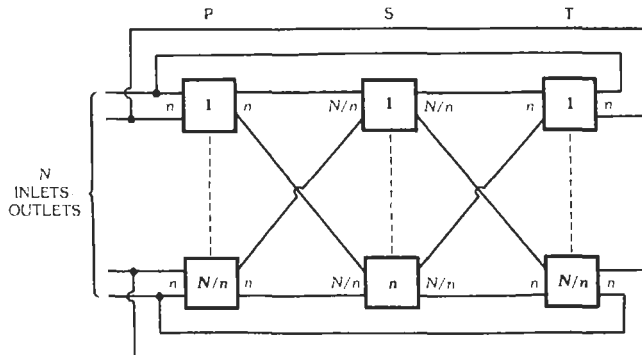


Fig. 11. Folded nonblocking three-stage network.

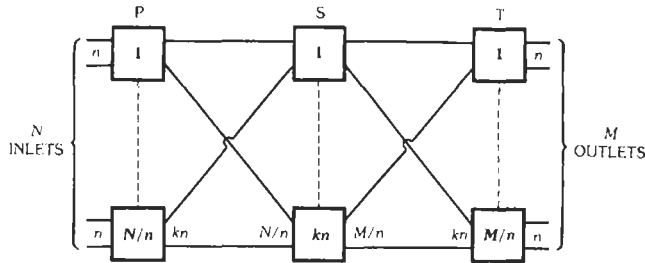


Fig. 12. Three-stage network with internal blocking.

stream, and the converse. This configuration is usual in the first stage of a time-division line-switching system.

The last two configurations require storage of signals—one store or register for each input or output time slot, whichever is larger. This type of time-division switching stage is known as *time slot interchange (TSI)*.

The first configuration is analogous to a set of crosspoints arranged to connect n inlets to n outlets but with each inlet having access to only one outlet. It is achieved, of course, with only one physical crosspoint. This crosspoint can be combined with others in a rectangular array to provide the equivalent of n rectangular switches in which all the same-positioned time slots of a number of inlet digital streams have access to like-positioned slots in a number of outlet streams.

The second and third configurations correspond to a rectangular switch in which any inlet, path, or time slot has access to any outlet.

Because digital transmission is generally unidirectional, two paths must be established, one for each direction of transmission. This may be done by pairing the transmitter and receiver for each inlet and outlet. A connection is set up through parallel networks from the transmitter of the inlet to the receiver of the outlet and from the receiver of the inlet to the transmitter of the outlet. Crosspoints in the parallel networks are controlled in common.

The above arrangement cannot be used on folded networks where the ability is needed to connect any inlet-outlet to any other, because connections would be made from transmitter to transmitter. It can, however, be implemented by placing all the transmitters on one side of a nonfolded network and all the receivers on the other. Independent paths are then set up for each direction of transmission.

All coordinate switching networks considered above are used in what has become known as *circuit switching*, where each path is assigned for the duration of a voice message or conversation. For data messages, a technique known as *packet switching* is used widely. Here data messages are usually divided into uniform-length segments. The address of the called station is prefixed onto each segment or packet as it is transmitted into an idle time slot.

For small networks, such as local area networks (LANs), the address is detected by the station to which

it is being sent, and the switching consists of removing the packet from the transmission medium. This is a distributed $1 \times N$ switching network.

As networks become larger or are interconnected, congestion may be encountered in finding idle time slots. In a packet switching network, the packets may be sent over any of several independent paths that might reach the called station. In such networks, each packet is assigned a number, and the called station checks that each is received and places them in order.

In emerging large data networks, packets of a given message are sent through the same network nodes over a virtual circuit. Buffer storage is used to hold the packets temporarily. If there is insufficient storage, packets may be lost.

Multistage coordinate switching networks are the basis for switching packets buffered for virtual circuit switching. Two-by-two matrices have the advantage that out of a possible 16 states only two are useful. Paths are sought through networks consisting of only 2×2 matrices, sometimes called *beta* elements, based upon the first binary digit (bit) address read from the beginning of each transmitted packet. When only 2×2 matrices are used, nonblocking cannot be achieved by the technique shown in Figs. 6 and 7. It can be shown that with $(2 \log_2 n) - 1$ stages of beta elements almost nonblocking can be achieved. The network divides into a mirror image as shown in Fig. 13. Half of the stages are known as a *baseline* network and the other half as the *inverse baseline*.

Appended to each packet as it leaves an input buffer are address bits used to control the state of a beta element at each successive stage of a switching network. A routing network, called a *Banyan* network (Fig. 14), allows the same address to be used to reach a particular output regardless of the input source.

It is possible that packets launched simultaneously from different inputs will collide or cause congestion at the output of any beta element. For this reason, steps must be taken to add intermediate buffers between network stages or to avoid the same routing for packets destined to the same output address.

Adding sorting stages to the network can prevent collisions within the network as packets move toward their destinations. Stages with this pattern, or *shuffle*, of link interconnection (Fig. 15) are known as *batcher* networks.

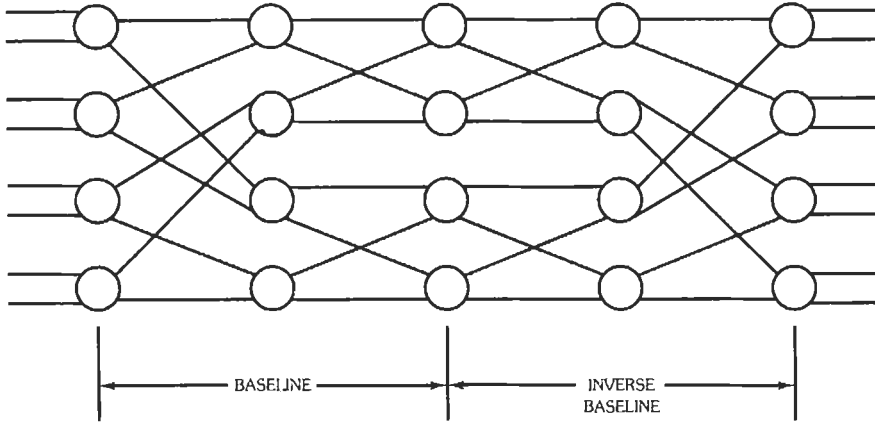


Fig. 13. Basic network of beta element.

Eventually, collisions will occur at the last stage. Various strategies to deal with such conflicts have been proposed. Most involve delaying all but one packet by additional buffers or feedback through delay stages.

The techniques described in this section are used in designing switching networks for serving a plurality of high-speed (fast) digital input lines carrying packets that arrive at random. (For this reason this technique has become known as *fast packet switching*.) Currently, there is a trend toward converting most digital messages, be they voice, images, video, or data, into fixed-sized packets, called *cells*. By the use of uniform-sized cells, each service may use its required portion of the capacity of the same switching network, and all services are switched ubiquitously.

Like packets, these cells arrive at the switch asynchronously or at random. When switching networks are used for this application, they are known as *asynchronous transfer mode* or *ATM* networks.

REFERENCES

Bellamy, J. *Digital Telephony*. New York: John Wiley & Sons, Inc., 1982.
 Hui, J. Y. *Switching and Traffic Theory for Integrated Broadband Networks*. Kluwer Academic Publishers, 1990.

Lee, C. Y. "Analysis of Switching Networks." *Bell Syst. T. J.*, Vol. 3, No. 6, Nov. 1955.
 Talley, D. *Basic Telephone Switching Systems*. Rochelle Park, N.J.: Hayden Book Co., Inc., 1979.

PART 2: TRAFFIC CONCEPTS

DEFINITIONS OF TERMS

Busy hour: The continuous one-hour period that, on consecutive days in the busy part of the year, contains the maximum average traffic intensity.

Call: A discrete engagement or occupation of a traffic path.

Calling rate: The average number of calls placed during the busy hour.

Call-second, call-minute, call-hour: Units of traffic quantity representing the occupation of a circuit for a second, minute, or hour.

Equated busy-hour call (EBHC): A European unit of traffic intensity equal to 1/30 of an erlang.

Erlang: The international dimensionless unit of traffic intensity. One erlang is the traffic intensity represent-

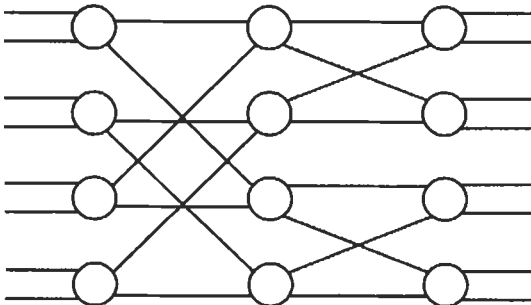


Fig. 14. Banyan routing network.

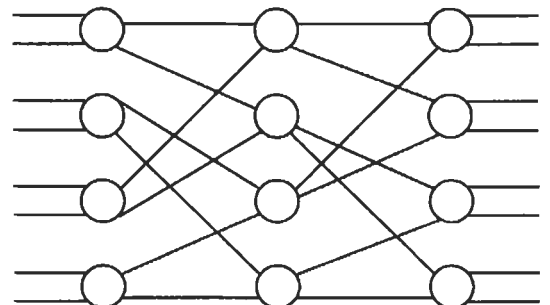


Fig. 15. Shuffle sorting network.

ed by an average of one circuit busy out of a group of circuits over some period of time.

Full-availability group: A group of traffic-carrying trunks or circuits in which every circuit is accessible to all the traffic sources.

Grade of service: A measure of the probability that, during a specified period of peak traffic, a call offered to a group of trunks or circuits will fail to find an idle circuit at the first attempt. Usually applied to the busy hour of traffic.

Holding time: The duration of occupancy of a traffic path by a call. Sometimes used to mean the average duration of occupancy of one or more paths by calls.

Hundred call-second per hour (CCS): A unit of traffic intensity equal to $\frac{1}{36}$ of an erlang.

Infinite sources: The assumption that the number of sources offering traffic to a group of trunks or circuits is large in comparison with the traffic load. A ratio of ten or higher is often considered "infinite."

Limited-access group: A group of traffic-carrying trunks or circuits in which only a fraction of the circuits is accessible to any one group of the traffic sources.

Limited sources: The assumption that the number of sources offering traffic to a group of trunks or circuits must be included in loss-probability calculations. Used in the binomial and Engset loss-probability equations.

Lost calls cleared: The assumption that calls not immediately satisfied at the first attempt are cleared from the system and do not reappear during the period under consideration. Used in the Erlang B loss-probability equation.

Lost calls delayed: The assumption that calls not immediately satisfied at the first attempt are held in the system until satisfied. Used in the Erlang C delay-probability equation.

Lost calls held: The assumption that calls not immediately satisfied at the first attempt are held in the system until served or abandoned. The sum of waiting and service time is assumed equal to what the service time would have been without delay. Used in the Poisson loss-probability equation. (When the service time is a negative exponential, this assumption has the same effect if holding and service times are independent and the distribution of the caller's holding time to abandon is a negative exponential with the mean equal to that of the service time.)

Occupancy: The traffic intensity per traffic path. One hundred percent occupancy implies all paths busy.

Traffic concentration: The average ratio of the traffic quantity during the busy hour to the traffic quantity during the day.

Traffic density: See Traffic intensity.

Traffic flow: See Traffic intensity.

Traffic intensity: The average number of calls present on a group of traffic paths over a period of time.

Traffic load: See Traffic intensity.

Traffic path: A channel, time slot, frequency band, line, trunk, switch, or circuit over which individual communications pass in sequence.

Traffic quantity: The aggregate engagement time or occupancy time of one or more traffic paths.

TRAFFIC MEASURES

In the terminology of traffic, there are three concepts that often cause confusion and need to be differentiated.

Traffic Quantity

Traffic quantity represents the total number of channel uses in units of time. It is of particular use for billing or other volume calculations but is not useful by itself. Where only this gross daily quantity is available, an estimate can be made of the busy-hour traffic by making an assumption about the proportion of daily traffic appearing in the busy hour. Direct busy-hour measurements are preferable, however, for calculating the traffic performance of a traffic-carrying facility. The latter is a function of traffic intensity, which is the quantity of traffic carried during a period of time.

Traffic Intensity

Traffic intensity, often called traffic load, is the average number of calls present during a period of time. It can be measured by averaging periodic counts of the number of calls present during the period or, alternatively, by summing the duration of calls within the period and dividing by the period length.

Call Intensity

For many traffic-carrying elements, the number of calls making up the total traffic load is immaterial; the load represented by two calls of ten minutes duration has the same impact as one call of twenty minutes duration. This is the case for switching-system processors and controllers, where the duration of calls on the traffic-carrying paths is irrelevant. Call intensity, or call rate, in calls per time unit, is frequently used, therefore, in estimating traffic loads on control equipment. Often the time unit is an hour and is implied, as in "a line calling rate of 3."

TRAFFIC EQUATIONS AND TABLES

Grade of Service

The overall "grade of service" of a switching system or trunk system refers to the ratio of calls that are not completed at the first attempt to the total number of

attempts to establish a connection through the system during a specific period of time, usually the busy hour. The traffic capacity of a switch, switching network, or trunk group is the traffic load that will, on the average, provide the grade of service that is chosen as the service objective.

For measurement purposes, it is sometimes easier to observe the proportion of time that all circuits are busy than to count the total of offered and carried calls (or carried and lost calls). For the usual assumptions of random calling, infinite sources, and a level value of offered traffic, the proportion of calls blocked and the proportion of time all circuits are busy are equal. This relation does not persist for assumptions of peaked traffic, smooth traffic, or limited sources.

Choice of Formula

The most important factors determining the choice of formula for estimating grade of service are as follows:

- A. The statistical nature of the call originations
 - Random from infinite sources
 - Random from finite sources
 - Peaked from infinite sources
 - Smooth from infinite sources
 - Periodic
- B. The probability distribution of the call holding time
- C. The availability of the circuits serving the traffic
 - Full access
 - Limited access in a graded multiple
 - Limited access through a switching network
- D. The behavior of calls when blocked
 - Lost calls cleared
 - Lost calls delayed
 - Lost calls held
 - Lost calls retried at a later time

Traffic Equations

The two most commonly used equations in trunk service estimation are the Erlang B and the Poisson. In general, the Erlang B is a good estimator of single-hour service, while the Poisson is a good estimator of the service given on the average over a busy season during which the offered load varies over a wider range than would be expected from the assumptions of random offered traffic. In general, the use of the Poisson equation is not recommended for grades of service that exceed 0.10 (10% blocking).

In small networks, it may be necessary to take limited sources into account. This results in the Engset and binomial equations as indicated below.

The Erlang B equation is given by

$$B(c,a) = (a^c/c!) / \sum_{x=0}^{x=c} a^x/x!$$

The corresponding Engset equation for limited sources

is solved only by iterative techniques. Tables are available for this purpose.*

The Poisson equation is given by

$$P(c,a) = e^{-a} \sum_{x=c}^{\infty} (a^x/x!)$$

The corresponding binomial equation is given by

$$P(c,n+1,r) = \sum_{x=c}^n \binom{n}{x} r^x(1-r)^{n-x}$$

Finally, the Erlang C equation is given by

$$C(c,a) = (a^c/c!) [c/(c-a)] /$$

$$\left\{ \sum_{x=0}^{c-1} a^x/x! + (a^c/c!) [c/(c-a)] \right\}$$

where,

- a* = offered load in erlangs,
- c* = number of circuits,
- n* = number of limited sources less one,
- r* = load offered by an idle source.

Although the above equations are derived on the assumption of a negative exponential distribution of holding times, they are reasonable approximations for most distributions encountered in practice.

The Erlang B equation is particularly useful in estimating the load that will be carried on and overflowing from a direct group of trunks to which random traffic is offered. Calls that find all direct circuits busy are directed to an alternate route. The traffic carried on the direct route will be $a[1 - B(c,a)]$, and the traffic overflowing will be $aB(c,a)$, where *a* = the random offered traffic. The overflow traffic is no longer random, and other approximations are necessary to estimate the grade of service given by an alternate trunk group to overflow traffic. For this purpose, it is necessary to estimate the variance of the overflow traffic, or more conveniently, the ratio of the variance to the mean. This latter quantity is known as the *peakedness* of the traffic. (Random traffic from an infinite source has a peakedness of one.) Peakedness, *z*, of traffic overflowing from *c* trunks offered *a* erlangs under the LCC assumptions is given by:

$$z = 1 - aB(c,a) + a/[c + 1 + aB(c,a) - a]$$

When a number of overflow loads from a number of direct groups are combined, the peakedness of the total is given closely by

$$z = \sum_{x=1}^n z_x a'_x / \sum_{x=1}^n a'_x$$

where *a'*_{*x*} is the overflow load from the *x*th source.

*Telephone Traffic Theory, Tables and Charts. Berlin-Munich: Siemens Aktiengesellschaft, 1970.

An approximation to the grade of service given by a group of trunks to peaked traffic may be made with the equation

$$P = B(c/z, y/z)$$

(See D. W. Hill and S. R. Neal, "Traffic Capacity of a Probability Engineered Trunk Group," *Bell Syst. T. J.*, Vol. 55, No. 7, September, 1976, for closer approximations.)

Delays

Many control systems operate on a delay basis, and it is convenient to have an equation for estimating de-

lays. Under the Erlang C assumptions, the delay has a mean of

$$d(c, a, h) = [c/(c - a)]C(c, a)h$$

If calls are served in order of arrival, the length of delay for delayed calls will be negative exponential.

Approximations to other service time distributions are given by

$$d(c, a, h) = [c/(c - a)]C(c, a)h(1 + V/h^2)/2$$

where,

V is the variance,

h is the average of the holding time distribution.

TABLE 4. OFFERED LOAD TO A TRAFFIC-CARRYING FACILITY FOR OBJECTIVE GRADE OF SERVICE

Trunks	Grade of Service											
	0.001		0.005		0.010		0.02		0.05		0.10	
	<i>B</i>	<i>P</i>	<i>B</i>	<i>P</i>	<i>B</i>	<i>P</i>	<i>B</i>	<i>P</i>	<i>B</i>	<i>P</i>	<i>B</i>	<i>P</i>
1	0.001	0.001	0.005	0.005	0.01	0.01	0.02	0.02	0.05	0.05	0.11	0.11
2	0.05	0.05	0.11	0.10	0.16	0.15	0.22	0.21	0.38	0.36	0.60	0.53
3	0.19	0.19	0.35	0.34	0.46	0.44	0.60	0.57	0.90	0.82	1.3	1.1
4	0.44	0.43	0.70	0.67	0.87	0.82	1.1	1.0	1.5	1.4	2.0	1.7
5	0.76	0.74	1.1	1.1	1.4	1.3	1.7	1.5	2.2	2.0	2.9	2.4
6	1.1	1.1	1.6	1.5	1.9	1.8	2.3	2.1	3.0	2.6	3.8	3.2
7	1.6	1.5	2.2	2.0	2.5	2.3	2.9	2.7	3.7	3.3	4.7	3.9
8	2.1	2.0	2.7	2.6	3.1	2.9	3.6	3.3	4.5	4.0	5.6	4.7
9	2.6	2.4	3.3	3.1	3.8	3.5	4.3	4.0	5.4	4.7	6.6	5.4
10	3.1	3.0	4.0	3.7	4.5	4.1	5.1	4.6	6.2	5.4	7.5	6.2
11	3.7	3.5	4.6	4.3	5.2	4.8	5.8	5.3	7.1	6.2	8.5	7.0
12	4.2	4.0	5.3	4.9	5.9	5.4	6.6	6.0	8.0	6.9	9.5	7.8
13	4.8	4.6	6.0	5.6	6.6	6.1	7.4	6.7	8.8	7.7	10.5	8.6
14	5.4	5.2	6.7	6.2	7.4	6.8	8.2	7.4	9.7	8.5	11.5	9.5
15	6.1	5.8	7.4	6.9	8.1	7.5	9.0	8.2	10.6	9.2	12.5	10.3
16	6.7	6.4	8.1	7.6	8.9	8.2	9.8	8.9	11.5	10.0	13.5	11.1
17	7.4	7.0	8.8	8.3	9.6	8.9	10.7	9.6	12.5	10.8	14.5	12.0
18	8.0	7.7	9.6	8.9	10.4	9.6	11.5	10.4	13.4	11.6	15.5	12.8
19	8.7	8.3	10.3	9.6	11.2	10.3	12.3	11.2	14.3	12.4	16.5	13.7
20	9.4	9.0	11.1	10.4	12.0	11.1	13.3	11.9	15.2	13.2	17.6	14.5
30	16.7	15.9	19.0	17.8	20.3	18.7	21.9	19.8	24.8	21.6	28.1	23.2
40	24.4	23.3	27.4	25.6	29.0	26.8	31.0	28.1	34.6	30.2	38.8	32.1
50	32.5	30.0	36.0	33.7	37.9	35.0	40.3	36.6	44.5	39.0	49.6	41.2
60	40.8	38.9	44.8	41.9	46.9	43.5	49.6	45.2	54.6	47.9	60.4	50.3
70	49.2	47.0	53.7	50.3	56.1	52.0	59.1	53.9	64.7	56.8	71.3	59.5
80	57.8	55.2	62.7	58.8	65.4	60.7	68.7	62.7	74.8	65.9	82.2	68.8
90	66.5	63.5	71.8	67.4	74.7	69.4	78.3	71.6	85.0	75.0	93.1	78.1
100	75.2	71.9	80.9	76.1	84.1	78.2	88.0	80.6	95.2	84.1	104.1	87.4

Legend:

All loads in erlangs

B = Lost calls cleared, Erlang B

P = Lost calls held, Poisson

For calls served in order of arrival, the negative exponential is an approximation to the delay distribution of delayed calls.

Table 4 gives values of loads for objective losses over a range of loads and trunks for the Poisson and Erlang B formulas. For computer use, advantage can be taken of the following recurrence relation for computing the infinite source loss formulas:

$$\text{Let } E(c + 1, a) = (c/a)[E(c, a) + 1]$$

$$\text{and } E(0, a) = 1$$

Then for lost calls cleared

$$B(c, a) = 1/E(c, a)$$

and for lost calls delayed

$$C(c, a) = 1/\{(c - a)/c[E(c, a) - 1] + 1\}$$

For lost calls held:

$$\text{Let } D(c, a) = [(a/c) + 1]D(c - 1, a) - (a/c)D(c - 2, a)$$

$$\text{and } D(-1, a) = 0$$

$$D(0, a) = 1$$

Then

$$P(c, a) = 1 - e^{-a} D(c - 1, a)$$

REFERENCES

- Telephone Traffic Theory, Tables and Charts*. Berlin-Munich: Siemens Aktiengesellschaft, 1970. (In English and German.)
- Bear, D. *Principles in Telecommunications—Traffic Engineering*. Peter Peregrinus, Ltd., 1976 (Repr. 1980).
- Mina, R. *Introduction to Teletraffic Engineering*. Telephony Publishing Co., 1974.

40

Electroacoustics

*Revised and Expanded by
Paul D. Schomer and George W. Swenson, Jr.*

Theory of Sound Waves	40-3
Wave Equation	
Sound in Gases	
Sound in Liquids	
Sound in Solids	
Sound Intensity	
Sound-Pressure Level (SPL)	
Electrical Analogies for Acoustical and/or Mechanical Systems	40-5
Hearing	40-6
Measurement of Sound	40-8
Sound Level Meters	
Measurement and/or Laboratory Microphones and Their Calibration	
Sound Level	
Sound Power Level	
Definitions	40-9
Environmental Noise	40-10
Assessment of Workplace-Industrial Noise	40-12
Sound in Enclosed Rooms	40-12
Standing Sound Waves	
Reverberation Time	
Measurement of Reverberation Time	
Calculation of Reverberation Time	
Absorption Coefficients	
Optimum Reverberation Time	

Indoor Design Noise Levels and Speech Communication	<i>40-15</i>
Preferred Speech Interference Levels/Speech Interference Levels	
Interior Design Noise Levels	
Sound Isolation	<i>40-17</i>
Walls	
Doors	
Resilient Channels	
Weighting of Components	
Flanking	
Leaks	
Sound-Reinforcement Systems	<i>40-23</i>

THEORY OF SOUND WAVES*

Sound (or a sound wave) is an alteration in pressure, stress, particle displacement, or particle velocity that is propagated in an elastic material, or the superposition of such propagated alterations. Sound (or sound sensation) is also the sensation produced through the ear by the above alterations.

Wave Equation

The behavior of small-amplitude sound waves is given by the wave equation

$$\nabla^2 p = (1/c^2)(\partial^2 p / \partial t^2) \quad (\text{Eq. 1})$$

where,

- p is the instantaneous pressure increment above and below a steady pressure (pascals),
- t is the time in seconds,
- c is the velocity of propagation in meters/second,
- ∇^2 is the Laplacian.

The quantity p is a function of time and of the three coordinates of space. For the particular case of rectangular coordinates x , y , and z (in meters), the Laplacian is given by

$$\nabla^2 = (\partial^2 / \partial x^2) + (\partial^2 / \partial y^2) + (\partial^2 / \partial z^2) \quad (\text{Eq. 2})$$

Plane Waves—For a plane wave of sound, where variations with respect to y and z are zero, $\nabla^2 p = \partial^2 p / \partial x^2 = d^2 p / dx^2$; the latter is approximately equal to the curvature of the plot of p versus x at some instant. Equation 1 states simply that, for variations in x only, the acceleration in pressure p (which is the second time derivative of p) is proportional to the curvature in p (which is the second space derivative of p).

Sinusoidal variations in time are usually of interest. For this case, the standard procedure is to put $p = (\text{real part of } \bar{p}e^{j\omega t})$, where the phasor \bar{p} now satisfies the equation

$$\nabla^2 \bar{p} + (\omega/c)^2 \bar{p} = 0 \quad (\text{Eq. 3})$$

The velocity phasor \bar{v} of the sound wave in the medium is related to the complex pressure phasor \bar{p} by

$$\bar{v} = -(1/j\omega\rho_0) \text{ grad } \bar{p} \quad (\text{Eq. 4})$$

The specific acoustic impedance \bar{Z} at any point in the medium is the ratio of the pressure phasor to the velocity phasor, or

$$\bar{Z} = \bar{p} / \bar{v} \quad (\text{Eq. 5})$$

Spherical Waves—The solutions of equations 1 and 3 take particularly simple and instructive forms for the case of one-dimensional plane and spherical waves in one direction. Table 1 summarizes the pertinent information.

For example, the acoustic impedance for spherical waves has an equivalent electrical circuit comprising a resistance shunted by an inductance. In this form, it is obvious that a small spherical source (r is small) cannot radiate efficiently since the radiation resistance $\rho_0 c$ is shunted by a small inductance $\rho_0 r$. Efficient radiation begins approximately at the frequency where the resistance $\rho_0 c$ equals the inductive (mass) reactance $\rho_0 r$. This is the frequency at which the period ($1/f$) equals the time required for the sound wave to travel the distance $2\pi r$.

Sound in Gases

The acoustic behavior of a medium is determined by its physical characteristics and, in the case of gases, by the density, pressure, temperature, specific heat, coefficients of viscosity, and the amount of heat exchange at the boundary surfaces.

The velocity of propagation in a gas is a function of the equation of state ($PV = RT$ plus higher-order terms), the molecular weight, and the specific heat.*

For small displacements relative to the wavelength of sound, the velocity is given by

$$c = (\gamma p_0 / \rho_0)^{1/2} \quad (\text{Eq. 6})$$

where,

- γ = ratio of the specific heat at constant pressure to that at constant volume,
- p_0 = the steady pressure of the gas in newtons/meter²,
- ρ_0 = the steady or average density of the gas in kilograms/meter³.

The values of the velocity in a few gases are given in Table 2 for 0 degrees Celsius and 760 millimeters of mercury barometric pressure.

The velocity of sound, c , in dry air is given by the experimentally verified equation

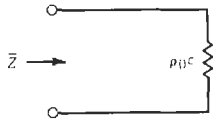
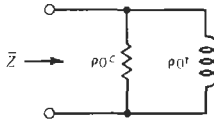
$$c = 331.45 \pm 0.05 \text{ meters/second} \\ = 1087.42 \pm 0.16 \text{ feet/second}$$

for the audible-frequency range, at 0 degrees Celsius and 760 millimeters of mercury with 0.03-mole-percent content of CO₂.

* Hardy, H. C., Telfair, D., and Piclemeier, W. H. "The Velocity of Sound in Air." *Journal of the Acoustical Society of America*, Vol. 13, January 1942, pp. 226-233. See also Beranek, L. *Acoustic Measurements*. New York: John Wiley & Sons, Inc., 1949, p. 46.

* Rayleigh, Lord. *Theory of Sound*. Vols. 1 and 2. New York: Dover Publications, 1945. Morse, P. M. *Vibration & Sound*. 2d ed. New York: McGraw-Hill Book Co., 1948.

TABLE 1. SOLUTIONS FOR VARIOUS PARAMETERS

Factor	Type of Sound Wave	
	Plane Wave	Spherical Wave
Equation for p	$\partial^2 p / \partial x^2 = (1/c^2)(\partial^2 p / \partial t^2)$	$(\partial^2 p / \partial x^2) + (2/r)(\partial p / \partial r) = (1/c^2)(\partial^2 p / \partial t^2)$
Equation for \bar{p}	$(d^2 \bar{p} / dx^2) + (\omega/c)^2 \bar{p} = 0$	$(d^2 \bar{p} / dx^2) + (2/r)(d\bar{p} / dr) + (\omega/c)^2 \bar{p} = 0$
Solution for p	$p = F[t - (x/c)]$	$p = (1/r) F[t - (x/c)]$
Solution for \bar{p}	$\bar{p} = \bar{P} \exp(-j\omega x/c + j\theta)$	$\bar{p} = (1/r) \bar{P} \exp(-j\omega r/c + j\theta)$
Solution for \bar{v}	$\bar{v} = (p/\rho_0 c) \exp(-j\omega x/c + j\theta)$	$\bar{v} = (p/\rho_0 c r) [1 + (c/j\omega r)] \exp(-j\omega r/c + j\theta)$
\bar{Z}	$\bar{Z} = \rho_0 c$	$\bar{Z} = \rho_0 c / [1 + (c/j\omega r)]$
Equivalent Electrical Circuit for \bar{Z}		
<p> p = excess pressure in pascals (newtons/meter²) \bar{p} = complex excess pressure in pascals (newtons/meter²) t = time in seconds x = space coordinate for plane wave in meters r = space coordinate for spherical wave in meters \bar{v} = complex velocity in meters/second \bar{Z} = specific acoustic impedance in rayls (newton-seconds/meter³) c = velocity of propagation in meters/second $\omega = 2\pi f$; f = frequency in hertz F = an arbitrary function θ = phase constant ρ_0 = density of medium in kilograms/meter³ \bar{P} = peak amplitude of the phasor (units dependent on wave type) </p>		

The velocity in air for a range of about 20 degrees Celsius change in temperature is given approximately by

$$c = 331.45 + 0.607T_c \text{ meters/second}$$

$$= 1052.03 + 1.106T_f \text{ feet/second}$$

where T_c is the temperature in degrees Celsius and T_f is the temperature in degrees Fahrenheit. For values of T_c greater than 20 degrees, the following equation may be used:

$$c = 331.45 \times (T_k/273)^{1/2} \text{ meters/second}$$

where T_k is the temperature in kelvins.

For other corrections, if extreme accuracy is desired, reference should be made to the literature.*

* Hardy, H. C., Telfair, D., and Pielemeier, W. H. "The Velocity of Sound in Air." *Journal of the Acoustical Society of America*, Vol. 13, January 1942, pp. 226-233.

From Equation 5 and Table 1, characteristic impedance is equal to the ratio of the sound pressure to the particle velocity.

$$\bar{Z} = \bar{p} / \bar{v} = \rho_0 c \cos \phi$$

For plane waves, $\phi = 0$ and $\cos \phi = 1$. For spherical waves, $\tan \phi = \lambda / 2\pi r$, where λ = wavelength of the acoustic wave, and r = distance from the sound source. For r greater than a few wavelengths, $\cos \phi \approx 1$.

Characteristic impedance $\rho_0 c$ in newton-seconds/meter³ for several gases at 0 degrees Celsius and 760 millimeters of mercury is given in Table 2.

Sound in Liquids

In liquids, the velocity of sound is given approximately by

$$c = (1/K\rho_0)^{1/2}$$

TABLE 2. GASES

Gas	Symbol	Density (g/L)	Velocity (m/s)	$\Delta v/\Delta t$ (m/s °C)	$\rho_0 c$ (N · s/m ³)
Air, dry		1.293	331.45	0.59	428.6
Argon	A	1.783	319	0.56	596
Carbon dioxide	CO ₂	1.977	259	0.4	511
Carbon monoxide	CO	1.25	338	0.6	421
Helium	He	0.178	965	0.8	173.2
Hydrogen	H ₂	0.0899	1284	2.2	114
Neon	Ne	0.900	435	0.8	383
Nitrogen	N ₂	1.251	334	0.6	418

where,

- c is the velocity in meters/second,
- K is the compressibility in meter-seconds²/kilogram and may be regarded as constant.

For most liquids,

$$K = (47 \times 10^{-8})/981$$

Figures for the velocity of sound in meters/second through some liquids are given in Table 3.

Sound in Solids

See Chapter 4, Tables 17, 18, and 19.

Sound Intensity

The sound intensity is the average rate of sound energy transmitted in a specified direction through a unit area normal to this direction at the point considered. Intensity is a vector quantity. In the case of a plane or spherical wave (several wavelengths from the source) in a fluid, the intensity in the direction of propagation is given by

$$I = p^2/\rho c \tag{Eq. 7}$$

The units of I are watts/meter².

Sound-Pressure Level (SPL)

It is sometimes convenient to use the decibel scale to express the ratio between any two sound intensities. Since sound pressure is usually proportional to the square root of the corresponding sound intensity, sound-pressure level (SPL) is defined as

$$SPL = 10 \log_{10} (p^2/p_0^2) \tag{Eq. 8}$$

where p_0 is 20 micropascals in air.*

At times, sound pressure is measured in units other than pascals. Table 4 lists the SPL for a quantity of one in these other units (e.g., 1 atmosphere).

ELECTRICAL ANALOGIES FOR ACOUSTICAL AND/OR MECHANICAL SYSTEMS**

Analysis of electromechanical-acoustical transducers, mechanical systems, or acoustical systems can be facilitated by the use of electrical analogies. Two forms of electrical analogies are possible: (1) force or pres-

* The reference of 1 micropascal is usually used for liquids and all other media.

** Swenson, George W. *Principles of Modern Acoustics*. New York: Van Nostrand, 1953.

TABLE 3. VELOCITY OF SOUND IN LIQUIDS

Liquid	Temperature in °C	Velocity in (m/s) × 10 ³
Alcohol, ethyl	12.5	1.24
	20	1.17
Mercury	20	1.45
Petroleum	15	1.33
Water, fresh	17	1.43
Water, sea (36 parts/thousand salinity)	15	1.505

TABLE 4. SOUND-PRESSURE LEVELS FOR SEVERAL UNITS USED TO MEASURE SOUND PRESSURE IN AIR (RE 20 MICROPASCALS)

SPL	Unit	
94	1 pascal	(Pa, N/m ²)
127.6	1 pound/foot ²	(psf)
170.8	1 pound/in ²	(psi)
194.1	1 atmosphere	(atm)
74	1 microbar	(μ bar)

sure-voltage analogue and (2) force or pressure-current analogue. Tables 5 and 6 tabulate the pertinent quantities in these two analogies.

In a specific problem in which it is desired to represent an acoustical or mechanical system by an electrical circuit, a choice must be made between the force-voltage and force-current analogies. If the system includes electrical as well as mechanical or acoustical elements, as in a motor or a loudspeaker, for example, the choice is dictated by the type of coupling between the electrical and mechanical parts. However, if the system is entirely mechanical or acoustical, the choice is arbitrary and depends upon personal preference. In any event, given the equivalent circuit derived on the basis of one analogy, the equivalent circuit corresponding to the other analogy can be derived by taking the dual of the given circuit.

In a system that combines acoustical and mechanical quantities, it is necessary to establish a common frame of reference. Usually, it is easiest to convert acoustical quantities to mechanical quantities, as for example, integrating the pressure over a piston (loudspeaker or microphone) to calculate the force, etc.

HEARING*

The auditory system consists of the periphery sensors, acoustic neurological transducers, the ears, the eighth cranial nerve leading to a programming and priority switching center at various levels of the brain stem, and finally to the auditory area of the cortex located near the Sylvian fissure of the frontal-lobe convolution.

The auditory system does much more than detect minute sounds. Among other functions, it preferentially places more weight on certain preprogrammed characteristic sounds, localizes the direction of most sounds by a variety of ingenious techniques, and initiates involuntary actions for visual acquisition of the source.

The hearing mechanism was probably evolved to help man survive in a hostile environment, and not for linguistic communication or for musical entertainment.

* Flanagan, J. L. *Speech Analysis, Synthesis and Perception*. New York: Academic Press, 1965. Also Richardson, E. G., ed. *Technical Aspects of Sound*. New York: Elsevier Press, 1953.

TABLE 5. FORCE/PRESSURE-VOLTAGE ANALOGY

Mechanical Quantity	Acoustical Quantity	Electrical Quantity
Force (newtons)	Sound Pressure (newtons meter ⁻²)	Voltage (volts)
Velocity (meters/second)	Volume Velocity (meters ³ second ⁻¹)	Current (ampercs)
Displacement (meters)	Volume Displacement (meters ³)	Charge (coulombs)
Mass (kilograms)	Acoustic Mass (kilograms meter ⁻⁴)	Inductance (henrys)
Stiffness (newtons/meter)	Acoustic Stiffness (kilograms meter ⁻⁴ second ⁻²)	Capacitance ⁻¹ (farads ⁻¹)
Viscous Friction (newtons/meter/second)	Acoustic Resistance (kilograms meter ⁻⁴ second ⁻¹)	Resistance (ohms)
Mechanical Impedance (mechanical ohms) or (force/velocity)	Acoustic Impedance (kilograms meter ⁻⁴ second ⁻¹)	Impedance (ohms)

TABLE 6. FORCE/PRESSURE-CURRENT ANALOGY

Mechanical Quantity	Acoustical Quantity	Electrical Quantity
Force (newtons)	Sound Pressure (newtons meter ⁻²)	Current (amperes)
Velocity (meters/second)	Volume Velocity (meters ³ second ⁻¹)	Voltage (volts)
Displacement (meters)	Volume Displacement (meters ³)	Impulse $\int v dt$ (volt-seconds)
Mass (kilograms)	Acoustic Mass (kilograms meter ⁻⁴)	Capacitance (farads)
Stiffness (newtons/meter)	Acoustic Stiffness (kilograms meter ⁻⁴ second ⁻²)	Inductance ⁻¹ (henrys ⁻¹)
Viscous Friction (newtons/meter/second)	Acoustic Resistance (kilograms meter ⁻⁴ second ⁻¹)	Conductance (mhos)
Mechanical Impedance (mechanical ohms) or (force/velocity)	Acoustic Impedance (kilograms meter ⁻⁴ second ⁻¹)	Admittance (mhos)

The apparent loudness attributed to a sound varies not only with the sound pressure, but also with the frequency (or pitch) of the sound. In addition, the way it

varies with frequency depends on the pressure. Fig. 1 illustrates experimentally determined equal loudness contours (SPL) as a function of frequency.

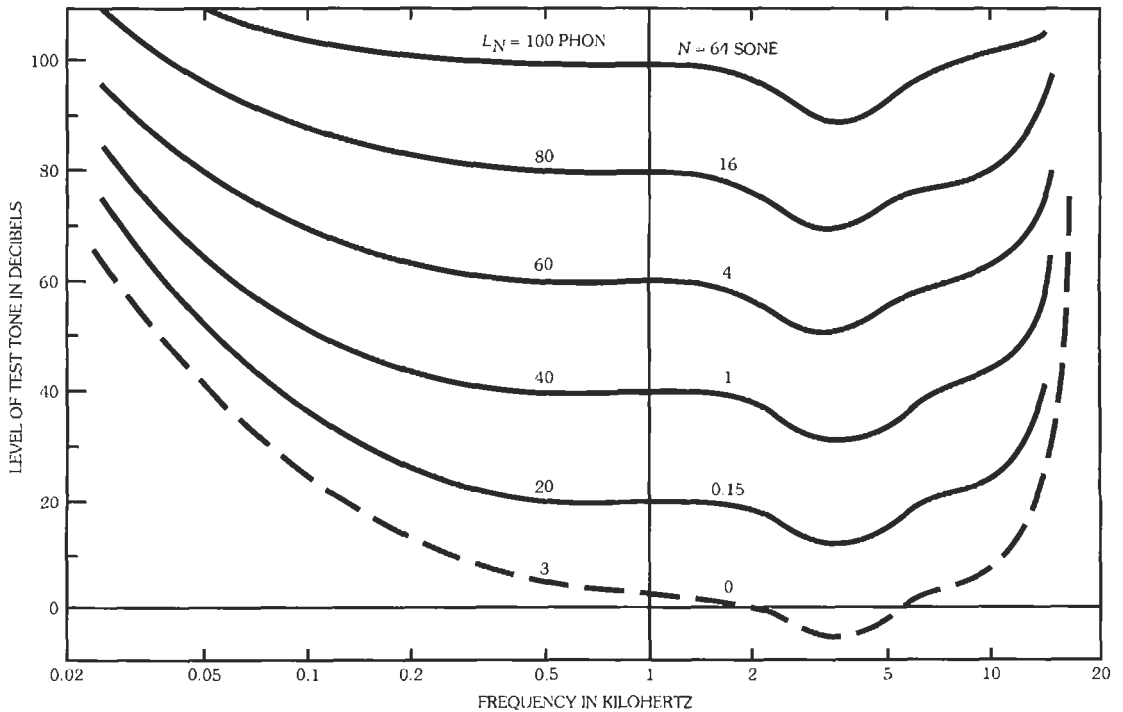


Fig. 1. Equal loudness contours. (After Zwicker and Fastl, Psychoacoustics—Facts and Models. Springer Verlag, 1990.)

MEASUREMENT OF SOUND

Sound Level Meters

A Sound Level Meter (SLM) is an instrument designed to measure directly sound-pressure levels while incorporating frequency weighting networks which approximate the inverse of equal loudness contours at various levels. The current American National Standard Specification for Sound Level Meters (ANSI S1.4A-1985) designates three alternate frequency weighting networks (Fig. 2) which correspond to the inverse of an experimentally determined moderate level "equal noisiness contour" and a "flat," or unweighted, measure. (In the latter case, the measurement bandwidth is determined by the microphone and electrical network characteristics.) Standard SLM also complies with International Electro-technical Commission (IEC) 651.*

An SLM measures the "quasi-root-mean-square" (rms) value of a time-varying sound pressure. Table 7 lists the rms detector time constant and corresponding label as given in ANSI S1.4A-1985. In addition, some SLMs incorporate a peak-hold (p) which captures the instantaneous largest sound pressure level.

An integrating-averaging SLM provides a "true-integrating" detector circuit in lieu of the quasi-rms detector. These meters are built in accordance with IEC-804. Typically, they measure time-average sound level (also called equivalent sound level) and sound exposure level.*

In text, complete description must be given for any measurement made (e.g., 55 dB B-weighted slow). In equations, SPLs are abbreviated as L_{AF} for A-weighted fast, L_{BF} for B-weighted impulse, L_{pC} for C-weighted

* Definitions for these terms are given in the next section.

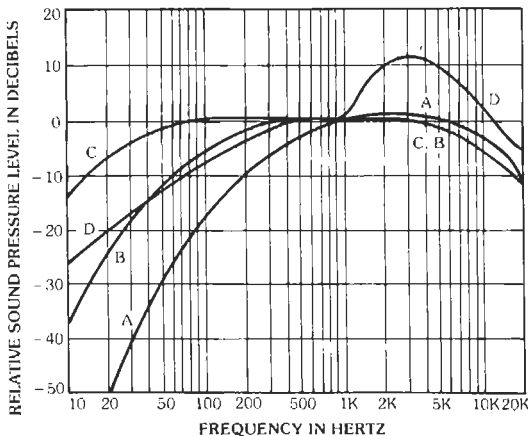


Fig. 2. IEC Standard A, B, and C weighting curves for Sound Level Meters. (EPA. Public Health and Welfare Criteria for Noise. 550/9-73-002. p. 2-2.)

TABLE 7. SOUND LEVEL METER TIME CONSTANTS

Label	RMS Detector Time Constant
Slow (S)	1 s
Fast (F)	0.1 s
Impulse (I)	35 ms

peak, etc. Maximum should not be confused with peak. Peak refers to an instantaneous measure, while maximum refers to a detected level [e.g., max A-weighted slow (L_{AFmax})].

Measurement and/or Laboratory Microphones and Their Calibration

Current measurement and/or laboratory microphones are generally capacitor, electret, or ceramic (piezoelectric) types. These are chosen because of their flat frequency response, low noise, and stability. One-inch and 1/2-inch microphones are the most common with smaller sizes also available. Capacitor microphones require a polarization voltage but are not potentially as susceptible to aging or temperature effects. All three are very high impedance devices requiring a FET follower close coupled for impedance transformation. For a given size, the ceramic microphone is on the order of 20 dB less sensitive than the capacitor or electret microphone.

Calibration is normally performed with a close-coupled single-frequency acoustical source. The acoustical signal may be generated mechanically (pistons driven by a cam and motor) or electromechanically. Typical frequencies are 1000 Hz and 250 Hz. Other frequencies, high pressure, and reciprocity calibrators are also available as needs dictate, but are generally not required.

Sound Level

In recent years it has become customary to measure most environmental or industrial work-place (occupational) noise with the A-weighting. In the United States, A-weighted SPL is the commonly designated sound level. (Special high-amplitude impulse sounds, such as sonic booms or explosions, which can noticeably vibrate structures are measured with C-weighting.)

Sound Power Level

One is usually concerned with the acoustic power radiated by some source of noise or radiated purposefully by a transducer such as a loudspeaker. Typically, these sources have a directivity pattern Q which is a function of angle. Radiated power is the normal component of the intensity integrated over an imaginary surface encompassing the source. (Sufficiently far from

TABLE 8. ACOUSTIC POWER AND SOUND POWER LEVELS OF TYPICAL NOISE SOURCES*
(A-WEIGHTED)

Total Power (A-weighted) (Watts)	Power Level (dB re 10^{-12} Watts)	Source
100 000	170	Ram jet Turbojet engine with afterburner
10 000	160	Turbojet engine, 7000 lb thrust
1000	150	4-propeller airliner
100	140	75-piece orchestra Pipe organ { Peak rms level in 1/8-second intervals
10	130	
3	125	Small aircraft engine
1.0	120	Large chipping hammer Piano BB ³ tuba { Peak rms level in 1/8-second intervals
0.1	110	Blaring radio Centrifugal ventilating fan (13 000 cfm)
0.01	100	4-foot loom Auto on highway
0.001	90	Vanaxial ventilating fan (1500 cfm) Voice—shouting (average long-time rms)
0.0001	80	
0.000 01	70	Voice—conversational level (average long-time rms)
0.000 001	60	
0.000 0001	50	
0.000 000 01	40	
0.000 000 001	30	Voice—very soft whisper

*Space average sound pressure level at 10 meters = Power level - 28 dB for a source radiating into a hemisphere.

the source, the intensity is proportional to the pressure squared.) Table 8 lists the PWL for some common noise sources.

DEFINITIONS

The following definitions are in accordance with ANSI S12.9-1988 and ANSI S12.7-1986.

Instantaneous sound pressure—Total instantaneous pressure at a point in the presence of a sound wave minus the static pressure at that point. Unit: pascal (Pa).

Sound power (level)—(A.) Sound power: The radiated acoustic power in a stated frequency band. Unit: watt (W). (B.) Sound power level: Ten times the common logarithm of the ratio of the sound power to the reference sound power, W_0 , of 10^{-12} watts. Unit: decibel (dB).

$$L_w = 10 \log W/W_0$$

where W is the radiated acoustic power in watts and W_0 is the reference power of 10^{-12} watts.

Sound pressure (level)—(A.) Sound pressure: Root mean square of the instantaneous sound pressures in a stated frequency band and during a specified time interval, unless another time-averaging process is indicated. Unit: pascal (Pa). (B.) Sound pressure level: Ten times the common logarithm of the square of the ratio of the sound pressure to the reference sound pressure, p_0 , of 20 micropascals. Unit: decibel (dB).

Sound (Level)—(A.) (Frequency-weighted) sound pressure: Root mean square of the instantaneous sound pressure which is filtered (frequency-weighted) with a standard frequency characteristic (e.g., A or C, etc.) and exponential time weighted in accordance with the

standardized characteristics slow (S), fast (F), impulse (I), or peak, both weightings as specified in ANSI S1.4A-1985. The time weighting should be specified. The frequency weighting should be specified; otherwise, A-weighting will be understood. Unit: pascal (Pa). (B.) Sound level: Ten times the common logarithm of the square of the ratio of the frequency-weighted (and time-averaged) sound pressure to the reference sound pressure, p_0 , of 20 micropascals. Unit: decibel (dB). NOTE: A, C, or octave band sound levels are all sound levels having different frequency weightings which must be specified.

Time-average sound (level)—(A.) Time-average (frequency-weighted) sound pressure: (1) square root of the quotient of the time integral of frequency-weighted squared instantaneous sound pressures divided by the time period of integration in seconds; or (2) square root of the quotient of the sound exposure, in pascal-squared seconds, in a specified time period, divided by the time period of integration in seconds. Unit: pascal (Pa). (B.) Time-average sound level: Ten times the common logarithm of the square of the ratio of time-average (frequency-weighted) sound pressure to the reference pressure, p_0 , of 20 micropascals. Unit: decibel (dB). NOTES: (1.) Time-average sound level is also termed equivalent-continuous sound level or equivalent-continuous frequency-weighted sound pressure level. (2.) A, C, or octave band time-average sound levels are all time-average sound levels having different frequency weightings which must be specified.

Peak sound pressure (level)—(A.) Peak sound pressure: For any specified time interval, the maximum absolute value of the instantaneous sound pressure in that interval. Unit: pascal (Pa). (B.) Peak sound pressure level: Ten times the common logarithm of the square of the ratio of the peak sound pressure to the reference sound pressure of 20 micropascals. Unit: decibel (dB).

Sound pressure (level)—(A.) Maximum sound pressure: For any specified time interval, the maximum value of the sound pressure in that interval. Unit: pascal (Pa). (B.) Maximum sound pressure level: Ten times the common logarithm of the square of the ratio of the maximum sound pressure to the reference sound pressure of 20 micropascals. Unit: decibel (dB).

Sound exposure (level)—(A.) Sound exposure: Time integral of squared, frequency-weighted instantaneous sound pressure. The time period of integration should be specified when the sound exposure of the background noise is a significant contributor to the total sound exposure, when the data may be useful for identifying a source, or when the time period of integration is otherwise useful. The frequency weighting should be specified; otherwise, A-weighting will be understood. Unit: pascal-squared second (Pa^2s). NOTES: (1.) One pascal-squared second is called a sound exposure unit. (2.) In principle, sound exposure is independent of the time period of integration. However, in practice, a measurement of sound exposure may be limited by the square of the equivalent sound

pressure of the background noise multiplied by the time period of integration. In such a case, the time period of integration should be specified. (3.) Some instruments for measurement of sound exposure include a threshold sound level below which the instrument does not accumulate contributions to the integral. If the threshold sound level is above the level of the background noise, then the time period of integration need not be stated, but the threshold should be stated. (4.) For most transient sounds such as an aircraft flyby or a vehicle pass-by, etc., the A-frequency weighting is preferred. For short-duration, high-energy impulsive sounds, the C-frequency weighting is preferred, and measurements are to be made in accordance with ANSI S12.4-1986. (B.) Sound exposure level: Ten times the common logarithm of the ratio of sound exposure to the reference sound exposure (E_0) of 400 micropascal-squared seconds [$400 (\mu\text{Pa})^2\text{s}$]. Unit: decibel (dB). NOTES: (1.) In a given measurement time period of T seconds, the sound exposure level (L_{ET}) is related to the time-average sound (pressure) level (L_{pT}) as follows:

$$L_{ET} = L_{pT} + 10 \log_{10}(T/t_0)$$

where t_0 is the reference duration of 1 second. (2.) A, C, or octave-band sound exposure levels are all sound exposure levels having different frequency weightings which must be specified.

ENVIRONMENTAL NOISE

Current methods to assess environmental noise embody several relatively simple concepts. A fundamental concept is loudness. The human ear is more sensitive to the middle frequencies than it is to very low or very high frequencies.* So for purposes of noise assessment, the sound measured by a microphone, like the sound "measured" by the human ear, is filtered to remove much of the low and high frequencies. The International Organization for Standardization (ISO) and the American National Standards Institute (ANSI) recommend the A frequency-weighting curve to accomplish this filtering, and in the USA (as in most countries), the Federal Environmental Protection Agency (USEPA), the Federal Aviation Administration (FAA), etc., have adopted the use of this A frequency-weighting curve, a curve which removes much of the low and high frequencies (references 2 and 3).†

* Middle frequencies are about middle C to four octaves above middle C on the piano.

† The Illinois property-line-noise-source regulations go further and divide the sound spectrum into octave bands so as to create an environment which is more acceptable to communities than would be an environment regulated only by the A-weighting curve. These regulations go a second step further and divide the sound spectrum into 1/3 octave bands in order to regulate "pure tones" which may be particularly bothersome.

The first rule of current environmental noise assessment makes use of the concept of loudness. In effect, it states: All else being equal, a louder noise is more annoying to a community than is a quieter noise. In other words, all else being equal, community annoyance grows as loudness increases.

Environmental noise regulations typically take one of two forms. Either they are specified as simple SPL limits at some measurement locations or they are specified by time-integrated quantities. The former are similar in concept to "speed limits" or "discharge limits" and are relatively easy to enforce. They may take the form of property line limits such as so many decibels (A-weighted, slow) or may even be expressed as octave band or 1/3-octave band limits. (Some SLMs include an octave band filter set; others have octave and 1/3-octave filter sets available as attachments. The filters are specified by ANSI S1.11-1966 and IEC 225.) For vehicular sources, the limits are usually specified in terms of a maximum A-weighted fast (or slow) level at a reference distance (typically 15 m).

Airports, highways, railways, and other large, distributed time-varying sources are usually regulated and assessed by the use of time-integrated measures. Extensive research over the past three decades has shown that these time-integrated measures better correlate with human and community response than do the single event limits described above, but they are more difficult to use for purposes of enforcement. Time-integrated measures incorporate the second and third rules of current noise assessment procedures.

The second rule of environmental noise assessment deals with the duration or repetition of a sound. Basically, it states: For the same loudness, a longer-duration sound is more annoying than a shorter-duration sound, and sounds repeated many times are more annoying than are sounds repeated only a few times.

The third rule of noise assessment deals with the time of day. It basically states: All else being equal, sounds during some period(s) of the day result in greater overall net community annoyance than do these same sounds during other periods of the day.

In terms of noise assessment these three rules imply:

- (1) We use a filter (or filters) proportional to hearing to measure the "loudness,"*
- (2) We add up all the sound an event or operation creates to account for the duration as well as the loudness,
- (3) We add in each event or operation separately to account for repetition, and,
- (4) We count each event or operation at night as equal to ten daytime events.

* The Illinois property-line-noise-source octave-band level spectra are very similar to the inverse of equal loudness spectra and thus yield results monotonically related to loudness.

To do all of this, we predict or measure the sound exposure of each event or operation. Sound exposure can be thought of as a measure of the loudness multiplied by the time it takes for the event or operation to occur, i.e., loudness times duration.

Next, we add up the sound exposure of each individual event or operation in order to develop the total sound exposure. In this fashion, we account for *all* of the noise. This point is very important and frequently misunderstood. Each event or operation in and of itself can be annoying, can interfere with communications, etc. But the noise impact to a community experiencing 10 or 100 of these events each day is far greater than is the impact to a community experiencing just one of these events.

Finally, we multiply the sound exposure of each event or operation which occurs at night by ten, or we divide its criterion sound exposure by ten before assessing the noise. In this way, we apply a nighttime adjustment. This, too, is a very important factor in assessing noise.

Two measures, defined below, are typically used to quantify noise impact. One is sound exposure (SE), and the other is time-average sound. These two are related by the time period of measurement. Sound exposure is proportional to the *total* sound energy during the time period, and time-average sound is the average sound during this same time period. Sound exposure and time-average sound level are usually expressed in decibels with reference to the standard reference pressure, p_0 , of $20 \mu\text{Pa}$ or the standard reference sound exposure, SE_0 , of $400 (\mu\text{Pa})^2\text{-s}$.

In the United States, with the help of the National Academy of Sciences (NAS), the US Environmental Protection Agency (USEPA) created the day-night sound level (DNL) as the descriptor of choice for composite environmental noise assessment (references 2 and 3).* This descriptor, DNL, is really a measure of total day-night sound exposure as described in ANSI S12.9-1988 (reference 4).† In reality, measures such as DNL model the relation: Average community response is proportional to total sound exposure received by that community during a typical day. The underlying measure in DNL is total day-night sound exposure (DNSE) which is given by:

$$\text{DNSE} = \text{Total } SE_{\text{Day}} + (10)(\text{Total } SE_{\text{Night}})$$

and

* The best summary of references 2 and 3 is: *Protective Noise Levels—Condensed Version of EPA Levels Document*. U.S. Environmental Protection Agency, EPA 550/9-79-100, November 1978.

† Sound exposure is exactly analogous to the concept of light exposure as used in photography. It is the product of the sound intensity times the duration of the sound. Sound exposure is discussed in Chapter 3 of reference 4, and a formal definition based on ANSI S12.9-1988 is given in Chapter 4 of reference 4.

$$DNL = 10 \log [DNSE/(T_0 SE_0)]$$

where

$$T_0 = 86\,400, \text{ the number of seconds in a day}$$

This is the actual way we calculate or measure DNL.

ASSESSMENT OF WORKPLACE-INDUSTRIAL NOISE

A Sound Level Meter (SLM) set on the A-scale can be used to assess industrial noise for hearing hazard. For enforcement purposes, the US Department of Labor (DOL) has set 90 dB (A) as the limit for eight hours of exposure per day. Most industrial noise is not uniform and continuous, and so the regulations of DOL include procedures for determining the equivalent exposure. These procedures incorporate a 5-dB trading rule in that halving the exposure time allows the noise limit to rise by 5 dB.* The present limits are given in Table 9.

As an alternative to using an SLM to measure industrial noise exposure, several personnel dosimeters are available that directly measure a worker's dose (usually as a percent of the eight-hour limit). These devices directly incorporate the 5-dB trading rule in their electronics.

Exposure to impulsive or impact noise should not exceed 140 dB peak sound pressure level.

* Most European countries use a 3-dB or equal energy rule, and the USAF uses a 4-dB rule.

TABLE 9. PERMISSIBLE NOISE EXPOSURES*

Duration per Day, Hours	Sound Level, dBA Slow Response
8	90
6	92
4	95
3	97
2	100
1½	102
1	105
½	110
¼ or less	115

* When the daily noise exposure is composed of two or more periods of noise exposure of different levels, their combined effect should be considered, rather than the individual effect of each. If the sum of the following fractions: $C_1/T_1 + C_2/T_2 + \dots + C_n/T_n$ exceeds unity, where C_n indicates the total time of exposure at the n th noise level and T_n indicates the total time of exposure permitted at that level, then the mixed exposure should be considered to exceed the limit value.

SOUND IN ENCLOSED ROOMS*

Indoors, acoustics is of concern for a variety of reasons. Auditoria, broadcast studios, etc., require low background noise levels and "good" acoustical design; offices require varying background levels and communications capabilities depending on use; and in the industrial workplace, conservation of hearing and ability to communicate are usually the primary concerns.

In any room, the reverberant SPL increases as the acoustic power entering the room increases and decreases as sound absorption in the room increases and/or as sound power flows from the room. Sources of sound may be internal to rooms, such as machinery or loudspeakers, or may be external. In the latter case, the sound flows through walls, windows, ventilation ducts, etc., from other building spaces or components, or from outdoors. In either case, the sound field in a room is composed of two parts, the direct field and the reverberant field. The direct field of a source is

$$p = \rho c k D Q^{1/2} / 2\pi^{1/2} r \quad (\text{Eq. 9})$$

where,

- $k = (2\pi f/c)$ (rad/m),
- $p = \text{pressure}$ (N/m²),
- $c = \text{speed of sound}$ (m/s),
- $\rho = \text{density of air}$ (kg/m³)
- $D = \text{source strength}$ (m³/s·rad),
- $Q = \text{directivity}$ (rad)†,
- $r = \text{distance from the source}$ (m).

Near a source, the sound field may be very complicated; the pressure and velocity may not be in phase, and hence the power radiated is not proportional to the pressure squared. This is particularly a problem when determining the noise power radiated by machinery in the industrial environment.‡

By substitution in equation 9, the direct field intensity (I_d) of a source can be written as

$$I_d = WQ/4\pi r^2 \quad (\text{Eq. 10})$$

where,

$W = \text{radiated source power}$ (watts),

* See, for example, reference 2.5, 2.6, or 2.8.

† The directivity factor of a source is the ratio of the sound pressure squared, at some fixed distance and specified direction (such as the axis of a loudspeaker), to the mean-square sound pressure at the same distance averaged over all directions from the source. The distance must be great enough so that the sound appears to diverge spherically from the effective acoustic center of the source.

‡ Sound-intensity meters are a valuable aid in these situations, as they indicate the direction of the intensity.

Q = directivity (rad),
 r = distance from the source (meters).

The total intensity (I) in an enclosed room (direct plus reverberant field) is

$$I = WQ/4\pi r^2 + 4W/R \quad (\text{Eq. 11})$$

The quantity R is the room constant defined by

$$R = S\bar{\alpha}/(1 - \bar{\alpha}) \quad (\text{Eq. 12})$$

where,

S = total surface area of the room,
 $\bar{\alpha}$ = average room absorption coefficient.

Fig. 3 illustrates the sound fields in a room. In terms of the sound pressure level,

$$\text{SPL} = \text{PWL} + 10 \log (Q/4\pi r^2 + 4/R) + K \quad (\text{Eq. 13})$$

where,

$K = 0.2$ dB if metric units are used,
 $K = 10.5$ dB if English units are used.

Standing Sound Waves

Resonant conditions in sound studios cause standing waves by reflections from opposing parallel surfaces, such as ceiling-floor and parallel walls, resulting in serious peaks in the reverberation-time/frequency curve. Standing sound waves in a room can be considered comparable to standing electrical waves in an improperly terminated transmission line where the transmitted power is not fully absorbed by the load.

In properly proportioned rooms, resonances can be effectively reduced and standing waves practically eliminated by introducing numerous absorbing surfaces disposed obliquely. Thus, large-order reflections can be

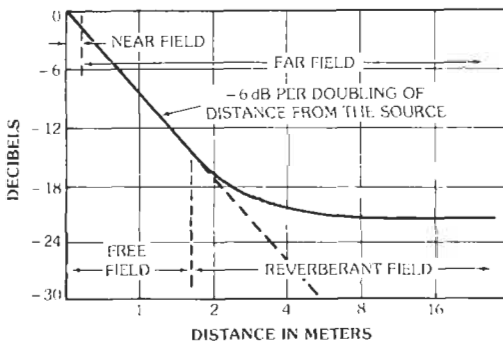


Fig. 3. Description of the sound field around a sound source in a reverberant room. (From Ginn, K. B. Architectural Acoustics. Brüel and Kjaer, 1978.)

avoided by breaking them up into numerous smaller reflections. The object is to prevent sound reflection back to the point of origin until after several absorbing reflections.

The most desirable ratios of dimensions for broadcast studios are given in Fig. 4.

Reverberation Time

Reverberation time varies with frequency and is measured by the time required for a sound, when suddenly interrupted, to die away or decay to a level 60 decibels below the original sound. Measurement of reverberation time is the easiest way to find the average absorption coefficient ($\bar{\alpha}$) and its associated value of reflectivity ($1 - \bar{\alpha}$) in an acoustic environment. Conversely, absorption coefficients and room sizes are used to predict the reverberation time for new construction or remodeling.

Measurement of Reverberation Time

The reverberation time of an enclosed space that already exists is an important quantity that is relatively easy to measure. The degree of accuracy required is determined by the use to which the data will be put. The majority of day-to-day reverberation measurements are taken with quite satisfactory results by using the interrupted-noise method (Fig. 5).

The Schroeder-Kuttruff method (Fig. 6) employs a 2.7-ms rectangular pulse that is used to excite a standard 1/3-octave bandpass filter. The resulting "ringing" of the filter provides a statistically reliable signal that allows highly repeatable decay recordings and excellent resolution of the EDT.

Calculation of Reverberation Time

The Sabine equation (equation 14 for room dimensions in feet, equation 14a for room dimensions in meters) is commonly used to calculate reverberation time.

$$RT_{60} = 0.049V/S \bar{\alpha} \quad (\text{Eq. 14})$$

$$RT_{60} = 0.161V/S \bar{\alpha} \quad (\text{Eq. 14a})$$

where,

RT_{60} is the time in seconds required for a sound to decay 60 dB,
 V is the volume of the room,
 S is the boundary surface area,
 $\bar{\alpha}$ is the average absorption coefficient.

The value of $\bar{\alpha}$ is:

$$\bar{\alpha} = (s_1\alpha_1 + s_2\alpha_2 + \dots + s_n\alpha_n)/S \quad (\text{Eq. 15})$$

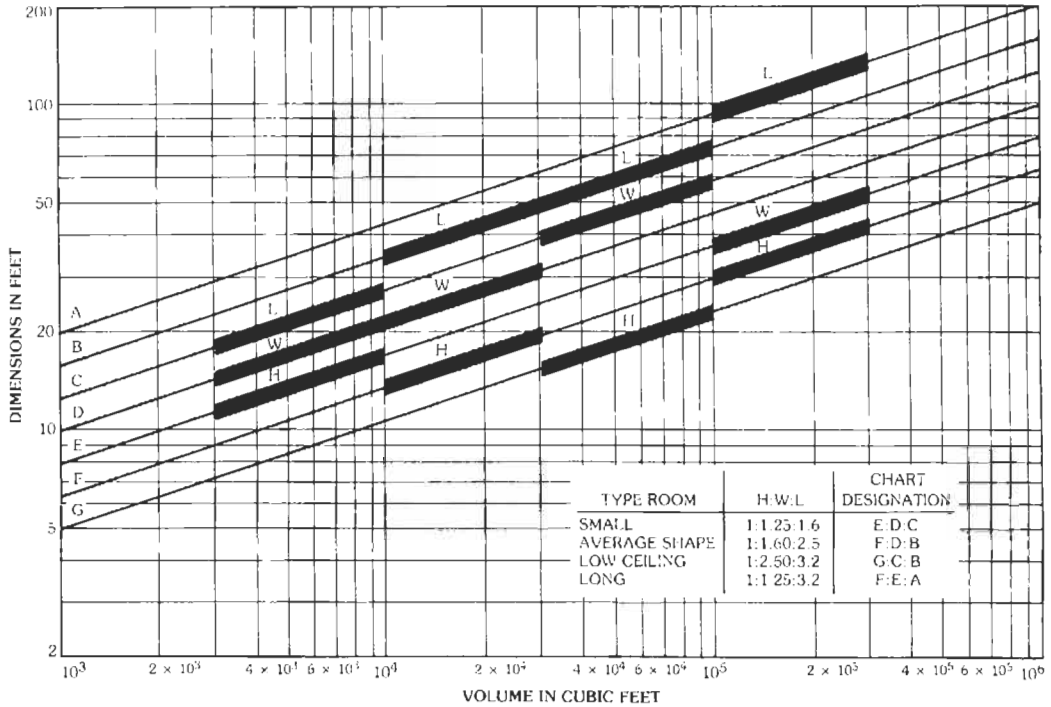


Fig. 4. Preferred room dimensions based on $2^{1/3}$ ratio. Permissible deviation is 5 percent. (Courtesy Acoustical Society of America and RCA.)

where,

- $s_1, s_2, \text{ etc.},$ are boundary surface areas,
- $\alpha_1, \alpha_2, \text{ etc.},$ are the absorption values for the boundary areas with which they are associated,
- $s_n \alpha_n$ is the total absorption of the people, furniture, etc., present in the room.

Note that $S\bar{\alpha}$ can be replaced by A , the total absorption in the room. This concept is useful when considering the effects of surfaces plus objects or people in the room.

In the limiting case, the Sabine equation predicts a finite reverberation time in a room with 100% absorption present, and for true absorption values in excess of 0.63, this equation can give α values in excess of 1.0 (100% absorption). The Norris-Eyring equation (equation 16 or 16a) gives $\bar{\alpha}$ values from 1.0 to 0 for true absorption values when calculated from actual RT_{60} measurements.

$$RT_{60} = 0.049V/[-S \ln(1 - \bar{\alpha})] \quad (\text{Eq. 16})$$

with room dimensions in feet or

$$RT_{60} = 0.161V/[-S \ln(1 - \bar{\alpha})] \quad (\text{Eq. 16a})$$

with room dimensions in meters.

If the value of RT_{60} is measured and the corresponding value of $\bar{\alpha}$ is calculated, insertion of this value of $\bar{\alpha}$ into the expression $-\ln(1 - \alpha)$ converts the $\bar{\alpha}$ value into a Sabine $\bar{\alpha}$. For example, $-\ln(1 - 0.63) = 0.99$.

Absorption Coefficients

When tables of absorption values are examined, it is important to know which formula was used in determining the numerical values. If the values were obtained from the Sabine equation, be sure to use the Sabine

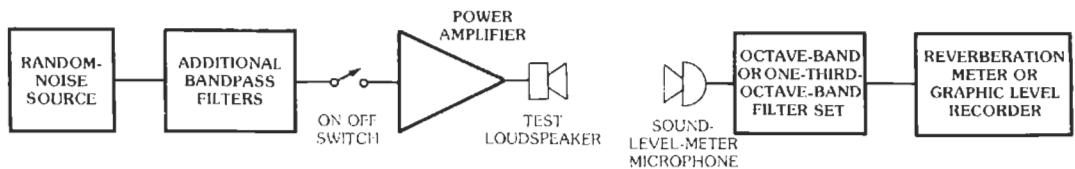


Fig. 5. The interrupted-noise method.

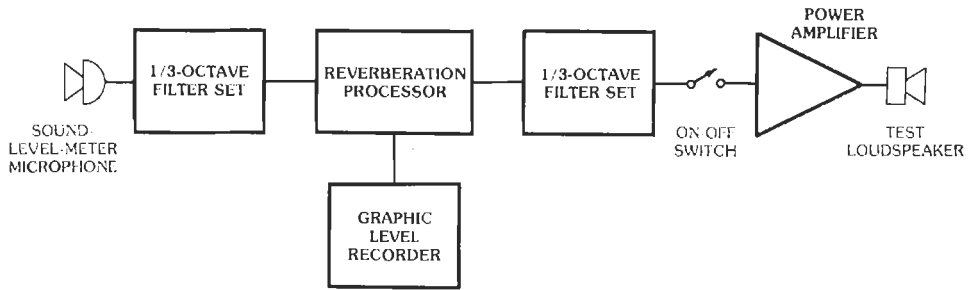


Fig. 6. The Schroeder-Kuttruff method.

equation variations consistently for any further manipulations of the data. If the Norris-Eyring equation was used, remain consistent in its use for any further manipulations of the data. The Sabine formula is used in measuring and specifying most building materials. Hence, the Sabinc equation must be used unless the absorption values are converted (see above). Table 10 gives absorption coefficients of some typical building materials.

Optimum Reverberation Time

Optimum, or most desirable, midfrequency (500–1000 Hz) reverberation time varies with room size and use, such as music, speech, etc. (see Fig. 7).

INDOOR DESIGN NOISE LEVELS AND SPEECH COMMUNICATION

Preferred Speech Interference Levels/ Speech Interference Levels

Noise interference with speech is usually a masking process. For many noises, the measurement and calculation of articulation index can be simplified further by the use of a three-band analysis.* The bands chosen are the octave bands centered on 500, 1000, and 2000 Hz. The arithmetic average of these three SPLs in these three bands gives the quantity called the three-band preferred octave speech-interference level (PSIL). One can use this level for determining when telephone use or speech communication is easy or difficult and what changes in levels are required to improve the situation (Figs. 8 and 9).

* Beranek, L. L. "The Design of Speech Communications Systems." *Proc. IEEE*, Vol. 35, No. 9, Sept. 1947, pp. 880–890.

For satisfactory intelligibility of difficult speech material, maximum permissible values of speech-interference levels for men with average voice strengths are given in Fig. 9, which is an extension by Webster of Beranek's work.* It is assumed in this chart that there are no reflecting surfaces nearby, that the speaker is facing the listener, and that the spoken material is not already familiar to the listener. For example, the speech-interference level in a factory might be 80 dB, which is high. The chart indicates that the two people must ordinarily be no more than two feet apart in order to be understood satisfactorily. If the words spoken are carefully selected and limited in number, intelligible speech will be possible at greater distances.

If a number of conversations are to be held in the same reverberant room, the procedure is more complicated. This chart cannot be used on the basis of the background-noise level before the conversations are in progress, because a given conversation will be subject to interference from the noise produced by all the other conversations. The general procedure for calculating a speech-interference level under those conditions has not been completely worked out.

Not only background, equipment, and other occupant noise affects speech intelligibility. In a live room (too much reverberation), speech syllables are smeared by reflected noise at the other end, even though the noise is not very great. That is, if the distance from speaker (or loudspeaker) to listener is too great or the room too reverberant, then the speaker's own sound energy may make the speech unintelligible. As a rule of thumb, the direct field from the speaker (or loudspeaker) must exceed his (its) own reverberant field by 6 dB in order to maintain good speech intelligibility and prevent feedback.

* Webster, J. C. "Effects of Noise on Speech Intelligibility." *Proc. Conference Noise as a Public Health Hazard*, Washington, D. C., June 13–14, 1968, ASHA Rpt. 4. Washington, DC: The American Speech and Hearing Association, 1969, pp. 49–73.

TABLE 10. COEFFICIENTS OF GENERAL BUILDING MATERIALS AND FURNISHINGS

Complete tables of coefficients of the various materials that normally constitute the interior finish of rooms may be found in the various books on architectural acoustics. This short list will be useful in making simple calculations.

Materials	Coefficients					
	125 Hz	250 Hz	500 Hz	1000 Hz	2000 Hz	4000 Hz
Brick, unglazed	0.03	0.03	0.03	0.01	0.05	0.07
Brick, unglazed, painted	0.01	0.01	0.02	0.02	0.02	0.03
Carpct, heavy, on concrete	0.02	0.06	0.14	0.37	0.60	0.65
Samc, on 40 oz hairfelt or foam rubber	0.08	0.24	0.57	0.69	0.71	0.73
Same, with impermeable latex backing on 40 oz hairfelt or foam rubber	0.08	0.27	0.39	0.34	0.48	0.63
Concrete block, coarse	0.36	0.44	0.31	0.29	0.39	0.25
Concrete block, painted	0.10	0.05	0.06	0.07	0.09	0.08
Fabrics (draperics)						
Light velour, 10 oz per sq yd, hung straight, in contact with wall	0.03	0.04	0.11	0.17	0.24	0.35
Medium velour, 14 oz per sq yd, draped to half area	0.07	0.31	0.49	0.75	0.70	0.60
Heavy velour, 18 oz per sq yd, draped to half area	0.14	0.35	0.55	0.72	0.70	0.65
Floors						
Concrete or terrazzo	0.01	0.01	0.015	0.02	0.02	0.02
Linoleum, asphalt, rubber, or cork tile on concrete	0.02	0.03	0.03	0.03	0.03	0.02
Wood	0.15	0.11	0.10	0.07	0.06	0.07
Wood parquet in asphalt on concrete	0.04	0.04	0.07	0.06	0.06	0.07
Glass						
Large panes of heavy plate glass	0.18	0.06	0.04	0.03	0.02	0.02
Ordinary window glass	0.35	0.25	0.18	0.12	0.07	0.04
Gypsum board, ½" nailed to 2×1's, 16" o.c.	0.29	0.10	0.05	0.04	0.07	0.09
Marble or glazed tile	0.01	0.01	0.01	0.01	0.02	0.02
Openings						
Stage, depending on furnishings				0.25-0.75		
Deep balcony, upholstered seats				0.50-1.00		
Grills, ventilating				0.15-0.50		
Plaster, gypsum or lime, smooth finish on tile or brick	0.013	0.015	0.02	0.03	0.04	0.05
Plaster, gypsum or lime, rough finish on lath	0.14	0.10	0.06	0.05	0.04	0.03
Same, with smooth finish	0.14	0.10	0.06	0.04	0.04	0.03
Plywood paneling, ¾" thick	0.28	0.22	0.17	0.09	0.10	0.11
Water surface, as in a swimming pool	0.008	0.008	0.013	0.015	0.020	0.025
Air, Sabins per 1000 cubic feet @ 50% RH				0.9	2.3	7.2

Absorption of Seats and Audience
Values given are in Sabins per square foot of seating area per unit.

	Absorption of Seats and Audience					
	125 Hz	250 Hz	500 Hz	1000 Hz	2000 Hz	4000 Hz
Audience, seated in upholstered seats, per sq ft of floor area	0.60	0.74	0.88	0.96	0.93	0.85
Unoccupied cloth-covered upholstered seats, per sq ft of floor area	0.49	0.66	0.80	0.88	0.82	0.70
Unoccupied leather-covered upholstered seats, per sq ft of floor area	0.44	0.54	0.60	0.62	0.58	0.50
Wooden pews, occupied, per sq ft of floor area	0.57	0.61	0.75	0.86	0.91	0.86
Chairs, metal or wood seats, each, unoccupied	0.15	0.19	0.22	0.39	0.38	0.30

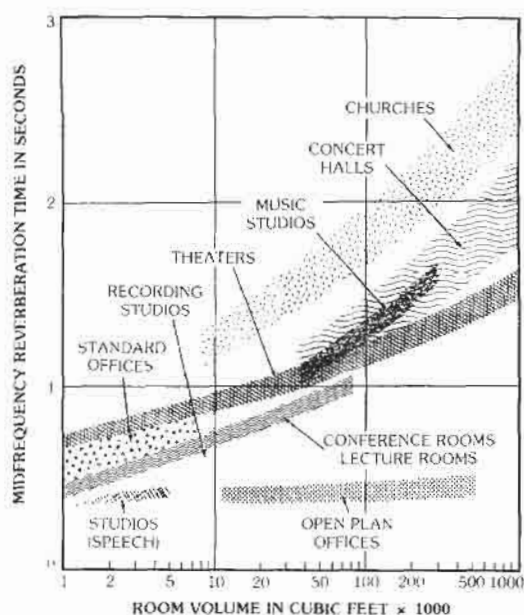


Fig. 7. Optimum midfrequency reverberation times.

Interior Design Noise Levels

Table 11 contains a list of planning levels for activities conducted in interior spaces.* The (interior) planning levels for *exterior noise and for interior equipment*

* In offices or office-type environments, Preferred Noise Criteria Curves (PNC) are frequently used to select and specify the ambient.

that is not continuously operated are given in terms of L_{eq} values. Continuous noise sources, for example ventilating systems or other mechanical equipment, emit steady-state noise that is measured in terms of L_s .* The sources must be considered separately. As indicated in the right-hand column of Table 11, permissible L_s values are 5 and 10 dB less than L_{eq} values for the same activity.

Interior steady-state noise levels more than 5–10 dB below the levels specified in Table 11 are not desirable. Annoyance will actually increase with the lowered background noise levels because individuals will hear intruding sounds that normally would be masked by the steady-state noise. Occasionally, where adequate noise insulation cannot be provided, increasing the continuous background noise levels over the values shown in Table 11 will provide better masking of intruding intermittent sounds. For such occasions, the characteristics of both the intruding noises and the background noise should be considered during the design of the facility.

SOUND ISOLATION†

Sound entering a room from external areas can be controlled by quieting the source, by adding absorption to the source room (if the source is indoors) to lower the

* L_s is the A-weighted noise level produced by the ventilation or mechanical systems (or other interior noise sources) which operate more or less continuously. The L_s value for design should be the noise level produced in the space during the time of occupancy while the equipment is at the typical mode of operation.

† See, for example, references 2.1 through 2.11.

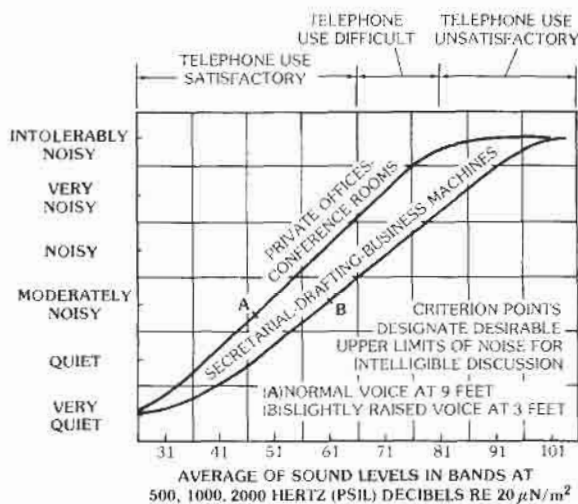


Fig. 8. Rating chart for office noises. Data were determined by an octave-band analysis and correlated with subjective tests. (From Peterson, A. and Gross, E. Handbook of Noise Measurement. General Radio, 1972. Modified from Beranek and Newman for preferred bands.)

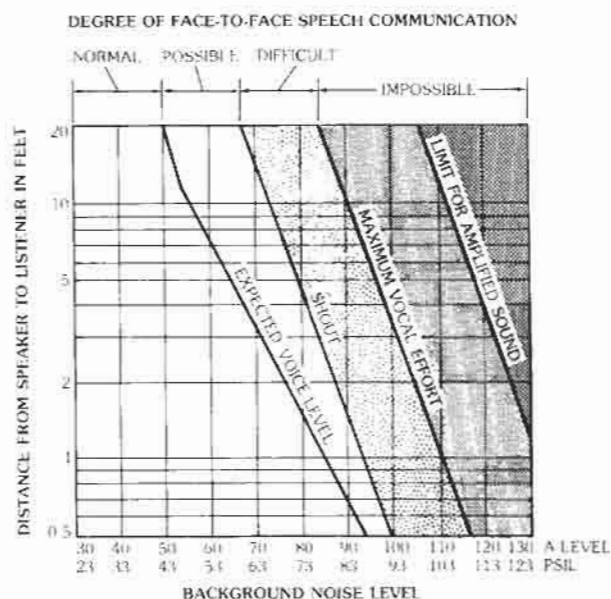


Fig. 9. Degree of face-to-face speech communication between speaker and listener as a function of distance and background noise level (A-levels or PSIL). (After Contractor Report to Army for TW-5-838-2.)

TABLE 11. INTERIOR NOISE ENVIRONMENT PLANNING LEVELS

Activity	All Noise Sources L_{eq} (dB)	Continuous Interior Sources* L_T (dB)**
Sleeping	45	40
Other residential activities (conversations, radio or TV listening, etc.)	50	40
Classrooms, libraries, churches, hospitals	45	40
Offices—private, conference	45	40
Offices/work spaces, telephone use satisfactory	55	45
Work spaces—occasional speech communication or telephone use	60	55
Work spaces—infrequent speech communication, telephone use infrequent	70	60
Concert halls, large auditoriums, theaters, or churches	30	25
Broadcast, TV, and recording studios	35	30
Small auditoriums, theaters, or churches	45	40

* Typically, ventilation systems and mechanical equipment in near-continuous operations.

** The L_T value is given in terms of A-weighted noise level. The approximate noise criteria (NC) curve values are 8 dB less than the A-level values.

ambient level in the source room and hence the energy striking the party* surface, and by increasing the transmission loss of the party surface. Sound travels from sources into receiving rooms directly through building elements such as walls, windows, or doors or indirectly through "sound leaks" such as cracks, HVAC, plumbing, electrical facilities, plenums above suspended ceilings, corridors, etc.

Sound isolation should be contrasted with absorption. Efficient sound absorbers such as glass fiber or open-cell polyurethane foam do not contain or isolate sound. Efficient sound isolators or containers such as heavy walls do not absorb sound.

In the United States, the sound isolating value of walls and other elements is frequently specified in terms of STC (see ASTM E413-87). More recently, it is also specified in terms of ΔL_A , the A-weighted sound reduction between typically furnished small rooms (see ASTM E597-87). Where usable, ΔL_A has the advantage of easier field measurement than STC.

Walls

Figs. 10, 11, and 12 illustrate many building details and their corresponding STC and ΔL_A values. The total noise reduction between a source room and a receiving room is a function of the area of the party surface and the absorption in the receiving room. Small amounts of energy entering a very reverberant room will build up to high levels. Large wall areas with adequate STC being used to isolate a recording studio from a machine room will likely fail. Thus, party boundaries between noise-critical and noisy spaces should be minimized or eliminated.

Doors

Solid-core wood doors and hollow metal doors filled with insulation are generally better sound isolators than hollow-core doors; however, the amount of air space around the edges of the door is usually the controlling sound path. (See also the section on sound leaks.) For maximum sound isolation there must be a soft gasketing or weather stripping around the door to provide an airtight seal. The gasket must not be so stiff that it is difficult to close the door.

Resilient Channels

Frequent mention is made of the use of *resilient channels* to improve the noise-isolation performance of typical constructions. Typically, these channels are made of light-gauge metal in a Z-section, with one flange of the channel firmly attached to a stud or joist, and the other flange of the channel providing support for one or two layers of gypsum wall board. Because of the

spring-like action of the metal between the two flanges, the gypsum board is able to vibrate independently of the wall supports.

In normal rigid constructions, sound energy from the source room passes directly from one surface (wall or floor), through the structural support (studs or joists), to the surface on the opposite side, causing this second surface to vibrate and thus reradiate sound energy from the source room. However, resilient channels reduce this structural continuity. Now, most of the sound energy from one side of the wall or floor will pass through the airspace between the studs or joists. This structural discontinuity greatly improves the sound isolation of the construction. Because the sound path in this case is through the wall or floor cavity, it is especially useful to have a sound-absorbing blanket of glass fiber or mineral wool insulation in the cavity air space. When resilient channels are not used, this insulation offers very little benefit, but when the resilient channels are used, this insulation should always be added.

It is very easy to short-circuit the isolation provided by the resilient channel. For example, wall-hung cabinets require rigid contact with the wall; they must be fastened directly to the studs, thus eliminating the isolation of the resilient channel. In general, use of the channels for wall constructions is not recommended; however, in special instances and for ceilings, which are out of direct contact and abuse, resilient channels can be very helpful.

For ceilings and walls, resilient channels might typically be attached 24 inches on center at right angles to the floor joists or studs. Gypsum board is then attached to the channels with screws (Fig. 13).

For wall installations, manufacturers often recommend that a 3-inch-wide continuous filler strip be applied at the base. This will provide good mechanical support for the baseboard and gypsum board against damage without sacrificing much resilience.

Weighting of Components

Each structural component (window, door, ceiling plenum, etc.) in a composite wall or sound path will contribute sound energy to the receiving room in proportion to the percentage of the total area occupied by that element. Therefore the ΔL_A value for the element must be weighted to correct for the area it occupies.

To weight the contribution of each element, use Table 12. For example, a wood solid-core door has a ΔL_A rating of 27 when fully gasketed; it occupies 20% of the wall area with a neighboring office. Therefore, its weighted contribution to the noise reduction between rooms is $27 + 7 = 34 \Delta L_A$.

By looking at the weighted ΔL_A value for each structural component or flanking path or leak (described later), it is readily apparent what the controlling noise path is in any isolation analysis. Consider the

* The party surface is the surface (e.g., wall or ceiling) common to two spaces.

1 2x4 WOOD STUDS, 16" ON CENTER WITH 1/2" GYPSUM BOARD ON EACH SIDE		ΔI_A /STC 35/33
INSULATION ADDED IN THE AIR SPACE		39/38
RESILIENT CHANNELS ADDED ON ONE SIDE		44/41
2 2x4 WOOD STUDS, STAGGERED ON 2x6 PLATE WITH TWO LAYERS OF 1/2" GYPSUM BOARD ON EACH SIDE		46/44
3 2x4 WOOD STUDS, DOUBLED, EACH ON SEPARATE PLATE WITH 3/4" INSULATION IN AIR SPACE		55/54

(A) Wood stud framing.

THE EFFECTS ON THE ΔI_A RATING FOR METAL-STUD PARTITIONS DUE TO VARIOUS ADDITIONS OR CHANGES IN THE COMPONENTS OF THE BASIC PARTITION ARE SUMMARIZED BELOW

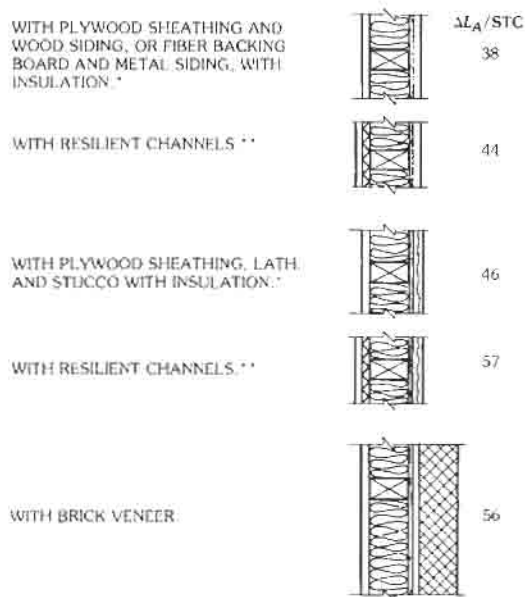
1. BASIC PARTITION. 2 1/2" STEEL SCREW STUDS, 24" ON CENTER SINGLE LAYER OF 1/2" GYPSUM BOARD ON EACH SIDE. (SUBSTITUTION OF 5/8" WALLBOARD WILL NOT SUBSTANTIALLY CHANGE THE ΔI_A RATING.)		ΔI_A /STC 37/38
2. ADD 1" OF GLASS FIBER INSULATION OR 1 LAYER OF WALLBOARD		ADD 3
3. ADD SECOND ITEM FROM 2 ABOVE.		ADD 3
4. ADD WALLBOARD TO OTHER SIDE.		ADD 3
5. ADD ADDITIONAL GLASS FIBER INSULATION, UP TO 3".		ADD 2 PER INCH
6. IF MORE THAN TWO ADDITIONS, SUBTRACT 1.		SUBTRACT 1
7. FOR 1-5/8" STUDS		SUBTRACT 1
FOR 3-5/8" STUDS		ADD 1

(B) Metal stud partitions.

1. 4" DENSE HOLLOW-CORE BLOCK, PAINTED		ΔI_A /STC 41/38
PLASTER ADDED ON BOTH SIDES.		45/42
2. 6" DENSE HOLLOW-CORE, PAINTED		44/43
3. 8" DENSE HOLLOW CORE WITH 4" BRICK ONE SIDE.		55/54
4. 8" LIGHTWEIGHT BLOCK, PAINTED		45/45
1x2 FURRING, 16" O.C. WITH 1/2" GYPSUM BOARD ADDED ON ONE SIDE.		50/50

(C) Concrete and concrete block.

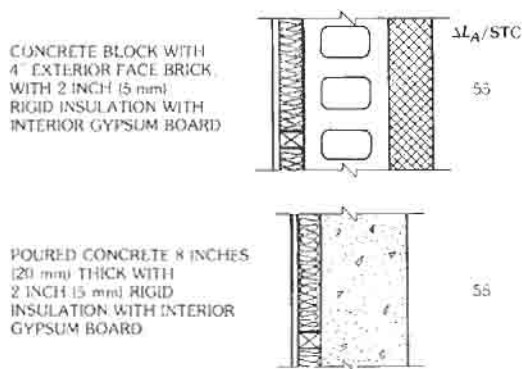
Fig. 10. Interior-wall sound isolation.



* It is assumed that all new wood frame construction would have glass fiber or mineral wool insulation in the stud cavity.

** Always use sound-absorptive insulation with resilient channels.

(A) Wood stud framing.

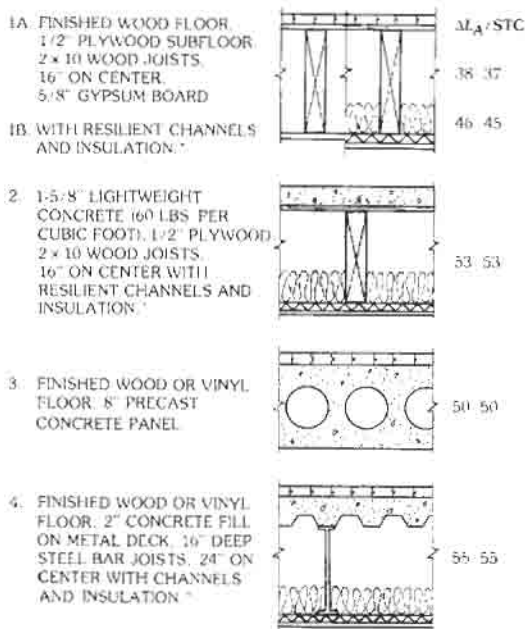


(B) Concrete and masonry.

Fig. 11. Exterior-wall sound isolation.

example in Table 13. The door in this example is the "weak link." Improvements to other components will be of no benefit unless the ΔL_A value of the door is upgraded or its area is reduced.

The acoustical contributions for all sound-transmitting structural components (wall, door, window, etc.) must be combined to determine the performance of the composite wall. Since ΔL_A values for each element represent a reduction in sound from one side of



Insulation and resilient channels should always be used together.

Fig. 12. Floor/ceiling-system sound isolation.

the sound path to the other, the overall, combined ΔL_A value will decrease when the elements are combined (Table 14). For example: To determine the composite ΔL_A for the wall example, first combine ΔL_A 36 and ΔL_A 34. These two values differ by 2 dB; therefore, subtract 2 dB from the lower ΔL_A value: ΔL_A 34 less 2 dB = ΔL_A 32. Next, combine ΔL_A 32 and ΔL_A 42. These two values differ by more than 9 dB; therefore, subtract nothing from the lower ΔL_A value: ΔL_A stays at 32.

Flanking

Flanking is the transfer of sound through secondary sound paths that go around the common wall barrier between two spaces. These indirect passages can seriously degrade the sound isolation performance of an otherwise acceptable construction.

Typical flanking paths, A, B, C, D, and E, are shown in Fig. 14. In addition, flanking paths F and G are structural connections, which are less suitable for analysis with specific ΔL_A sound-isolation performance ratings. In these cases, a sound from a source room will be transmitted to a receiver room by way of a common wall or floor surface.

It is most important to break the structural paths that connect two spaces. Floors, rafters, and joists should not span between separate dwelling units. Interior

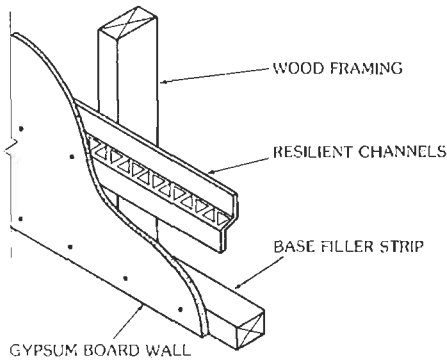


Fig. 13. Installation of resilient channels.

wallboard skins should not span between noise-sensitive spaces.

Leaks

Leaks are sound paths through open air passages caused by poor detailing, inadequate field supervision, and sloppy workmanship. They represent a major cause for poor performance of partitions and can be easily resolved by proper caulking and sealing of all openings and penetrations between spaces. Typical locations of leaks are shown in Fig. 15.

The noise isolation value of an untreated hole in a wall is, of course, $\Delta L_A = 0$. Depending on the size of the leak, the noise isolation performance of an entire wall or other structural component can be severely degraded. Fig. 16 plots this effect.

For example, consider a perfectly good composite construction with $\Delta L_A = 50$ and a typical area of 100 square feet (9 square meters). A crack only 0.01 inch (2.5 mm) high along the base of a wall 12 feet (4 meters) long would be a total opening of about 1.4

TABLE 12. WEIGHTING FACTORS FOR COMBINING PARTY SURFACE COMPONENTS

Percent of Total Area Occupied by Component	Weighting Factor for ΔL_A /STC
100	Same value
50	Add 3 points
33	Add 5 points
20	Add 7 points
10	Add 10 points
5	Add 13 points
2	Add 17 points
1	Add 20 points
0.1	Add 30 points

TABLE 13. WALL EXAMPLE

Component	Intrinsic STC/ ΔL_A Value	Percent of Area	Weighted STC/ ΔL_A Value
Window	26	10	36
Door	27	20	34
Wall	40	70	42
Flanking (floor)	50	—	50

TABLE 14. COMBINED ΔL_A VALUE

Where two ΔL_A values differ by	Subtract the following amount from the lower value ΔL_A
0-1	3
2-3	2
4-8	1
9 or more	0



- A OPEN PLENUM ABOVE WALLS
- B UNBAFFLED DUCTS
- C OUTDOOR PATH, WINDOW TO WINDOW
- D CONTINUOUS UNBAFFLED INDUCTOR UNITS
- E CORRIDOR PATH
- F COMMON WALL
- G COMMON FLOOR

Fig. 14. Flanking paths.

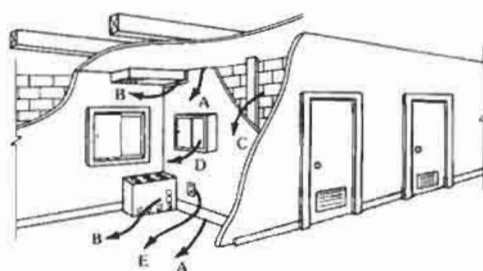
square inches. The ΔL_A performance would be degraded from 50 to 39. Also consider that a typical properly installed door with a 1-inch undercut for rugs or ventilation has an opening of about 40 square inches.

SOUND-REINFORCEMENT SYSTEMS*

Listeners should be almost unaware of the operation of a sound-reinforcement system. There should be adequate level, high intelligibility, and natural sound.

The purpose of a sound system is to boost the message or program sufficiently above the ambient level. Ideally, the amplification at the receiver should be

* Yerges, Lyle F. *Sound, Noise, and Vibration Control*. New York: Van Nostrand-Reinhold Co., 1969.

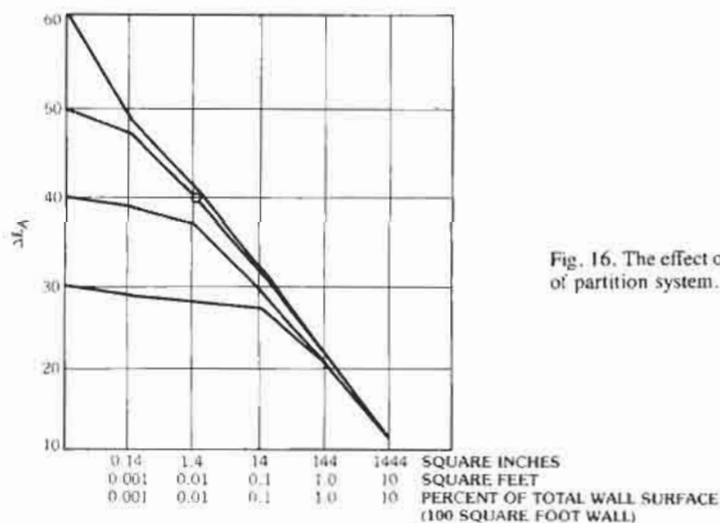


- A. POOR SEAL AT EDGES OF CEILING, FLOOR, AND WALLS
- B. POOR SEAL AROUND DUCT PENETRATIONS, GAPS AND HOLES AT PIPE PENETRATIONS
- C. POOR MORTAR JOINTS AND POROUS MASONRY
- D. BACK TO BACK CABINETS (AND BATH TUBS) WITHOUT WALL BEHIND
- E. BACK TO BACK ELECTRICAL OUTLETS

Fig. 15. Sound leaks.

unity. That is, the sound level at a receiver's ear should be about 65–70 dB, or the level generated by a speaker at a listening distance of 1 meter. In most auditoria or lecture rooms, the background level is not a problem, but in sports arenas, spectator noise can raise the ambient such that the direct sound amplification at the receiver must be much greater than unity. It has been found that about 9 dB of amplification at the receiver is the most that can go undetected—the sports-arena situation frequently requires more.

In general, a sound-reinforcement system may contain a single loudspeaker cluster or a distribution of loudspeakers. If the sound originates from a stage or performance arena, then a single cluster should be used if the sound is to be natural. An arrangement of two loudspeakers flanking the stage is inappropriate; the two speakers create interference patterns and an unnatural sound. If one speaker cluster cannot maintain sufficient level all the way to the rear of an audience, then

Fig. 16. The effect of leaks on ΔL_A value of partition system.

additional speakers can be added closer to the rear with appropriate time-delay hardware so that the sound remains realistic and in correct time phase with the source. Distributed loudspeaker systems are used when there is no identifiable source. Situations requiring distributed systems include paging, such as in offices, and background music, such as in food stores.

A rule of thumb for placing loudspeakers in a distributed system is:

$$D = 2(H - 4)$$

where,

D = distance between speaker centers in feet,
 H = ceiling height in feet.

or

$$D = 2(H - 1)$$

where D and H are in meters.

The center of a single speaker should be aimed at the most distant centrally located receiver. The directivity pattern (vertically and horizontally) should be chosen to provide a uniform sound field to the audience given the elevation and plan of the room.*

REFERENCES ON ENVIRONMENTAL NOISE

- Schultz, Theodore J., "Synthesis of social surveys on noise annoyance." *Journal of the Acoustical Society of America*, 64(2): August 1978, pp. 337-405.
- Public Health and Welfare Criteria for Noise*. US Environmental Protection Agency Office of Noise Abatement and Control. Report Number 550/9-73-002, July 27, 1973.†
- Information on Levels of Environmental Noise Requisite to Protect Public Health and Welfare with an Adequate Margin of Safety*. U.S. Environmental Protection Agency Office of Noise Abatement and Control, Report Number 550/9-74-004, March 1974.†
- American National Standard Quantities and Procedures for Description and Measurement of Environmental Sound. Part 1*. American National Standards Institute (ANSI) Standard S12.9-1988.
- Guidelines for Considering Noise In Land Use Planning and Control*. Federal Interagency Committee on Urban Noise, June 1980.
- Army Regulation 200-1, 1990. *Environmental Protection and Enhancement, Chapter 7. Environmental Noise Abatement Program*. Department of the Army, Washington D.C., 23 April 1990.
- Annoyance, Loudness, and Measurement of Repetitive Type Impulsive Noise Sources*. US Environmental Protection Agency Office of Noise Abatement and Control, Report Number 550/9-79-103, November 1979.

SELECTED BIBLIOGRAPHY

1. General Acoustics

- Kinsler, Frey, Coppens, and Sanders. *Fundamentals of Acoustics*. 3rd ed. New York: John Wiley and Sons, Inc., 1982.
- Morse and Ingard. *Theoretical Acoustics*. New York: McGraw-Hill Book Co., 1968.
- Rayleigh, J. *The Theory of Sound*. New York: Dover Publications, 1945 (from the 1894 edition).
- Swenson, George W. *Principles of Modern Acoustics*. New York: Van Nostrand, 1953.
- Pierce, Allan D. *Acoustics*. Woodbury, NY: Acoustical Society of America, 1989.

2. Building Acoustics

- Department of Housing and Urban Development (HUD), *A Guide to Airborne, Impact and Structure Borne Noise Control in Multi-Family Dwellings*, FT/TS-24, 1967.
- Dept. of Commerce, *Acoustical and Thermal Performance of Exterior Residential Walls, Doors and Windows*, NBS Building Science Series No. 77, Nov., 1975.
- Acoustical Manual: Apartment and Home Construction*. Rockville, Md.: National Association of House Builders (NAHB), Research Foundation, No. 315.04, June, 1971.
- Guide and Data Book*. New York: American Society of Heating, Refrigeration and Air-Conditioning Engineers (ASHRAE).
- Ginn, K. *Architectural Acoustics*. Denmark: Brüel and Kjaer, 1978.
- Knudsen and Harris. *Acoustical Design in Architecture*. New York: John Wiley and Sons, Inc., 1950.
- Newman, R., et al. *Acoustics, Time Saver Standard*. New York: McGraw-Hill Book Co., 1973.
- Dept. of Commerce, *Quieting: A Practical Guide to Noise Control*, NBS Handbook No. 119, July, 1971.
- Sound Research Laboratories, Ltd. *Practical Building Acoustics*. London: E. and F. N. Spon, Ltd.; or New York: Halsted Press, a division of John Wiley and Sons, Inc.; 1976.
- Sweet's Architectural Catalog File*. New York: Sweet's Division of McGraw-Hill Information Systems Co.

* Sound systems in large and/or critical areas (e.g., playhouses are dealt with in reference 4.1).

† The best summary of references 2 and 3 is: *Protective Noise Levels—Condensed Version of EPA Levels Document*. US Environmental Protection Agency, EPA 550/9-79-100, November 1978.

- 2.11 Yerges, Lyle F. *Sound, Noise, and Vibration Control*. New York: Van Nostrand Reinhold Co., 1969.

3. Environmental Noise Assessment and Land Use Planning With Respect to Noise

- 3.1 Beranck, Leo L., ed. *Noise and Vibration Control*. New York: McGraw-Hill Book Co., 1971.
- 3.2 Dept. of Housing and Urban Development (HUD), Office of Policy Development and Research, *Criteria for Acceptable Outdoor and Indoor Acoustical Environments for Schools, Hospitals, Nursing and Convalescent Homes, Churches, Business Offices and Other Community Facilities*, June, 1976.
- 3.3 Dept of Commerce (NBS), *Design Guide for Reducing Transportation Noise in and Around Buildings*, Building Science Series 84, April, 1978.
- 3.4 Dept. of the Army, *Environmental Protection: Planning in the Noise Environment*, Technical Manual (TM) 5-803-2, 15 June 1978.
- 3.5 *Environmental Assessment of High Energy Impulsive Sounds*, Report of Working Group 84, Committee on Hearing, Bioacoustics and Biomechanics (CHBAB), Assembly of Behavioral and Social Sciences, National Research Council, National Academy of Science, 1981.
- 3.6 Dept. of Transportation, Federal Interagency Committee on Urban Noise, *Guidelines for Considering Noise in Land Use Planning and Control*, June, 1980.
- 3.7 *Guidelines for Preparing Environmental Impact Statements on Noise*, Committee on Hearing, Bioacoustics and Biomechanics, (CHBAB), Assembly of Behavioral and Social Sciences, National Research Council, National Academy of Science, 1977.
- 3.8 EPA, *Information on Levels of Environmental Noise Requisite to Protect Public Health and Welfare with an Adequate Margin of Safety*, 550/9-74-004, March, 1974.
- 3.9 Kryter, Karl. *The Effects of Noise on Man*. New York: Academic Press, 1970.
- 3.10 von Gierke, H. "Noise—How Much Is Too Much." *Noise Control Engineering*, Vol. 5, No. 1, 1975.
- 3.11 Suter, Alice H. "Noise Sources and Effects—A New Look." *Sound and Vibration*, 26, January 1992, pp. 18–38. (This article contains an extensive bibliography on noise and its human effects.)

4. Noise Measurement and Instrumentation

- 4.1 Davis, Don and Carolyn. *Sound System Engineering*, 2nd ed. Indianapolis: Howard W. Sams & Co., Inc., 1987.

- 4.2 Hassall, J. R., and Zaveri, K. *Acoustic Noise Measurement*, Denmark: Brüel and Kjaer, 1979.
- 4.3 Peterson, A. G., and Gross, E. E., Jr. *Handbook of Noise Measurement*. Concord, Mass.: General Radio Co., 1972.
- 4.4 Harris, Cyril M., ed. *Handbook of Acoustical Measurements and Noise Control*. 3rd ed. New York: McGraw-Hill, 1991.

5. Selected Acoustical Standards

AMERICAN NATIONAL STANDARDS INSTITUTE 1430
Broadway, New York, NY 10018

- ANSI S1.4-1983, Sound Level Meters, Specification for
- ANSI S1.4a-1985, Sound Level Meters
- ANSI S1.6-1984(R1990), Preferred Frequencies, Frequency Levels, and Band Numbers for Acoustical Measurements
- ANSI S1.8-1989, Reference Quantities for Acoustical Levels
- ANSI S1.10-1966(R1986), Calibration of Microphones, Method for the
- ANSI S1.11-1986, Octave-band and Fractional Octave-Band Analog and Digital Filters, Specifications for
- ANSI S1.12-1967(R1986), Laboratory Standard Microphones, Specifications for
- ANSI S1.13-1971(R1986), Sound Pressure Levels, Methods for the Measurement of
- ANSI S1.20-1988, Procedures for Calibration of Underwater Electroacoustic Transducers
- ANSI S1.26-1978(R1989), Absorption of Sound by the Atmosphere, Method for the Calculation of the
- ANSI S1.40-1984(R1990), Specifications for Acoustical Calibrators
- ANSI S1.42-1986, Weighting Networks for Acoustical Measurements, Design Response of
- ANSI S3.1-1977(R1986), Permissible Ambient Noise During Audiometric Testing, Criteria for
- ANSI S3.2-1989, Method for Measuring the Intelligibility of Speech over Communications System
- ANSI S3.3-1960(R1990), Methods for Measurement of Electroacoustical Characteristics of Hearing Aids
- ANSI S3.4-1980(R1986), Procedure for the Computation of Loudness of Noise
- ANSI S3.5-1969(R1986), Methods for the Calculation of the Articulation Index
- ANSI S3.6-1989, Specifications for Audiometers
- ANSI S3.7-1973(R1986), Coupler Calibration of Earphones, Method for
- ANSI S3.13-1987, Mechanical Coupler for Measurement of Bone Vibrators
- ANSI S3.14-1977(R1986), Rating Noise with Respect to Speech Interference
- ANSI S3.17-1975(R1980), Rating the Sound Power Spectra of Small Stationary Noise Sources, Method for

- ANSI S3.18-1979(R1986), Whole-Body Vibration, Guide for the Evaluation of Human Exposure to
- ANSI S3.19-1974(R1990), Method for Measurement of Real-Ear Protection of Hearing Protectors and Physical Attenuation of Earmuffs
- ANSI S3.20-1973(R1978), Psychoacoustical Terminology
- ANSI S3.21-1978(R1986), Manual Pure-Tone Threshold Audiometry, Method for
- ANSI S3.22-1987, Hearing Aid Characteristics, Specification of
- ANSI S3.36-1985(R1990), Specification for Manikin for Simulated in Situ Airborne Acoustic Measurements
- ANSI S3.37-1987, Preferred Earhook Nozzle Thread for Postauricular Hearing Aids
- ANSI S3.39-1987, Specifications for Instruments to Measure Aural Acoustic Impedance and Admittance (Aural Acoustic Immittance)
- ANSI S12.7-1986, Methods for Measurement of Impulse Noise
- ANSI S12.8-1987, Methods for Determination of Insertion Loss of Outdoor Noise Barriers
- ANSI S12.9-1988, Quantities and Procedures for Description and Measurement of Environmental Sound, Part 1
- ANSI S12.10-1985(R1990), Methods for the Measurement and Designation of Noise Emitted by Computer and Business Equipment
- ANSI S12.1-1983, Preparation of Standard Procedures to Determine the Noise Emission from Sources, Guidelines for the
- ANSI S12.3-1985(R1990), Statistical Methods for Determining and Verifying Stated Noise Emission Values of Machinery and Equipment
- ANSI S12.4-1986, Method for Assessment of High-Energy Impulsive Sounds with Respect to Residential Communities
- ANSI S12.5-1990, Reference Sound Sources, Requirements for the Performance and Calibration of
- ANSI S12.6-1984(R1990), Method for the Measurement of the Real-Ear Attenuation of Hearing Protectors
- ANSI S12.11-1987, Methods for the Measurement of Noise Emitted by Small Air-Moving Devices
- ANSI S12.23-1989, Method for the Designation of Sound Power Emitted by Machinery and Equipment
- ANSI S12.30-1990, Sound Power Standards and the Preparation of Noise Test Codes, Guidelines for the Use of
- ANSI S12.31-1990, Broad-Band Noise Sources in Reverberation Rooms, Precision Methods for the Determination of Sound Power Levels of
- ANSI S12.32-1990, Discrete-Frequency and Narrow-Band Noise Sources in Reverberation Rooms, Precision Methods for the Determination of Sound Power Levels of
- ANSI S12.33-1990, Sound Power Levels of Noise Sources in a Special Reverberation Test Room, Engineering Methods for the Determination of
- ANSI S12.34-1988, Free-Field Conditions over a Reflecting Plane, Engineering Methods for the Determination of Sound Power Levels of Noise Sources for Essentially
- ANSI S12.35-1990, Sound Power Levels of Noise Sources in Anechoic and Hemi-Anechoic Rooms, Determination of
- ANSI S12.36-1990, Sound Power Levels of Noise Sources, Survey Methods for the Determination of (redesignation of ANSI S1.36-1979)
- ANSI S12.40-1990, Compatible Land Use, Sound Level Descriptors for Determination of
- ANSI PH22.202M-1984, B Chain Electro-Acoustic Response—Control Rooms and Indoor Theaters
- INTERNATIONAL STANDARDS ORGANIZATION (ISO)
Central Secretariat, Case Postale 56, CH-1211, Geneva 20, Switzerland
- ISO 362:1981, Acoustics—Measurement of noise emitted by accelerating road vehicles—Engineering method
Ed. 1 4 p. Code C TC 43 HB 4
Amendment 01-1985
Ed. 1 1 p. Code A TC 43
- ISO 1680-1:1986, Acoustics—Test code for the measurement of airborne noise emitted by rotating electrical machinery—Part 1: Engineering method for free-field conditions over a reflecting plane
Ed. 1 15 p. Code J TC 38
- ISO 1680-2:1986, Acoustics—Test code for the measurement of airborne noise emitted by rotating electrical machinery—Part 2: Survey method
Ed. 1 14 p. Code H TC 43
- ISO 2151:1972, Measurement of airborne noise emitted by compressor/primemover-units intended for outdoor use (To be replaced by future ISO 3989 of TC 43)
Ed. 1 7 p. Code D TC 118 HB 4
- ISO 2922:1975, Acoustics—Measurement of noise emitted by vessels on inland water-ways and harbours
Ed. 1 4 p. Code C TC 43 HB 4
- ISO 2923:1975, Acoustics—Measurement of noise on board vessels
Ed. 1 4 p. Code C TC 43 HB 4
- ISO 3095:1975, Acoustics—Measurement of noise emitted by railbound vehicles
Ed. 1 6 p. Code D TC 43 HB 4
- ISO 3381:1976, Acoustics—Measurement of noise inside railbound vehicles
Ed. 1 4 p. Code C TC 43 HB 4
- ISO 3891:1978, Acoustics—Procedure for describing aircraft noise heard on the ground
Ed. 1 24 p. Code N TC 43 HB 4
- ISO 4412-1:1979, Hydraulic fluid power—Test code for the determination of airborne noise levels—Part 1: Pumps
Ed. 1 6 p. Code D TC 131 HB 25
- ISO 4412-2:1984, Hydraulic fluid power—Test code

- for the determination of airborne noise levels—Part 2: Motors
Ed. 1 6 p. Code D TC 131 HB 25
- ISO 4871:1984, Acoustics—Noise labelling of machinery and equipment
Ed. 1 5 p. Code D TC 43 HB 4
- ISO 4872:1978, Acoustics—Measurement of airborne noise emitted by construction equipment intended for outdoor use—Method for checking compliance with noise limits
Ed. 1 11 p. Code G TC 43 HB 4
- ISO 5128:1980, Acoustics—Measurement of noise inside motor vehicles
Ed. 1 6 p. Code D TC 43 HB 4
- ISO 5129:1987, Acoustics—Measurement of noise inside aircraft
Ed. 2 4 p. Code C TC 43
- ISO 5130:1982, Acoustics—Measurement of noise emitted by stationary road vehicles—Survey method
Ed. 1 6 p. Code D TC 43 HB 4
- ISO 5131:1982, Acoustics—Tractors and machinery for agriculture and forestry—Measurement of noise at the operator's position—Survey method
Ed. 1 8 p. Code E TC 43 HB 4
- ISO 5135:1984, Acoustics—Determination of sound power levels of noise from air terminal devices, high/low velocity/pressure assemblies, dampers and valves by measurement in a reverberation room
Ed. 1 11 p. Code G TC 43 HB 4
- ISO 6081:1986, Acoustics—Noise emitted by machinery and equipment—Guidelines for the preparation of test codes of engineering grade requiring noise measurements at the operator's or bystander's position
Ed. 1 9 p. Code F TC 43
- ISO 6190:1988, Acoustics—Measurement of sound pressure levels of gas turbine installations for evaluating environmental noise—Survey method
Ed. 1 6 p. Code D TC 43
- ISO 6393:1985, Acoustics—Measurement of airborne noise emitted by earth-moving machinery—Method for determining compliance with limits for exterior noise—Stationary test condition
Ed. 1 7 p. Code E TC 43 TC 127 HB 4
- ISO 6394:1985, Acoustics—Measurement of airborne noise emitted by earth-moving machinery—Operator's position—Stationary test condition
Ed. 1 5 p. Code D TC 43 TC 127 HB 4
- ISO 6395:1988, Acoustics—Measurement of exterior noise emitted by earth-moving machinery—Dynamic test conditions
Ed. 1 13 p. Code H TC 43 TC 127
- ISO 7182:1984, Acoustics—Measurement at the operator's position of airborne noise emitted by chain saws
Ed. 1 3 p. Code C TC 43 TC 23 HB 4
- ISO 7188:1985, Acoustics—Measurement of noise emitted by passenger cars under conditions representative of urban driving
Ed. 1 5 p. Code D TC 43 HB 11
- ISO 7574-1:1985, Acoustics—Statistical methods for determining and verifying stated noise emission values of machinery and equipment—Part 1: General considerations and definitions
Ed. 1 4 p. Code C TC 43
- ISO 7574-2:1985, Acoustics—Statistical methods for determining and verifying stated noise emission values of machinery and equipment—Part 2: Methods for stated values for individual machines
Ed. 1 2 p. Code B TC 43
- ISO 7574-3:1985, Acoustics—Statistical methods for determining and verifying stated noise emission values of machinery and equipment—Part 3: Simple (transition) method for stated values for batches of machines
Ed. 1 2 p. Code B TC 43
- ISO 7574-4:1985, Acoustics—Statistical methods for determining and verifying stated noise emission values of machinery and equipment—Part 4: Methods for stated values for batches of machines
Ed. 1 14 p. Code H TC 43
- ISO 7779:1988, Acoustics—Measurement of airborne noise emitted by computer and business equipment
Ed. 1 39 p. Code S TC 43
- ISO/TR 7849:1987, Acoustics—Estimation of airborne noise emitted by machinery using vibration measurement
Ed. 1 20 p. Code L TC 43
- ISO 7917:1987, Acoustics—Measurement at the operator's position of airborne noise emitted by brush saws
Ed. 1 4 p. Code C TC 43 TC 23
- ISO 9295:1988, Acoustics—Measurement of high-frequency noise emitted by computer and business equipment
Ed. 1 11 p. Code G TC 43
- ISO 9296:1988, Acoustics—Declared noise emission values of computer and business equipment
Ed. 1 7 p. Code E TC 43
- ISO Handbook 4—Acoustics, vibration and shock
 Standards used in examining the physical and environmental effects of sound, mechanical vibration and shock, and in the balancing of machinery. The standards included are under: basic documents for noise measurements and limits; noise measurement of various sources; building acoustics; general acoustics; terminology in vibration and shock; balancing and balancing machines; human exposure to vibration and shock.
 Price group XZ. 850 pages. (2nd ed. 1985)
 ISBN 92-67-10106-4
- AMERICAN SOCIETY OF TESTING AND MATERIALS (ASTM)
 1916 Race St., Philadelphia, PA 19103-1187
- C 367-78 (Reapproved 89), Standard Test Methods for Strength Properties of Prefabricated Architectural Acoustical Tile of Lay-In Ceiling Panels

- C 384-90a, Standard Test Method for Impedance and Absorption of Acoustical Materials by the Impedance Tube Method
- C 423-90a, Standard Test Method for Sound Absorption and Sound Absorption Coefficients by the Reverberation Room Method
- C 522-87, Standard Test Method for Airflow Resistance of Acoustical Materials
- C 634-89, Standard Terminology Relating to Environmental Acoustics
- C 635-91, Standard Specification for the Manufacture, Performance, and Testing of Metal Suspension Systems for Acoustical Tile and Lay-In Panel Ceilings.
- C 636-91, Standard Practice for Installation of Metal Ceiling Suspension Systems for Acoustical Tile and Lay-In Panels
- C 800-89, Standard Specification for Glass Fiber Blanket Insulation (Aircraft Type)
- C 919-84 (Reapproved 88), Standard Practice for Use of Sealants in Acoustical Applications
- C 971-82, Standard Guide for Selection and Application of Insulation Systems for Heating, Ventilating, and Air Conditioning Duct Work
- C 1014-88, Standard Specification for Spray-Applied Mineral Fiber Thermal of Acoustical Insulation
- C 1071-86, Standard Specification for Thermal and Acoustical Insulation (Mineral Fiber, Duct Lining Material)
- C 1149-90, Standard Specification for Self-Supported Spray Applied Cellulosic Thermal/Acoustical Insulation
- D 1779-65 (Reapproved 83), Standard Specification for Adhesive for Acoustical Materials
- E 90-90, Standard Test Method for Laboratory Measurement of Airborne Sound Transmission Loss of Building Partitions
- E 477-90, Standard Test Method for Measuring Acoustical and Airflow Performance of Duct Liner Materials and Prefabricated Silencers
- E 492-90, Standard Test Method for Laboratory Measurement of Impact Sound Transmission Through Floor-Ceiling Assemblies Using the Tapping Machine
- E 497-89, Standard Practice for Installing Sound-Isolating Lightweight Partitions
- E 557-77 (Reapproved 90), Standard Practice for Architectural Application and Installation of Operable Partitions
- E 580-91, Standard Practice for Application of Ceiling Suspension Systems for Acoustical Tile and Lay-In Panels in Areas Requiring Seismic Restraint
- E 596-90, Standard Test Method for Laboratory Measurement of the Noise Reduction of Sound-Isolating Enclosures
- E 597-91 (Reapproved 87), Standard Practice for Determining a Single-Number Rating of Airborne Sound Isolation for Use in Multiunit Building Specifications
- E 717-84 (Reapproved 89), Standard Guide for Preparation of the Accreditation Annex of Acoustical Test Standards
- E 989-89, Standard Classification for Determination of Impact Insulation Class (IIC)
- E 1041-85 (Reapproved 90), Standard Guide for Measurement of Masking Sound in Open Offices
- E 1042-85, Standard Classification for Acoustically Absorptive Materials Applied by Trowel or Spray
- E 1050-90, Standard Test Method for Impedance and Absorption of Acoustical Materials Using a Tube, Two Microphones, and a Digital Frequency Analysis System
- E 1110-86 (Reapproved 90), Standard Classification for Determination of Articulation Class
- E 1111-88, Standard Test Method for Measuring the Interzone Attenuation of Ceiling Systems
- E 1123-86 (Reapproved 90), Standard Practice for Mounting Test Specimens for Sound Transmission Loss Testing of Naval and Marine Ship Bulkhead Treatment Materials
- E 1124-86, Standard Test Method for Field Measurement of Sound Power Level by the Two-Surface Method
- E 1130-90, Standard Test Method for Objective Measurement of Speech Privacy in Open Offices Using Articulation Index
- E 1179-87, Standard Specification for Sound Sources Used for Testing Open Office Components and Systems
- E 1224-88, Standard Guide for Categorizing Fields of Testing for Laboratory Accreditation Purposes
- E 1264-90, Standard Classification for Acoustical Ceiling Products
- E 1265-90, Standard Test Method for Measuring Insertion Loss of Pneumatic Exhaust Silencers
- E 1289-91, Standard Specification for Reference Specimen for Sound Transmission Loss
- E 1374-90, Standard Guide for Open Office Acoustics and Applicable ASTM Standards
- E 1375-90, Standard Test Method for Measuring the Interzone Attenuation of Furniture Panels Used as Acoustical Barriers
- E 1376-90, Standard Test Method for Measuring the Interzone Attenuation of Furniture Panels Used as Acoustical Barriers
- E 1408-91, Standard Test Method for Laboratory Measurement of the Sound Transmission Loss of Door Panels and Door Systems
- E 1414-91a, Standard Test Method for Airborne Sound Attenuation Between Rooms Sharing a Common Ceiling Plenum
- E 1425-91, Standard Practice for Determining the Acoustical Performance of Exterior Windows and Doors

41 Lasers

Joseph T. Verdeyen and Thomas A. DeTemple

Elementary Relationships	41-2
The Line Shape Function	41-3
Amplification in a Quantum System	41-3
Laser Threshold	41-4
CW Oscillation	
Transient Operation	
Modelocking	
Optical Cavities	41-6
Hermite-Gaussian Beam Modes	41-8
Common Lasers	41-8
Semiconductor Lasers	41-9
Nonlinear Optics	41-9
Second-Harmonic Generation	41-12
Pockels Effect	
Stimulated Raman Emission	

The term "laser" is an acronym for *Light Amplification by Stimulated Emission of Radiation*. The first demonstration of an electronic amplifier using stimulated emission was in the microwave frequency domain where the acronym "maser" is used. Stimulated emission has been observed in amorphous and single crystal solids, semiconductors, and dyes and in atomic, molecular, and ionized gases. Fig. 1 shows the spectral region of operation of the various types of lasers, a smoothed version of atmospheric transmission, common detector classes, and transmissive materials for various optical functions. Since the operation of lasers is dependent on quantum interactions, they are known as quantum electronic devices.

ELEMENTARY RELATIONSHIPS

In order for a quantum electronic device to amplify (or attenuate) at a frequency ν_{21} , there must be a set of allowed energy levels (or bands) E_2 or E_1 in the material separated by the Planck spacing.

$$E_2 - E_1 = h\nu_{21} \quad (\text{Planck's law}) \quad (\text{Eq. 1})$$

An atom can make discontinuous jumps between these energy levels, accompanied by either the emission

or absorption of a photon by one of three processes (Fig. 2).

- A. *Stimulated emission* occurs when an external applied wave causes an atom to yield its potential energy given by Eq. 1 to the inducing wave. This added energy (photon) propagates in the same direction, at the same frequency and phase, and in the same sense of polarization as the stimulating field. This results in amplification of the wave.
- B. *Absorption* occurs when the atom increases its potential energy by Eq. 1, thereby decreasing the wave energy by the same amount, and results in attenuation.
- C. *Spontaneous emission* is the name for a $2 \rightarrow 1$ radiative transition for which there is no apparent external influence and which yields a photon at the frequency given by Eq. 1 in any of 4π directions. It is the cause of noise in quantum electronic amplifiers; however, discharge lamps produce most of their radiation by this process.

The rate of production of electromagnetic energy by the above processes is given by

$$\text{Spontaneous: } h\nu_{21}A_{21}N_2 \quad (\text{Eq. 2})$$

$$\text{Stimulated: } h\nu_{21}B_{21}N_2g(\nu)I_\nu/c' \quad (\text{Eq. 3})$$

$$\text{Absorption: } h\nu_{21}B_{12}N_1g(\nu)I_\nu/c' \quad (\text{Eq. 4})$$

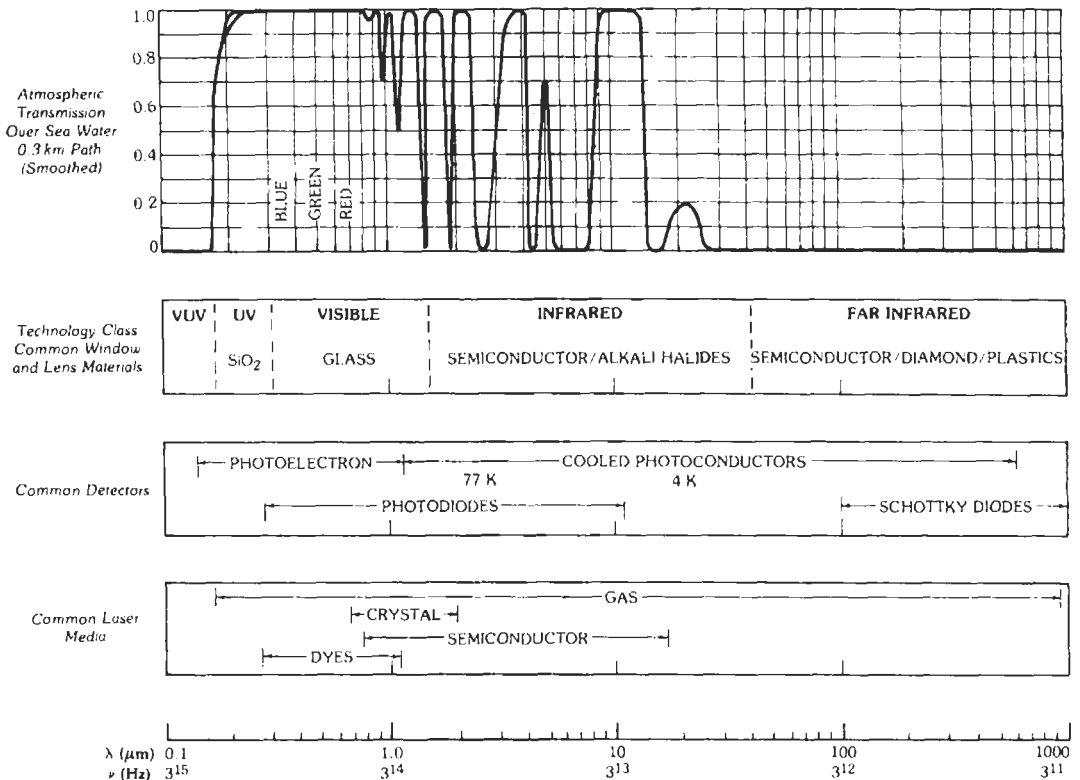


Fig. 1. Regions of operation of various types of laser and examples of the technology required in the various domains.

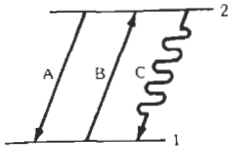


Fig. 2. Radiative processes between two levels in an atom. (The letters A, B, and C refer to stimulated emission, absorption, and spontaneous emission of radiation.)

where,

- $N_{2,(1)}$ is the density of atoms in energy state 2 (1),
- A_{21} , B_{21} , and B_{12} are the Einstein coefficients,
- $g(\nu)$ is the line shape function,
- I_ν is the intensity (watts/area) of the externally applied wave,
- $c' = c/n$ is the velocity of propagation of the wave,
- n is the index of refraction of the material.

The Einstein coefficients are interrelated.

$$A_{21}/B_{21} = 8\pi n^3 \nu^3 / c^3 \quad B_{21}/B_{12} = g_1/g_2 \quad (\text{Eq. 5})$$

where $g_{2(1)}$ is the degeneracy of the state, i.e., number of quantum levels with the same energy.

At thermal equilibrium, the ratio of the population of the two energy states, E_2 and E_1 , which are separated by Eq. 1 is given by the Boltzmann distribution:

$$N_2/N_1 = (g_2/g_1) \exp [-h\nu_2/kT] \quad (\text{Eq. 6})$$

THE LINE SHAPE FUNCTION

The function $g(\nu)d\nu$ has the following physical interpretations:

- A. It is the fraction of the spontaneous emission from state 2 to state 1 that appears in the frequency interval ν to $\nu + d\nu$.
- B. It is the relative strength of interaction between the medium and the wave for stimulated emission or absorption.
- C. It is the line profile of an optically thin transition.

Since it is a probability function, it obeys a normalization condition

$$\int_0^\infty g(\nu) d\nu = 1 \quad (\text{Eq. 7})$$

Transitions in the near IR through the UV in low-pressure gases (a few Torr) are broadened by the thermal motion of the atoms, which leads to a Doppler broadened (Gaussian) line shape:

$$g_D(\nu) = [(4 \ln 2)/\pi]^{1/2} (1/\Delta\nu_D) \exp \{(-4 \ln 2) [(\nu - \nu_{21})/\Delta\nu_D]^2\} \quad (\text{Eq. 8})$$

with

$$\Delta\nu_D = (8kT/Mc^2)^{1/2} \nu_{21} \quad (\text{Eq. 9})$$

where,

- M is the mass of the atom (or molecule),
- T is the temperature,
- $\Delta\nu_D$ is the full-width-at-half maximum (FWHM) of the transition centered at ν_{21} .

At higher pressures and/or longer wavelengths (lower frequency), elastic gas collisions broaden the transition, and the homogeneous (Lorentzian) line shape function becomes:

$$g(\nu) = \Delta\nu_h/2\pi [(\nu - \nu_{21})^2 + (\Delta\nu_h/2)^2] \quad (\text{Eq. 10})$$

with

$$\Delta\nu_h = (1/2\pi) [A_2 + k_2 + A_1 + k_1 + 2\nu_{col.}] \quad (\text{Eq. 11})$$

In Eq. 11, the A 's are the total radiative decay rates of levels 2 and 1, the k 's are total quenching rates of the two states, and $\nu_{col.}$ is the elastic collision rate of the active atoms or molecules with the other gas constituents. Usually, the elastic-collision term dominates, and one refers to the resulting line profile as "pressure broadened."

Seldom does one attempt to describe the line shape of a transition in a solid or band-to-band transition in semiconductors by elementary functions such as those used in Eqs. 8 and 10. Rather, one depends on the experimental determination of the line profile (e.g., C above) and often approximates $g(\nu)$ by

$$g(\nu) \approx 1/\Delta\nu \quad (\text{Eq. 12})$$

where $\Delta\nu$ is the FWHM of the spontaneous emission of an optically thin line.

AMPLIFICATION IN A QUANTUM SYSTEM

The net rate of increase of an externally applied electromagnetic wave $I\nu$ (watts/area) per unit of length of active material is the difference between the stimulated rate of Eq. 3 and the absorptive rate given by Eq. 4:

$$dI_\nu/dz = h\nu_{21} B_{21} [N_2 - (g_2/g_1) N_1] g(\nu) I_\nu/c' \quad (\text{Eq. 13})$$

If $N_2 < (g_2/g_1) N_1$, one has a "normal" population with a positive "temperature" specified by Eq. 6; such a medium will attenuate an electromagnetic wave. If $N_2 > (g_2/g_1) N_1$, one has an "inverted" population with a "negative" temperature, and this medium will amplify an electromagnetic wave. By combining Eqs. 5 and 13, one obtains the *small-signal gain coefficient* given by

$$\begin{aligned} \gamma_o(\nu) &\triangleq (1/I_\nu)(dI_\nu/dz) \\ &= A_{21} (\lambda^2/8\pi) g(\nu) [N_2 - (g_2/g_1) N_1] \end{aligned} \quad (\text{Eq. 14})$$

where $\lambda = \lambda_o/n$. The collection of terms multiplying N_2 is often referred to as the *stimulated emission cross section* (i.e., "area")

$$\sigma_{sc}(v) = A_{21} (\lambda^2/8\pi) g(v_{21}) \quad (\text{Eq. 15})$$

The *absorption cross-section* refers to the factors multiplying N_1 in Eq. 14.

$$\sigma_{abs}(v) = (g_2/g_1) \sigma_{sc}(v) \quad (\text{Eq. 16})$$

The small-signal gain of the medium of length l is given by the expression

$$G_o = \exp \gamma_o(v)l \quad (\text{Eq. 17})$$

The subscript "o" in Eqs. 14 and 17 implies that the intensity of the stimulating wave is sufficiently small so as to create negligible perturbation on the density of atoms in the two states.

If the intensity is large enough, the populations, $N_{2,1}$, in Eq. 14 are functions of the intensity, I_v , of the stimulating wave. For homogeneously broadened CW systems, the *saturated gain coefficient* is given by

$$\gamma(v) = \gamma_o(v)/[1 + (I_v/I_s) \bar{g}(v)] \triangleq (1/I_v) dI_v/dz \quad (\text{Eq. 18})$$

where,

- $\gamma_o(v)$ is the small-signal value,
- $\bar{g}(v)$ is the line shape function normalized to unity at line center,
- I_s is the saturation intensity given by

$$I_s = h\nu_{21}/\tau_2 \sigma_{sc}(v_{21}) \quad (\text{Eq. 19})$$

In Eq. 19, τ_2^{-1} is the decay rate of state 2 due to all causes—radiation plus quenching, and $\sigma_{sc}(v_{21})$ is the stimulated-emission cross section given by Eq. 15 evaluated at line center.

The overall gain, G , of the amplifier is thus a function of the input intensity and is found by integrating Eq. 18 from 0 to l .

$$\ln G + (I_{in}/I_s) \bar{g}(v) (G - 1) = \gamma_o(v)l \quad (\text{Eq. 20})$$

where $G = I_{out}/I_{in}$.

The maximum intensity that can be extracted from an amplifier by stimulated emission is given by

$$(I_{out} - I_{in}) \Big|_{\max} = \gamma_o(v_{21})I_s l \quad (\text{Eq. 21})$$

For inhomogeneously broadened transitions the saturated law is modified:

$$\gamma(v) = \gamma_o(v)/(1 + I_v/I_s)^{1/2} \quad (\text{Eq. 22})$$

where all factors are as defined previously, and Eqs. 20 and 21 must be suitably modified.

LASER THRESHOLD

Oscillation with an inverted population is possible provided sufficient feedback is provided, usually in the form of an open cavity. In general, the necessary condition for oscillation is that net round-trip gain must exceed 1 so that an initially small wave can grow in amplitude as it passes around the feedback path. For a simple system such as that shown in Fig. 3, the threshold condition is

$$R_1 R_2 T_a^2 T_b^2 \exp [2 \gamma_o(v)l] \geq 1$$

or

$$\gamma_o \geq (1/l) \ln (1/\sqrt{R_1 R_2}) + (1/l) \ln (1/T_a T_b) \triangleq \alpha \quad (\text{Eq. 23})$$

In Fig. 3, the power transmission coefficients for the windows at the ends of the laser cell are T_a and T_b , and those of the mirrors are $T_1 = 1 - R_1$ and $T_2 = 1 - R_2$. Implicit in Eq. 23 is the assumption that the population inversion exists long enough to allow many transits back and forth between mirrors.

If the transition is homogeneously broadened, then oscillation will occur on a cavity mode that has the highest gain-to-loss coefficient ratio γ_o/α . For the simple laser shown in Figs. 3 and 4, this would correspond to the cavity mode nearest line center. If the transition is inhomogeneously broadened, multimode oscillation is possible on those cavity modes at frequencies ν_q such that $\gamma_o(\nu_q)/\alpha > 1$.

CW Oscillation

For the simple laser shown in Fig. 3, there are three types of losses: internal discrete losses represented by the imperfect transmission by the windows, a residual loss due to imperfections in the medium, and external or coupling losses by virtue of the transmission through the mirrors. The photon flux inside the cavity builds up to a large intensity so that the gain coefficient (Eq. 18 or 21) saturates at the threshold value given by Eq. 23. If the mirrors are highly reflecting, little power escapes the cavity; if the mirrors have too large a transmission coefficient, the laser fails to oscillate. The transmission coefficient of M_2 that provides maximum output power from the laser illustrated in Fig. 3 is

$$T_2 = -L + (g_o L)^{1/2}; \quad (T_1 = 0) \quad (\text{Eq. 24})$$

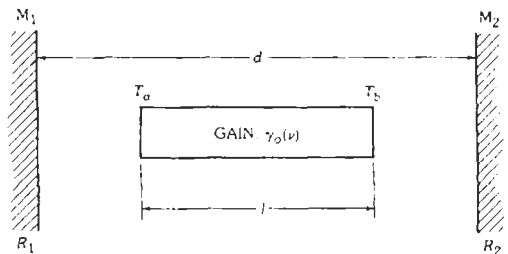


Fig. 3. A simple laser.

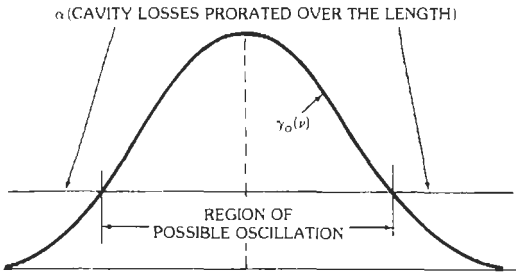


Fig. 4. Spectral bandwidth available for oscillation.

where,

$g_o = \gamma_o \cdot 2l$ is the gain coefficient integrated over the double pass through the active medium,
 $L = \alpha_{int} \cdot 2l$ is the integrated internal losses given by

$$\alpha_{int} = (1/2l) [\ln(1/R_1) + \ln(1/T_a^2 T_b^2)] + \alpha_{residual} \quad (\text{Eq. 25})$$

The output intensity (through M_2) is given by

$$I_{out} = (I_s/2) \cdot [g_o/(L + T_2) - 1] T_2 \quad (\text{Eq. 26})$$

For T_2 given by Eq. 24, the output intensity is given by the expression

$$I_{out} = (I_s/2) \cdot [\sqrt{g_o} - \sqrt{L}]^2 \quad (\text{Eq. 27})$$

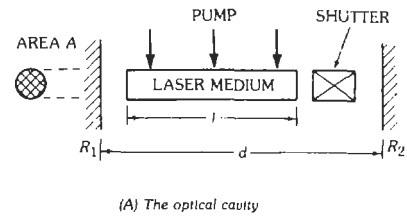
If $\sqrt{g_o} \gg \sqrt{L}$, then $I_{out} \approx (I_s g_o)/2 \approx \gamma_o I_s$.

Transient Operation

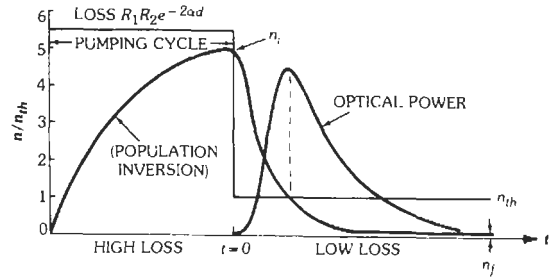
The previous discussion presumes steady-state or CW laser oscillation. However, some lasers are operated in a pulsed mode with rather spectacular results, as discussed below.

If the upper laser level has a lifetime that is long compared to the pumping duration, then it is a candidate for Q -switching. The scheme is similar to that used for radar modulators in which energy is slowly stored in a capacitor, transmission line, or PFN and then is rapidly discharged into the magnetron load. In the laser scheme, the energy is stored in the upper quantum state and is then extracted by stimulated emission.

The scheme is illustrated in Fig. 5. The Q of the laser cavity is spoiled (intentionally) during the time the upper state is being populated by the external pump, thus preventing lasing. When the upper state has reached an equilibrium with the pump, the cavity Q is restored to its high value. The population inversion is larger than threshold for CW oscillation, and thus stimulated emission rapidly amplifies the photons initially present in the cavity. This buildup occurs so rapidly that one can neglect any pumping on the time scale of the output pulse; in fact, it is so rapid and intense that the population inversion is driven below threshold for CW oscillation.



(A) The optical cavity



(B) Time evolution of the population difference and the optical pulse.

Fig. 5. Q -switching of a simple laser. (From J. T. Verdeyen, Laser Electronics. Englewood Cliffs, N. J.: Prentice-Hall, Inc., 1981; Fig. 9-5.)

The parameters that characterize the Q -switched pulse are:

- n_i = initial number of inverted atoms or molecules in the cavity
- n_{th} = the inversion required for CW oscillation in the high- Q cavity
- n_f = the final inversion in population (see Fig. 5)
- $T = t/\tau_c$, time divided by the passive cavity lifetime

The peak power produced by such a laser is given by:

$$P_p = k(h\nu/\tau_c) [(n_i - n_{th}/2) - (n_{th}/2) \ln(n_i/n_{th})] \quad (\text{Eq. 28})$$

where k is the coupling factor on the order of 1/2.

The output energy of the Q -switch pulse is given by

$$W_{out} = k \cdot [(n_i - n_f)/2]h\nu \quad (\text{Eq. 29})$$

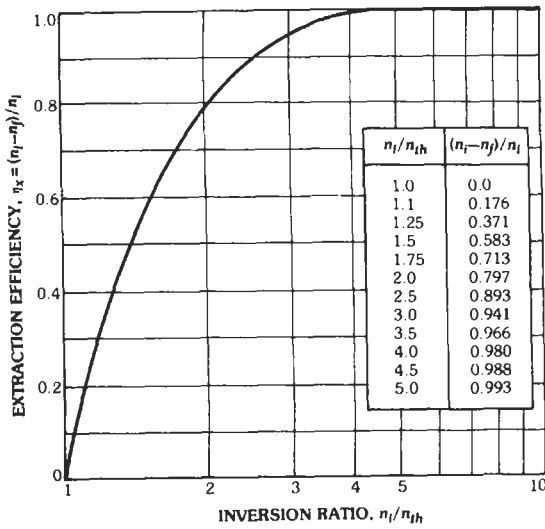
and the final inversion is found from a solution to

$$n_f/n_i = \exp [-(n_i - n_f)/n_{th}] \quad (\text{Eq. 30})$$

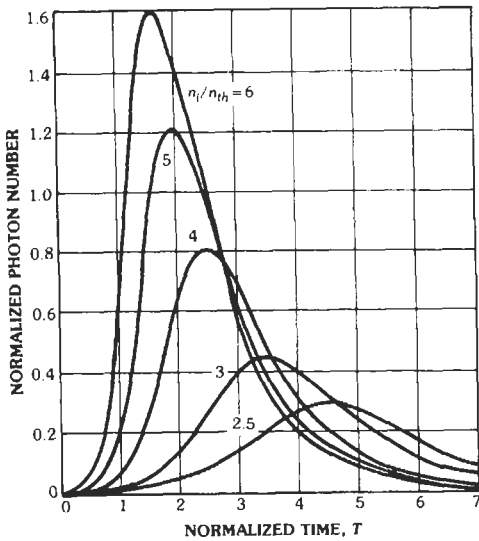
Inasmuch as the maximum energy that could be extracted is $n_i \cdot h\nu/2$, Eq. 29 can be written as:

$$W_{out} = k \cdot [(n_i - n_f)/n_i] \cdot (n_i h\nu/2) \quad (\text{Eq. 31})$$

The quantity $(n_i - n_f)/n_i$ is called the energy extraction efficiency and is plotted in Fig. 6A. The time evolution of this Q -switch pulse requires a numerical solution of the coupled differential equations. A typical solution is shown in Fig. 6B for various values of the ratio n_i/n_{th} . A



(A) The extraction efficiency for a Q-switched laser.



(B) The time evolution of the Q-switched pulse, time normalized to the passive cavity lifetime.

Fig. 6. Characteristic parameters of Q-switching. (From J. T. Verdeyen, Laser Electronics. Englewood Cliffs, N.J.: Prentice-Hall, Inc., 1981; Figs. 9-7 and 9-8.)

rough estimate for the FWHM of the pulse can be found by combining Eqs. 28 and 31:

$$\Delta t_{1/2} = W_{out}/P_p \quad (\text{Eq. 32})$$

Modelocking

If the system is inhomogeneously broadened, the laser may oscillate on many cavity modes within the

spectral bandwidth available for oscillation (Fig. 4). If these modes can be locked together in a definite phase sequence, then the time-domain representation of the output consists of a series of pulses at a repetition rate equal to the round-trip transit time of a photon in the cavity.

Each individual pulse is very short; typically, 0.1 to 2 ns is common, with 30 fs (30×10^{-15} s) being the shortest presently recorded. The time evolution of these pulses is related to the spectral distribution of modes by the Fourier transform pair (see Chapter 7). Rough rules of thumb for the pulse characteristics are:

$$\text{Peak power} = N \times \text{average power} \quad (\text{Eq. 33a})$$

$$\text{Pulse repetition rate} = c'/2d \quad (\text{Eq. 33b})$$

$$\text{Pulse width} = (2d/c')/N \approx 1/\Delta\nu \quad (\text{Eq. 33c})$$

where,

N is the number of modes locked,
 $2d/c'$ is the round-trip transit time,
 $\Delta\nu$ is the FWHM of the transition.

OPTICAL CAVITIES

Optical cavities perform the essential function of providing positive feedback for the quantum electronic amplifier so as to obtain oscillation. The spontaneous emission from the upper state of the quantum system is reflected by the mirrors and amplified by the inverted population. The cavity also selects the oscillation frequency within the spectral bandwidth defined in Fig. 4. The laser depends upon stimulated emission for amplification, and that rate is a maximum at the resonant frequencies of the cavity since the electromagnetic energy is a maximum under these conditions.

The resonant buildup of electromagnetic energy in any cavity is due to the near-perfect phase addition of the various partial waves bouncing back and forth between the mirrors. In general then, the resonant frequency of any cavity is determined by

Round-trip phase shift

$$= \text{Integral multiple of } 2\pi \text{ radians} \quad (\text{Eq. 34})$$

If one assumes a uniform plane wave for the normal mode of the simple Fabry-Perot cavity shown in Fig. 7A, then Eq. 34 reduces to

$$v_q = q c'/2d \quad (\text{Eq. 35})$$

where,

c' is the velocity of light in the medium between the mirrors,
 q is an integer.

Typically $d \approx 0.1$ to 1.0 m for solid-state and gas lasers, and thus $c'/2d \approx 150$ MHz to 1.5 GHz,

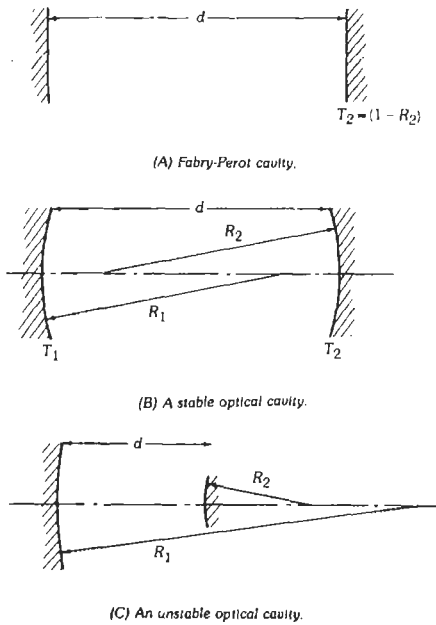


Fig. 7. Typical cavities used with lasers.

considerably smaller than the line width of many transitions (approximately 5 to 100 GHz).

Most optical cavities use one or more focusing elements in the manner indicated in Figs. 7B and 7C. The cavity shown in Fig. 7B is *stable* in the sense that the maximum excursion of a paraxial ray from the axis of the cavity is bounded. However, a ray will “walk off” the surfaces of the feedback mirrors of Fig. 7C; hence, this cavity is classified as *unstable*. Simple cavities such as those shown in Fig. 7 are stable if

$$0 \leq g_1 g_2 \leq 1 \quad (\text{Eq. 36})$$

where,

$$g_{1,2} = 1 - d/R_{1,2} \quad (\text{Eq. 37})$$

If the equality is not satisfied, the cavity is unstable. The stability criterion given by Eq. 36 is quite often embodied in a graphical format such as shown in Fig. 8. The dashed curves illustrate the equi-loss contours for the unstable cavity. The quantity Γ is related to the mean fractional power lost per pass by $\Gamma = 1 - L$. For stable cavities, the diffraction loss (i.e., loss around edges of the mirrors) is very small; most of the loss of electromagnetic energy is transmission through the mirrors.

Unstable cavities have a very large mode volume and thus can utilize the stimulated emission from a large number of atoms. The losses from around the edges of the mirrors are usually large compared to the loss from direct coupling—usually by design. One can relate Γ , the mean fraction of the power that survives a single pass through the cavity, to the cavity g parameters by

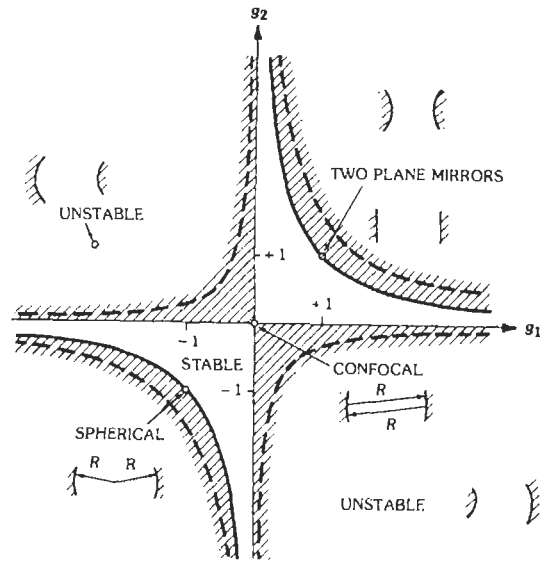


Fig. 8. Graphical display of the stability of optical cavities.

$$\Gamma = \pm \frac{1 - [1 - (g_1 g_2)^{-1}]^{1/2}}{1 + [1 - (g_1 g_2)^{-1}]^{1/2}} \quad (\text{Eq. 38})$$

where the choice of the sign is determined by the requirement that $0 < \Gamma < 1$. This leads to the equi-loss contours shown in Fig. 8.

If all external sources of electromagnetic energy for the resonant mode of a cavity are suddenly removed, the energy will decay with time constant τ_p , called the photon lifetime:

$$dW/dt = -W/\tau_p \quad (\text{Eq. 39a})$$

$$\tau_p = \frac{\text{Time for a round trip}}{\text{Fraction of energy lost per round trip}} \quad (\text{Eq. 39b})$$

where W is the energy stored in the cavity.

For the simple cavity shown in Fig. 7B, this leads to the expression

$$\tau_p = (2d/c)/(1 - R_1 R_2) \quad (\text{Eq. 40})$$

where all losses except those due to imperfect mirror reflectivities have been neglected.

The photon lifetime is directly related to the cavity Q by the equation

$$Q = 2\pi \times \frac{\text{Energy stored}}{\text{Energy lost in a cycle}} \quad (\text{Eq. 41a})$$

or

$$Q = \omega_o W/(-dW/dt) \quad (\text{Eq. 41b})$$

where ω_o is the (angular) resonant frequency of the

cavity. The differential equation of Eq. 41b is the same as Eq. 39a and relates Q and τ_p by

$$\tau_p = Q/\omega_0 \quad (\text{Eq. 42a})$$

The full-width at half maximum (FWHM) of the cavity resonance is given by the familiar formula

$$\Delta\omega = \omega_0/Q \quad \text{or} \quad \Delta\nu = \nu_0/Q \quad (\text{Eq. 42b})$$

Hence

$$\Delta\omega\tau_p = 1 \quad (\text{Eq. 42c})$$

HERMITE-GAUSSIAN BEAM MODES

Stable cavities such as that shown in Fig. 7B produce a beam whose electric field is described by

$$E_{m,n}(x,y,z) = E_0(w_0/w(z)) H_m(\sqrt{2x}/w) H_n(\sqrt{2y}/w) \cdot \exp[-(x^2 + y^2)/w^2] \cdot \exp[-jk(x^2 + y^2)/2R(z)] \cdot \exp\{-j[kz - (1 + m + n)\tan^{-1} z/z_0]\} \quad (\text{Eq. 43})$$

where H_m is a Hermite polynomial and $k = \omega/c$. The mode indices, m and n , refer to the variation of the TEM field in the x and y directions similar to the convention used for waveguides. The interrelationships between the parameters in Eq. 43 are listed below.

$$w^2(z) = w_0^2 [1 + (z/z_0)^2] \quad (\text{Eq. 44a})$$

$$R(z) = z [1 + (z_0/z)^2] \quad (\text{Eq. 44b})$$

$$z_0 = \pi w_0^2/\lambda \quad (\text{Eq. 44c})$$

$$H_n(u) = (-1)^n (\exp u^2) (d^n/du^n) \exp(-u^2) \quad (\text{Eq. 44d})$$

The plane $z = 0$ is defined to be the location where the spot size, w , is a minimum and the wave-front curvature becomes planar ($R = \infty$). The spot size of the beam will expand by $\sqrt{2}$ when it has propagated a distance z_0 from $z = 0$. The interrelationships are plotted in Fig. 9 for the TEM_{0,0} mode.

Notice that in the far field the beam expands as if it were a spherical wave of limited angular extent and originating at $z = 0$. Angle θ is the angular divergence of the beam given by

$$\theta = 2\lambda/\pi w_0 \quad (\text{Eq. 45})$$

and is an important laser characteristic.

The characteristic parameter, z_0 , and the location of the plane $z = 0$ are related to the radius of curvature of

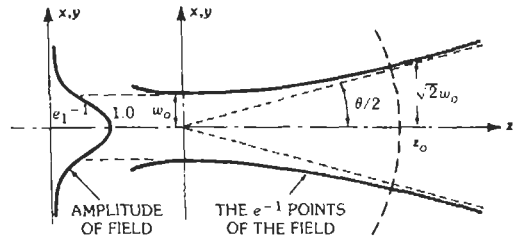


Fig. 9. The expansion of a TEM_{0,0} mode.

the mirrors and the spacing for the cavity shown in Fig. 7B by:

$$z_0^2 = (\pi w_0^2/\lambda)^2 = \frac{d(R_1 - d)(R_2 - d)(R_1 + R_2 - d)}{(R_1 + R_2 - 2d)^2} \quad (\text{Eq. 46a})$$

$$z_1 = d(R_2 - d)/(R_1 + R_2 - 2d) \quad (\text{Eq. 46b})$$

$$z_2 = d(R_1 - d)/(R_1 + R_2 - 2d) \quad (\text{Eq. 46c})$$

The resonant frequencies of the TEM_{m,n} modes differ slightly from those predicted by Eq. 35:

$$\nu_{m,n,q} = (c/2d) \{q + [(1 + m + n)/\pi] \cos^{-1} (g_1 g_2)^{1/2}\} \quad (\text{Eq. 47})$$

COMMON LASERS

Gases, impurity doping in crystalline and amorphous solids, direct bandgap semiconductors, and dyes dissolved in various solvents are the common materials used for lasers. All have a common theoretical framework, but the details of the excitation route differ from one type to another as illustrated in Table 1.

The excitation can be either pulsed or CW, resulting in the corresponding temporal characteristic of the laser radiation. With the exception of the semiconductors, the wall-plug efficiency of most lasers is typically less than 1%; the CW CO₂ gas laser system is an exception with efficiencies of 10–20% being common. Therefore, many very high power lasers are operated in a pulsed mode so as to dissipate the waste heat associated with the excitation. Many lasers, such as N₂ at 0.337 μm, require pulsed operation since only a transient inversion can be created.

Table 2 lists some common high-power pulsed lasers; typical pulse widths and the optical energy produced are shown. Table 3 lists common CW lasers. It should be noted here that even a milliwatt (10⁻³ W) coherent laser is potentially hazardous to the eye. Consequently, wavelength-selective eyeglasses are recommended, since even the backscatter of high-power lasers from rough surfaces can exceed this level.

TABLE I. COMMON EXCITATION MECHANISMS FOR VARIOUS LASERS

Material	Excitation Route
Gases	Electric discharge. Optical pumping, chemical reactions.
Impurity doping in solids	Optical pumping usually with incoherent flash lamps, halogen lamps, or semiconductor LEDs.
Semiconductor	Carrier injection.
Dyes	Coaxial flash lamps, N ₂ laser, excimer, or harmonics of YAG (0.532 μm , 0.355 μm) for pulsed operation; argon or krypton ion lasers for CW operation.

The output of various dyes depends on the solvent, the concentration of the dye, and the source of excitation. This is illustrated in Fig. 10, where the relative energy versus wavelength is plotted for the common sources of excitation. Flash-lamp-pumped dyes produce the highest energy per pulse, approximately 0.2 J, and the longest pulse width, approximately 0.2 μs , but typically at a rather slow repetition rate of a few pulses per second. The N₂ laser at 0.337 μm or the second and third harmonics of Nd:YAG are used to produce high power (tens of kilowatts) and short pulse widths (5 ns) at moderate repetition rates of 10–100 pps. Continuous-wave dye lasers are usually excited with the lines from an argon or krypton ion laser. While the conversion efficiency of the dye is quite reasonable (10–25%), the electrical efficiency of the pump laser is usually less than 0.1%.

Titanium-doped sapphire lasers, pumped by argon ion lasers, have tuning characteristics which are similar to those of near-infrared dye lasers, but with a higher output power. Typical tuning curves for this laser are shown in Fig. 11.

Many of the far-infrared lasers, using various gases such as CH₃OH, are pumped by the CO₂ laser. Fig. 12 illustrates the range of powers and frequencies available when such a pumping scheme is used.

SEMICONDUCTOR LASERS

Population inversion, and thus optical gain in the sense of Eq. 14, may be obtained in a direct bandgap semiconductor exhibiting radiative recombination of electrons and holes provided there is, simultaneously in space and time, a large density of electrons in the conduction band and holes in the valence band such as that provided by the junction region of a diode under heavy forward bias where large numbers of electrons and holes exist.

Such a diode is shown in Fig. 13 for a homojunction laser consisting of, for example, heavily doped regions of GaAs to form a pn junction. Under forward bias, electrons are injected into the *p* region, leading to an "inverted" population of width 1 μm or so perpendicular to the junction. The light (*L*) output is directly proportional to the recombining carriers and thus is directly related to the injected current (*I*). At low

currents, the device acts as a LED. Beyond a threshold current, stimulated emission of recombination takes place, and the slope of light output versus current shows a dramatic increase. Due to the short recombination lifetime (typically 10⁻⁹ second), one is able to amplitude modulate the light output by simple current modulation.

Semiconductor lasers that involve only one material (e.g., GaAs) are classified as *homojunctions*; those that utilize a different material on one side of the active region are called single heterostructures (SH); if both sides use different materials, the device is called a double heterostructure (DH). Usually, the different materials used to obtain a heterostructure have a different bandgap and index of refraction than that of the active region. A typical configuration for a double heterostructure is shown in Fig. 14. The higher bandgap confines the carriers, and change in index of refraction guides the electromagnetic mode. Due to the small size of the active region, the divergence angle of the mode may be on the order of 10–20° and is different in the planes parallel or perpendicular to the junction.

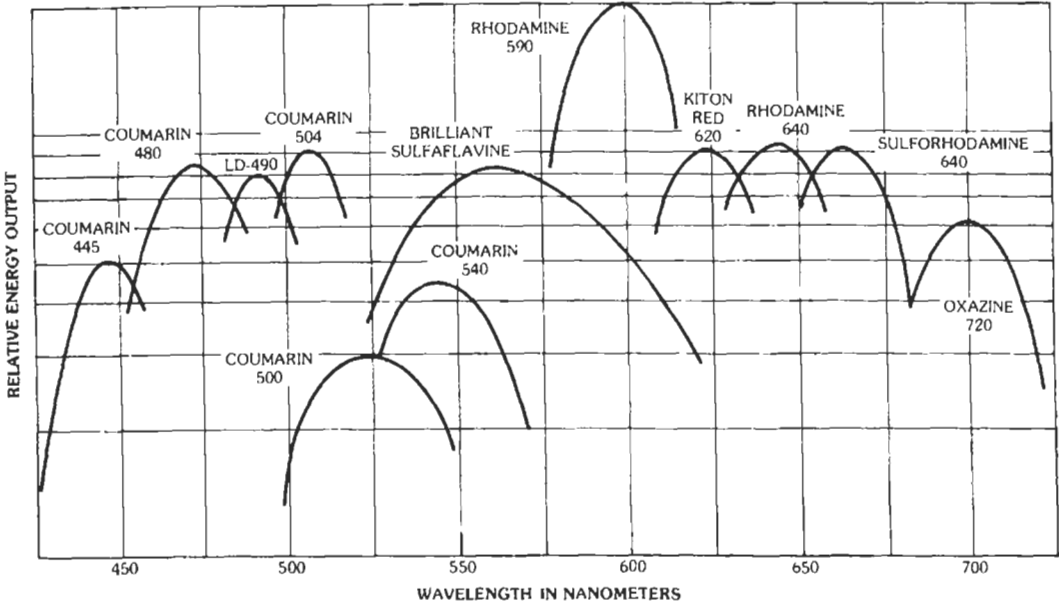
Fig. 15 shows the emission ranges of semiconductor lasers obtained with various combinations of III-V and IV-VI materials. For a given material system, the emission range corresponds to different crystal compositions and not the wavelength tuning of a single laser. Typical threshold currents for heterostructure lasers are in the range from 1 to 50 mA; output powers range from 1 to 50 mW. The very-long-wavelength diodes require cryogenic cooling.

NONLINEAR OPTICS

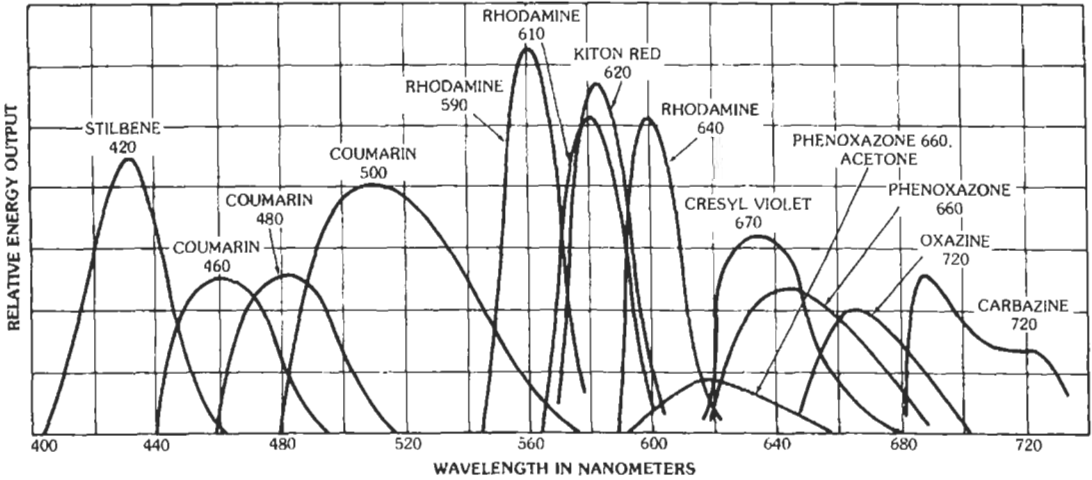
As in the case of other electronic devices, under suitable conditions the ordinary linear behavior of transparent materials becomes nonlinear. These nonlinearities, which are caused by strong optical or low-frequency electric fields, are used in the frequency and amplitude control of light.

The various effects are described by a simple Taylor series expansion of the microscopic polarization

$$P = \epsilon_0 \chi^{(1)} E + \epsilon_0 \chi^{(2)} EE + \epsilon_0 \chi^{(3)} EEE + \dots \quad (\text{Eq. 48})$$



(A) Coaxial flash lamp pumped dyes (Candela Corp.).



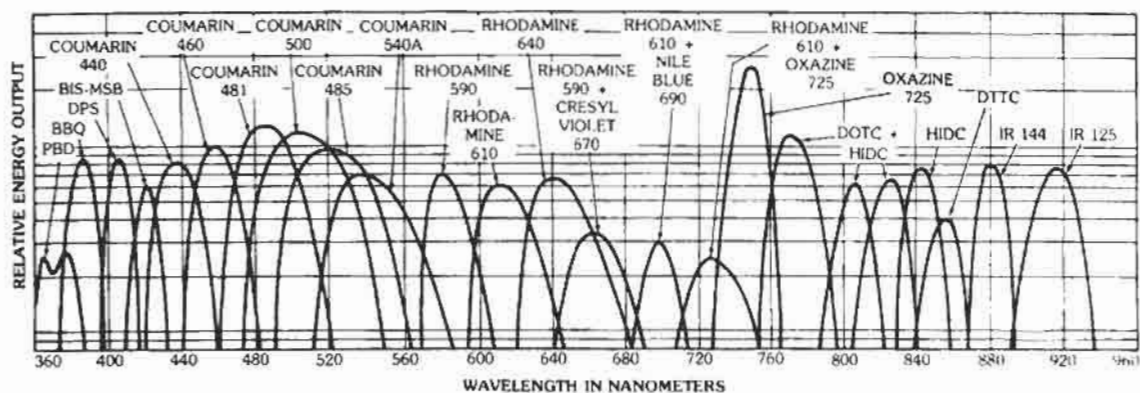
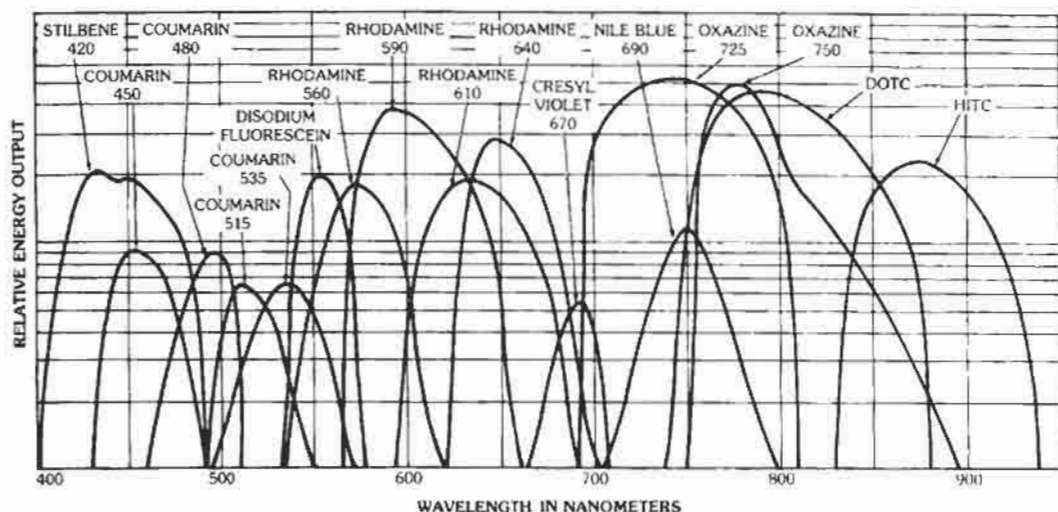
(B) Nd:YAG pumped dyes (Quanta-Ray).

Fig. 10. Typical tuning curves for various dyes when

so that the normal index of refraction is given by $n^2 = 1 + \chi^{(1)}$.

If the symmetry property of the material is such that a reversal of the sign of all coordinates causes no change in the spatial arrangement of the atoms, the material is said to be centrosymmetric. Applying the reversal to P and E , one sees that the $\chi^{(2)}$ coefficient must be zero in centrosymmetric systems but that $\chi^{(1)}$ and $\chi^{(3)}$ are not. Accordingly, the lowest order or quadratic nonlinearity is seen in noncentrosymmetric systems which comprise 21 of the 32 known crystal classes.

In general, $\chi^{(1)}$ has nine components. If one chooses a natural coordinate system that reflects the symmetry axis of the material, $\chi^{(1)}$ becomes diagonal with three components only. This is referred to as the principal axis frame. In the principal axis frame $\chi^{(2)}$ has at most 27 components while $\chi^{(3)}$ has 3^4 components. Not all components are separately distinguishable, so it is possible to reduce the independent number down to 18 for $\chi^{(2)}$. This is done by noting that the i th component of the quadratic part of the polarization is $P_i = \sum_k \sum_l$

(C) N_2 pumped dyes (Molelectron)(D) Ar^* and Kr^* CW pumped dyes (Spectra-Physics)

excited by the indicated sources. (Courtesy Exciton Corp.)

$\epsilon_0 \chi_{ijk}^{(2)} E_j E_k$, so a new parameter is accordingly defined as $d_{iq} = \epsilon_0 \chi_{i(jk)}^{(2)}$ where q is a contracted index for the jk combination in parentheses (the Voigt notation). Principal axes along x , y , and z are labeled 1, 2, and 3 so that q ranges from 1 to 6 according to the jk combinations 11, 22, 33, 23, or 32; 13 or 31; and 12 or 21, respectively. So, for example, $\chi_{i(23)}^{(2)} = \chi_{i(32)}^{(2)}$ contributes to d_{i4} . With this, d is a 3×6 matrix with at most 18 elements.

The symmetry of the individual system is then used to deduce which elements are nonzero. Since $\chi^{(2)}$ is the

same order tensor as that used to describe the piezoelectric effect, the nonzero elements are well known and are compiled in Chapter 4. Similarly, the individual values for d_{iq} for various materials and optical frequencies are tabulated elsewhere. The same kind of symmetry-based reductions can be applied to $\chi^{(3)}$ for crystals, but none are used for amorphous materials.

A specific example is for the crystal class $\bar{4}2m$, which is appropriate to the common nonlinear material potassium dihydrogen phosphate (KDP). For this class, $\chi_1^{(1)} = \chi_2^{(1)} \neq \chi_3^{(1)}$, and

TABLE 2. PULSED LASERS—HIGH POWER

λ^*	Atom or Molecule	Materials	Pulse Width (ns)	Energy (mJ)	Rate (pps)	Comments
157	F ₂	g	6	10	50	RGH excimers: mixtures of rare gas and halide donors; 10 ⁵ –10 ⁶ laser pulses per gas fill. Can also be used with N ₂ , HF, CO ₂ in a typical laser system.
193	ArF	g	14	200	50	
222	KrCl	g	9	30	50	
249	KrF	g	16	250	50	
282	XeBr	g	8	17	50	
308	XeCl	g	6	150	50	
351	XeF	g	14	400	50	
266	Nd: YAG, Glass	s	4	50	0.02–20	Fourth harmonic of 1.06 μm Nd
337	N ₂	g	6	16	100	
347	Cr:Al ₂ O ₃	s	25	100	0.1	Second harmonic of 694.3 nm (ruby)
355	Nd: YAG	s	5	100	0.1–20	Third harmonic of 1.06 μm Nd
502	HgBr	g	50	100	5–100	High temperature (> 150°C) req.
532	Nd: YAG, Glass	s	7	200	20	Second harmonic of 1.06 μm Nd
510.6 } 578.2 }	Cu	g	30	2.5	6 kHz	High temperature, high average power
694.3	Cr:Al ₂ O ₃	s	20	1–10 J	0.02	Q Switched
850.0	GaAs	Semi	100	0.01	1 kHz	Semiconductor diode array
1.06 μm	Nd: YAG	s	15	0.1–1 J	10	Other transitions are possible.
1.06	Nd: Glass	s	20	1–20 J	0.03	Other transitions are possible.
2.64–3.01 μm	HF	g	500	300	2	Chemical lasers, discrete line spectra associated with VR bands
5–6 μm	CO	g	1000	8 J	5	Discrete ~ 50 line spectra associated with VR structure
9.4 } 10.6 μm }	CO ₂	g	100	10 ² J	0.1–10 ²	Discrete ~ 50 line spectra associated with VR structure
12–13 μm	NH ₃	g	100	0.1 J	1–10	Optically pumped with CO ₂
385 μm	D ₂ O	g	100	0.1 J	1–10	Optically pumped with CO ₂
496.1 μm	CH ₃ F	g	100	10	1–10	Optically pumped with CO ₂

*Nanometers unless noted otherwise.

$$d = \begin{bmatrix} 0 & 0 & 0 & d_{14} & 0 & 0 \\ 0 & 0 & 0 & 0 & d_{25} & 0 \\ 0 & 0 & 0 & 0 & 0 & d_{36} \end{bmatrix} \quad (\text{Eq. 49})$$

with $d_{14} = d_{25}$.

SECOND-HARMONIC GENERATION

With the presence of an intense optical field, the quadratic nonlinearity leads to a polarization at twice

TABLE 3. COMMON CW LASERS

λ (nm)	Atom or Molecule	Materials	Power	Comments
325.0	Cd ⁺	g	8 mW	Small compact
441.6			40 mW	
290-305	Dye	l	50 mW	SHG of R6G dye
351.1, 363.8	Ar ⁺⁺	g	2.5 W	1. Many other lines possible 2. Very intense discharge 3. Dominant lines denoted by * 4. Water cooling usually required
454.5			1.1 W	
457.9			1.35 W	
476.5			2.7 W	
488.0*			6.5 W	
496.5			2.5 W	
501.7			1.5 W	
514.5*			7.5 W	
528.7			1.0 W	
405-805	Dye	l	0.05-2.0 W	Optically pumped by Ar ⁺ or K ⁺ lasers (See Fig. 10 for more detail.)
413.1	Kr ⁺⁺	g	0.5 W	1. Many other lines possible 2. Very intense discharge 3. Dominant lines denoted by * 4. Water cooling required
468.0			0.3 W	
476.2			0.40 W	
482.5			0.40 W	
520.8			0.7 W	
530.9*			1.5 W	
568.2*			1.1 W	
647.1*			3.5 W	
676.4*			0.9 W	
752.5			1.2 W	
799.0*	0.3 W			
632.8	Ne	g	1-50 mW	Common "red" laser
670-1130	Ti ³⁺ (Al ₂ O ₃)	s	100 mW	Optically pumped by Ar ⁺
1152.3	Ne	g	10 mW	First gas laser
3391.3	Ne	g	10 mW	Frequency standard
750-900	Al _x Ga _{1-x} As	Semi	10 mW	Typical drive current = 150 mA
864-904	GaAs	Semi	1-1000 mW	λ is temperature dependent.
1.06 μ m	Nd ³⁺ (YAG)	s	1-10 W	Optically pumped (flash lamp)
1.1-1.3 μ m	In _{1-x} Ga _x P _{1-r}	Semi	0.1-10 mW	
2.3-3.3 μ m	F-center	s	1-10 mW	Color center
2.6-3 μ m	HF	g	1-50 W	Chemical laser, gas flow required
3.6-4.1 μ m	DF	g		

Continued on next page.

TABLE 3 (CONT). COMMON CW LASERS

λ (μm)	Atom or Molecule	Materials	Power	Comments
5-6.5 μm	CO	g	10-20 W	Line selectable
9.2-11.2 μm	CO ₂	g	4 W	Line selectable, sealed
10.6 μm	CO ₂	g	50 W-8.5 kW	Flowing gas
28, 78, 118 μm	H ₂ O	g	10 mW	Electrically excited
118.1 μm	CH ₃ OH	g	100 mW	Optically pumped (See Fig. 12.)
311, 377 μm	HCN	g	100 mW	Electrically excited
496.1 μm	CH ₃ F	g	25 mW	Optically pumped with CO ₂

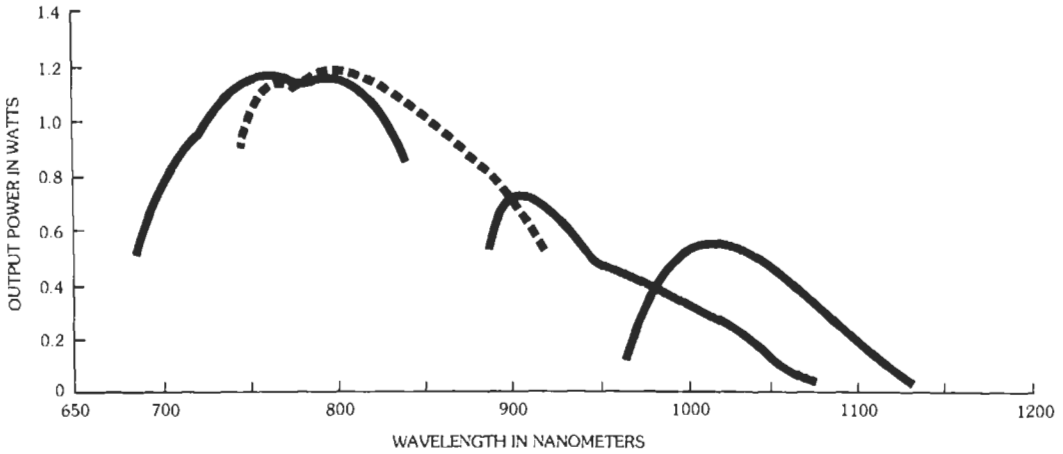


Fig. 11. Output power versus wavelength for a titanium-doped sapphire laser pumped with an argon laser. The curves refer to different mirror sets with a 5-watt pump laser, except for the far right curve which is for a 7.5-watt pump. (Schwartz Electro-Optics, Inc.)

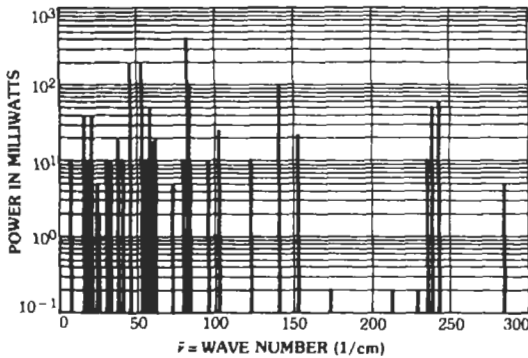


Fig. 12. Optically pumped far-infrared laser lines.

the optical frequency, which may act as a driving term for a new field component at 2ω . Since the growth tends to occur on spatial scales long compared with the wavelengths, the full set of Maxwell's equations can be reduced to a first-order equation for the evolution of the field amplitude. The growth of the harmonic is governed by

$$d(E_{2\omega})/dz = A E_{\omega}^2 \exp(i\Delta kz) \quad (\text{Eq. 50})$$

where,

$$\Delta k = k(2\omega) - 2k(\omega) = 2k_o[n(2\omega) - n(\omega)],$$

$$k_o(\omega) = \omega/c,$$

A = a constant proportional to the d_{iq} elements, the exact form of which depends on the field orientations and the d matrix.

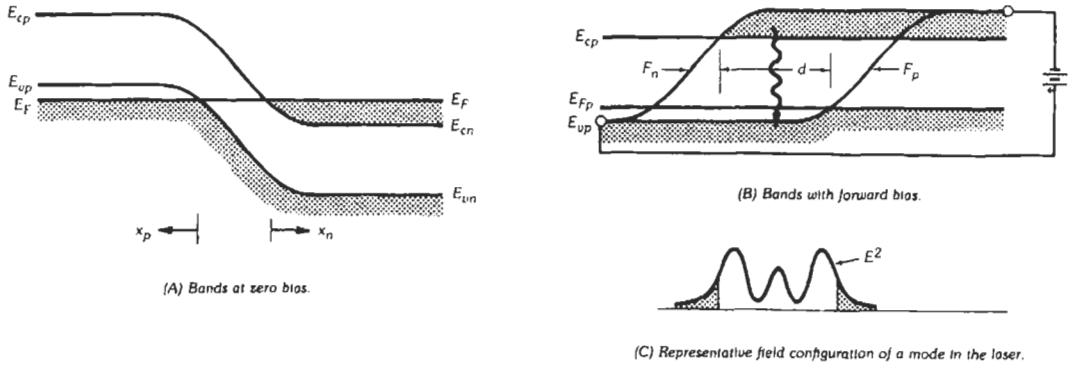


Fig. 13. Semiconductor injection laser. (From J. T. Verdeyen, *Laser Electronics*, Englewood Cliffs, N.J.: Prentice-Hall, Inc., 1981; Fig. 10-19.)

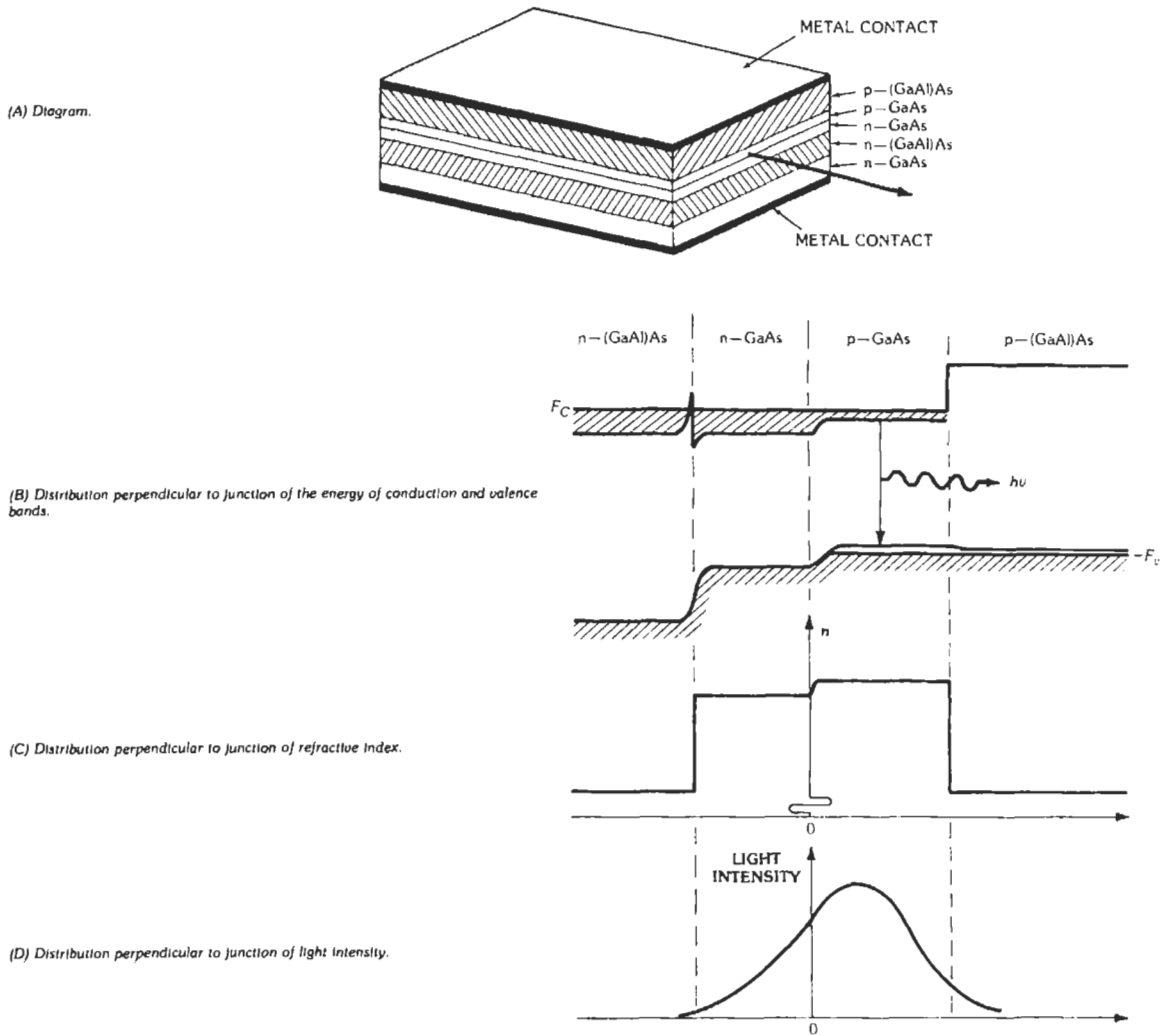


Fig. 14. Four-layer separate confinement heterostructure. (From G. H. B. Thompson, *Physics of Semiconductor Laser Devices*, New York: John Wiley & Sons, Inc., 1980; Figs. 3.16 and 3.17.)

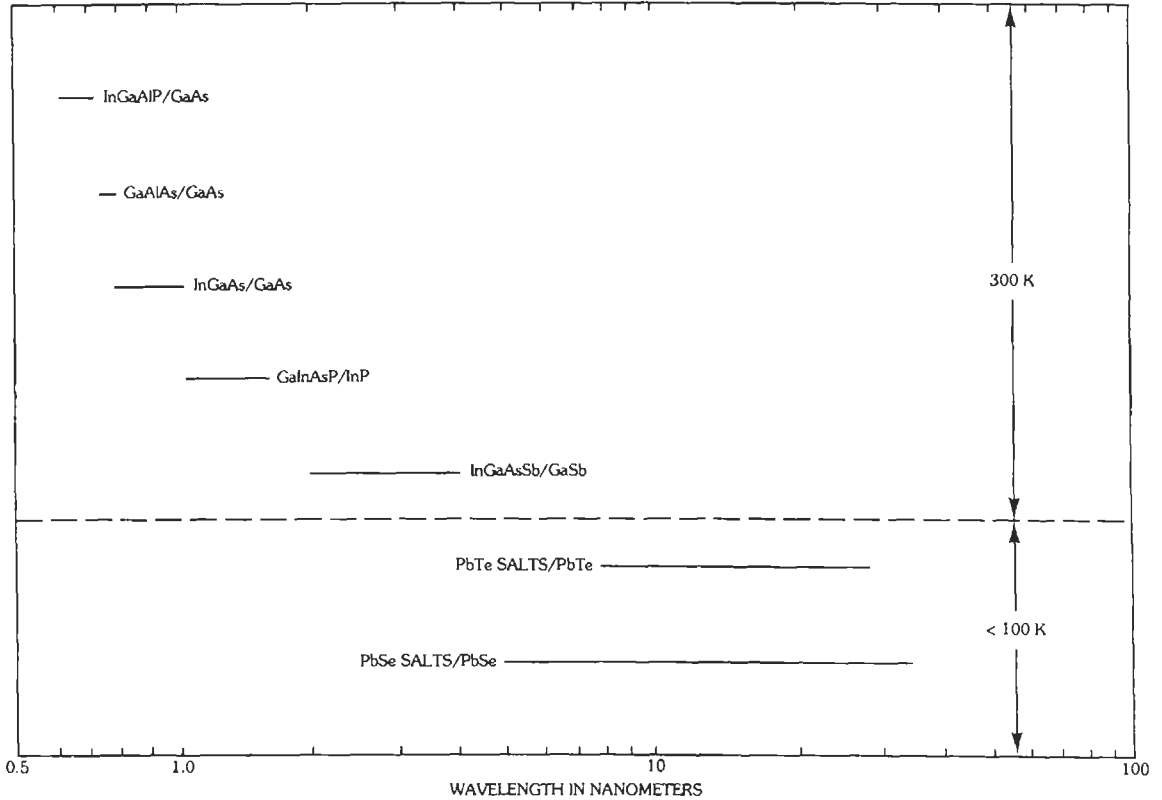


Fig. 15. Common CW compound semiconductor lasers grouped according to the active-layer material and substrate material. The wavelength range refers to different chemical compositions.

If E_ω is approximately constant, $E_{2\omega}$ is seen to be harmonic in z if $\Delta k \neq 0$. If $\Delta k = 0$, then $E_{2\omega} \approx z$, and the growth increases with distance. This latter case is clearly favorable and requires $\Delta k = 0$; it is referred to as a phase-matched condition.

Since most materials are normally dispersive, $n(2\omega) > n(\omega)$, the phase-matched condition is not automatic. It can be achieved by using the anisotropic properties of the index of refraction. An important linear optical class of materials used as harmonic generators is uniaxial, for which $n_1 = n_2 \neq n_3$. Other classes are isotropic, $n_1 = n_2 = n_3$, and biaxial, $n_1 < n_2 < n_3$. For the uniaxial class, $n_1 = n_2 = n_o$, the ordinary index of refraction, and $n_3 = n_e$, the extraordinary index of refraction. If $n_e > (<)n_o$, the material is said to be positive (negative) uniaxial.

A wave polarized in the x - y or 1-2 plane (ordinary ray) has an index of refraction equal to n_o . A wave polarized with a field component along the z or 3 axis (optic axis) is an extraordinary ray and has an effective index of refraction, n_{eff} , given by

$$1/n_{eff}^2 = \cos^2 \theta/n_o^2 + \sin^2 \theta/n_e^2 \quad (\text{Eq. 51})$$

where θ is the angle between the optic axis and the

direction of propagation. A common representation of these is a polar plot of $n(\theta)$, which is referred to as an index surface. Fig. 16A shows the index surface for a negative uniaxial crystal.

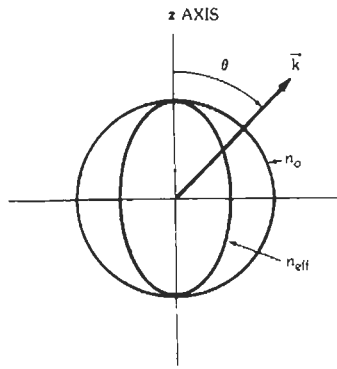
Phase matching may be achieved if the index-surface curves at ω and 2ω cross as shown in Fig. 16B. For $n_e(2\omega) < n_o(\omega)$, there is an angle θ_m for which $n_{eff}(2\omega) = n_o(\omega)$. Since the angular width of the phase matching ($\Delta kL = \pi$) is comparable to the divergence angle of a laser, the phase matching of this type is referred to as critical phase matching. In some cases, the temperature dependence of n_o and n_e is sufficiently different that phase matching can be achieved at elevated temperatures.

Under fully phase-matched conditions, the harmonic intensity is given by

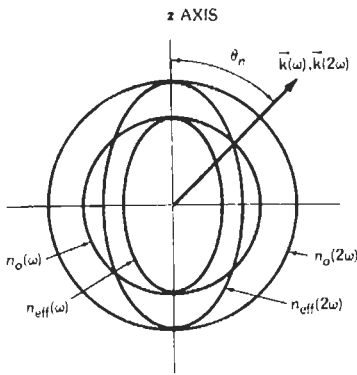
$$I_{2\omega}(z) = I_\omega(0) \tanh^2 [AE_\omega(0)z] \xrightarrow{z \rightarrow \infty} I_\omega(0) \quad (\text{Eq. 52})$$

which shows the full conversion capability of the process, which has an initial growth of

$$I_{2\omega}(z) = A^2 E_\omega^2(0) I_\omega(0) z^2 \cong I^2(0) z^2 \quad (\text{Eq. 53})$$



(A) Negative uniaxial crystal.



(B) Phase matching.

Fig. 16. Wave vectors in a uniaxial crystal.

The term $AE_\omega(0)$ represents an inverse conversion distance, l_{conv} , such that at $z = l_{conv}$, 58% of the fundamental has been converted to a harmonic. For practically all materials with large d values, l_{conv} on the order of centimeters requires initial pumping intensities of $10\text{--}100 \text{ MW/cm}^2$. Although the d_{iq} values tend to be larger at lower frequencies, the equivalent A terms tend to be numerically about the same. Hence, harmonic generation by this method tends to be efficient, 50–80%, but requires large powers. Representative nonlinear materials and coefficients are listed in Table 4.

Pockels Effect

A second effect associated with $\chi^{(2)}$ is the Pockels, or linear electro-optic, effect used to amplitude modulate light. Here, the nonlinear polarization is at the optical frequency and is proportional to $E_\omega E^0$, where E^0 is a nonoptical field. The effect, which has the same basic d matrix as for the case of second-harmonic generation, leads only to a shift in the index of refraction caused by E^0 . The effect is described by a three-dimensional version of the index surface called the index ellipsoid, which relates n_{eff} to the displacement vector rather than to the direction of propagation. With no applied field, the index ellipsoid is of the general form

$$1 = D_1^2/n_1^2 + D_2^2/n_2^2 + D_3^2/n_3^2 \quad (\text{Eq. 54})$$

where D_i are the normalized displacement vectors along the principal axis. The effective index of refraction is simply the length of this vector.

The effect of E^0 is to distort the index ellipsoid. The distortion is expressed in terms of a new matrix, r , which structurally is the transpose of the d matrix. The generalized index ellipsoid is then

$$1 = \sum_{i=1}^3 \left[D_i^2/n_i^2 + D_i^2 \sum_{j=1}^3 \Delta B_{(ij)} \right] + \sum_{i=1}^3 \sum_{j \neq i}^3 D_i D_j \sum_k \Delta B_{(ij)k} \quad (\text{Eq. 55})$$

where $\Delta B_{(ij)k} = r_{(ij)k} E_k^0$, and the indices in parentheses are contracted according to the Voigt notation. Continuing with the example of $\bar{4}2m$ crystals,

$$r = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ r_{41} & 0 & 0 \\ 0 & r_{52} & 0 \\ 0 & 0 & r_{63} \end{bmatrix}$$

and $r_{41} = r_{52}$. The resulting index ellipsoid is

TABLE 4. NONLINEAR OPTICAL COEFFICIENTS

Material	Class	λ (μm)	r_{qk} (10^{-12} m/V)	d_{iq}/ϵ_o (10^{-12} m/V)	Class
KDP	$\bar{4}2m$	1	$r_{63} = -9$	$d_{14} = 0.49$	Uniaxial
			$r_{41} = 8.6$	$d_{36} = 0.47$	
GaAs	$\bar{4}3m$	10	$r_{41} = 1.6$	$d_{14} = 134$	Isotropic
CdTe	$\bar{4}3m$	10	$r_{41} = 6.8$	$d_{14} = 16.7$	Isotropic

$$\begin{aligned}
 1 &= D_1^2/n_o^2 + D_2^2/n_e^2 + D_3^2/n_e^2 \\
 &+ 2D_2D_3r_{41}E_1^0 \\
 &+ 2D_1D_3r_{52}E_2^0 \\
 &+ 2D_1D_2r_{63}E_3^0
 \end{aligned} \tag{Eq. 56}$$

The problem then reduces to finding a new principal axis frame for which the index ellipsoid is diagonal.

Two special cases exist. For longitudinal modulators, the optical field is perpendicular to E^0 , whereas for transverse modulators, it is parallel to E^0 . Fig. 17A illustrates the former, for which the index ellipsoid is diagonal in the indicated primed axis for $\sqrt{2}m$. The index ellipsoid is then

$$\begin{aligned}
 1 &= D_1'^2[1/n_o^2 - r_{63}E_3^0] + \\
 &D_2'^2[1/n_o^2 + r_{63}E_3^0] + D_3'^2/n_e^2 \\
 &= D_1'^2/n_1'^2 + D_2'^2/n_2'^2 + D_3'^2/n_e^2 \tag{Eq. 57}
 \end{aligned}$$

For all materials, the product $r_{63}E_3^0$ is small compared with the indices of refraction, so in the primed axis the new indices of refraction are approximately $n_1' \cong n_o + (1/2)n_o^3r_{63}E_3^0$ and $n_2' \cong n_o - (1/2)n_o^3r_{63}E_3^0$. A wave originally polarized along the 1 axis will have equal field components along the 1' and 2' axes, each of which propagates with a different phase velocity. After exiting from the crystal, the net optical polarization will be rotated from the incident polarization by an amount proportional to $n_1' - n_2' = n_o^3r_{63}E_3^0$. Amplitude modulation is accomplished by placing a polarizer after the crystal. For this orientation, the transmission, T , through the crystal-polarizer combination is given by

$$T = \sin^2(\Gamma/2) \tag{Eq. 58}$$

where,

$$\begin{aligned}
 \Gamma &\text{ is the retardation given by } \omega(n_1' - n_2')L/c \\
 &= \omega n_o^3 r_{63} E_3^0 L/c = \pi V/V_{1/2}, \\
 V &= E_3^0 L, \\
 V_{1/2} &= \pi c/\omega n_o^3 r_{63}, \text{ the half-wave voltage.}
 \end{aligned}$$

The half-wave voltage is the most important parameter in specifying electro-optic modulators.

For the transverse case shown in Fig. 17B, the retardation becomes

$$\begin{aligned}
 \Gamma &= (\omega/c)(n_o - n_e - \frac{1}{2}n_o^3r_{63}E_3^0)L \tag{Eq. 59a} \\
 &= \Gamma_m - \pi V/V_{1/2} \tag{Eq. 59b}
 \end{aligned}$$

where,

$$\begin{aligned}
 V &= E_3^0 t, \\
 V_{1/2} &= 2\pi c t/\omega n_o^3 r_{63} L, \\
 \Gamma_m &= \text{a fixed retardation due to the natural birefringence.}
 \end{aligned}$$

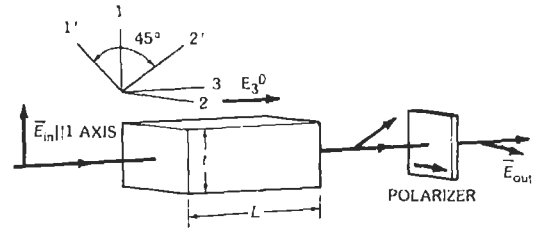
In practice, $t/L \ll 1$, so the half-wave voltage in the transverse case is considerably smaller than in the longitudinal case.

Half-wave voltages of 1–10 kV are not uncommon so that for high-frequency modulation, operational voltages are typically less than $V_{1/2}$. Because of the \sin^2 dependence, small modulating voltages lead to a distortion, since $\sin^2\theta = \theta^2$ for small θ . This can be eliminated by optically biasing the system to the 50% transmission point by using a quarter-wave plate between the polarizer and crystal as shown in Fig. 17C. With this arrangement, the transmission becomes

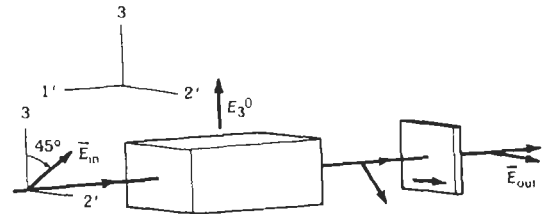
$$T = (1 + \sin\Gamma)/2 \tag{Eq. 60}$$

and results in a change in transmission linear in $V(\Gamma)$ for small modulation. Hence, the standard amplitude modulator is a three-element device. For certain types of crystals, a reduction of the half-wave voltage can be achieved by using multiple crystals connected optically in series but electrically in parallel.

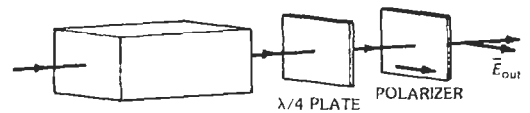
The intrinsic material response is almost instantaneous, so the system response time tends to be determined simply by circuit effects. If these are made small, a second frequency limitation occurs because of the finite optical transit time through the crystal. Operation above 1 GHz generally requires a traveling-wave modulation wherein a certain type of phase matching is used.



(A) Longitudinal modulator.



(B) Transverse modulator.



(C) Bias with 1/4 plate.

Fig. 17. Configurations for amplitude modulation of light.

If the incident optical field is oriented along the primed principal axis, no rotation of the polarization occurs, but modulation of the phase velocity exists. By this method, phase modulation may be achieved.

Stimulated Raman Emission

There are many other different effects associated with $\chi^{(2)}$ and $\chi^{(3)}$ nonlinearities depending on the number of different fields present. Frequency tripling into the deep ultraviolet is commonly done with gases for which $\chi^{(3)}$ is significant. Similarly, another third-order effect leads to a self-induced index of refraction change and causes focusing and defocusing of a beam by virtue of its radial intensity variation.

One important $\chi^{(3)}$ effect used to frequency shift light is the stimulated Raman effect. Here, a polarization and hence a gain is created at frequency ω' due to a strong field at ω where $\hbar\omega - \hbar\omega' = \Delta E$, a characteristic atomic or molecular energy, typically a vibrational energy in a molecule. The interaction can also be viewed as an inelastic scattering process wherein the shift is determined by the characteristic energy. It is possible to have stimulated rotational, vibrational, and electronic Raman scattering with shifts ranging from a few to a few thousand wave numbers. Representative Raman materials and shifts are listed in Table 5.

In the usual Raman situation, the gain is large, so the conversion can be very efficient in terms of photons. In terms of power, the conversion efficiency is limited by the Manley-Rowe limit of $\hbar\omega'/\hbar\omega$. Typically, the requirements for pump intensities lie in the range 10–100 kW/cm².

The normal Raman effect occurs when $\omega' < \omega$ and is called a Stokes wave. The growing first Stokes can reach a sufficient intensity to act as a pump for a second Stokes, and so on. With the presence of a strong pump

and first Stokes, an anti-Stokes wave may be generated at $\omega'' = \omega + \Delta E/\hbar$. Hence, under normal conditions a number of Stokes and anti-Stokes waves may be created. A Raman cell is usually constructed locally because of its simplicity.

REFERENCES

General Laser Theory

Elementary Laser Theory

Siegman, A. *An Introduction to Lasers and Masers*. New York: McGraw-Hill Book Co., 1976.

Verdeyen, J. T. *Laser Electronics*. 2nd ed. Englewood Cliffs, N.J.: Prentice-Hall, Inc., 1989.

Yariv, A. *Optical Electronics*. 4th ed. New York: Holt, Rinehart & Winston, Inc., 1991.

Advanced Laser Theory

Pantell, R., and Puthoff, H. *Fundamentals of Quantum Electronics*. New York: John Wiley & Sons, Inc., 1969.

Sargent, M. III, Scully, M. O., and Lamb, W. E., Jr. *Laser Physics*. Reading, Mass.: Addison-Wesley Publishing Co., Inc., 1974.

Yariv, A. *Quantum Electronics*. 3rd ed. New York: John Wiley & Sons, Inc., 1989.

Specific Lasers

Solid State and Gas Lasers

Brown, D. C., ed. *High-Peak-Power ND: Glass Laser Systems*. Springer Series in Optical Sciences. Berlin: Springer-Verlag, 1979.

TABLE 5. RAMAN MATERIALS AND SHIFTS

Substance	States	Stokes Shift (cm ⁻¹)
H ₂	Vibrational	4155
N ₂	Vibrational	2331
D ₂	Vibrational	2991
Carbon Tetrachloride	Vibrational	460
Carbon Disulfide	Vibrational	656
Benzene	Vibrational	992
Nitrobenzene	Vibrational	1344
Acetone	Vibrational	2921
Water	Vibrational	3651
H ₂	Rotational	365
Lead	Electronic	10 600

- Pressley, R. J., ed. *Handbook of Lasers With Selected Data on Optical Technology*. Cleveland, O.: Chemical Rubber Co., 1971.
- Weber, M. J., ed. *Handbook of Laser Science and Technology*. Vol. I, *Lasers and Masers*; Vol. II, *Gas Lasers*. Cleveland, O.: Chemical Rubber Co., 1982.
- Willett, C. S. *Introduction to Gas Lasers*. New York: Pergamon Press, Inc., 1974.

Semiconductor Lasers

- Agrawal, G. P., and Dutta, N. K. *Long Wavelength Semiconductor Lasers*. New York: Van Nostrand Reinhold, 1986.
- Casey, H. C., Jr., and Panish, M. B. *Heterostructure Lasers: Part A and Part B*. New York: Academic, 1978.
- Kressel, H., ed. *Semiconductor Devices for Optical Communications*. Vol. 39 in *Topics of Applied Physics*. Berlin: Springer-Verlag, 1980.
- Kressel, H., and Butler, J. K. *Semiconductor Lasers and Heterojunction LED's*. New York: Academic Press, Inc., 1977.
- Thompson, G. H. B. *Physics of Semiconductor Laser Devices*. New York: John Wiley & Sons, Inc., 1980.

Dye Lasers

- Hansch, T. W. "Repetitively Pulsed Tunable Dye Laser for High Resolution Spectroscopy." *Applied Optics*, 4, 1972, pp. 895-898.
- Wallenstein, R. "Pulse Dye Lasers," in *Laser Handbook*, Vol. 3. M. L. Stitch, ed. New York: North-Holland Publishing Co., 1979.

Nonlinear Optics

- Bloembergen, N. *Nonlinear Optics*. New York: W. A. Benjamin Inc., 1965.

- Hellwege, K. H. and Hellwege, A. M., eds. *Landolt-Bornstein New Series Group III Vol. 2*. Berlin: Springer-Verlag, 1979. (See also A. Yariv, *Quantum Electronics*, Chapters 14-19.)
- Shen, Y. R. *The Principles of Nonlinear Optics*. New York: John Wiley, 1984.

Pulse Behavior

Modelocking

- Haus, H. A. "Theory of Mode Locking With a Slow Saturable Absorber." *IEEE J. of Quant. Electr.*, QE-11, 1975, pp. 736-746.
- IEEE J. of Quant. Electr.*, QE-19, 1983. Special issue on Picosecond Phenomena and references cited therein.
- Kuzenga, D. J., and Siegman, A. E. "FM and AM Mode Locking of the Homogeneous Laser." Part I, Theory. *IEEE J. of Quant. Electr.*, QE-6, 1970, p. 694.
- Laubereau, A., and Eisenthal, K. B., eds. *Picosecond Phenomena III*. Berlin: Springer-Verlag, 1982.
- Shank, C. V., Ippen, E. P., and Shapiro, S. L., eds. *Picosecond Phenomena*. Berlin: Springer-Verlag, 1978.
- Shapiro, S. L., ed. *Ultrashort Light Pulses*. Berlin: Springer-Verlag, 1979.

Q-Switching

- Hellwarth, R. W. "Control of Fluorescent Pulsations," in *Advances in Quantum Electronics*, J. R. Singer, ed. New York: Columbia University Press, 1961.
- Hellwarth, R. W. "A Modulation of Lasers," in *Lasers I*, A. K. Levine, ed. New York: Marcel Dekker Inc., 1966.
- Wagner, W. G., and Lengyel, B. A. "Evolution of a Giant Pulse in a Laser." *J. of Appl. Phys.*, 34, 1963, p. 2044.

42

Computer Organization and Programming

John Wakerly

Introduction	42-3
Hardware	
Software	
Computer Data Types	
Number Systems and Arithmetic	42-4
Decimal and Binary Positional Number Systems	
Octal and Hexadecimal Numbers	
General Positional Number System Conversions	
Addition and Subtraction of Nondecimal Numbers	
Representation of Negative Numbers	
Twos-Complement Addition and Subtraction	
Binary Multiplication	
Binary Division	
Fixed-Point and Floating-Point Representations	
BCD Representation	
Character Codes	
Basic Computer Organization	42-10
Processors	
Memory	
Input/Output	
Processors and Programming	42-14
Organization of a General-Register Processor	
Instruction Formats	
Instruction Set	
Addressing Modes	
Operands and Operations	
A Machine Language Program	42-22
Assembly Language	

Operation of a Simple Program 42-25

 Indirect Addressing

 Subroutines

Input/Output 42-30

 Buses

 Devices and Interfaces

 Ports

 I/O Programming

INTRODUCTION*

Computers come in all shapes and sizes, from large corporate data-processing machines that fill a room to tiny microcontrollers buried deep in automobiles, stereo systems, and toasters. Even though computers span a wide range of sizes, capabilities, and cost, they share a great many characteristics and operating principles. The goal of this chapter is to describe general principles of computer organization and programming that apply to computers of any category.

Taken as a whole, a computer is an incredibly complex system, with many more levels of detail than, say, the noncomputer part of a radio. Fortunately, a computer system *is* structured into many levels, so that it is easy to understand if taken one level at a time; see Fig. 1.

Hardware

Central to the organization, or indeed existence, of a computer is its *architecture*. The architecture of a computer is really an abstract concept—it is just a definition of the functionality of the computer as seen by a machine-language programmer. A user of a particular computer architecture does not care how its functions are achieved, only that they match their definitions. An example of a particular computer architecture, the Motorola 68000, is given later in the chapter.

At the next level down is *implementation*, the logic structures used to accomplish the functions defined by the architecture. An implementation consists of a block diagram of the system hardware modules, their interconnections, and further refinements of the internal logic structure of each module. Note that implementation is still a somewhat abstract concept in that an implementation can exist on paper without the machine actually having been built. Chapter 43 discusses logic design—the tools and techniques for designing hardware modules.

The lowest level of detail in a computer is its *hardware realization*—the electronic circuits and other components from which the computer is built. Realization encompasses integrated-circuit technologies, such as those discussed in Chapter 20, as well as packaging, interconnections, and all other physical details of the machine.

The bottom three levels of a computer just described are often grouped together and simply called “the hardware.” The higher levels in Fig. 1 are usually grouped together as “software.”

Software

Computer *software* consists of the instructions and data that the computer hardware manipulates to perform useful work. A sequence of instructions for a computer is called a *program*. The data manipulated by a program is called a *data base*, a *file*, *input*, or simply *data*, depending on its nature and extent.

The most primitive instructions that can be given to a computer are those interpreted directly by hardware, in the *machine language* of the computer. Machine-language instructions are encoded as strings of bits in the computer memory, often one instruction per memory location. The processor fetches machine instructions from memory and executes them one by one.

Since it is difficult for humans to read and recognize strings of bits, machine-language programs are written in *assembly language* and translated into bit strings by an *assembler*. Assembly language represents machine instructions by mnemonic names and allows memory addresses and other constants to be represented by symbols rather than bit strings.

Most programs are written in *high-level languages* that allow common operations such as expression evaluation, repetition, assignment, and conditional action to be invoked in a single high-level *statement*. Popular high-level languages include BASIC, FORTRAN, and Pascal.

Few computers execute a high-level language directly. Therefore, a *compiler* is needed to translate a high-level-language program into a sequence of machine instructions that performs the desired task.

Assemblers and compilers are not the only *software tools* that a programmer may encounter. Other useful tools related to program development are interpreters, simulators, and on-line debuggers. Like a compiler, an *interpreter* processes a high-level-language program. Unlike a compiler, an interpreter actually executes the high-level-language program one statement at a time, rather than translating each statement into a sequence of machine instructions to be run later. Most BASIC environments use an interpreter.

A *simulator* is a program that simulates individual machine instructions, usually on a machine other than the one being simulated. A typical use of a simulator is to test programs to be run on a processor before the processor hardware is available. An *on-line debugger* executes a program on a machine one or a few instructions at a time, allowing the programmer to see the effects of small pieces of the program and thereby isolate programming errors (*bugs*).

Text editors are used to enter and edit text in a general-purpose computer, whether the text is a letter, a report, or a computer program. *Text formatters* read text with imbedded formatting commands and produce formatted documents such as this book. Text editors and formatters belong to the area of computing known as *word processing*.

In a medium to large computer system, cooperating

* The material in this chapter is adapted, with permission, from portions of Chapters 1, 4, 5, and 10 of *Microcomputer Architecture and Programming: The 68000 Family* by John F. Wakerly, © 1989 by John Wiley & Sons, Inc., New York.

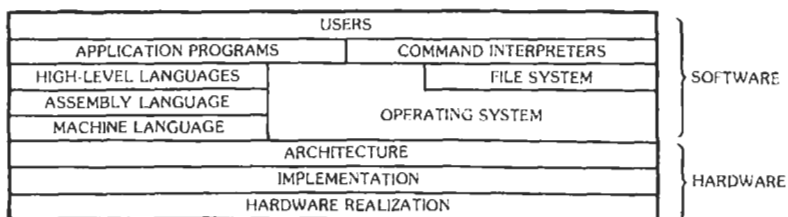


Fig. 1. Levels of detail of a computer.

programs run under the control of an *operating system*. An operating system schedules programs for execution, controls the use of I/O devices, and provides utility functions for all of the programs that run on the computer. Programs and text stored on disks and other mass-storage devices are managed by a *file system*, a collection of programs for reading, writing, and structuring such information in "files." The operating systems in most computers include file systems. Even a very small computer with no mass-storage or file system has a simple operating system, at least to monitor inputs and accept commands from the outside world.

Computer Data Types

The basic unit of information storage in a digital computer is the *bit*, which has a value of either 0 or 1. Obviously a single-bit data type is not very useful for numeric computation, since it only allows us to count from 0 to 1.

By assembling two or more bits into a string, we can represent more than two values or conditions. The bits in a string of n bits can take on 2^n different combinations of values.

Strings of eight bits are usually referred to as *bytes*. The name "byte" was invented at IBM in the early days of electronic computers. The name for a string of four bits is fancifully derived from the byte—the *nibble*.

A bit string manipulated by a computer in one operation is usually called a *word*. Some computers have a *word length* as short as four bits, others as long as 64 bits or more. For many minicomputers and microcomputers, the nomenclature has been standardized to use "word" to describe 16-bit strings and "double word" or "long word" to describe 32-bit strings.

Bits, nibbles, bytes, and words are easy data types to classify because they require differing amounts of storage in the computer memory. There are other data types that are classified not by how much storage they take, but by how they are interpreted and used by the computer hardware and software. For example, the following three data types, all using an eight-bit byte, might be defined for a microcomputer:

Unsigned integer: The byte represents an unsigned integer between 0 and 255.

Signed integer: The byte represents a signed integer between -128 and $+127$.

Character: The byte represents an ASCII character.

Even though all three data types are stored exactly the same way in the computer memory, they may be manipulated differently by the hardware and software.

NUMBER SYSTEMS AND ARITHMETIC

Decimal and Binary Positional Number Systems

Positional number systems are used in all computers and almost all day-to-day business of people. In a *positional number system*, a number is represented by a string of digits in which each digit position has an associated weight. For example, the value D of a 4-digit decimal number, $d_3d_2d_1d_0$, is

$$D = d_3 \cdot 10^3 + d_2 \cdot 10^2 + d_1 \cdot 10^1 + d_0 \cdot 10^0$$

Each digit d_i has a weight of 10^i . Thus, the value of 6851 is computed as follows:

$$6851 = 6 \cdot 1000 + 8 \cdot 100 + 5 \cdot 10 + 1 \cdot 1$$

A decimal point is used to allow negative as well as positive powers of 10 in a decimal number representation. Thus, $d_1d_0 \cdot d_{-1}d_{-2}$ has the value

$$D = d_1 \cdot 10^1 + d_0 \cdot 10^0 + d_{-1} \cdot 10^{-1} + d_{-2} \cdot 10^{-2}$$

In a general positional number system, each digit position has an associated weight of b^i , where b is called the *base*, or *radix*, of the number system. The general form of a number in such a system is

$$d_{p-1}d_{p-2} \cdots d_1d_0 \cdot d_{-1}d_{-2} \cdots d_{-n}$$

where there are p digits to the left of the point and n digits to the right of the point, called the *radix point*.

The value of the number is

$$D = \sum_{p-1 \geq i \geq -n} d_i \cdot b^i$$

the summation of each digit times the corresponding power of the radix. If the radix point is missing, it is assumed to be to the right of the rightmost digit. Except for possible leading and trailing zeros, the representation of a number in a positional number system is unique. (Obviously, 34.85 equals 034.85000, and so on.)

The *binary radix* is used in almost all computers. The allowable digits, 0 and 1, are called *bits*, and each bit d_i has weight 2^i . In the examples below, subscripts distinguish between radix-2 and radix-10 numbers.

$$10001_2 = 1 \cdot 16 + 0 \cdot 8 + 0 \cdot 4 + 0 \cdot 2 + 1 \cdot 1 = 17_{10}$$

$$110.011_2 = 1 \cdot 4 + 1 \cdot 2 + 0 \cdot 1 + 0 \cdot 0.5 + 1 \cdot 0.25 + 1 \cdot 0.125 = 6.375_{10}$$

The leftmost bit of a binary number is called the *high-order* or *most significant bit (MSB)*; the rightmost is the *low-order* or *least significant bit (LSB)*.

Octal and Hexadecimal Numbers

Most computer software uses either the *octal number system* (radix 8) or the *hexadecimal number system* (radix 16) to represent binary numbers. Table 1 shows the binary numbers from 0 through 1111 and their octal, decimal, and hexadecimal equivalents. The octal system uses eight digits, 0-7. The hexadecimal system

requires 16 digits, so the letters A-F are used in addition to the digits 0-9 of the decimal system.

The octal and hexadecimal number systems are useful for representing binary numbers because their radices are powers of two. Since a string of three bits can take on eight different combinations, it follows that each 3-bit string is uniquely represented by one octal digit, according to the third and fourth columns of Table 1. Likewise, a 4-bit string is represented by one hexadecimal digit according to the fifth and sixth columns of the table.

Thus, it is very easy to convert a binary integer to octal (or hexadecimal). Starting at the binary point and working left, we simply separate the bits into groups of three (or four) and replace each group with the corresponding octal (or hexadecimal) digit:

$$101011000110_2 = 101\ 011\ 000\ 110_2 = 5306_8 = 1010\ 1100\ 0110_2 = AC6_{16}$$

$$11011001110101001_2 = 011\ 011\ 001\ 110\ 101\ 001_2 = 331651_8 = 0001\ 1011\ 0011\ 1010\ 1001_2 = 1B3A9_{16}$$

Conversion from octal or hexadecimal to binary is easy, too. We simply replace each octal or hexadecimal digit with the corresponding 3- or 4-bit string. And to convert from octal to hexadecimal or vice versa, we first convert to binary:

$$1573_8 = 001\ 101\ 111\ 011_2 = 0011\ 0111\ 1011_2 = 37B_{16}$$

TABLE 1. BINARY, DECIMAL, OCTAL, AND HEXADECIMAL NUMBERS

Binary	Decimal	Octal	3-Bit String	Hexadecimal	4-Bit String
0	0	0	000	0	0000
1	1	1	001	1	0001
10	2	2	010	2	0010
11	3	3	011	3	0011
100	4	4	100	4	0100
101	5	5	101	5	0101
110	6	6	110	6	0110
111	7	7	111	7	0111
1000	8	10	—	8	1000
1001	9	11	—	9	1001
1010	10	12	—	A	1010
1011	11	13	—	B	1011
1100	12	14	—	C	1100
1101	13	15	—	D	1101
1110	14	16	—	E	1110
1111	15	17	—	F	1111

General Positional Number System Conversions

In the more general case, conversion between two bases cannot be done by simple substitutions; base-10 arithmetic operations are generally required.

The base-10 value of a number in any base can be found from the expansion formula given earlier:

$$D = \sum_{p-1 \geq i \geq -n} d_i \cdot b^i$$

We convert each digit d_i of the number to its base-10 equivalent and expand the formula using base-10 arithmetic. For example:

$$\begin{aligned} 1BE8_{16} &= 1 \cdot 16^3 + 11 \cdot 16^2 + 14 \cdot 16^1 + 8 \cdot 16^0 \\ &= 7144_{10} \end{aligned}$$

$$\begin{aligned} 3176_8 &= 3 \cdot 8^3 + 1 \cdot 8^2 + 7 \cdot 8^1 + 6 \cdot 8^0 \\ &= 1662_{10} \end{aligned}$$

In the opposite direction, to convert a decimal number D to base b , we perform a series of divisions by b . In the first step, we divide D by b ; the remainder is used as the least significant digit (LSD) of the desired result, and the quotient is used as the input to the next step. At each succeeding step, we divide the input number by b , use the remainder as the next-higher-order digit of the result, and use the quotient as the input to the next step. Since a remainder is always less than b , it is always representable as a single base- b digit. The process eventually terminates when a zero quotient is produced. Examples are given below.

$$89 \div 2 = 44 \text{ remainder } 1 \text{ (LSB)}$$

$$\div 2 = 22 \text{ remainder } 0$$

$$\div 2 = 11 \text{ remainder } 0$$

$$\div 2 = 5 \text{ remainder } 1$$

$$\div 2 = 2 \text{ remainder } 1$$

$$\div 2 = 1 \text{ remainder } 0$$

$$\div 2 = 0 \text{ remainder } 1 \text{ (MSB)}$$

$$89_{10} = 1011001_2$$

$$3417 \div 16 = 213 \text{ remainder } 9 \text{ (LSD)}$$

$$\div 16 = 13 \text{ remainder } 5$$

$$\div 16 = 0 \text{ remainder } 13 \text{ (MSD)}$$

$$3417_{10} = D59_{16}$$

Addition and Subtraction of Nondecimal Numbers

Addition and subtraction of nondecimal numbers by hand uses the same technique that we learned in grammar school for decimal numbers; the only catch is that the addition and subtraction tables are different. Table 2 gives the addition and subtraction tables for binary numbers. To add two binary numbers X and Y , we add together the least significant bits with an initial carry (c_{in}) of 0, producing sum ($x + y + c_{in}$) and carry (c_{out}) bits according to the table. We continue processing bits from right to left, including the carry out of each column in the sum for the next column. Two examples of decimal additions and the corresponding binary additions are shown below, with the carries shown as a bit string, C .

C		101111000
X	190	10111110
Y	+ 141	+ 10001101
$X + Y$	331	101001011

C		001011000
X	173	10101101
Y	+ 44	+ 00101100
$X + Y$	217	11011001

TABLE 2. BINARY ADDITION AND SUBTRACTION TABLE

c_{in}	0	0	0	0	1	1	1	1
x	0	0	1	1	0	0	1	1
y	0	1	0	1	0	1	0	1
$x+y+c_{in}$	0	1	1	0	1	0	0	1
c_{out}	0	0	0	1	0	1	1	1
b_{in}	0	0	0	0	1	1	1	1
x	0	0	1	1	0	0	1	1
y	0	1	0	1	0	1	0	1
$x-y-b_{in}$	0	1	1	0	1	0	0	1
b_{out}	0	1	0	0	1	1	0	1

Subtraction is performed similarly, except with borrows (b_{in} and b_{out}) instead of carries:

B		001111100
X	229	11100101
Y	- 46	-00101110
$X - Y$	183	10110111

Addition and subtraction tables can be developed for octal and hexadecimal numbers, or any other desired base. However, rather than memorize tables, many computer engineers use hand-held "hex calculators" to perform arithmetic on binary, octal, and hexadecimal numbers.

Representation of Negative Numbers

So far, we have dealt only with positive numbers. There are many ways to represent negative numbers.

The representation of decimal numbers used in everyday business is called the *signed-magnitude system*. In this system, a number consists of a magnitude and a symbol indicating whether the magnitude is positive or negative. The signed-magnitude system can be applied to binary numbers quite easily by using an extra bit position to represent the sign. Traditionally, the most significant bit (MSB) is used to represent the sign (0 = plus, 1 = minus), and the lower-order bits contain the magnitude. In this system, there are two representations of zero, "plus zero" and "minus zero."

To add two signed-magnitude numbers, we follow the rules that we learned in grammar school. If the signs are the same, we add the magnitudes and give the result the same sign. If the signs are different, we subtract the smaller magnitude from the larger and give the result the sign of the larger. To subtract signed-magnitude numbers, we change the sign of the subtrahend and proceed as in addition.

Despite the conceptual simplicity of the signed-magnitude system, for performance and efficiency reasons negative numbers in a computer are usually represented in *complement number systems*. In such a system, a number is negated not by inverting a sign bit, but rather by "taking the complement."

Complement number systems can be defined for any base or radix, but we shall concentrate on the base-2 radix complement, commonly called the *twos-complement* system. In this system, we assume that a binary number, B , has n bits,

$$B = b_{n-1}b_{n-2} \cdots b_1b_0,$$

so that the radix point is on the right and the number is an integer. If any operation produces a result that requires more than n bits, we throw away the extra high-order bit(s).

By definition, the two's complement of an n -bit number, B , is obtained by subtracting B from 2^n . If B is between 1 and $2^n - 1$, this subtraction will result in another number between 1 and $2^n - 1$. If B is 0, the result of the subtraction is b^n ; this has the form $100 \cdots 00$, where there are a total of $n + 1$ bits. We throw away the extra high-order bit and get the result $00 \cdots 00 = 0$. Thus there is only one representation of zero in the two's-complement system.

From the foregoing, it might seem that a subtraction operation is needed to form the two's complement of a number B . However, this subtraction is avoided by rewriting 2^n as $(2^n - 1) + 1$ and $b^n - B$ as $((2^n - 1) - B) + 1$. Since $2^n - 1$ has the form $11 \cdots 11$, any n -bit number B may be subtracted from it by simply complementing each bit of B . Thus, the two's complement of B is obtained by complementing the individual bits of B and adding 1. Some examples of 8-bit numbers and their two's complements are shown below.

$$\begin{array}{r} 17_{10} = \quad 00010001_2 \\ \quad \quad 11101110 \quad (\text{complement bits}) \\ \quad \quad \quad \quad + 1 \\ \hline \quad \quad 11101111_2 = -17_{10} \end{array}$$

$$\begin{array}{r} 119_{10} = \quad 01110111_2 \\ \quad \quad 10001000 \quad (\text{complement bits}) \\ \quad \quad \quad \quad + 1 \\ \hline \quad \quad 10001001_2 = -119_{10} \end{array}$$

$$\begin{array}{r} 0_{10} = \quad 00000000_2 \\ \quad \quad 11111111 \quad (\text{complement bits}) \\ \quad \quad \quad \quad + 1 \\ \hline 1| 00000000_2 = \quad 0_{10} \end{array}$$

$$\begin{array}{r} -99_{10} = \quad 10011101_2 \\ \quad \quad 01100010 \quad (\text{complement bits}) \\ \quad \quad \quad \quad + 1 \\ \hline \quad \quad 01100011_2 = 99_{10} \end{array}$$

The decimal value of a two's-complement number can be obtained from the same expansion formula that we showed for unsigned numbers, except that the weight of the MSB is -2^{n-1} instead of $+2^{n-1}$ (for example, $-119 = -128 + 8 + 1$ above).

Some computers use the *ones-complement* system, in which the complement of an n -bit number, B , is defined to be $2^n - 1 - B$. Thus, the ones complement of a number is obtained by simply inverting its bits. Although taking the complement is faster, addition may be slower than in the two's-complement system, requiring an "end-around carry" operation. Also, there are two representations of zero, $00 \cdots 00$ and $11 \cdots 11$.

Twos-Complement Addition and Subtraction

Twos-complement numbers, regardless of sign, are added by ordinary binary addition, ignoring any carries beyond the MSB. The result is always the correct sum as long as the range of the number system is not exceeded. Some examples of decimal addition and the corresponding 4-bit twos-complement additions illustrate this:

+3	0011	-2	1110
+ +4	+0100	+ -6	+ 1010
+7	0111	-8	1 1000
+6	0110	+4	0100
+ -3	+1101	+ -7	+1001
+3	1 0011	-3	1101

If an addition operation produces a result that exceeds the range of the number system, *overflow* is said to occur. Addition of two numbers with different signs can never produce overflow, but addition of two numbers of like sign can, as shown by the following examples:

-3	1101	+5	0101
+ -6	+1010	+ +6	0110
-9	1 0111 = +7	+11	1011 = -5
-8	1000	+7	0111
+ -8	+1000	+ +7	+0111
-16	1 0000 = +0	+14	1110 = -2

Fortunately, there is a simple rule for detecting overflow in addition: an addition overflows if the signs of the addends are the same and the sign of the sum is different. Most computers have built-in hardware for detecting overflow.

Two numbers may be subtracted by negating the subtrahend and adding the result to the minuend. A twos-complement number is negated by taking its twos complement. Thus, negating the subtrahend and adding the minuend can be accomplished with only one addition operation as follows: Perform a bit-by-bit complement of the subtrahend, and add the complemented subtrahend to the minuend with an initial carry of 1 instead of 0. Examples are the following:

+4	0100	0100	1—initial carry
- +3	-0011	+1100	
+1	1 0001		

+3	0011	0011	1—initial carry
- +4	-0100	+1011	
-1		1111	

+3	0011	0011	1—initial carry
- -4	-1100	+0011	
+7		0111	

Overflow in subtraction can be detected by examining the signs of the minuend and the *complemented* subtrahend, using the same rule as in addition.

Since twos-complement numbers are added and subtracted by the same basic binary addition and subtraction algorithms as unsigned numbers, most computers have only one type of addition or subtraction instruction that handles both signed and unsigned numbers. However, a program must interpret the results of such an instruction differently, depending on whether it thinks it is dealing with signed numbers (e.g., -8 through +7) or unsigned numbers (e.g., 0 through 15).

Thus, most computers have an *overflow* flag to indicate that the range of signed numbers has been exceeded, and a *carry* or *borrow* flag to indicate that the range of unsigned numbers has been exceeded. Addition and subtraction instructions affect both flags, but a program looks only at the flag for the number system it is using. Also, a program interprets the MSB of a number as a sign bit only when dealing with signed numbers.

Binary Multiplication

In grammar school, we learned to multiply by adding a list of shifted multiplicands computed according to the digits of the multiplier. The same technique can be used to obtain the product of two unsigned binary numbers as shown below:

11	
× 13	
33	
11	
143	
1011	multiplicand
× 1101	multiplier
1011	
0000	shifted multiplicands
1011	
1011	
10001111	product

Forming the shifted multiplicands is trivial in binary multiplication, since the only possible values of the multiplier digits are 0 and 1.

Instead of listing all the shifted multiplicands and then adding, in a computer it is more convenient to add each shifted multiplicand as it is created to a *partial product*. The previous example is repeated in Chart 1, calculated with this technique.

When we multiply an n -bit word and an m -bit word in a computer, the resulting product will take at most $n + m$ bits to express. Therefore, a typical multiplication instruction multiplies two n -bit words and produces a $2n$ -bit double-word product.

Multiplication of signed numbers can be accomplished by using unsigned multiplication and the usual grammar-school rules: Perform an unsigned multiplication of the magnitudes, and make the product positive if the operands had the same sign, negative if they had different signs.

Binary Division

The simplest binary division algorithm is also based on the technique we learned in grammar school, as shown in Chart 2. In both the decimal and binary cases, we mentally compare the reduced dividend with multiples of the divisor to determine which multiple of the shifted divisor to subtract. In the decimal case, we first pick 11 as the smallest multiple of 11 less than 21, and then pick 99 as the smallest multiple less than 107. In the binary case, the choice is somewhat simpler, since the only two choices are zero and the divisor itself. Still, a comparison operation is needed to pick the proper shifted divisor.

Unsigned division in a computer processor is complementary to multiplication. A typical division instruction accepts a double-word dividend and a single-word divisor, and produces single-word quotient and remainder. Such a division *overflows* if the divisor is zero or the quotient would take more than one word to express. The second situation occurs only if the divisor is less than or equal to the high-order word of the dividend.

CHART 2. EXAMPLE OF LONG DIVISION

19	10011	quotient
11)217	1011)11011001	dividend
11	1011	shifted divisor
107	0101	reduced dividend
99	0000	shifted divisor
8	1010	reduced dividend
	0000	shifted divisor
	10100	reduced dividend
	1011	shifted divisor
	10011	reduced dividend
	1011	shifted divisor
	1000	remainder

Fixed-Point and Floating-Point Representations

All of the number systems discussed so far fix the binary point to the right of the rightmost bit. Thus, a 16-bit unsigned number lies in the range 0 through +65 535. This type of number system is most appropriate for programs that count objects or otherwise deal with integer quantities. However, many programs must deal with fractional quantities. For example, in a scientific program that computes a table of positive sines, it might be convenient to fix the binary point to the right of the *leftmost* bit of a 16-bit number, so that numbers range between 0 and 1.999 969 482 421 875 in increments of 0.000 030 517 578 125 (2^{-15}).

In either case, the binary point is fixed in a particular location for all numbers in the system. Thus, the system is called a *fixed-point representation*. Addition and subtraction may be performed directly on fixed-point numbers without regard to the location of the binary point; however, scaling may be required after multiplication and division operations.

To avoid error-prone bookkeeping and to represent a large range of numbers with relatively few bits, a *floating-point representation* can be used to encode explicitly a scale factor in each number. Similar to numbers in scientific notation, a floating-point number has two variable components, the mantissa (M) and the exponent (E). The value of the number is $M \cdot R^E$, where R is a fixed *radix*. In decimal scientific notation, the radix is 10, while in the floating-point representations of most computers, the radix is 2.

Within a computer, a floating-point number is represented by a bit string containing explicit fields for the mantissa and the exponent; the radix is implicit. The number of mantissa bits determines the precision of the number system, while the number of exponent bits determines the dynamic range of the system (the ratio between the largest and smallest nonzero numbers). Both the mantissa and the exponent are signed numbers, allowing representation of a large range of positive and negative numbers with absolute values both greater than and less than 1. Floating-point formats in typical

CHART 1. MULTIPLICATION WITH PARTIAL PRODUCTS

1011	multiplier
× 1101	multiplier
00000000	partial product
1011	shifted multiplicand
00001011	partial product
0000	shifted multiplicand
00001011	partial product
1011	shifted multiplicand
00110111	partial product
1011	shifted multiplicand
10001111	product

computers use a total of 32 bits or more for the mantissa, exponent, and sign fields.

Many computers use a "single precision" floating-point format developed by the IEEE Standards Committee. This format uses a 24-bit mantissa and an 8-bit exponent to represent signed numbers with absolute values in the range of about 2^{-127} to 2^{+127} (10^{-38} to 10^{+38}). The IEEE standard also specifies a 64-bit ("double precision") format and gives guidelines for extended formats (e.g., 128-bit "quad precision" format).

Computer systems may use either of two methods to provide floating-point operations. Some computers have built-in hardware and instructions to process floating-point numbers in a prescribed format. Others have software programs that operate on floating-point numbers using sequences of primitive arithmetic instructions. The hardware approach typically performs operations 100 times faster than the software approach.

A computer that "supports" floating-point numbers is one that provides a minimum set of floating-point instructions including addition, subtraction, multiplication, division, and comparison of floating-point numbers, and conversion to and from integer fixed-point format. Special algorithms are needed to perform these instructions, including scaling or "normalization" of operands and results, as well as arithmetic manipulations of mantissas and exponents.

BCD Representation

The *binary-coded decimal (BCD)* number system encodes the digits 0 through 9 by their 4-bit unsigned binary representations, 0000 through 1001. The codes 1010 through 1111 are not used. Conversions between BCD and decimal representations are trivial, a direct substitution of four bits for each decimal digit. Two BCD digits may be packed into one byte; thus one byte may represent the values from 0 to 99 as opposed to 0 to 255 for a normal unsigned 8-bit binary number. Binary-coded decimal numbers with any desired number of digits may be obtained by using a string of bytes, one byte for each two digits.

Computers "support" BCD representation by providing instructions that perform arithmetic on packed-BCD numbers. Some computers have a "decimal adjust" instruction that corrects the result of an ordinary 8-bit add instruction, assuming that the operands were packed-BCD bytes. Others have decimal add, subtract, multiply, and divide instructions that operate directly on packed-BCD bytes, words, or arbitrary-length strings of bytes.

Character Codes

A string of bits in a computer need not represent a number; in fact most input and output of contemporary computers is nonnumeric. The most common type of nonnumeric data is *text*, strings of characters from

some character set. Each character is represented in the computer by a bit string according to an established convention.

The most commonly used character code is ASCII, the American Standard Code for Information Interchange. This code represents each of 128 different characters by a 7-bit string, as shown in Table 3. Thus, the text string "Yecch!" is represented by a rather innocent looking list of seven 7-bit numbers:

```
1011001 1100101 1100011 1100011 1100011
                    1101000 0100001
```

Some of the 7-bit strings in ASCII denote device control functions instead of "printing" characters. For example, CR (0001101) returns the print head or cursor on a printer or display to the first column, and LF (0001010) advances to the next line. Most of the other control characters are intended for use by data communication links, but different computer systems may use these characters for different functions.

Most computers manipulate an 8-bit quantity as a single unit, a byte, and store one character in each byte. The disposition of the extra bit when 7-bit ASCII is used depends on the system or program. Sometimes this bit is set to a particular value, sometimes it is ignored, and sometimes it is used to encode an additional 128 non-ASCII characters.

An important feature of ASCII is that the bit strings for letters and digits form a reasonable numerical sequence, so that text strings can be sorted by computer instructions that compare numerical values.

BASIC COMPUTER ORGANIZATION

A computer system consists of three major subsystems: processor, memory, and input/output (I/O), as shown in Fig. 2.

Processors

The *processor* (or *central processing unit, CPU*) is the heart of the computer. Shown in Fig. 3, a simple processor contains control circuits for fetching and executing instructions, an arithmetic logic unit (ALU) for manipulating data, and registers for storing processor status and a small amount of data. It also has interface circuits for controlling and communicating with the memory and I/O subsystems.

Different processors have different organizations, both in their internal implementations and in their architectures as seen by programmers. Typical organizations can be roughly classified as accumulator-based processors, general-register processors, and stack machines. Recently, a number of new general-register processor architectures have employed principles of *reduced instruction-set computing (RISC)* to obtain performance advantages over older processor architec-

TABLE 3. AMERICAN STANDARD CODE FOR INFORMATION INTERCHANGE (ASCII), STANDARD NO. X3.4-1968 OF THE AMERICAN NATIONAL STANDARDS INSTITUTE

		$b_6b_5b_4$ (Column)							
$b_3b_2b_1b_0$	Row (hex)	000 0	001 1	010 2	011 3	100 4	101 5	110 6	111 7
0000	0	NUL	DLE	SP	0	@	P		p
0001	1	SOH	DC1	!	1	A	Q	a	q
0010	2	STX	DC2	"	2	B	R	b	r
0011	3	ETX	DC3	#	3	C	S	c	s
0100	4	EOT	DC4	\$	4	D	T	d	t
0101	5	ENQ	NAK	%	5	E	U	e	u
0110	6	ACK	SYN	&	6	F	V	f	v
0111	7	BEL	ETB	'	7	G	W	g	w
1000	8	BS	CAN	(8	H	X	h	x
1001	9	HT	EM)	9	I	Y	i	y
1010	A	LF	SUB	*	:	J	Z	j	z
1011	B	VT	ESC	+	;	K	[k	{
1100	C	FF	FS	.	<	L	\	l	
1101	D	CR	GS	-	=	M]	m	}
1110	E	SO	RS	_	>	N	^	n	~
1111	F	SI	US	/	?	O	-	o	DEL

Control Codes			
NUL	Null	DLE	Data link escape
SOH	Start of heading	DC1	Device control 1
STX	Start of text	DC2	Device control 2
ETX	End of text	DC3	Device control 3
EOT	End of transmission	DC4	Device control 4
ENQ	Enquiry	NAK	Negative acknowledge
ACK	Acknowledge	SYN	Synchronize
BEL	Bell	ETB	End transmitted block
BS	Backspace	CAN	Cancel
HT	Horizontal tab	EM	End of medium
LF	Line feed	SUB	Substitute
VT	Vertical tab	ESC	Escape
FF	Form feed	FS	File separator
CR	Carriage return	GS	Group separator
SO	Shift out	RS	Record separator
SI	Shift in	US	Unit separator
SP	Space	DEL	Delete or rubout

tures. Processors are also classified by the size of the data units that they can manipulate in a single instruction, typically 8, 16, 32, or 64 bits.

The simplest processor architecture organization is *accumulator based*. Such a processor has only one or two registers, called *accumulators*, in which arithmetic and logical operations and data transfers take place. Popular 8-bit microprocessors for embedded-control applications are accumulator based. For example, the Intel 8051 and the Zilog Z80 are single-accumulator processors; the Motorola 68HC11 is a two-accumulator processor. Accumulator-based processors usually have other special-purpose registers, in addition to the accumulator(s), for address manipulation and other operations.

Processors with more than two registers for arithmetic and logical operations are classified as *general-register processors*. In such processors, any operation can be performed using any (or almost any) register.

The IBM 370, the DEC VAX, and most other mainframes and minicomputers have general-register processors. The Motorola 68000, 68020, 68030, and 68040 microprocessors, used in Apple's Macintosh computers, are the most widely used 32-bit general-register microprocessors. Later in this chapter, we use a subset

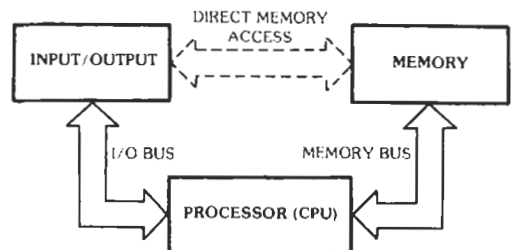


Fig. 2. Block diagram of a typical computer.

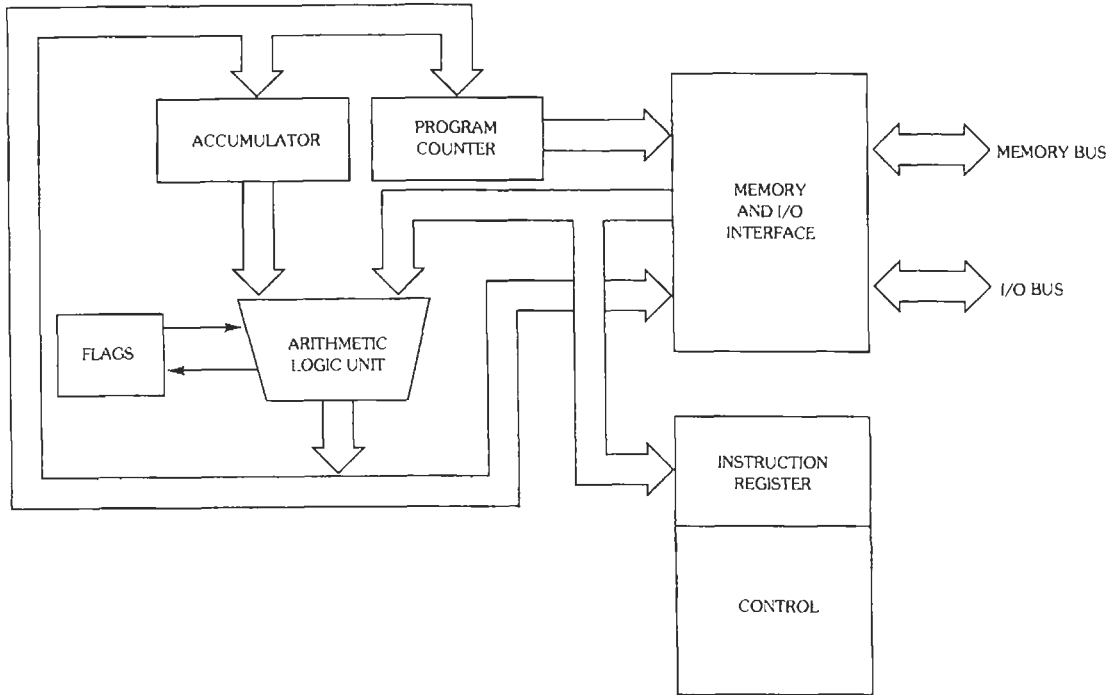


Fig. 3. Block diagram of a simple processor.

of the 68000 architecture to illustrate microprocessor organization and programming.

RISC processors are general-register processors whose instructions have been reduced to a simple, uniform set in order to streamline the implementation of the processor, thereby increasing the speed at which individual instructions are executed. RISC implementations are also characterized by the use of *pipeline* hardware, which allows three or more instructions to be in various stages of execution at any moment, and *caches*, which give the processor instant access to frequently-used instructions and data without having to wait for relatively slow main memory. As a result of these and other features, a typical RISC processor can execute an average of close to one instruction per clock cycle. Advanced RISC processors use *superscalar* and *superpipelining* techniques to execute an average of one to two or more instructions per clock cycle. Popular 32-bit RISC architectures include the Sun SPARC, the Silicon Graphics MIPS, and the Intel 960. The DEC Alpha and the IBM, Motorola, and Apple POWER have 64-bit architectures; the SPARC and MIPS architectures have recently been extended to 64 bits.

Intel's original 16-bit microprocessor, the 8086, is best classified as a single-accumulator processor with many special-purpose registers. Intel's 32-bit extensions to the 8086 architecture, including the 386 and the 486 families, superimposed a general-register structure on the original special-purpose registers, giving it more of the flavor of a general-register architecture. Recent

high-performance implementations of the 386/486 architecture, including the so-called 586 family, have incorporated many of the same implementation methods first applied to RISC and to older mainframe architectures, in order to approach the same high levels of performance that are possible with "pure" RISC architectures.

A *stack machine* has neither general registers nor accumulators, only a stack pointer, SP, that contains the address of a pushdown stack in memory; all operations are performed on the stack. Many of Hewlett-Packard's desktop and hand-held calculators, for example, are stack machines. Also, most accumulator and general-register processors have at least a few stack-oriented instructions.

Typical RISC processors do not typically have stack-oriented instructions, but their programs can implement stacks using appropriate combinations of simple instructions. Likewise, several high-level programming languages, including FORTH and PostScript, have a stack-based view of the world that can be achieved by appropriate combinations of instructions in any processor, stack machine or not.

Memory

The *memory* (or *main memory*) of a computer contains storage for instructions and data, and is tied to the processor via the memory bus in Fig. 2. A *bus* is

simply a bundle of wires or any other physical medium for transferring information. A computer memory has some number of *locations*, each of which stores a *b*-bit quantity. In most contemporary computers, each location stores an 8-bit *byte*. Associated with each location in the memory is a unique binary number called the *address*. If there are *n* locations, then the addresses range from 0 to *n* - 1.

The key feature that distinguishes main memory from other forms of storage in a computer is *random access*—the processor has equally fast access to every location in memory. Random access memory is analogous to a wall of post office boxes; a postal clerk can deposit mail in any box with equal ease. Compare this with the *serial access* method of a letter carrier who visits locations sequentially, in the order of the route. Magnetic tapes provide serial-access memory in computer systems.

Fig. 4 shows how the processor accesses main memory in a typical small computer system. The memory is an array of *n* locations of *b* bits each. To read the data stored at location *X*, the processor places the number *X* on the address bus, and activates a read control signal; the memory responds by placing the contents of location *X* on the data bus. To write a value *V* at location *X*, the processor places *X* on the address bus and *V* on the data bus and activates Write; the memory immediately writes the value *V* in the specified location. Subsequent reading of address *X* will now return the value *V*.

Fig. 5A provides a conceptual view of a memory with 2^{16} bytes; a 16-bit address uniquely specifies any byte in this memory. Larger memories use more address bits.

Fig. 5B shows the physical memory organization used in 32-bit Motorola 68000-family microprocessors. Notice that although the memory is 32 bits (4 bytes) wide, each byte has its own address. Also, more bytes (2^{32}) can be addressed, since addresses are 32 bits wide. In the 68000 memory, a 16-bit *word* is defined as two consecutively numbered bytes, and a 32-bit *long word* is four consecutive bytes. For example, "word 2" consists

of bytes 2 and 3, and "long word 4" consists of bytes 4 through 7. A word whose address is a multiple of two, or a long word whose address is a multiple of 4, is said to be *aligned*.

Fig. 5C shows the physical memory organization used in 32-bit Intel 80x86 microprocessors (the 80386 and its successors). The memory organization is similar to that of the 68000 family, except that the byte numbering within a word is reversed. (Somehow, it is befitting that the 68x and x86 families should have their byte order reversed!)

A 4-byte quantity in the 80x86 is called a *double word*. The 80x86 allow words and double words to be *nonaligned*. For example, "word 7" in Fig. 5C consists of bytes 7 and 8, with byte 7 occupying the most significant position in the word. However, aligned words can usually be accessed in physical memory faster than nonaligned ones.

Note that "word" doesn't mean the same thing in all processor architectures. For example, in most RISC architectures, "word" refers to a 32-bit quantity, and a 16-bit quantity is called a *halfword*.

Several types of memories are commonly used in computer main memory systems. With a *read/write memory (RWM)*, a processor can store data at any address and read back the stored information at any time. Most semiconductor RWMs are *volatile*, in that the stored information disappears when power is removed.

With a *read-only memory (ROM)*, a processor can read the contents of any address at any time, but data values can be stored only once, when the ROM is manufactured. The key advantage of a ROM is that it is a *nonvolatile*; that is, its contents are preserved even if no power is applied.

A *programmable read-only memory (PROM)* is similar to a ROM, except that the customer may store data values (i.e., "program the PROM") using a *PROM programmer*. A PROM chip is manufactured with all its bits at a particular value, typically 1. The PROM programmer can be used to set desired bits to the opposite value. In PROMs that use bipolar technology,

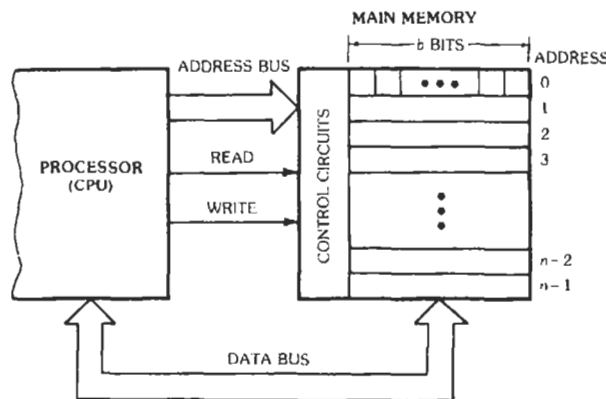


Fig. 4. Main memory in a typical computer.

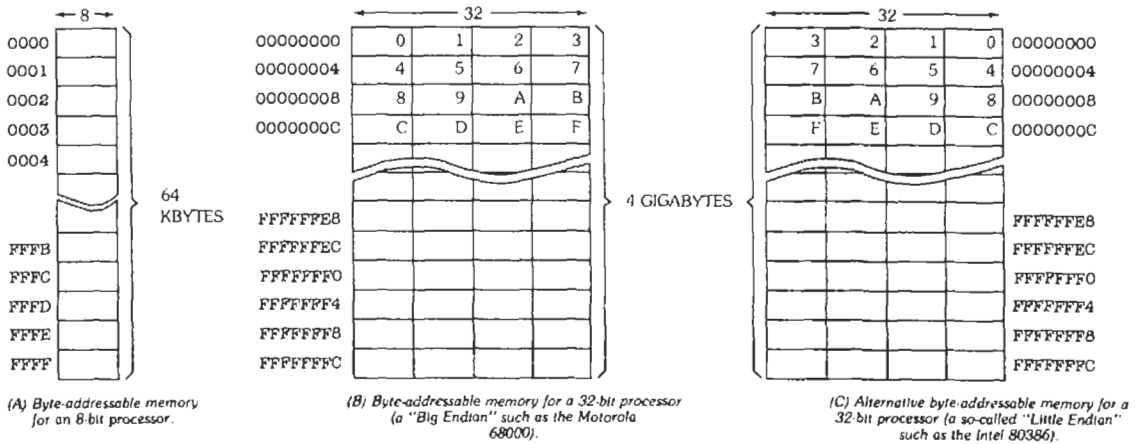


Fig. 5. Memory organizations.

this is typically done by vaporizing tiny fuses inside the PROM corresponding to each bit. In MOS PROMS, this is typically done by forcing electrons into tiny insulated wells corresponding to each 0 bit, where they can remain trapped for 10 years or more.

An erasable programmable read-only memory (EPROM) is similar to a PROM, except that the EPROM can be "erased" to the all-1's state by exposing it to ultraviolet light. An electrically erasable programmable read-only memory (EEPROM or E²PROM) is similar to an EPROM, except that it may be erased electrically.

Every computer has RWM for storing variable data. Depending on the computer and the application, the programs may also be stored in RWM, or they may be stored in ROM, PROM, EPROM, or EEPROM. Storing programs in EEPROM is especially advantageous in embedded-control applications, since software can be updated (for example, via modem) without requiring the chip to be physically replaced.

All of the memories described above are random-access memories, because the processor has equally fast access to every memory location. However, computer jargon has developed so that the acronym "RAM" most commonly refers to read/write memory only.

Input/Output

The input/output (I/O) subsystem of a computer contains peripheral devices for communicating with, observing, and controlling the world outside the computer. Peripheral devices include terminals, printers, communication devices, and mechanical sensors and actuators. Also included in the I/O subsystem are mass storage devices such as magnetic tapes and disks. These devices are used to store information not needed in the main memory at all times, such as applications programs and text files. Not all computers have mass storage devices, but all useful computers have at least one peripheral device, since by definition a peripheral is

the only means for the computer to communicate with the outside world.

The processor writes and reads information to and from peripherals by means of I/O instructions that place commands and data on the I/O bus. In many computers, both main memory and peripherals share the same physical bus. Going one step further, some processors communicate with their peripherals using registers that masquerade as memory locations. In such computer systems, the processor makes no distinction between accessing main memory and accessing peripherals.

In simple computer systems, there is no direct path from peripherals to main memory; the only way to transfer data between a peripheral and memory is for the processor to read it from the peripheral and store it in memory, or vice versa. However, systems requiring a higher data transfer rate incorporate direct memory access (DMA), a link between a special peripheral controller and memory that allows a peripheral to read and write memory without processor intervention.

PROCESSORS AND PROGRAMMING

The Motorola 68000 microprocessor family, used in Apple's Macintosh computers, contains several different microprocessor chips that share a common base architecture. The 68000 architecture is an excellent example of a general-register processor, and we'll use it to introduce general concepts of processor architecture and assembly-language programming in this section.

Organization of a General-Register Processor

Fig. 6 shows the conceptual organization of a typical 68000-family processor. This processor accesses a

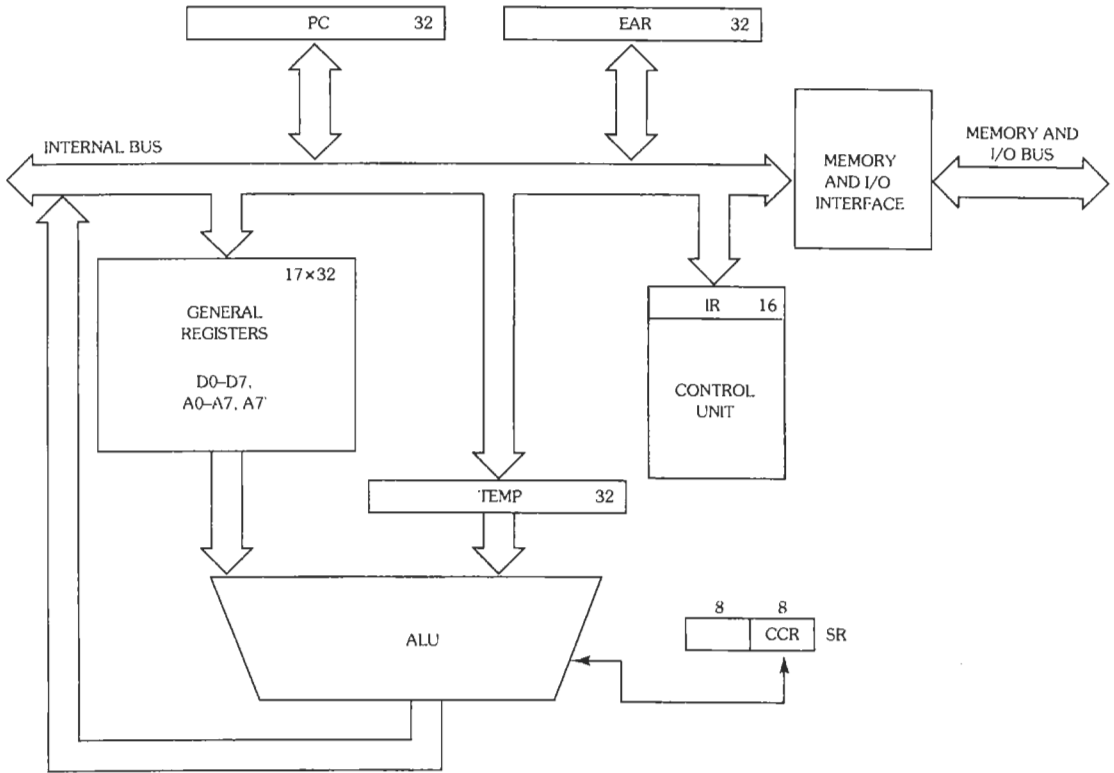


Fig. 6. Organization of a 68000 processor.

memory of up to 2^{32} bytes (16 megabytes) arranged as in Fig. 5B. Thus, an address is specified by a 32-bit number (4 bytes, or a longword).

The operation of the 68000 (or almost any other computer processor) consists of endless repetition of two steps: read the next instruction from memory (the *fetch cycle*) and perform the actions it requires (the *execution cycle*). Several registers and functional units are involved in instruction fetching and execution:

- **Program counter (PC):** A 32-bit register that holds the memory address of the next instruction to be executed.
- **Instruction register (IR):** A 16-bit register that holds the first word of the currently executing instruction.
- **Effective address register (EAR):** A 32-bit register that, when required, holds a memory address at which the processor reads or writes data during the execution of an instruction.
- **General registers (D0–D7, A0–A7, A7’):** A set of seventeen 32-bit registers containing data and addresses. Registers D0–D7 are called *data registers*, and A0–A7 and A7’ are called *address registers*. Register A7 is the *user stack pointer (SP or USP)*, and A7’ is the *supervisor stack pointer (SSP)*.

- **Temporary register (TEMP):** A 32-bit register that holds operands or intermediate results during the execution of an instruction.
- **Condition code register (CCR):** A set of 1-bit flags that the processor sets or clears during the execution of each data manipulation instruction.
- **Status register (SR):** A 16-bit register whose high-order byte contains “processor control” bits and whose low-order byte is the CCR.
- **Arithmetic and logic unit (ALU):** Operates on one or two 8-, 16-, or 32-bit quantities, producing a like-sized result.
- **Control unit:** Decodes instructions and controls the other blocks to fetch and execute instructions.
- **Memory and I/O interface:** Reads and writes memory and communicates with I/O (Input/Output) devices as commanded by the control unit.

Although all of the blocks above are essential to the internal operation of the 68000, only the registers PC, D0–D7, A0–A7, A7’, and SR are explicitly manipulated by instructions and have values that are meaningful after the execution of each instruction. Such registers comprise the *processor state*, and may be shown in a

programming model for the processor, as in Fig. 7. Only these registers are of concern to a programmer.

Instruction Formats

For the purposes of illustration, we introduce a subset of the complete 68000 instruction set. Instructions in the subset are one to four words long.

The first word of each instruction is called the *opcode word*. It contains an operation code (*opcode*) and zero, one, or two *effective-address (EA) fields*. The opcode uniquely specifies the operation to be performed, while each EA field specifies the register or memory location that contains an *operand* to be manipulated by the instruction.

If an operand is a variable stored in memory, then an additional instruction word may be used in computing the address of the memory location. If the operand is an "immediate" constant, then one or two additional instruction words may contain the constant value.

As shown in Fig. 8, the instructions in the subset use six different formats for the opcode word, as required by different operations and operand types. There are many other formats used by the full 68000 instruction set.

To understand why a particular set of formats is used in a given processor architecture, one must understand

the goals, optimizations, and compromises of the architects, which is beyond the scope of this chapter. When you look at the formats in Fig. 8, you will often ask, "Why did they do *that*?" This question is easier to answer in view of the full 68000 instruction set, because there are many more instructions to be encoded, but even the full set of 68000 instructions is encoded somewhat irregularly.

In RISC architectures, instruction formats are quite a contrast from those of the 68000. A hallmark of RISC architecture is that all instructions have the same, fixed width (one word), and there are very few distinct instruction formats. This greatly simplifies and speeds up instruction fetching and decoding.

Returning to the 68000, most of its data-manipulation instructions can operate on bytes, words, or longwords. The operand size is encoded as part of the opcode in MOVE instructions (Fig. 8A) and in the size or op-mode field in other instructions (Figs. 8B and 8C).

Instructions that reference memory contain a 6-bit EA field that specifies the location of the operand; MOVE contains two such fields. In all of these instructions, the EA field can specify one of several "addressing modes" as discussed shortly. Some addressing modes require an additional word of addressing information to be appended to the instruction, thus giving rise to two- to four-word instructions.

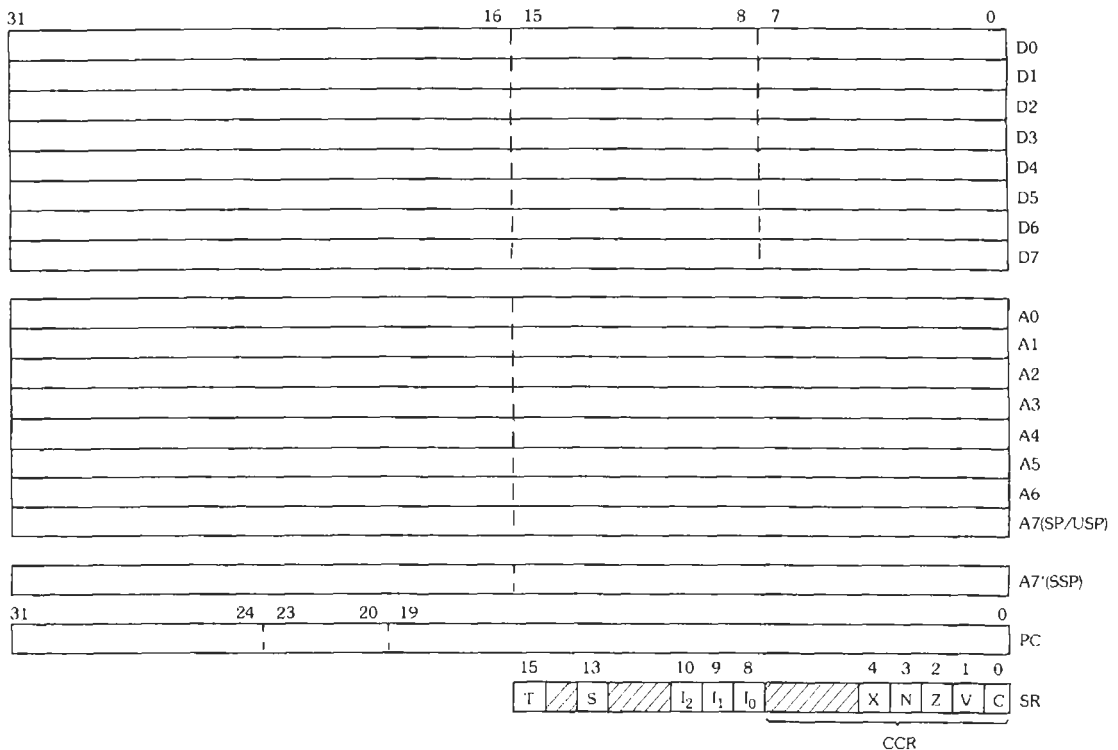


Fig. 7. Programming model for the 68000.

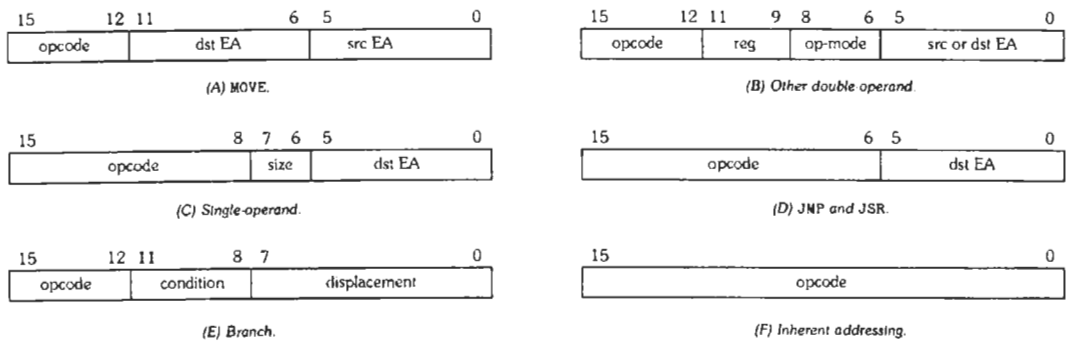


Fig. 8. Format of the opcode word of typical 68000 instructions.

Instruction Set

Table 4 lists a subset of the 68000 instruction set; the full 68000 has many more instructions. There are several different formats for the opcode word, as was shown in Fig. 8. Associated with each opcode word is an alphabetic *mnemonic* that we can use to name and recognize the operation performed by that instruction.

Note that our subset of the 68000 uses only four of the eight possible binary values in the 3-bit *op-mode* field; the complete 68000 uses them all, as listed in Table 5. In both the subset and the complete 68000, there are still more variations in each instruction according to the addressing modes specified for the source (*src*) and destination (*dst*) fields, as discussed in the next subsection. After that, we will describe the actual operations that are performed by each instruction.

Addressing Modes

In the 68000, an operand is specified by a 6-bit source or destination EA field in the instruction, consisting of a 3-bit addressing mode designator and a 3-bit register number, as shown in Fig. 9. In conjunction with the *reg* field, the eight mode combinations are used as follows:

- 0 *reg* specifies a data register (D0–D7) that contains the operand.
- 1 *reg* specifies an address register (A0–A7) that contains the operand.
- 2–6 *reg* specifies an address register that is used, possibly in combination with other values, to determine the effective address of an operand in memory.
- 7 *reg* does not specify a register at all; instead, the three *reg* bits specify one of five addressing modes that do not use a general register; three combinations are left over for future expansion.

A full treatment of 68000 addressing modes is beyond the scope of this chapter, but we shall discuss here the seven most important ones, which are summarized in Table 6.

In *register-direct addressing*, the operand is contained in one of the data registers D0–D7 or in one of the address registers A0–A7. As suggested by the dash lines in Fig. 7, word operations on the data registers use the low-order 16 bits, while byte operations use only the 8 low-order bits; the high-order bits in both cases are unaffected.

Word operations on the address registers also use the low-order 16 bits, but a write operation into an address register extends the most significant bit of the low-order word into the high-order word. This is consistent with the *short-addressing model* of the 68000 memory shown in Fig. 10, in which a 16-bit number can specify a memory address in the highest or lowest 32-kbyte portion of a 4-gigabyte (32-bit) address space.

Byte operations are not allowed on address registers, since an “address” should normally be at least a 16-bit quantity. The complete 68000 instruction set provides long-word operations on address registers, in order to handle addresses up to 32 bits long.

Register-direct addressing is not allowed with the 68000 JMP and JSR instructions, since it does not make sense to jump to a register.

In *immediate addressing*, a 1-, 2-, or 4-byte operand is contained in the second and possibly third word of the instruction, as appropriate for the size of the operation; with 1-byte operands, the high-order byte of the immediate word is unused. Thus, to place the value 1234₁₆ in register D5, we could use the instruction, MOVE.W #1234, D5. By convention, the notation “#” denotes an immediate operand, while the dollar sign (\$) denotes a hexadecimal value. The machine language instruction occupies two words as shown in Fig. 11A.

With immediate addressing, the operand is *constant* value. Immediate mode cannot be used for a “destination” operand, since it would require the instruction to modify itself. Thus, MOVE.W D5, #1234 is not allowed.

In *absolute short addressing*, the instruction contains a 16-bit absolute memory address for the operand. The

TABLE 4. A SUBSET OF THE INSTRUCTION SET OF THE 68000

Mnemonic	Operands	Format	Opcode Word	XNZVC	Description
MOVE.B	asrc, ddst	a	0001ddddddsssssss	-**00	Copy asrc to ddst (byte)
MOVE.W	asrc, ddst	a	0011ddddddsssssss	-**00	Copy asrc to ddst (word)
MOVE.L	asrc, ddst	a	0010ddddddsssssss	-**00	Copy asrc to ddst (long)
MOVEA.W	asrc, An	a	0011rrr001sssssss	-----	Copy asrc to An (word)
MOVEA.L	asrc, An	a	0010rrr001sssssss	-----	Copy asrc to An (long)
LEA.L	csrc, An	b	0100aaa111sssssss	-----	Load An with address of csrc
ADD.B	asrc, Dn	b	1101rrr000sssssss	*	Add asrc to Dn (byte)
ADD.W	asrc, Dn	b	1101rrr001sssssss	*	Add asrc to Dn (word)
ADD.L	asrc, Dn	b	1101rrr010sssssss	*	Add asrc to Dn (long)
ADDA.W	asrc, An	b	1101rrr011sssssss	-----	Add asrc to An (word)
ADDA.L	asrc, An	b	1101rrr111sssssss	-----	Add asrc to An (long)
AND.B	dsrc, Dn	b	1100rrr000sssssss	-**00	AND dsrc to Dn (byte)
AND.W	dsrc, Dn	b	1100rrr001sssssss	-**00	AND dsrc to Dn (word)
AND.L	dsrc, Dn	b	1100rrr010sssssss	-**00	AND dsrc to Dn (long)
CMP.B	asrc, Dn	b	1011rrr000sssssss	-	Set CCR using Dn—asrc (byte)
CMP.W	asrc, Dn	b	1011rrr001sssssss	-	Set CCR using Dn—asrc (word)
CMP.L	asrc, Dn	b	1011rrr010sssssss	-	Set CCR using Dn—asrc (long)
CMPA.W	asrc, An	b	1011rrr011sssssss	-	Set CCR using An—asrc (word)
CMPA.L	asrc, An	b	1011rrr111sssssss	-	Set CCR using An—asrc (long)
CLR.B	ddst	c	0100001000dddddd	-0100	Set ddst to 0 (byte)
CLR.W	ddst	c	0100001001dddddd	-0100	Set ddst to 0 (word)
CLR.L	ddst	c	0100001010dddddd	-0100	Set ddst to 0 (long)
NOT.B	ddst	c	0100011000dddddd	-**00	Complement bits of ddst (byte)
NOT.W	ddst	c	0100011001dddddd	-**00	Complement bits of ddst (word)
NOT.L	ddst	c	0100011010dddddd	-**00	Complement bits of ddst (long)
BNE	offset	e	0110011000000000	-----	Branch if result nonzero (Z = 0)
BEQ	offset	e	0110011100000000	-----	Branch if result zero (Z = 1)
JMP	cdst	d	0100111011dddddd	-----	Jump to cdst address
JSR	cdst	d	0100111010dddddd	-----	Jump to subroutine at cdst address
RTS		f	0100111001110101	-----	Return from subroutine

Notes: asrc = any source operand, cannot be an address register in byte instructions.
 csrc, cdst = control source or destination operand, cannot use register, autoincrement, autodecrement, or immediate addressing modes.

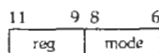
dsrc = data source operand, cannot be an address register.
 ddst = data destination operand, cannot use address register direct, immediate, or relative addressing.
 Dn = one of the data registers, D0–D7. An = one of the address registers, A0–A7.
 offset = 8-bit signed integer added to PC if branch is taken.
 s = bit in src field; d = bit in dst field; r = bit in reg (An or Dn) field; o = bit in offset field.

Condition Code Bits:

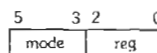
- N—most significant bit of result was 1.
- Z—result was zero.
- V—operation caused twos-complement overflow.
- C—addition or subtraction caused a carry or borrow from most significant bit, or shift caused a 1 bit to “fall off” the end.
- X—when affected, same as C (but not always affected).

Effects of instructions on condition code bits:

- : not affected;
- 1 : always set to 1;
- 0 : always cleared to 0;
- * : set according to operation result as detailed above.



(A) dst operand of MOVE instructions



(B) All other EA fields.

Fig. 9. Operand effective-address (EA) fields in the 68000.

TABLE 5. TYPICAL USE OF THE OP-MODE FIELD IN THE MOTOROLA 68000

Value (binary)	Operand Size	Register Type	Operation
000	byte	data	reg op src → reg
001	word	data	reg op src → reg
010	long	data	reg op src → reg
011	word	addr	reg op src → reg
100	byte	data	dst op reg → dst
101	word	data	dst op reg → dst
110	long	data	dst op reg → dst
111	long	addr	reg op src → reg

most significant bit of this 16-bit address is extended to produce a 32-bit address in the top or bottom 32 K-bytes of the address space, consistent with Fig. 10. Thus, the instruction `MOVE.W $4322, D5` Fig. 11B reads the 16-bit value currently stored at memory location `$4322` and copies (loads) it into register `D5`; the contents of memory are not disturbed. By convention, the absence of the modifier “#” denotes an absolute operand.

In the other direction, `MOVE.W D6, $5432` Fig. 11C stores the value of `D6` into memory location `$5432`, without disturbing `D6`. Thus, an operand with absolute addressing is typically a *variable*, since the program may store different values into the specified memory location at different times.

Since `MOVE` instructions have two EA fields, they can perform memory-to-memory moves. For example, `MOVE.W $2222, $4444` is a three-word long instruction that copies the contents of memory word `$2222` to memory word `$4444`. Immediate-to-memory moves are also possible; for example, `MOVE.W #$1234, $4444` loads memory word `$4444` with the value `$1234`.

In *address-register indirect addressing*, the specified address register (An) contains the address of the operand. In this mode, not only is the operand a variable,

but so is its address. Each time that the instruction is executed, the address register may “point to” a different memory location. This mode is useful for dealing with arrays, lists, and other data structures; we will give an example in a later section on Indirect Addressing.

In *based addressing*, the specified address register (An) contains a *base address* that is “near” the desired operand. The instruction contains a 16-bit, signed displacement value, `displ6`, which is added to the base address to obtain the effective address of the operand. Note that this sum is *not* written back into the address register, which remains unchanged.

Quite often, the base address is the starting address of a block of memory that has been reserved to store data. The displacement value `displ6` specifies the relative position of a particular datum in that block.

While absolute addressing requires a particular datum to be stored at the same absolute address every time the program is executed, based addressing allows addresses to be changed. For example, before giving control to a user’s program, an operating system could reserve a block of memory for the user program to store variables, and pass its base address to the user program in an address register, say `A5`. The user program would refer to all variables using based addressing with fixed displacements (`displ6` values) from the base address

TABLE 6. SOME ADDRESSING MODES OF THE 68000

Name	mode	reg	Notation	Operand	Extra Words
Data-register direct	0	0-7	Dn	Dn	0
Address-register direct	1	0-7	An	An	0
Immediate	7	4	# data	data	1 or 2
Absolute short	7	0	addr	MEM[addr]	1
Address-register indirect	2	0-7	(An)	MEM[An]	0
Based	5	0-7	displ6(An)	MEM[An+displ6]	1
Relative	7	2	raddr16(PC)	MEM[(raddr16-PLC)+PC]	1

Notes: Dn denotes a data register: D0-D7.

An denotes an address register: A0-A7 or SP (same as A7).

data is an 8-, 16-, or 32-bit value as appropriate for the size of the operation.

addr is the 16-bit absolute memory address of the operand.

displ6 is a 16-bit value that is sign extended to 32 bits before being combined with the address register.

raddr16 is a relative address within 32 768 bytes of the instruction. The extension word contains a 16-bit displacement (raddr16-PLC, where PLC is the address of the extension word) that is sign extended to 32 bits and added to the PC to obtain the effective address of the operand when the instruction is executed.

MEM[x] is the 8-, 16-, or 32-bit value beginning at memory address x, as appropriate for the size of the operation.

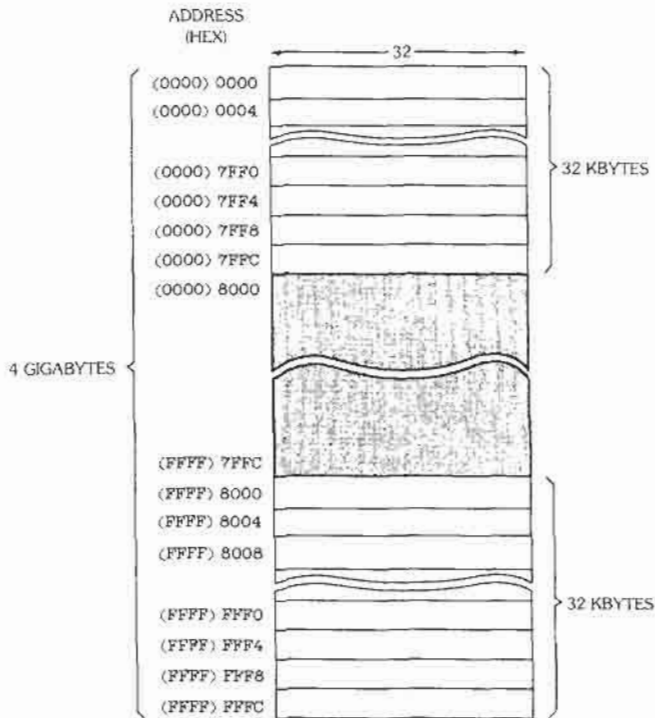


Fig. 10. Short-addressing model for memory in the 68000 family.

contained in A5. This technique is used extensively in programs written for the Macintosh computer.

Relative addressing is used to access constant data that is stored in the same block of memory as the program instructions themselves. Instead of specifying the absolute address of the operand, relative addressing specifies a *relative displacement*—the difference between the instruction address and the operand address. When the program runs, the instruction address (contained in the PC) and the relative displacement (contained in the extension word) are added to form the effective address of the operand. Instructions and operands that use relative addressing are *position independent*—the difference between instruction and the operand addresses remains constant even if the program and its constant data are picked up and moved to a different block of memory. Like based addressing, relative addressing is used extensively in position-independent programs written for the Macintosh and other computers.

The full 68000 processor and other "complex instruction-set computer" (CISC) processors provide many more addressing modes than shown here, in an attempt to provide efficient access to the data structures that are most commonly used in programs. For example, the full 68000 has "autoincrement" and "auto-decrement" modes that are often used to access push-down stacks, and based-indexed and memory indirect modes that are used to access arrays of records and records of arrays.

In contrast, RISC processors are normally characterized by the dearth of addressing modes that they support. In fact, typical RISC processors support only two of the addressing modes described in detail above—immediate and register indirect. Using just these two modes, a RISC processor can emulate any of the more complex addressing modes using a sequence of simple instructions that combine and adjust address values in whatever way is desired.

Operands and Operations

The first part of Table 4 contains double-operand instructions. Most instructions have three variations, indicated by ".B," ".W," and ".L" suffixes, for operating on byte, word, and longword data.

Each double-operand instruction has a "source" operand and a "destination" operand. The 68000 restricts the addressing modes that can be used to specify operands depending on the type of the operand (source or destination), the size of the operand (byte, word, or longword), and the instruction (e.g., MOVE vs MOVEA). In Table 4, *asrc* denotes the most general source operand, which may use any addressing mode (unless the operand is a byte, in which case address register direct addressing is not allowed). A *dsrc* operand may use any mode except address register direct. Similar restrictions apply to the destination operand types *adst* and *dst*, with the additional

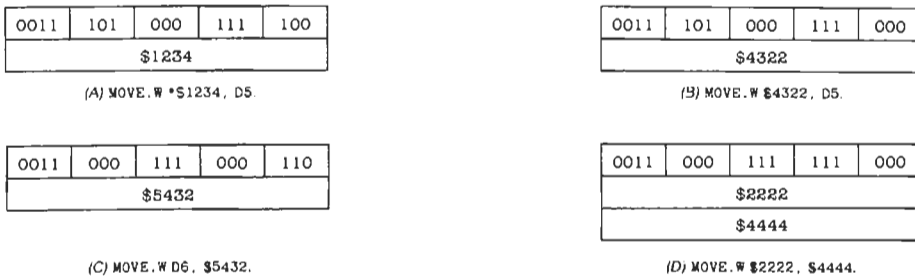


Fig. 11. Encodings of various machine instructions.

restriction that immediate and relative addressing are not allowed. In some cases, operands may be only an address register or a data register; these are explicitly indicated in the table as An or Dn, respectively.

Special cases like the ones described above are the sort of thing that drives compiler writers crazy, and historically motivated the architecture of RISC microprocessors. In a typical RISC microprocessor, any register may be used as the source or destination of any operation, and any operation can use any register. However, memory can be accessed only with two basic, essential instructions—LOAD (from memory) and STORE (into memory). As a result, a compiler has a much better chance of using all of the resources of the processor in an efficient manner.

Now we return to the operations of the 68000. The MOVE instructions simply copy data from one place to another—the operand specified by the src field is copied to the place specified by the dst field, without disturbing the src data. Since src and dst may use different addressing modes, a MOVE instruction can copy from register to register, register to memory, memory to register, or memory to memory. The MOVES are the only instructions that provide all these addressing options; other instructions do not provide memory-to-memory operation or register-to-memory operation (except arithmetic on data registers).

The MOVE instructions affect four of the condition code bits as a side effect. The condition code bits store information about the results of most operations. For example, when affected the zero bit (Z) is set to 1 if the instruction produces a zero result (all bits 0), and to 0 if the instruction produces a nonzero result. An instruction that affects the Z bit can be followed by a conditional branch instruction to perform one of two possible actions depending on the value of Z, as we will describe at the end of this section.

MOVEA is a special version of the MOVE instruction used only when the dst operand is an address register. Its encoding is consistent with that of MOVE, and it differs only in that it does not affect the condition codes; this reflects the 68000 designers' philosophy that the "arithmetic" value of an address should have no significance, and that address registers should contain only addresses. In keeping with this philosophy, note also that instructions with single-byte operands cannot

have an address register as an src or dst operand, since in the 68000 addresses are never as short as one byte.

The LEA instruction is sort of a half-hearted MOVE. Rather than actually read its operand, LEA merely computes the effective address of the operand and deposits it in an address register for later use. This instruction is a result of the 68000 designers' philosophy that, while 32-bit addresses are "expensive" (it would waste a lot of program memory to carry one or two addresses in every instruction), address registers are "cheap" (there are eight of them available). Therefore, if a memory operand is going to be referenced more than once, it makes sense to compute its address only once and keep the address lying around in an address register where it can be quickly and inexpensively used later as needed. We will address this subject in more detail in a later section on Indirect Addressing. But before leaving LEA, also notice that in keeping with the 68000 designers' philosophy about the nature of addresses, the LEA instruction does not affect the condition bits.

The ADD instructions perform a twos-complement addition of a source operand and a data register, and store the result in the data register. The condition codes are affected, so that an ADD may be followed by a conditional branch to test the value of the result. The AND instructions are similar to the ADDs, except that they perform a bit-by-bit logical AND of their operands.

The ADDA instructions add to an address register; there is no byte version of this instruction since addresses are never only 8 bits wide. An ADDA instruction is typically used to compute the address of a particular item in a complex data structure. Like other instructions with address-register destinations, ADDA does not affect the condition bits.

The CMP and CMPA instructions subtract an src operand from a register, and set the condition codes according to the result, without affecting either operand. Thus, CMP and CMPA may be used to compare two operands for equality or other relationships, without disturbing the current values of the operands.

The second part of Table 4 contains instructions that have just one destination operand. The CLR instructions set their operand to 0, obviously useful for initializing variables. The NOT instructions complement each bit of

their operand, useful in conjunction with AND instructions for isolating bits in a single byte or word where different bits or fields contain unrelated data. Note that these instructions do not work on address registers, once again because of the philosophy that addresses have no arithmetic or bitwise significance.

The third part of Table 4 contains "program control" instructions that can affect the normally sequential fetching and execution of instructions. The BNE and BEQ instructions add an offset value to the PC if the Z bit is 0 or 1, respectively. Otherwise, execution continues with the next instruction.

The 68000 processor interprets the offset field of a branch instruction as a signed, two-complement integer in the range -128 to +127. Since instructions must start on even addresses (word boundaries), the range is actually restricted to -128 to +126. At the time that the processor tests the branch condition, the PC is already pointing to the next instruction, and so the branch offset is limited to -128 to +126 bytes from the next instruction.

According to this convention, if the branch offset is 0, the next instruction should be executed whether the condition is true or false. However, since such an instruction performs no useful work and would therefore never be used, the designers of the 68000 used the offset=0 case to specify another instruction type in

the full instruction set. Also, since instructions are aligned on word boundaries, odd offsets such as -127 are not used; the 68020 uses this to encode yet another instruction type.

The JMP instruction in the 68000 is equivalent to the GOTO instruction in C, Pascal, and most other high-level languages. It unconditionally transfers control to an instruction starting at the destination address. It does this by simply replacing the current value of the PC with the destination address.

The JSR and RTS instructions are used to call and return from subroutines, as we will discuss in detail in a later section on Subroutines.

A MACHINE LANGUAGE PROGRAM

Table 7 shows the values stored in memory for a sequence of instructions and data that form a program for multiplying 123 by 456. A list of machine instructions stored in memory, as defined by the two left-hand columns of the table, is called a *machine language program*. Even though they don't mean much to a human reader, these two columns completely specify the operations to be performed by the computer. The remaining columns of Table 7 are an *assembly language*

TABLE 7. MEMORY CONTENTS FOR A SEQUENCE OF INSTRUCTIONS AND DATA

Machine Language		Assembly Language			
Address (hex)	Contents (hex)	Label (sym)	Opcode (mnem)	Operand (sym)	Comments
00000000		*			Program to multiply MCND by MPY.
00000000			ORG	\$1000	
00001000		*			D0 will accumulate product.
00001000	4240	MULT	CLR.W	D0	
00001002	3238		MOVE.W	MPY,D1	D1 holds loop count (multiplier).
00001004	101C				
00001006	670A		BEQ	DONE	Done if count is down to zero. Else add MCND to product
00001008	D078	LOOP	ADD.W	MCND,D0	
0000100A	101E				and do loop MPY (D1) times.
0000100C	D27C		ADD.W	#-1,D1	
0000100E	FFFF				BNE
00001010	66F6			LOOP	
00001012	31C0	DONE	MOVE.W	D0,PROD	Save product.
00001014	101A				
00001016	4EF8		JMP	MON	Return to debug monitor.
00001018	7008				
0000101A	????	PROD	DS.W	1	Storage for PROD.
0000101C	007B	MPY	DC.W	123	
0000101E	01C8	MCND	DC.W	456	Multiplier value. Multiplicand value.
00001020		MON	EQU	\$7008	
00001020			END	MULT	Address of debug monitor.

Symbol Table

DONE	00001012	LOOP	00001008	MCND	0000101E	MON	00007008
MPY	0000101C	MULT	00001000	PROD	0000101A		

Notes: hex = hexadecimal; sym = symbolic; mnem = mnemonic.

program that gives the same information in a much more understandable symbolic form. In the next two subsections, we will discuss assembly language and then describe the actual behavior of the machine-language program.

Assembly Language

Machine language and assembly language are not the same thing. The *machine language* of a computer is the set of bit strings recognized as instructions; the actions performed by each machine instruction are defined by the computer hardware. *Assembly language* is a software tool, a symbolic language that can be directly translated into machine language by a system program called an *assembler*. The output of an assembler is an *object module* containing the bit strings that make up the machine language program, and information that tells a *loader* program where to place these bit strings in the computer memory.

The process of developing programs in assembly language is illustrated in Fig. 12. A typical programmer uses a *text editor* to create a text file containing an assembly language program. The names *source file* and *source program* are often used for assembly language and high-level language text files. The assembler accepts a source program as input, checks for format errors, and produces an object module containing the machine language program. A loader then loads the object module into the memory of the target machine. There the machine language program is run, possibly with the aid of a *debugger*.

The loader, debugger, and machine language program described above *must* run on the target machine; the text editor and assembler may run there or on a different machine. An assembler that runs on one machine and produces object modules for another is called a *cross assembler*. For example, it is possible to create an assembly language program by using a text editor on a laptop computer, transfer the source file by

using a serial data link to a workstation that runs a cross assembler, and transfer (or *download*) the object module to a loader in an embedded microcomputer by way of another serial data link.

Although every computer manufacturer defines a standard assembly language for a new machine when it is introduced, other vendors may define different assembly languages for the same machine. While the effect of each machine language instruction is fixed in hardware, the person who defines an assembly language is free to specify:

- A mnemonic for each machine language instruction
- A standard format for the lines of an assembly language program
- Formats for specifying addressing modes and other instruction variations
- Formats for specifying character and integer constants in different bases
- Mechanisms for associating symbolic names with addresses and other numeric values
- Mechanisms for defining constant data to be stored in memory along with the instructions when the program is loaded
- Directives that specify how the program is to be assembled

In a typical assembly language for the 68000 or any other processor, each line has four fields arranged as shown below:

LABEL OPCODE OPERANDS COMMENTS

The LABEL field is optional. A *label* is simply an identifier (or *symbol* in assembler parlance), that is, a sequence of letters and digits beginning with a letter. Every symbol in an assembly language program is assigned a value at the time that it is defined; the assembler program keeps track of labels and their values

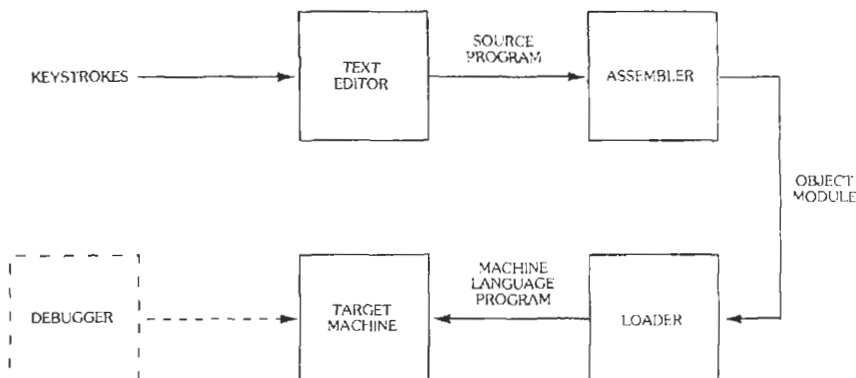


Fig. 12. Assembly language program development.

by an internal *symbol table*. For example, in the statement below,

```
START MOVE.W XX,DO      DO := con-
contents of memory location XX.
```

the value of the symbol *START* equals the memory address at which the *MOVE.W* instruction is stored. In general, the value of a symbol is the memory address at which the corresponding instruction or data value is stored (except with the *EQU* pseudo-operation, described later). Each symbol may be defined only once, but may be referenced as often as needed.

Symbols only have values at the time a program is being assembled, that is, at *assembly time*. Symbols exist in the source text of the program, and in the environment of the assembler, but do not appear in the machine-language program. For instance, when our example program is loaded into memory (at *load time*) or is actually run (at *run time*), there is no way to discover, simply by looking at the object module or at memory, that the programmer had associated the symbol *MULT* with address \$00001000, or even that there had been a symbol *MULT* at all. (However, *symbolic debuggers* may access this information in the run-time environment by referring to a copy of the source code and the symbol table that were used by the assembler.)

The *OPCODE* field contains the mnemonic of either a machine instruction or a pseudo-operation or assembler directive. In 68000 assembly language, many opcode mnemonics have a *size suffix* of *.B*, *.W*, or *.L* to indicate the size of the operands and operation. The size suffix is optional for operations that have only one size (such as *LEA*).

Depending on the contents of the *OPCODE* field, the *OPERAND* field specifies zero or more operands separated by commas. An operand is an expression consisting of symbols, constants, and operators such as *+* and *-*. The simplest expression consists of a single symbol or constant.

In 68000 assembly language, a sequence of decimal digits denotes a *decimal constant*, and a sequence of hexadecimal digits preceded by \$ denotes a *hexadecimal constant*. *Character constants* are surrounded by single quotes (e.g., 'A') and have the corresponding ASCII value.

The *COMMENTS* field is ignored by the assembler, but it is used by the programmer to convey a high-level explanation of the program to human readers. It is also possible to use an entire line as a comment by placing and asterisk (*) at the beginning of the line.

Four different *assembler directives* appear in Table 7; their functions are described below:

- *ORG* (origin). The operand is the address at which the next instruction or datum is to be deposited when the program is loaded into memory. Subse-

quent instructions and data are deposited in successive memory addresses.

- *DS.W* (define storage—word). The operand is a number of memory words to be skipped without storing any instructions or data, thereby reserving space to be used by variables in a program.
- *DC.W* (define constant—word). The specified word value is stored into memory when the program is first loaded into memory, thereby establishing a constant value that may be accessed when the program is run.
- *EQU* (equate). The identifier in the label field is assigned the value in the operand field, instead of the value in the address field. This makes the identifier a synonym for a constant value for the duration of the assembly process.
- *END* (end assembly). This pseudo-operation denotes the end of the text to be assembled. Its operand, if present, is the address of the first executable instruction of the program.

Of course, there are many other directives in a full assembly language, and different assembly languages may use different names and conventions for the same thing.

The input of an assembler is the text file containing the source code, while the output consists of an object module (or "object file" or "binary file") and a listing similar to Table 7.

To understand a little more about how an assembler works, it is useful to examine the *Address* column in the listing. The *Address* column indicates the address at which the assembler "thinks" that it should assemble the instruction on the current line. By default, this address is initialized to zero.

The *ORG* directive tells the assembler to start assembling subsequent instructions starting from a new address, and hence changes the *Address* column on the *next* line. Some directives, such as *EQU*, do not assemble any instructions or data and therefore do not affect the *Address* column at all. Others, such as *DS.W*, advance the *Address* value but do not cause anything to be loaded into the corresponding memory location at load time.

For its own and the programmer's use, the assembler creates a *symbol table* that lists the numeric value associated with each symbol. As evident in Table 7, in most cases the symbol receives the value of the *Address* column in the line in which it is defined, except in the case of *EQU* statements, which assign the symbol the value in the *Operand* column. In all cases, the symbols and their values exist only at assembly time.

Many instructions in the 68000 are more than one word long. It is therefore convenient to compress the program listing, showing all words associated with the same instruction on one line as in Table 8.

OPERATION OF A SIMPLE PROGRAM

We are now ready to explain the program in Table 8. It multiplies MCND by MPY by initializing the product to 0 and then adding MCND to it MPY times.

Notice that the values of the multiplier and multiplcand are "passed" to the program in fixed memory locations, \$101A and \$101C, that are initialized to 123 and 456 when the program is loaded. If, after loading the program, we start the program at location \$1000, it will indeed compute the product of 123 and 456. Alternatively, if after loading the program we place different numbers in locations \$101A and \$101C, then the program will compute the product of the new numbers. In any case, the program "returns" its result by placing the result in memory location \$1018, which may be examined by whoever ran the program in the first place.

The program uses registers to accumulate the product (in D0) and to keep track of the multiplier count (in D1), because registers can be accessed faster than main memory in most computers, and hence the program will run faster. In fact, we could optimize the program to run even faster by placing MCND into a register (say D2) before entering the multiplication loop.

The program execution is traced by Table 9, which shows the values of relevant registers and memory *after* each instruction is executed. Only the low-order words of PC, D0, and D1 are shown; the high-order words are not affected by this program. During its execution, the program executes 374 instructions, including three instructions in a loop that is executed 123 times.

So, the program begins by setting D0 to zero and copying MPY into D1. The MOVE instructions set the

condition bits in CCR according to the value stored. In particular, the Z bit is set to 1 if a zero value (0000₁₆) was stored, or else Z is cleared to 0. Therefore, if MPY is zero, the Z bit is 1 when the BEQ DONE instruction is first reached, and the branch will be taken, leaving zero in D0. The purpose of this instruction is to exit the program with a product of 0 if the multiplier happens to be zero. Notice that the branch instruction gives a *relative offset* of \$0A, so that either the next instruction (at address \$1008) or the instruction at the branch offset (address \$1008+\$0A = \$1012) is executed.

If the BEQ branch is not taken, then the two ADD instructions add MCND to the product in D0 and subtract 1 from the loop count in D1. Like a MOVE instruction, an ADD instruction sets Z to 1 if its result is zero. Thus, if D1 has *not* been reduced to zero, the BNE instruction will branch back to address \$1008. (The 8-bit signed offset \$F6, interpreted in twos-complement, equals -\$0A; hence the branch is to address \$1012-\$0A = \$1008.)

Eventually, after executing the loop 123 times, the program reaches the MOVE instruction at address \$1012, which saves the accumulated product in memory. Finally, the program unconditionally jumps to address \$7008, where we assume there exists a debug monitor that prints out or gives the user some way to examine the computational result stored at location \$1018.

Indirect Addressing

The program in Table 8 manipulated only simple variables and constants. More complicated data structures such as arrays, stacks, queues, records, and lists are used in almost all programs. Consider the problem

TABLE 8. COMPRESSED PROGRAM LISTING

Machine Language		Assembly Language			
Address (hex)	Contents (hex)	Label (sym)	Opcode (mnem)	Operand (sym)	Comments
00000000		*			Program to multiply MCND by MPY.
00000000			ORG	\$1000	
00001000		*			
00001000	4240	MULT	CLR.W	D0	D0 will accumulate product.
00001002	3238 101C		MOVE.W	MPY,D1	D1 holds loop count (multiplier).
00001006	670A		BEQ	DONE	Done if count is down to zero.
00001008	D078 101E	LOOP	ADD.W	MCND,D0	Else add MCND to product
0000100C	D27C FFFF		ADD.W	#-1,D1	and do loop MPY (D1) times.
00001010	66F6		BNE	LOOP	
00001012	31C0 101A	DONE	MOVE.W	D0,PROD	Save product.
00001016	4EF8 7008		JMP	\$7008	Return to debug monitor.
0000101A	????	PROD	DS.W	1	Storage for PROD.
0000101C	007B	MPY	DC.W	123	Multiplier value.
0000101E	01C8	MCND	DC.W	456	Multiplicand value.
00001020		MON	EQU	\$7008	Address of debug monitor.
00001020			END	MULT	

TABLE 9. REGISTER AND MEMORY CONTENTS AFTER EXECUTING INSTRUCTIONS IN MULTIPLICATION PROGRAM

Step	Instruction	PC	D0	D1	Z	MEM[1018] (PROD)
0	...	1000	????	????	?	????
1	CLR.W D0	1002	0000	????	1	????
2	MOVE.W MPY,D1	1006	0000	007B	0	????
3	BEQ DONE	1008	0000	007B	0	????
4	ADD.W MCND,DO	100C	01C8	007B	0	????
5	ADD.W #-1,D1	1010	01C8	007A	0	????
6	BNE LOOP	1008	01C8	007A	0	????
7	ADD.W MCND,DO	100C	0390	007A	0	????
8	ADD.W #-1,D1	1010	0390	0079	0	????
9	BNE LOOP	1008	0390	0079	0	????
10	ADD.W MCND,DO	100C	0558	0079	0	????
11	ADD.W #-1,D1	1010	0558	0078	0	????
12	BNE LOOP	1008	0558	0078	0	????
	...					
367	ADD.W MCND,DO	100C	D950	0002	0	????
368	ADD.W #-1,D1	1010	D950	0001	0	????
369	BNE LOOP	1008	D950	0001	0	????
370	ADD.W MCND,DO	100C	DB18	0001	0	????
371	ADD.W #-1,D1	1010	DB18	0000	1	????
372	BNE LOOP	1008	DB18	0000	1	????
373	MOVE.W D0,PROD	1016	DB18	0000	0	DB18
374	JMP \$7008	7008	DB18	0000	0	DB18

of initializing the components of an array of five bytes to zero. An assembly language solution is shown in Table 10. Note that the operand expressions $Q+1$, $Q+2$, and so on are evaluated at assembly time. The choice of a five-component array above was very judicious—the corresponding program for a 100-component array would have 101 instructions! *Indirect addressing* avoids this problem by taking the address from an address register at run time. Thus we can write a loop to initialize an array, in which an address register points to a different array element on each iteration of the loop.

The program in Table 11 solves the array initialization problem by using indirect addressing. As shown in Fig. 13, it initializes A0 to point to the first component of Q, and A1 to point just past the last component of Q. Then it executes a loop that clears successive components of Q, incrementing A0 once per iteration of the loop. The CMPA instruction compares the contents of A0 with A1, which was initialized to contain the address just past the last array component. CMPA sets Z to 1 if A0 and A1 are equal, which causes the BNE instruction to terminate the loop.

Notice the use of the LEA instruction in Table 11, in particular the second one, which initializes A1. Performing this initialization outside the main loop reduces the length and execution time of the instructions inside the loop. An obviously less desirable alternative would have been to put the LEA instruction inside the loop, right before CMPA.

Not only does the program in Table 11 occupy fewer bytes than the one in Table 10, but it also stays the same length for an array of any size. The program is easily

modified to work on a different length array by changing the occurrences of the length "5" to the desired length.

Many other variations on indirect addressing are found in the full 68000 and most other CISC processors. RISC processors, on the other hand, typically have register indirect addressing as their one and only memory addressing mode.

Subroutines

A *subroutine* is the machine language equivalent of a procedure or function in C, Pascal, and other high-level languages: a sequence of instructions, defined and stored only once, that may be invoked (or *called*) from many places. In order to use subroutines, we need instructions to save the current value of the PC (the *return address*) each time the subroutine is called, and restore it when the subroutine is finished. In the 68000, these instructions save and restore PC by use of a stack.

A *pushdown stack* (or simply a *stack*) is a one-dimensional data structure in which values are entered and removed one item at a time at one end, called the *top of stack*. A register called the *stack pointer* (SP) points to the top of stack. An item is entered by a *push* operation that advances the stack pointer to the next available memory location and then stores the item at the top of stack. An item is removed by a *pop* operation that removes the item at the top of stack and then backs up the stack pointer.

In the 68000, the JSR and RTS instructions perform subroutine calls and returns using address register A7 (also called SP) as a stack pointer. Any 68000 program

TABLE 10. INITIALIZING AN ARRAY THE HARD WAY

Machine Language			Assembly Language		
Address (hex)	Contents (hex)	Label (sym)	Opcode (mnem)	Operand (sym)	Comments
00000000			ORG	\$3000	Set components of Q to zero.
00003000	4238 3100	INIT	CLR.B	Q	First component.
00003004	4238 3101		CLR.B	Q+1	Second component.
00003008	4238 3102		CLR.B	Q+2	Third component.
0000300C	4238 3103		CLR.B	Q+3	Fourth component.
00003010	4238 3104		CLR.B	Q+4	Fifth component.
00003014	4EF8 7008		JMP	MON	Return to debug monitor.
00003018			ORG	\$3100	
00003100	???? ????	Q	DS.B	5	Reserve 5 bytes for array.
00003105		MON	EQU	\$7008	Address of debug monitor.
00003105			END	INIT	

that uses subroutines is required to reserve a small area of memory for a pushdown stack for return addresses. At the beginning of such a program, the stack pointer SP must be initialized to point at this area by a MOVE.W #addr, SP instruction.

As shown in Fig. 14, SP points to the top item in the stack, or just past the stack area if the stack is empty. SP is decremented by 2 before each word is stored on the stack, and incremented by 2 after each word is popped from the stack.

Since a return address may be up to 32 bits wide, two words on the stack are used for each return address. As usual, the high-order word of the 32-bit return address is stored in the lower-numbered address.

The JSR addr instruction saves the address of the next instruction by pushing it onto the stack and then jumps to the instruction at location addr, the first instruction of the subroutine. At the end of the subroutine, RTS pops an address from the stack into PC, effecting a return to the original program sequence.

A stack is the most appropriate data structure for saving return addresses, because it can store more than one return address when subroutines are nested, that is, when one subroutine calls another. The number of levels of nesting is limited only by the size of the memory area reserved by the programmer for the stack.

A detailed example program using subroutines is presented in Table 12. Before describing it, we should

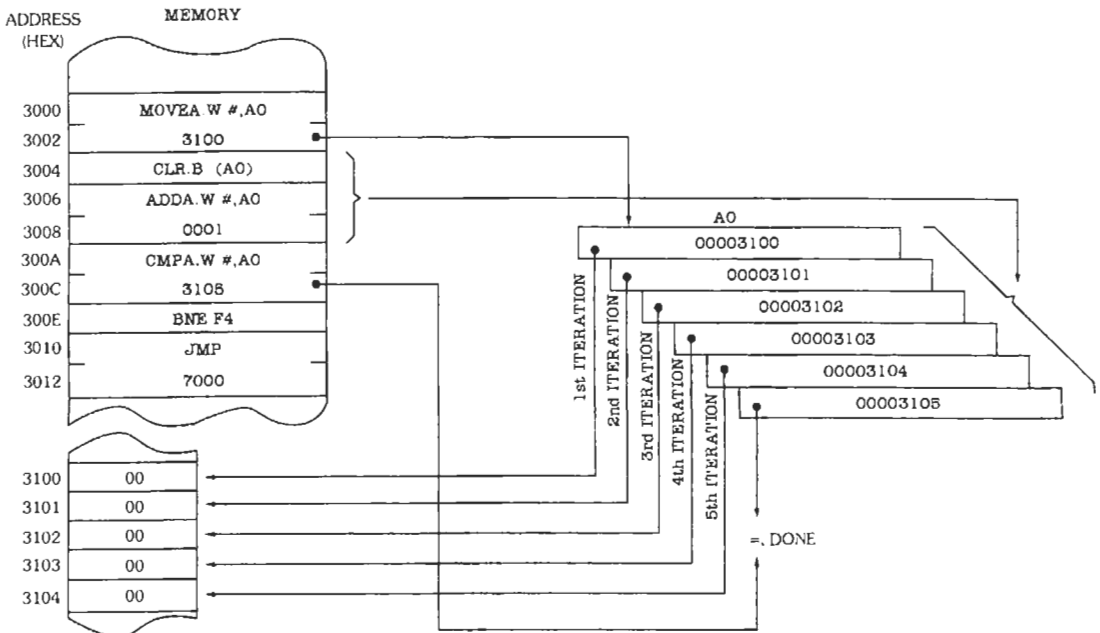


Fig. 13. Effects of indirect addressing.

TABLE 11. INITIALIZING AN ARRAY BY USING INDIRECT ADDRESSING

Machine Language		Assembly Language			
Address (hex)	Contents (hex)	Label (sym)	Opcod (mnem)	Operand (sym)	Comments
00000000			ORG	\$3000	Set components of Q to zero.
00003000	41F8 3100	INIT	LEA.L	Q, A0	A0 points to first component.
00003004	43F8 3105		LEA.L	QEND, A1	A1 points just past array end.
00003008	4210	ILUP	CLR.B	(A0)	Clear byte that A0 points to.
0000300A	D0FC 0001		ADDA.W	#1, A0	Point to next byte.
0000300E	B1C9		CMPL.A	A1, A0	Past the end?
00003010	66F6		BNE	ILUP	Nope, keep clearing.
00003012	4EF8 7008		JMP	MON	Yep, return to debug monitor.
00003016			ORG	\$3100	
00003100	???? ???? Q	Q	DS.B	5	Reserve 5 bytes for array.
00003105		QEND	EQU	Q+5	Address of byte just past Q.
00003105		MON	EQU	\$7008	Address of debug monitor.
00003105			END	INIT	

point out a few additional assembly language pseudo-operations and features that it uses:

- DS.B (define storage—byte). The operand specifies a number of memory bytes to be skipped without storing any instructions or data, thereby reserving space to be used by variables in a program.
- DC.B (define constant—byte). The specified byte values are stored into memory when the program is first loaded into memory, thereby establishing a constant value that may be accessed when the program is run.
- * (program location counter). When used in an expression, the symbol "*" denotes the current address at which assembly is taking place.

The program in Table 12 contains a main program and two subroutines. The main program initializes SP to point to an 8-word stack. A stack of four words would have been sufficient for the two nested subroutines in this program, but it is a good practice to provide "headroom" in case programming errors, modifications, or interrupts increase the space required.

The main program loads a 16-bit word into the low-order word of D0 and calls a subroutine WCNT1S that counts the number of 1 bits. WCNT1S splits the low-order word of D0 into two bytes and calls a subroutine BCNT1S to count 1's in each byte. The usefulness of subroutines is evidenced by the fact that BCNT1S can be called more than once and with a different byte to be processed each time.

Fig. 15 shows the state of the stack after each instruction that affects it. When WCNT1S returns to the

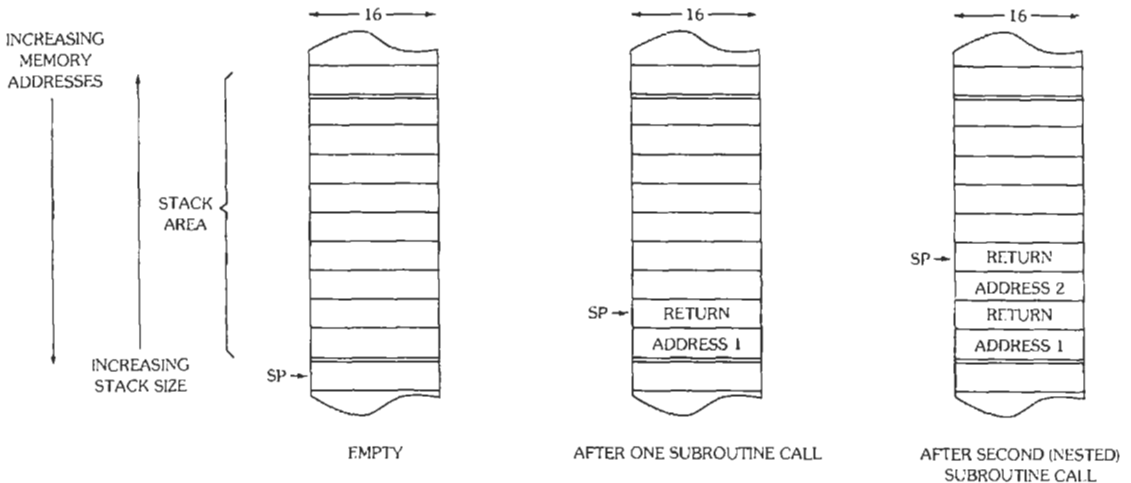


Fig. 14. A 68000 return-address stack.

TABLE 12. PROGRAM THAT USES SUBROUTINES TO COUNT THE NUMBER OF '1' BITS IN A WORD

Machine Language		Assembly Language			
Address (hex)	Contents (hex)	Label (sym)	Opcode (mnem)	Operand (sym)	Comments
00000000			ORG	\$2000	
00002000		MON	EQU	\$7008	Debug-monitor address.
00002000	4FF8 2022	MAIN	LEA.L	STKE.SP	Initialize SP (A7).
00002004	3038 2010		MOVE.W	TWORD,D0	Get test word.
00002008	4EB8 2022		JSR	WCNT1S	Count number of 1s in it.
0000200C	4EF8 7008		JMP	MON	Return to debug monitor.
00002010	5B29	TWORD	DC.W	\$5B29	Test-word to count 1s.
00002012	???? ????	STK	DS.W	8	Space for 4 return addr's.
00002022		STKE	EQU	*	Initialization addr for SP.
00002022		*			Count the number of '1' bits in a word.
00002022		*			Enter with word in D0, exit with count in D1.
00002022	31C0 203C	WCNT1S	MOVE.W	D0,CWORD	Save input word.
00002026	1038 203C		MOVE.B	CWORDH,D0	Get high-order byte.
0000202A	4EB8 203E		JSR	BCNT1S	Count 1s.
0000202E	3E01		MOVE.W	D1,D7	Save '1' count.
00002030	1038 203D		MOVE.B	CWORDL,D0	Get low-order byte.
00002034	4EB8 203E		JSR	BCNT1S	Count 1s.
00002038	D247		ADD.W	D7,D1	Add high-order count.
0000203A	4E75		RTS		Done, return.
0000203C	????	CWORD	DS.W	1	Temporary word.
0000203E		CWORDL	EQU	CWORD+1	Low-order byte address.
0000203E		CWORDH	EQU	CWORD	High-order byte address.
0000203E		*			Count number of '1' bits in a byte.
0000203E		*			Enter with byte in D0, exit with count in D1.
0000203E	4241	BCNT1S	CLR.W	D1	Initialize '1' count.
00002040	41F8 205C		LEA.L	MASKS,A0	Point to 1-bit masks
00002044	43F8 2064		LEA.L	MASKE,A1	and end of masks.
00002048	1410	BLOOP	MOVE.B	(A0),D2	Get next bit mask.
0000204A	C400		AND.B	D0,D2	Is there a '1' there?
0000204C	6704		BEQ	BN01	Skip if not.
0000204E	0641 0001		ADD.W	#1,D1	Else increment '1' count.
00002052	D0FC 0001	BN01	ADDA.W	#1,A0	Point to next mask.
00002056	B1C9		CMPA.L	A1,A0	Past last mask?
00002058	66EE		BNE	BLOOP	Continue if not.
0000205A	4E75		RTS		Return.
0000205C		*			Define 1-bit masks to test bits of byte.
0000205C	0102 0408	MASKS	DC.B	\$1,\$2,\$4,\$8,\$10,\$20,\$40,\$80	
00002060	1020 4080				
00002064		MASKE	EQU	*	Address just after table.
00002064			END	MAIN	

Symbol Table

BCNT1S	0000203E	BLOOP	00002048	BN01	00002052	CWORD	0000203C
CWORDH	0000203C	CWORDL	0000203D	MAIN	00002000	MASKE	00002064
MASKS	0000205C	MON	00007008				
STK	00002012	STKE	00002022				
TWORD	00002010	WCNT1S	00002022				

main program, the stack is again empty and D1 contains the 1st count. The main program terminates by jumping to the operating system.

The individual subroutines in Table 12 are worth discussing. On entry, WCNT1S expects the input word to be in D0, which it writes into a memory word that is read back as two individual bytes passed to BCNT1S. The memory write and reads are used because there are

no instructions to directly break out the high-order byte of a register.

When BCNT1S is entered, it expects the input byte to be in D0. The subroutine checks each bit position for a 1 and maintains a count in D1 accordingly. It uses a table of eight "mask bytes," each having a 1 in a different bit position. At each iteration of the loop, A0 contains the address of a mask byte, and a correspond-

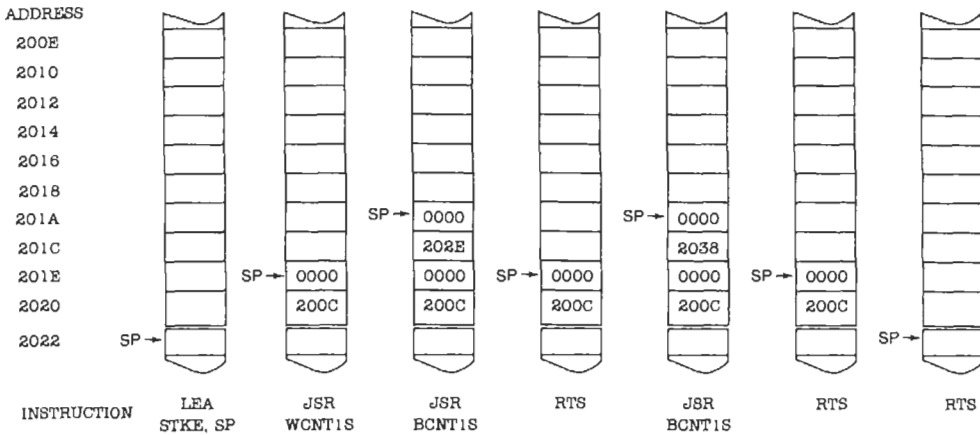


Fig. 15. Stack contents after the execution of instructions in Table 12.

ing bit of D0 is tested by the AND.B D0, D2 instruction. AND.B produces a nonzero result if and only if the tested bit of D0 is 1.

Like procedures and functions in high-level languages, subroutines are the key to structuring assembly language programs. A typical program is divided into many "modules," each of which is a subroutine with inputs, outputs, and local data.

INPUT/OUTPUT

Of the three major subsystems of a computer, input/output has experienced the biggest evolution, because of the explosion of computer applications and the hundreds of different devices that are now part of computer systems. In the 1950s, few computer architects would have predicted that some day more computers would be used in automobiles than in any other application, and that one of the most common output devices would be a fuel-injected carburetor.

Despite the proliferation of devices, fairly standard techniques are still used to connect typical devices to a computer system, as described in this section.

Buses

Fig. 2 showed the basic organization of a computer as consisting of processor, memory, and input/output (I/O). The processor communicated with the I/O subsystem by means of an I/O bus. In Fig. 16, we expand our view. Like a memory bus, the I/O bus in Fig. 16 contains data, address, and control lines. The address lines allow a program to select among different I/O devices connected to the system, while the data lines carry the actual data being transferred.

Devices and Interfaces

The I/O subsystem in Fig. 16 contains both devices and interfaces. A *peripheral device* (or *I/O device*)

performs some function for the computer. An *I/O interface* (or *device interface*) controls the operation of a peripheral device according to commands from the computer processor; it also converts computer data into whatever format is required by the device and vice versa. Also as shown in the figure, a peripheral device is often housed separately from the processor, while the interface is almost always packaged together with the processor and memory in one "CPU box."

There are many different peripheral devices that convert computer data into forms that are useful in the world outside the computer; such devices include displays, printers, plotters, digital-to-analog converters, mechanical relays, and fuel-injected carburetors. Many other devices convert data from the outside world into forms usable by the computer; examples include keyboards, text scanners, joysticks, analog-to-digital converters, mechanical switches, and crash detectors. The sole purpose of some devices is simply to store data for later retrieval; these are called *mass storage devices* and include magnetic disks and tapes.

Sometimes the dividing line between an interface and the device it controls is fuzzy. For example, Fig. 17 shows the circuitry associated with a simple mechanical keyboard. The encoder circuit converts a mechanical switch depression into a 7-bit number in the ASCII code. The bus interface can place this number on the I/O bus on demand by the processor. So it seems that the device interface consists of the Encoder and Bus Interface blocks. However, in a typical system the encoder is packaged with the keyboard; then only a small number of wires are needed between the encoder (in the keyboard package) and the bus interface (in the CPU box). Most computer designers would say that the Encoder block is part of the keyboard and the device interface consists of the Bus Interface block alone. Fortunately, the dividing line is unimportant to I/O programs that deal with the keyboard. More important is the "I/O programming model" that an I/O program sees, as discussed next.

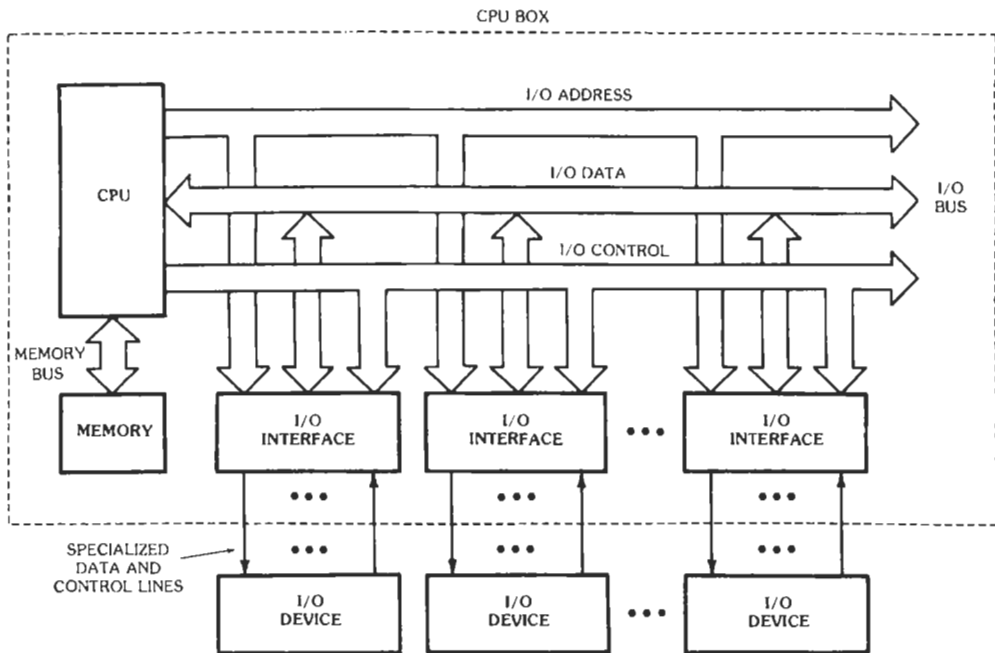


Fig. 16. Input/output (I/O) subsystem.

Ports

An *I/O port* (or *I/O register*) is a part of a device interface, a group of bits accessed by the processor during I/O operations. The "I/O programming model" of the keyboard in Fig. 17 contains one 8-bit I/O port named KBDATA, as shown in Fig. 18. The high-order bit of KBDATA is always 0. The low-order bits of KBDATA contain the output of the Encoder block in Fig. 17, that is, the 7-bit ASCII code for the key that is currently being depressed, or 0000000 if no key is depressed.

In order to read data from the keyboard in Fig. 17, a

program must execute an instruction that transfers the contents of KBDATA into one of the processor registers. Once the data is in the processor, it can be manipulated like any other data. Although the interface "writes" keyboard data into KBDATA, the port is read-only from the point of view of the processor; any attempt by the processor to write data into KBDATA has no effect. Therefore, we can call KBDATA an "input port."

Fig. 19A shows a very simple output device that interprets an 8-bit byte as two 4-bit BCD digits and displays the digits on two seven-segment displays. The I/O programming model consists of one 8-bit port, DIGOUT, shown in Fig. 19B. In order to display two

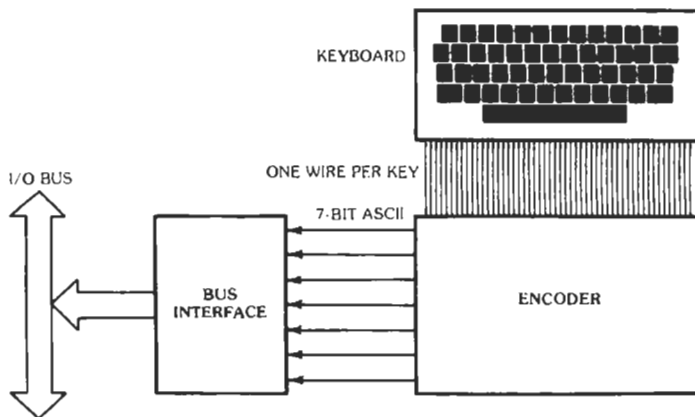


Fig. 17. Keyboard and interface.



Fig. 18. Programming model for a keyboard.

digits, the processor must transfer an 8-bit value into DIGOUT. In this case, DIGOUT is an "output port" and is write-only from the point of view of the processor; an attempt to read it produces an undefined value.

I/O Programming

So far, nothing has been said about how I/O-port data is transferred to and from processor registers. This subsection discusses two techniques that are used in different processors to perform I/O transfers.

In *isolated I/O*, the ports are accessed by special I/O instructions. For example, 80x86 processors have two instructions for transferring the contents of an I/O port to and from registers:

```
IN  rn, pn  REGISTER[ rn ] ←
                INPUTPORT[ pn ];

OUT pn, rn  OUTPUTPORT[ pn ] ←
                REGISTER[ rn ] ;
```

The IN and OUT instructions perform simple data transfers, as do load and store instructions, except that they access an array of I/O ports instead of an array of memory bytes. Since the main memory and the I/O ports are on different buses, the "address spaces" accessed by memory reference and I/O instructions are different, even though they both may happen to use 16- or 32-bit addresses.

It has been observed that I/O buses are very similar to memory buses, and that I/O instructions are similar to load and store instructions on memory. *Memory-mapped I/O* takes advantage of the similarity by eliminating the I/O bus and I/O instructions.

Fig. 20 shows the hardware organization of a computer with memory-mapped I/O. Both the main memory and all I/O ports communicate with the processor by way of a shared Memory and I/O Bus. Each I/O port has an address in the main memory address space of the processor. An input port responds to any instruction that reads at its address; an output port responds to any instruction that writes at its address. Typically, the system designer reserves a portion of the total address space for I/O ports, for example, the top 4 kilobytes. However, theoretically a hardware designer can locate an I/O port at any address, as long as there is no memory at that address also.

Memory-mapped I/O is a necessity in processors that have no special I/O instructions. The PDP-11 was the first minicomputer to require memory-mapped I/O;

the 68000 and all RISC processors also require it. Memory-mapped I/O has a number of advantages:

- No opcodes or processor circuits are used up for I/O instructions.
- All memory reference instructions, not just loads and stores, may be used to manipulate I/O ports.
- The number of available I/O port addresses is virtually unlimited.
- The hardware bus structure is simplified.

There are also disadvantages:

- Part of the memory address space is used up. Interfaces may need more circuitry to recognize longer addresses.
- Memory reference instructions may be longer or slower than optimized I/O instructions.

A simple Motorola 68000 program using memory-mapped I/O is shown in Chart 3.

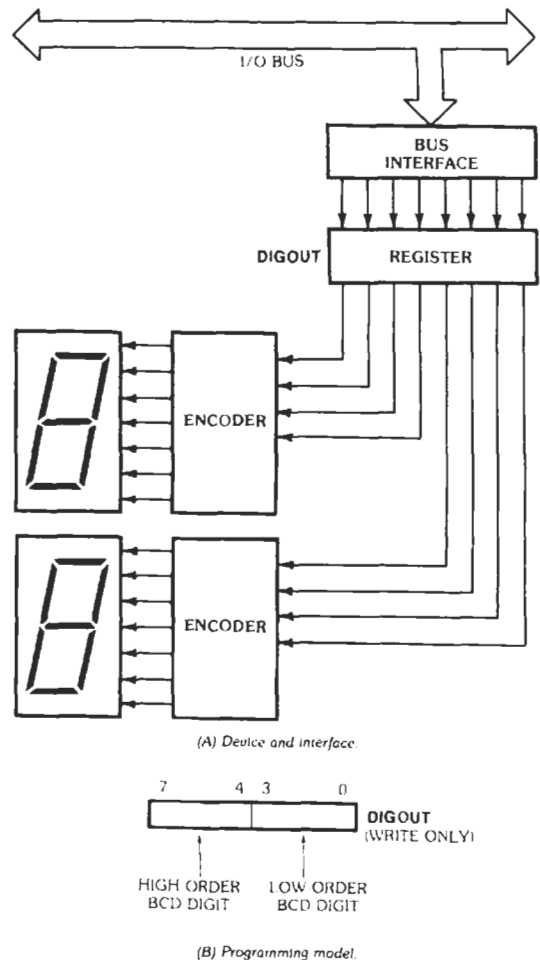


Fig. 19. Seven-segment display.

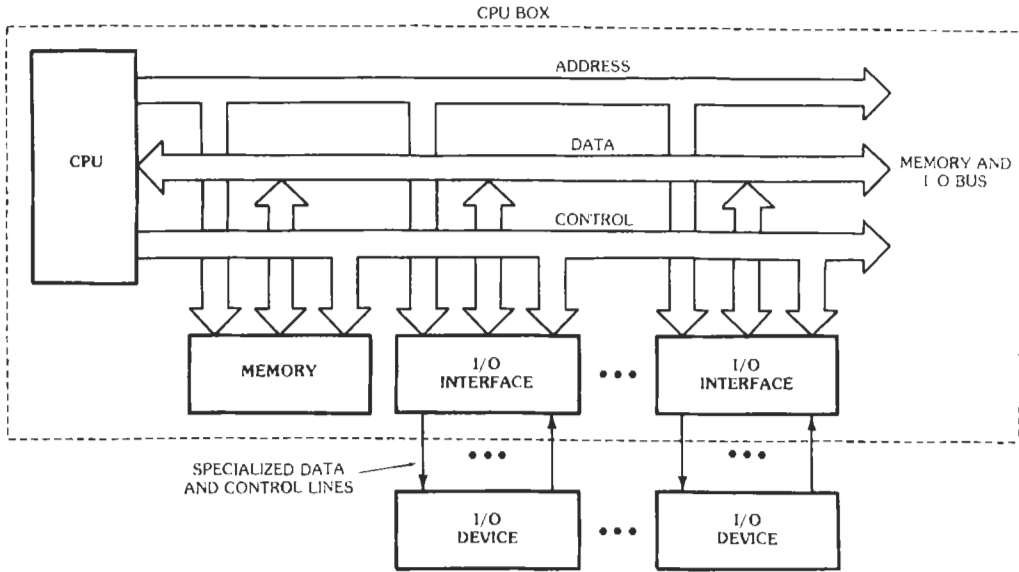


Fig. 20. Memory-mapped I/O structure.

CHART 3. KEYBOARD INPUT AND DISPLAY OUTPUT FOR THE 68000

```

*
*      Read two decimal digits from the keyboard and display them on the
*      seven-segment display. Ignore illegal characters.
KBDATA EQU      $F004      Keyboard input port address.
DIGOUT EQU      $F008      Seven-segment display port address.
*
*      First, a subroutine to read one decimal digit, convert it to 4-bit bcd,
*      and return the result in register D0.
RDDIG  MOVE.B    KBDATA,D0      Read current character.
      BEQ      RDDIG           Wait for a key to be pressed.
WAITUP MOVE.B    KBDATA,D1      Now the character is in D0,
      BNE      WAITUP          wait for the key to be released.
      CMP.B   #$30,D0          Was it a valid decimal digit?
      BLT      RDDIG           Not if it is less than ASCII '0' ...
      CMP.B   #$39,D0          ... or greater than ASCII '9'.
      BGT      RDDIG           Convert ASCII '0'-'9' to 0-9.
      ADD.B   D0,#-$30         Done, return.
      RTS
*
*      Now, the main program.
DIGDSP JSR      RDDIG           Read high-order decimal digit into D0.
      ASL.B   #4,D0           Shift left 4 bits
      MOVE.B  D0,D2           and save in D2.
      JSR    RDDIG           Read low-order decimal digit.
      OR.B   D2,D0           Merge (logical OR) with high-order digit.
      MOVE.B  D0,DIGOUT      Send to seven-segment display.
      JMP    DIGDSP         Do another pair of digits.
      END    DIGDSP
    
```

43 Logic Design

Edward J. McCluskey

Introduction	43-2
Postulates	43-2
Analysis	43-3
Synthesis	43-4
Canonical Expressions	
Networks	
Theorems	43-8
Single-Variable Theorems	
Two- and Three-Variable Theorems	
Example	
<i>n</i> -Variable Theorems	
General Gate Networks	43-10
The Map Method	43-12
Maps for Two, Three, and Four Variables	
Prime Implicants	
Formation of Minimal Sums	
Sequential Circuits	43-16

INTRODUCTION

The distinguishing feature of the circuits to be discussed in this chapter is the use of two-valued, or binary, signals. There will be some deviation of the signals from their nominal values, but within certain limits this variation will not affect the performance of the circuit. If the variations exceed these limits, the circuit will not behave properly, and steps must be taken to confine the signals to the proper ranges. When the statement is made that the signals are two-valued, what is really meant is that the value of each signal is within one of two (nonoverlapping) continuous ranges. Since the operation of the circuit does not depend on exactly which value within a given range the signal takes on, a particular value is chosen to represent the range, and the signal is said to be equal to this value. The exact numerical value of the signal is not important. It is possible to have two circuits perform the same function and have completely different values for their signals. In order to avoid any possible confusion that might arise because of this situation and to simplify the design procedures, it is customary to carry out the logic design without specifying the actual values of the signals. Once the logic design has been completed, the actual values must be assigned to the signals in the course of designing the detailed electrical circuit. For the purposes of the logic design, arbitrary symbols are chosen to represent the two values to which the signals are to be restricted. An algebra* using these symbols is then developed as the basis for formal design techniques. The development of such an algebra will be described next.

POSTULATES

The two symbols most commonly chosen to represent the two logic values taken on by binary signals are 0 and 1. It should be emphasized that these logic symbols have no numerical significance. For an electronic circuit that has its signals equal to either 0 or 5 volts, logical 1 usually represents 5 volts and logical 0 represents 0 volts. This choice—in which 1 corresponds to the more positive physical value—is called *positive logic*. *Negative logic* uses logic 1 for the more negative physical value; it is much less common than positive logic, but is natural when the active signal value corresponds to the lower physical value (as in PMOS, for example). We will use positive logic here.

Other symbols such as T and F, H and L, or + and - are also used. Logicians like T and F. IC data books favor H and L.

Switching variables are used in logic networks to represent the signals present at gate inputs and outputs.

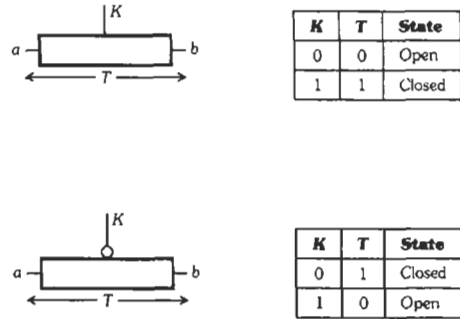


Fig. 1. The transmission, *T*, of a switch.

Switching variables can also represent the state of a switch. Fig. 1 shows this. The switch operation is defined as the two switch terminals (*a* and *b*) being connected together if and only if the *control variable*, *K*, equals 1. When *K* = 0, there is an open circuit between the two switch terminals. There is a variable, *T*, associated with the switch that equals 1 when the terminals are connected together and that equals 0 when there is an open circuit between the terminals. The variable, *T*, is called the *transmission* of the switch. (It is also possible to associate with the switch a variable that equals 1 only when the switch is open. Such a variable is called the *switch hindrance*. This was used in the very early papers on switching theory in connection with contact networks. The transmission concept is the standard usage at present.)

The first postulate of the switching algebra can now be presented. This is merely a formal statement of the fact that switching variables are always equal to either 0 or 1. (In the statements of postulates and theorems that follow, the symbols *X*, *Y*, *Z*, *X*₁, *X*₂ . . . , *X*_{*n*} will be used to represent switching variables.)

$$(P1) X = 0 \text{ if } X \neq 1 \quad (P1') X = 1 \text{ if } X \neq 0$$

To implement general switching networks, it is necessary to be able to obtain, for any signal representing a switching variable, a signal that has the opposite value. A circuit for realizing this function is called an *inverter*. Fig. 2 shows two of the most common inverter symbols.

If *X* represents a switching variable, the symbol *X'* is used to represent the signal having the opposite value. This notation is specified formally in the second switching-algebra postulate:

$$(P2) \text{ If } X = 0, \text{ then } X' = 1$$

$$(P2') \text{ if } X = 1, \text{ then } X' = 0$$



Fig. 2. Inverter symbols.

*This algebra will be called *switching algebra*. It is identical with a Boolean algebra and was originally applied to switching circuits (reference 1) by reinterpreting Boolean algebra in terms of switching circuits rather than by developing a switching algebra directly, as will be done here.

terminals. Only two-terminal networks will be considered for the present. In a gate network, an output function specifies those input conditions for which the voltage of an output node will be at the high level.

Fig. 5 shows a very simple analysis example. In an acyclic network, there will be at least one gate, such as the OR gate in Fig. 5, that has only network inputs connected to it. The output of such gates can be written down directly. There will then be other gates that have either input variables or expressions on all their inputs; the AND gate in Fig. 5 is such a gate. It is then possible to write expressions for the output of all these gates. By repeating this procedure, an expression will eventually be obtained for the network output.

A procedure for using these output functions to determine the circuit performance is to substitute a set of values for the input variables and then use the postulates to simplify the resulting expressions. A value for the output function will be obtained that specifies the output for the particular input combination chosen. For example, if the input combination $X_1 = 1, X_2 = 1, X_3 = 0, X_4 = 1$, is chosen, the output function of Fig. 5 becomes

$$f = X_1 X_2 (X_3 + X_4) \\ = 1 \cdot 1(0 + 1) = 1(1) = 1$$

By carrying out this procedure for all possible input combinations, it is possible to form a table that lists the output for each input combination. Such a table describes completely the circuit performance and is called a *table of combinations*.^{*} Table I is the table of combinations for the circuit of Fig. 5.

SYNTHESIS

To design a combinational circuit, the procedure just described must be reversed. The designer writes a table of combinations that specifies the desired circuit performance. From this table, an algebraic function is written from which a circuit can be derived. A concise means of specifying the table of combinations, called a *decimal specification*, is to list the numbers of the rows for which the output is to equal 1. For Table I, this specification is:

$$f(X_1, X_2, X_3, X_4) = \Sigma(13, 14, 15)$$

^{*}In logic, the table is called a *truth table*, and some writers use this term when discussing switching circuits.

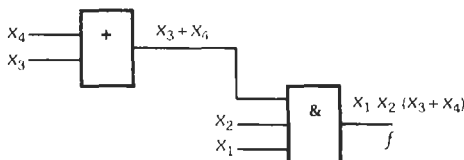


Fig. 5. Simple example of analysis of a gate network.

TABLE I. TABLE OF COMBINATIONS FOR $f = X_1 X_2 (X_3 + X_4)$

	X_1	X_2	X_3	X_4	f
0	0	0	0	0	0
1	0	0	0	1	0
2	0	0	1	0	0
3	0	0	1	1	0
4	0	1	0	0	0
5	0	1	0	1	0
6	0	1	1	0	0
7	0	1	1	1	0
8	1	0	0	0	0
9	1	0	0	1	0
10	1	0	1	0	0
11	1	0	1	1	0
12	1	1	0	0	0
13	1	1	0	1	1
14	1	1	1	0	1
15	1	1	1	1	1

where the Σ signifies that the rows for which the function equals 1 are being listed. It is also possible to list the rows for which the function equals 0, such a list being preceded by the symbol Π to indicate that it is the zero rows which are listed. This specification for the preceding table is

$$f(x_1, x_2, x_3, x_4) = \\ \Pi(0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12).$$

In order to avoid any ambiguity in these specifications, it is necessary to adopt some rule for numbering the rows of the table of combinations. The usual procedure is to regard each row of the table as a binary number and then use the decimal equivalent of this binary number as the row number. The output-column entries are not included in forming the binary row numbers. There is nothing special about decimal numbers other than the fact that they are the most familiar; any other number base, such as octal, could be used. The reason for using a number system other than binary is simply that binary numbers take too much space to write down.

It has been pointed out that the table of combinations is a complete specification for a combinational circuit. The first step in designing a circuit is to formulate such a table. There are no general formal techniques for doing this. When a sequential circuit is being designed, it is customary to reduce the sequential-design problem to (several) combinational problems, and formal techniques exist for doing this. However, when a combinational circuit is being designed, no formal techniques are available, and it is necessary to rely on common sense. This is not too surprising, since any formal technique must start with a formal statement of the problem, and this is precisely what the table of combinations is. As an example of how this is done, the table of combinations for a circuit to check binary-coded decimal digits for the 8, 4, 2, 1 code is shown in Table

TABLE 2. TABLE OF COMBINATIONS FOR CIRCUIT TO CHECK BINARY-CODED-DECIMAL DIGITS

	b_8	b_4	b_2	b_1	f
0	0	0	0	0	0
1	0	0	0	1	0
2	0	0	1	0	0
3	0	0	1	1	0
4	0	1	0	0	0
5	0	1	0	1	0
6	0	1	1	0	0
7	0	1	1	1	0
8	1	0	0	0	0
9	1	0	0	1	0
10	1	0	1	0	1
11	1	0	1	1	1
12	1	1	0	0	1
13	1	1	0	1	1
14	1	1	1	0	1
15	1	1	1	1	1

} Invalid code words

2. This circuit is to deliver an output whenever a digit having an invalid combination of bits is received.

In forming a table of combinations, there very often are rows for which it is unimportant whether the function equals 0 or 1. The usual reason for this situation is that the combination of inputs corresponding to these rows can never occur (when the circuit is functioning properly). As an example of this, consider a circuit to translate from the 8, 4, 2, 1 BCD code to a Gray (cyclic binary) code. When the circuit is working correctly, the input combinations represented by rows 10 through 15 of the table of combinations cannot occur. Therefore, the output need not be specified for these rows. The symbol d will be used to indicate the output condition for such a situation.* The output conditions so denoted are called *don't-care conditions* (see Table 3). It is possible to include the d rows in the decimal specification of a function by listing them after the symbol d . Thus, the decimal specification for g_1 of Table 3 would be

$$\begin{aligned}
 g_1(b_8, b_4, b_2, b_1) &= \Sigma(1, 2, 5, 6, 9) + \\
 &\quad d(10, 11, 12, 13, 14, 15) \\
 &= \Pi(0, 3, 4, 7, 8) + \\
 &\quad d(10, 11, 12, 13, 14, 15)
 \end{aligned}$$

Canonical Expressions

After the table of combinations has been formed, the next step in designing a circuit is to write an algebraic expression for the output function. The simplest output functions to write are those that equal 1 for only one row of the table of combinations or those that equal 0 for only one row. Two functions correspond to each row: one that equals 1 only for the row and one that

equals 0 only for the row (see Table 4). These functions are called *fundamental products* or *minterms* and *fundamental sums* or *maxterms*, respectively. Each fundamental product or sum contains all the input variables. The rule for forming the fundamental product for a given row is to prime any variables that equal 0 for the row and leave unprimed variables that equal 1 for the row. The fundamental product equals the product of the literals so formed. The fundamental sum is formed by a completely reverse, or dual, procedure. Each variable that equals 0 for the row is left unprimed, and each variable that equals 1 for the row is primed. The fundamental sum is the sum of the literals obtained by this process. The algebraic expression for any table for which the output is equal to 1 (or 0) for only one row can be written down directly by choosing the proper fundamental product (or sum). For example, the output function specified by $f(x_1, x_2, x_3) = \Sigma(6)$ is written algebraically as $f = x_1x_2x_3'$, and the output function $f(x_1, x_2, x_3) = \Pi(6)$ is written as $f = x_1' + x_2' + x_3$. The fundamental product corresponding to row i of the table of combinations will be denoted by p_i , and the fundamental sum corresponding to row i will be denoted by s_i .

The algebraic expression that equals 1 (or 0) for more than one row of the table of combinations can be written directly as a sum of fundamental products or as a product of fundamental sums. A function f that equals 1 for two rows, i and j , of the table of combinations can be expressed as a sum of the two fundamental products p_i and p_j : $f = p_i + p_j$. When the inputs correspond to row i , $p_i = 1$ and $p_j = 0$, so that $f = 1 + 0 = 1$. When the inputs correspond to row j , $p_i = 0$, $p_j = 1$, and $f = 0 + 1 = 1$. When the inputs correspond to any other row, $p_i = 0$, $p_j = 0$, $f = 0 + 0 = 0$. This shows that the function $f = p_i + p_j$ does equal 1 only for row i and j . This argument can be extended to functions which equal 1 for any number of input combinations—they can be represented algebraically as a sum of the corresponding fundamental products (see Table 5). An algebraic expression that is a sum of fundamental products is called a *canonical sum*. An arbitrary function can also be expressed as a product of fundamental sums. This form is called *canonical product*. The canonical product for a function that is equal to 0 only for rows i and j of the table of combinations is given by $f = s_i \cdot s_j$. For row i , $s_i = 0$ so that $f = 0$. For row j , $s_j = 0$ so that $f = 0$, and for any other row $s_i = s_j = 1$ so that $f = 1$. In general, the canonical product is equal to the product of all fundamental sums that correspond to input conditions for which the function is to equal 0.

It is possible to write a general expression for the canonical sum by making use of the following theorems*:

$$(T2') x \cdot 0 = 0$$

$$(T1) x + 0 = x \quad (T1') x \cdot 1 = x$$

*In the literature, the symbol ϕ is also used.

*The theorems are presented formally in the next section.

TABLE 3. TABLE OF COMBINATIONS FOR CIRCUIT TO TRANSLATE FROM BCD 8, 4, 2, 1 CODE TO GRAY CODE.

	BCD-Code Inputs				Gray-Code Outputs			
	b_8	b_4	b_2	b_1	g_4	g_3	g_2	g_1
0	0	0	0	0	0	0	0	0
1	0	0	0	1	0	0	0	1
2	0	0	1	0	0	0	1	1
3	0	0	1	1	0	0	1	0
4	0	1	0	0	0	1	1	0
5	0	1	0	1	0	1	1	1
6	0	1	1	0	0	1	0	1
7	0	1	1	1	0	1	0	0
8	1	0	0	0	1	1	0	0
9	1	0	0	1	1	1	0	1
10	1	0	1	0	d	d	d	d
11	1	0	1	1	d	d	d	d
12	1	1	0	0	d	d	d	d
13	1	1	0	1	d	d	d	d
14	1	1	1	0	d	d	d	d
15	1	1	1	1	d	d	d	d

If the value of the function $f(x_1, x_2, \dots, x_n)$ for the i -th row of the table of combinations is f_i ($f_i = 0$ or 1), then the canonical sum is given by

$$f(x_1, x_2, \dots, x_n) = f_0p_0 + f_1p_1 + \dots + f_{(2^n-1)}p_{(2^n-1)}$$

$$= \sum_{i=0}^{2^n-1} f_i p_i$$

For the function $f(x_1, x_2) = \Sigma(0, 2)$ the values of the f_i are $f_0 = f_2 = 1, f_1 = f_3 = 0$ so that

$$f(x_1, x_2) = 1 \cdot p_0 + 0 \cdot p_1 + 1 \cdot p_2 + 0 \cdot p_3$$

$$= p_0 + 0 + p_2 + 0$$

$$= p_0 + p_2$$

$$= x'_1 x'_2 + x_1 x_2'$$

In a similar fashion, a general expression for the canonical product can be obtained by using the theorems

$$(T2) \quad x + 1 = 1$$

$$(T1) \quad x + 0 = x \quad (T1') \quad x \cdot 1 = x$$

The resulting expression is

$$f(x_1, x_2, \dots, x_n) = (f_0 + s_0)(f_1 + s_1) \dots (f_{(2^n-1)} + s_{(2^n-1)})$$

$$= \prod_{i=0}^{2^n-1} (f_i + s_i)$$

For $f(x_1, x_2) = \Sigma(0, 2) = \Pi(1, 3)$ this becomes

$$f(x_1, x_2) = (1 + s_0)(0 + s_1)(1 + s_2)(0 + p_3)$$

$$= 1 \cdot s_1 \cdot 1 \cdot s_3$$

$$= s_1 \cdot s_3$$

$$= (x_1 x'_2)(x'_1 + x_2')$$

Networks

A technique for obtaining an algebraic expression from a table of combinations has just been described. A circuit can be drawn directly from this expression by reversing the analysis procedures. The circuit for a single fundamental product is just an AND gate with appropriate inputs or a series connection of switches. For a canonical sum involving more than one fundamental product, the circuit consists of a number of

TABLE 4. TABLE OF COMBINATIONS SHOWING FUNDAMENTAL PRODUCTS AND FUNDAMENTAL SUMS

	x_1	x_2	x_3	Fundamental Product	Fundamental Sum
0	0	0	0	$x'_1 x'_2 x'_3$	$x_1 + x_2 + x_3$
1	0	0	1	$x'_1 x'_2 x_3$	$x_1 + x_2 + x'_3$
2	0	1	0	$x'_1 x_2 x'_3$	$x_1 + x'_2 + x_3$
3	0	1	1	$x'_1 x_2 x_3$	$x_1 + x'_2 + x'_3$
4	1	0	0	$x_1 x'_2 x'_3$	$x'_1 + x_2 + x_3$
5	1	0	1	$x_1 x'_2 x_3$	$x'_1 + x_2 + x'_3$
6	1	1	0	$x_1 x_2 x'_3$	$x'_1 + x'_2 + x_3$
7	1	1	1	$x_1 x_2 x_3$	$x'_1 + x'_2 + x'_3$

TABLE 5. $f(x_1, x_2, x_3) = \Sigma(1, 2, 3, 4)$
(A) Table of Combinations

	x_1	x_2	x_3	f
0	0	0	0	0
1	0	0	1	1
2	0	1	0	1
3	0	1	1	1
4	1	0	0	1
5	1	0	1	0
6	1	1	0	0
7	1	1	1	0

(B) Canonical Sum

$$f = x_1'x_2'x_3 + x_1'x_2x_3 + x_1x_2x_3 + x_1x_2'x_3$$

(C) Canonical Product

$$f = (x_1 + x_2 + x_3)(x_1' + x_2 + x_3') \\ (x_1' + x_2' + x_3)(x_1' + x_2' + x_3')$$

parallel subnetworks, each subnetwork corresponding to one fundamental product, or a number of AND gates with their outputs connected as inputs to an OR gate. This is shown in Figs. 6A and 7A. Similarly, the switch network corresponding to a canonical product consists of a number of subnetworks in series, each subnetwork corresponding to one fundamental sum and consisting of switches in parallel (Fig. 6B). The gate network corresponding to a canonical product consists of a number of OR gates with their outputs connected as the inputs of an AND gate (Fig. 7B).

It follows that, for any arbitrary table of combinations, a network whose performance corresponds to the table of combinations can be constructed of*:

*This assumes there is no limit on the number of elements used.

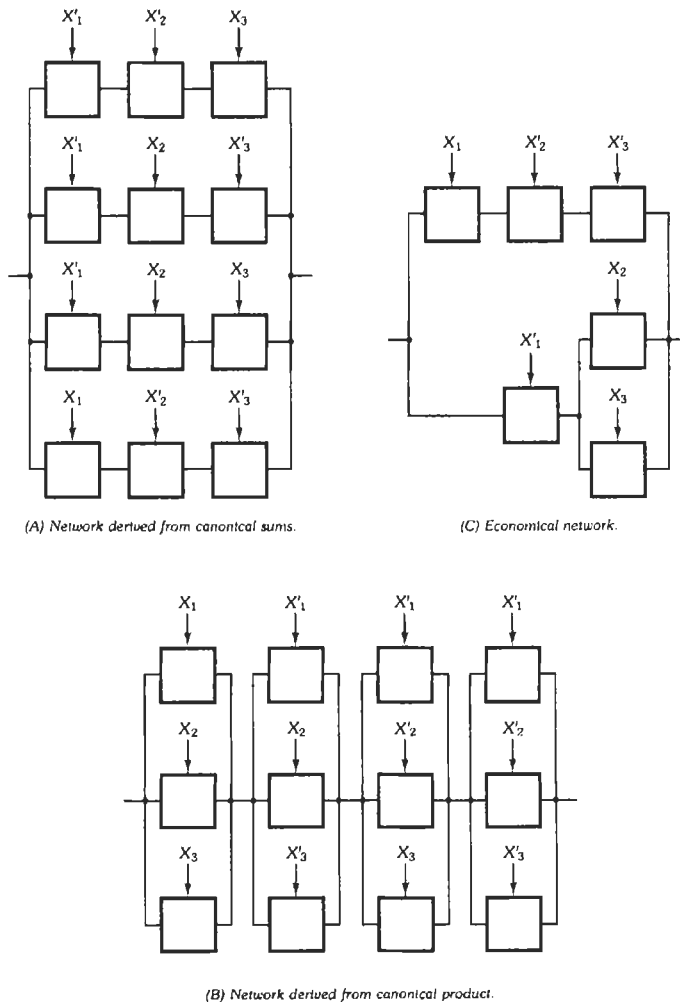


Fig. 6. Switch networks for $\Sigma(1, 2, 3, 4)$.

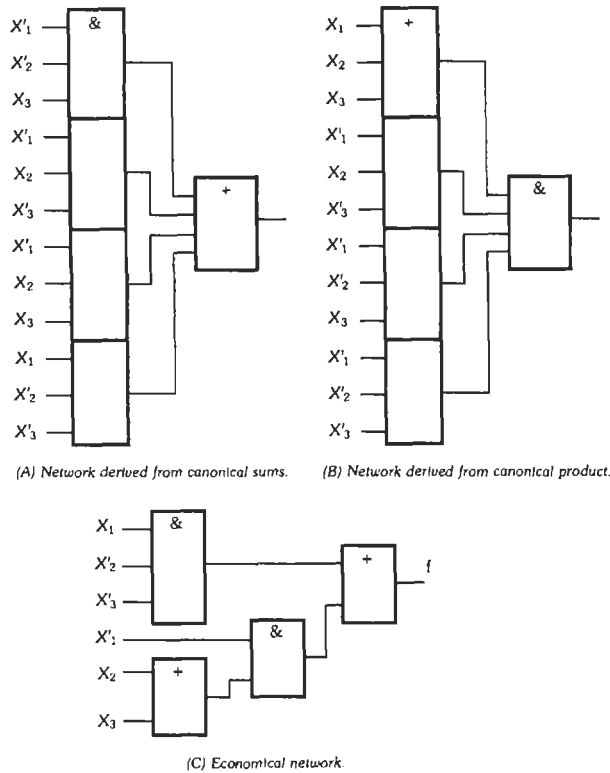


Fig. 7. Gate networks for $\Sigma (1, 2, 3, 4)$.

1. AND gates, OR gates, and inverters, or
2. Switches

In a certain sense, this completes the design procedure. A method has been presented for going from an arbitrary table of combinations to a circuit diagram. However, the canonical circuits so designed are usually very uneconomical and therefore unsatisfactory. An example of this can be seen by comparing the circuits of Figs. 6A, 6B, 7A, and 7B with those of Figs. 6C and 7C. In order to design satisfactory circuits, it is necessary to have procedures for simplifying them so that they correspond to simpler circuits.

THEOREMS

To design an efficient switching circuit, one formulates the desired circuit performance as an algebraic expression and then manipulates the expression into a form from which a desirable circuit can be arrived at directly. The manipulations are carried out by means of the theorems that will be presented in this section.

Single-Variable Theorems

The switching-algebra theorems that involve only a single variable are shown in Chart 1. Note the three theorems, T2, T3, and T3', that are false for ordinary

algebra. Theorems T3 and T3' can be extended to $X + X + X + \dots + X = nX = X$ and $X \cdot X \cdot X \cdot \dots \cdot X = X^n = X$.

Two- and Three-Variable Theorems

Chart 2 lists the switching-algebra theorems that involve two or three variables. Theorems T7 and T7' are used to eliminate terms from algebraic expressions so that the corresponding elements can be removed from the corresponding networks. Theorem T10' is not true for ordinary algebra even though its dual, T10, is.

In reducing algebraic expressions, Theorems T11 and T11' are very important and are used frequently, as is illustrated in the example that follows. This example shows how the theorems are used to manipulate a given algebraic expression into some other form. Very frequently, the form desired is one that has as few literals

CHART 1. SWITCHING-ALGEBRA THEOREMS INVOLVING ONE VARIABLE

(T1) $X + 0 = X$	(T1') $X \cdot 1 = X$	(Identities)
(T2) $X + 1 = 1$	(T2') $X \cdot 0 = 0$	(Null elements)
(T3) $X + X = X$	(T3') $X \cdot X = X$	(Idempotency)
(T4) $(X')' = X$		(Involution)
(T5) $X + X' = 1$	(T5') $X \cdot X' = 0$	(Complements)

CHART 2. SWITCHING-ALGEBRA THEOREMS INVOLVING TWO OR THREE VARIABLES

(T6)	$X + Y = Y + X$	(T6')	$XY = YX$	(Commutative)
(T7)	$X + XY = X$	(T7')	$X(X + Y) = X$	(Absorption)
(T8)	$(X + Y')Y = XY$	(T8')	$XY' + Y = X + Y$	
(T9)	$(X + Y) + Z = X + (Y + Z) = X + Y + Z$			
(T9')	$(XY)Z = X(YZ) = XYZ$			(Associative)
(T10)	$XY + XZ = X(Y + Z)$			
(T10')	$(X + Y)(X + Z) = X + YZ$			(Distributive)
(T11)	$(X + Y)(X' + Z)(Y + Z) = (X + Y)(X' + Z)$			
(T11')	$XY + X'Z + YZ = XY + X'Z$			(Consensus)
(T12)	$(X + Y)(X' + Z) = XZ + X'Y$			

occurring as possible. The number of literal occurrences corresponds roughly to the number of gate inputs in a gate network.

Theorem T12 has no dual (T12'), for the dual would be identical with the original theorem.

Example

By use of the theorems, the expression

$$(c' + abd + b'd + a'b)(c + ab + bd)$$

is to be shown equal to

$$b(a + c)(a' + c') + d(b + c)$$

$$(c' + abd + b'd + a'b)(c + ab + bd)$$

- (T12) $c'(ab + bd) + c(abd + b'd + a'b)$
- (T10) $abc' + bc'd + abcd + b'cd + a'bc$
- (T6) $abc' + abcd + bc'd + a'bc + b'cd$
- (T10) $ab(c' + cd) + bc'd + a'bc + b'cd$
- (T8')
- (T11')
- (T10) $abc' + abd + bc'd + a'bc + a'bd + b'cd$
- (T6) $(abd + a'bd) + abc' + bc'd + a'bc + b'cd$
- (T10) $bd(a + a') + abc' + bc'd + a'bc + b'cd$
- (T5) $bd(1) + abc' + bc'd + a'bc + b'cd$
- (T1')
- (T6) $bd + bc'd + abc' + a'bc + b'cd$

- (T7) $bd + abc' + a'bc + b'cd$
- (T6),(T10) $d(b + b'c) + abc' + a'bc$
- (T8')
- (T10) $d(b + c) + b(ac' + a'c)$
- (T12) $d(b + c) + b(a + c)(a' + c')$

n-Variable Theorems

The switching-variable theorems that involve an arbitrary number of variables are shown in Chart 3. Three of these theorems (T13, T13', and T14) cannot be proved by perfect induction. For these theorems, the proofs require the use of finite induction.* Theorems T13 and T13' are proved by first letting $n = 2$ and using perfect induction to prove their validity for this special case. It is then assumed that the theorems are true for $n = k$, and this is shown to imply that they must then be true for $n = k + 1$. This completes the proof, the details of which are given in the reference. Theorem T14 is proved by using Theorems T13 and T13' along with the fact that every function can be split into the sum of several functions or the product of several functions. By successively splitting the function into subfunctions and using T13 and T13', it is possible to prove T14.

Theorem T14, which is a generalization of T13, forms the basis of a method for constructing complementary networks. Two networks having outputs T_1 and T_2 are said to be *complementary* if $T_1 = T_2'$. The complementary network for any given network can be designed by writing the output, T_1 , for the first network, then forming T_1' by means of T14, and then designing a network having output T_1' . For example, if $T_1 = (x + y)[w(y' + z) + xy]$, then $T_1' = x'y' + (w' + yz')(x' + y')$.

It was pointed out in connection with Figs. 6 and 7 that the canonical networks are generally uneconomical. By manipulating the canonical sum or product with the aid of the theorems just presented, it is usually possible to obtain algebraic expressions that correspond to more economical networks than the canonical networks. The following example shows how this is done

*See reference 2.

CHART 3. SWITCHING-VARIABLE THEOREMS INVOLVING n VARIABLES

	(DeMorgan's theorems)
(T13)	$(X_1 + X_2 + \dots + X_n)' = X_1' X_2' \dots X_n'$
(T13')	$(X_1 X_2 \dots X_n)' = X_1' + X_2' + \dots + X_n'$
	(Generalized DeMorgan's theorem)
(T14)	$f(X_1, X_2, \dots, X_n, +, \cdot)' = f(X_1', X_2', \dots, X_n', \cdot, +)$
	(Expansion theorem)
(T15)	$f(X_1, X_2, \dots, X_n) = X_1 f(1, X_2, \dots, X_n) + X_1' f(0, X_2, \dots, X_n)$
(T15')	$f(X_1, X_2, \dots, X_n) = [X_1 + f(0, X_2, \dots, X_n)][X_1' + f(1, X_2, \dots, X_n)]$

for the networks of Figs. 6 and 7. The final expressions correspond to the networks of Figs. 6C and 7C.

$$\begin{aligned}
 f &= X_1' X_2' X_3 + X_1' X_2 X_3' && + X_1' X_2 X_3 + X_1 X_2' X_3' \\
 f &= X_1' X_2' X_3 + X_1' X_2 X_3' && + X_1' X_2 X_3' + X_1 X_2' X_3' \\
 f &= X_1' X_2' X_3 + X_1' X_2 X_3' && + X_1' X_2 X_3 + X_1' X_2 X_3' + X_1 X_2' X_3' \\
 f &= X_1' X_3 (X_2' + X_2) && + X_1' X_2 (X_3 + X_3') + X_1 X_2' X_3' \\
 f &= X_1' X_3 (1) && + X_1' X_2 (1) + X_1 X_2' X_3' \\
 f &= X_1' X_3 && + X_1' X_2 + X_1 X_2' X_3' \\
 f &= X_1' (X_3 + X_2) && + X_1 X_2' X_3'
 \end{aligned}$$

Many of the theorems of ordinary algebra are also valid for switching algebra. One that is not is the cancellation law. In ordinary algebra, it follows that $X = Z$ if $X + Y = Y + Z$. In switching algebra, this is not true. For example, it is generally true that $X + XY = X + 0$, but it is not necessarily true that $XY = 0$. This can be easily verified by writing out the tables of combinations for $f_1(X, Y) = X + XY$, $f_2(X, Y) = X + 0$, and $f_3(X, Y) = XY$. Similar remarks apply to the situation in which $XY = XZ$ does not imply that $Y = Z$.

GENERAL GATE NETWORKS

The previous discussion of gate networks in this chapter has been concerned solely with networks constructed of AND gates and OR gates. This can be considered only an introduction to the topic of gate networks, for other types of gates are equally important. In this section, networks with other types of gates will be considered.

Any arbitrary switching function can be realized by a network of AND gates, OR gates, and inverters. A natural question to ask in this connection is whether all three types of elements are necessary. Inverters are required if the inputs to the network consist of signals representing the input variables but not of signals representing the complements of the input variables. The situation when signals representing the complements are available is called *double-rail logic*, and when the complements are not available, the term *single-rail* is used. Both techniques are employed, but for the

purposes of the present discussion it will be assumed that complements are not directly available (single-rail logic). The function $f(x) = x'$ cannot be realized by a network of AND gates and OR gates only.

Arbitrary functions can require inverters, but maybe AND gates and inverters can realize any function. That OR gates are not necessary is easily demonstrated, for it is possible to construct a network having the function of an OR gate and using only AND gates and inverters. This is done by making use of DeMorgan's theorem— $X + Y = (X'Y)'$ —as illustrated in Fig. 8. Thus, any network consisting of AND gates, OR gates, and inverters can be changed into a network containing only AND gates and inverters by using the replacement shown in Fig. 8 to remove the OR gates. By duality, a similar technique can be used to remove the AND gates instead.

Since it is not possible to use only inverters to realize arbitrary functions, a minimum set of elements has now been determined. Because it is possible to construct a network containing only AND gates and inverters for any arbitrary function, the AND gate and inverter are said to form a *complete gate set*. Similarly, the OR gate and inverter form a complete gate set.

The two operations of the AND function and the complement can be combined in a single gate, the NAND gate, shown in Fig. 9. This is a very common integrated-circuit gate. It comprises a complete gate set in one gate since: (1) an inverter is obtained if all inputs are connected to the same source, as in Fig. 9C; (2) an AND gate is formed by combining two NAND gates, as in Fig. 9D.

Two symbols for the NAND gate are shown in Fig. 9A. This is because the basic NAND-gate function $(XY)'$ can also be written as $X' + Y'$ by using DeMorgan's Theorem (T13). Use of the two symbols facilitates analysis and synthesis using these gates. The small circles ("bubbles") in Fig. 9 indicate inversion, and from a logical standpoint each bubble can be replaced by an inverter.

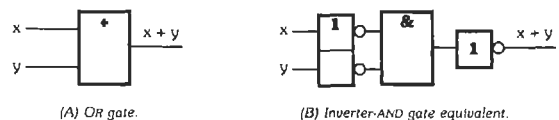
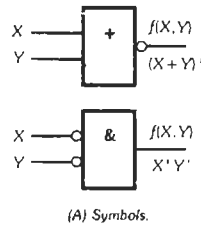
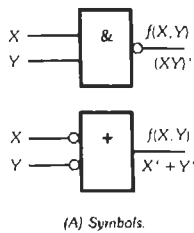


Fig. 8. Realization of an OR gate by means of an AND gate and inverters.



X	Y	NOR $f(X, Y)$
0	0	1
0	1	0
1	0	0
1	1	0

(B) Table of combinations.

X	Y	NAND $f(X, Y)$
0	0	1
0	1	1
1	0	1
1	1	0

(B) Table of combinations.

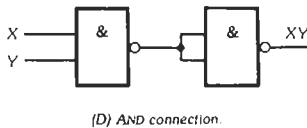
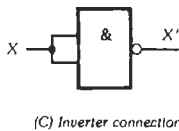


Fig. 9. NAND gate.

Another very common integrated-circuit gate is the NOR gate shown in Fig. 10. It is the dual of the NAND gate and is also a single gate that is sufficient to implement a network for any arbitrary switching function.

The other important IC gate type is the *sum modulo two*, or *EXCLUSIVE-OR (XOR) gate*,* which has a high output only when an odd number of its inputs are high. The table of combinations for a two-input XOR gate is given in Table 6. This table shows that the output of an

*The term XOR is used since a two-input gate has a 1 output if one *but not both* of the inputs is equal to 1. The OR gate (sometimes called *INCLUSIVE OR*) has a 1 output if one or both inputs are equal to 1.

TABLE 6. TABLE OF COMBINATIONS FOR A TWO-INPUT XOR GATE

X	Y	XOR $f(X, Y)$
0	0	0
0	1	1
1	0	1
1	1	0

Fig. 10. NOR gate.

XOR gate with input x_1 and x_2 is given by $x_1 \oplus x_2 = x_1x_2' + x_1'x_2$. It is easily demonstrated that the XOR operation is commutative and associative; that is,

$$x \oplus y = y \oplus x \quad \text{and} \quad (x \oplus y) \oplus z = x \oplus (y \oplus z)$$

If both inputs of an XOR gate are connected together to x , the output is given by $x \oplus x = xx' + x'x = 0$. The fact that this operation is associative means that any network composed only of XOR gates is equivalent to a single XOR gate with all inputs connected directly to it. This is illustrated in Fig. 11. Because of the facts that $x \oplus x = 0$ and $x \oplus 0 = x$, it is possible to remove any duplicated inputs so that the resulting network contains only one gate for which no input appears more than once. Since this gate, called the *reduced gate*, cannot function as an inverter, it follows that no network containing only XOR gates can be equivalent to an inverter. Thus the XOR gate is not a complete set by itself.

The XOR gate can perform as an inverter if a signal representing a constant 1 is available, since $x \oplus 1 = x1' + x'1 = x'$ (see Fig. 12). It is still not possible to construct an OR gate by using XOR gates and a 1 signal, for the output of the reduced gate will still be of the form $X \oplus Y \oplus 1$, which does not equal 1 when either or both of X and Y are equal to 1. A complete set can be formed by using both AND gates and XOR gates. As demonstrated above, any arbitrary function can be realized by a network containing only AND gates and inverters. By using XOR gates as in Fig. 12 to replace the inverters, it is possible to obtain a network containing only AND gates and XOR gates.

Canonical expressions involving AND, OR, and inverter operations have been derived. It is possible to obtain similar canonical expressions for any complete set. The canonical expression using AND and XOR for two-variable functions is

$$f(x, y) = g_0 \oplus g_1 \cdot x \oplus g_2 \cdot y \oplus g_3 \cdot x \cdot y$$

where,

$$g_0 = f_0$$

$$g_1 = f_0 \oplus f_2$$

$$g_2 = f_0 \oplus f_1$$

$$g_3 = f_0 \oplus f_1 \oplus f_2 \oplus f_3$$

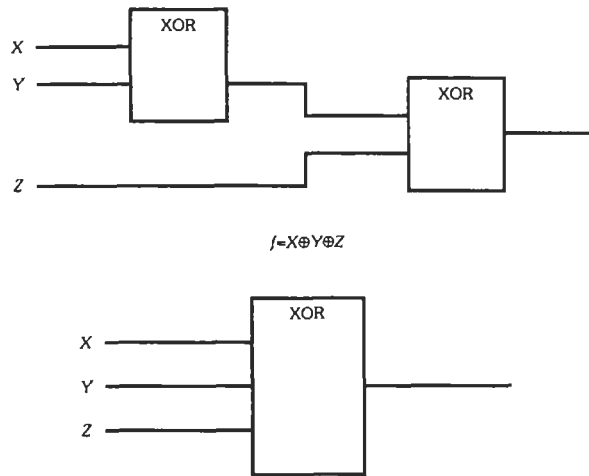


Fig. 11. Equivalent networks of XOR gates.

Thus for $f(x,y) = \Sigma(0,3) = x'y' + xy$

$$g_0 = 1$$

$$g_1 = 1$$

$$g_2 = 1$$

$$g_3 = 0$$

so that

$$f(x,y) = 1 \oplus x \oplus y$$

THE MAP METHOD

One important aspect of logic design is concerned with obtaining very efficient circuits. It is often desirable to minimize the total number of gates used for a specified maximum propagation delay through the circuit. No general algorithms are known for obtaining this objective. The best that can be done in an algorithmic fashion is to solve this minimization problem for networks that are limited to having at most two gates between any input and any output. This problem is important because: (1) the solution corresponds directly to a minimum PLA (programmable logic array) realization, (2) it is a minimum-delay solution, and (3) the resulting expressions are good starting places from which to derive efficient multistage networks. This section presents an introduction to the topic of two-stage minimization.

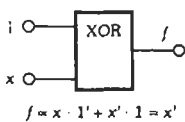


Fig. 12. Use of a constant 1 signal to form an inverter from an XOR gate.

The simplest sum-of-product-terms form of a function will be called a *minimal sum*.^{*} The sum-of-products form that has the fewest terms will be taken as the minimal sum. If there is more than one sum-of-products form having the minimum number of terms, and if these forms do not all contain the same total number of literals, then only the form(s) with the fewest literals will be called the minimal sum(s). For example, the function $f = x'yz + xyz + xyz'$ can be written as $f = yz + xyz'$, $f = x'yz + xy$, and $f = yz + xy$. Each of these forms contains two terms, but only the third form is a minimal sum, since it contains four literals, while the other two forms contain five literals each. The minimal sum corresponds to a gate circuit in which two-rail circuit inputs are connected to AND gates and the outputs of the AND gates form the inputs to an OR gate whose output is the circuit output. Such a circuit is called a *two-stage circuit*, since there are two gates connected in series between the circuit inputs and output. It is also possible to have two-stage circuits in which the circuit inputs are connected to OR gates and the circuit output is obtained from an AND gate. The minimal sum just defined corresponds to the two-stage circuit in which the output is derived from an OR gate and which contains the minimum number of gates.

The basic method for obtaining the minimal sum is to apply the theorem $XY + X'Y = Y$ to as many terms as possible and then to use the theorem $XY + X'Z + YZ = XY + X'Z$ to eliminate as many terms as possible.

Example:

$$f = x'y'z' + x'y'z + xy'z + xyz$$

$$x'y'z' + x'y'z = x'y'$$

$$xy'z + xyz = xz$$

$$f = x'y' + xz \quad \text{Minimal sum.}$$

^{*}This is called "minimal" rather than "minimum," since there may be more than one such form.

Example:

$$\begin{aligned}
 f &= w'x'y'z + w'x'yz + w'xy'z + w'xyz + wxy'z' \\
 &\quad + wxy'z + wx'y'z' + wx'y'z \\
 w'x'y'z + w'x'yz &= w'x'z \\
 w'xy'z + w'xyz &= w'xz \\
 wxy'z' + wxy'z &= wxy' \\
 wx'y'z' + wx'y'z &= wx'y' \\
 w'x'z + w'xz &= w'z \\
 wxy' + wx'y' &= wy'
 \end{aligned}$$

$$f = w'z + wy' \quad \text{Minimal sum.}$$

Example:

$$\begin{aligned}
 f &= xyz + x'yz + xy'z \\
 xyz &= xyz + xyz
 \end{aligned}$$

$$f = (xyz + x'yz) + (xyz + xy'z)$$

$$f = yz + xz \quad \text{Minimal sum.}$$

The last example illustrates the fact that it may be necessary to apply the theorem $XY + X'Y = Y$ several times, the number of literals in the terms being reduced each time. A single term may be paired with more than one other term, as shown in this example. The process of comparing pairs of terms to determine whether or not the theorem $XY + X'Y = Y$ applies can become very tedious for large functions. This comparison process can be simplified by using an "n-cube map," which is called a *Karnaugh map* when used for minimization.

Maps for Two, Three, and Four Variables

A map for a function of two variables, as shown in Fig. 13, is a square of four cells, or a 2-cube map. The value 0 or 1 that the function is to equal when $x = 1, y = 0$ (the entry in the 10 location or 2 row of the table of combinations) is placed in the cell having coordinates $x = 1, y = 0$. In general, the scheme for filling in the map is to place a 1 in all cells whose coordinates form a binary number that corresponds to one of the fundamental products included in the function and to place a 0 in all cells whose binary numbers correspond to fundamental products not included in the function. This is done very simply by writing a 1 in each cell whose decimal designation (decimal equivalent of the binary



(A) General form. (B) Map for $f = x'y' + xy' = \Sigma(0,2) - \Pi(1,3)$.

Fig. 13. Two-variable map.

number formed by the coordinates) occurs in the decimal specification of the function and writing 0s in the remaining cells.

The maps for functions of three and four variables are direct extensions of the two-variable map and are shown in Figs. 14 and 15. Maps for more than four variables are possible, but they are much more difficult to use and will not be discussed here.

Prime Implicants

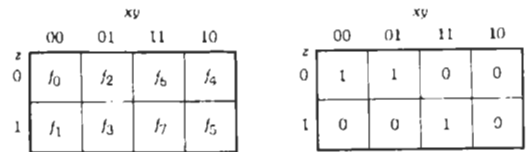
Two fundamental products can be "combined" by means of the theorem $XY + XY' = X$ if their corresponding binary numbers differ in only 1 bit. For the fundamental products $wxyz$ and $wxyz'$,

$$wxyz + wxyz' = wxy$$

The corresponding binary numbers are 1111 and 1110, which differ only in the lowest-order bit position. The fundamental products $wxyz$ and $w'xyz'$ cannot combine, and their corresponding numbers, 1111 and 0110, differ in the first and last bit positions.

Fundamental products that can be combined correspond to adjacent cells on the n-cube map. Thus cells that represent fundamental products that can be combined can be determined very quickly by inspection. In carrying out this inspection process, it must be remembered that cells such as f_4 and f_6 or f_1 and f_9 in Fig. 15 must be considered to be adjacent.

In a four-variable map, each cell is adjacent to four other cells, corresponding to the four bit positions in which two binary numbers can differ. In inspecting a map to determine which fundamental products can be combined, only cells with 1 entries (*1 cells*) need be considered, since these correspond to the fundamental products included in the function. Fig. 16 shows a four-variable map with adjacent 1 cells encircled. Notice that the 0111 cell is adjacent to two 1 cells. The rule for writing down the algebraic expression corresponding to a map is that there will be one product term for each pair of adjacent 1 cells and a fundamental product for each 1 cell that is not adjacent to any other 1 cell. The fundamental products are written down according to the rule: Any variable corresponding to a 0 in the binary number formed by the coordinates of the corresponding 1 cell is primed; the variables corresponding to 1s are left unprimed. The product terms corresponding to pairs of adjacent 1 cells are obtained by the same rule, with the exception that one variable is



(A) General form.

(B) Map for $f = x'y'z' + x'y'z + xyz$; $f = \Sigma(0,2,7) = \Pi(1,3,4,5,6)$.

Fig. 14. Three-variable map.

	wx			
	00	01	11	10
yz				
00	f_0	f_4	f_{12}	f_8
01	f_1	f_5	f_{13}	f_9
11	f_3	f_7	f_{15}	f_{11}
10	f_2	f_6	f_{14}	f_{10}

(A) General form.

	wx			
	00	01	11	10
yz				
00	1	0	0	0
01	1	0	0	0
11	0	1	1	0
10	0	1	1	0

(B) Map for $f = w'x'y'z' + w'x'y'z + w'xy'z' + w'xy'z + wxyz' + wxyz$
 $f = \Sigma (0,1,6,7,14,15) = \Pi (2,3,4,5,8,9,10,11,12,13)$.

Fig. 15. Four-variable map.

not included in the product. The variable excluded is that corresponding to the bit position in which the coordinates of the two 1 cells differ (see Fig. 16).

The situation in which it is possible to combine two of the terms obtained from pairs of the fundamental products must be considered next. In such a situation, four of the fundamental products can be combined into a single product term by successive applications of the $XY + XY' = X$ theorem. A function that is the sum of four such fundamental products is $f = wxyz + wxyz' + wxy'z + wxy'z'$. Application of the theorem to this function yields

$$f = (wxyz + wxyz') + (wxy'z + wxy'z')$$

$$= wxy + wxy' = wx$$

The characteristic property of four fundamental products that can be combined in this fashion is that all but two of the variables are the same (either primed or unprimed) in all four terms. The corresponding four binary numbers are identical in all but two bit positions. The corresponding cells on a map form "squares" (Fig. 17A) or "lines" (Fig. 17B) of four adjacent cells. For such a group of four cells on the map of a function, the corresponding products term is written just as for two adjacent cells, except that two variables corresponding to the two bit positions for which the cell coordinates change must be omitted.

It is also possible that eight of the fundamental products can be combined. In this case, all but three of the variables are identical (either primed or unprimed) in all eight terms.

	wx			
	00	01	11	10
yz				
00	1			1
01			1	
11			1	1
10				

Fig. 16. Four-variable map with adjacent 1 cells encircled.
 $f = \Sigma (0,5,6,8,15) = x'y'z' + w'xz + xy$.

The general rule is that, if in 2^i fundamental products all but i of the variables are identical (primed or unprimed), then the 2^i products can be combined and the i variables that change can be dropped.

In searching for a minimal sum for a function by means of a map, the first step is to encircle all sets of cells corresponding to fundamental products that can be combined (see Fig. 18). If one such set is contained in a larger set, only the larger set is encircled.* In Fig. 18 the set (0101, 0111) is not encircled. The encircled sets and the corresponding product terms will be called *prime implicants*.† These are exactly the terms that would result from repeated applications of the theorem $XY + XY' = X$. The terms appearing in the minimal sum will be some or all of the prime implicants.

Formation of Minimal Sums

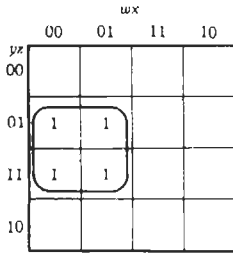
Once the prime implicants have been formed, the minimal sum can be determined directly from the map. The rule that must be followed in choosing the prime implicants that are to correspond to terms of the minimal sum is: *each 1 cell must be included in at least one of the chosen prime implicants*. The problem of obtaining a minimal sum is equivalent to that of selecting the fewest prime implicants. This rule is based on the fact that, for each combination of values of the input variables for which the function is to equal 1, the minimal sum must equal 1, and therefore at least one of its terms must equal 1. More simply, the map corresponding to the minimal sum must have the same 1 cells as the map of the original function.

A procedure for determining the minimal sum is first to determine whether any 1 cells are included in only one prime implicant. In Fig. 19, an asterisk has been placed in each 1 cell that is included in only one prime implicant. A 1 cell that is included in only one prime implicant is called a *distinguished 1 cell*.

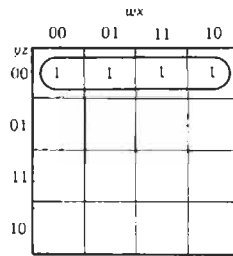
A prime implicant that includes a 1 cell that is not included in any other prime implicant is called an

*This corresponds to using the theorem $X + XY = X$.

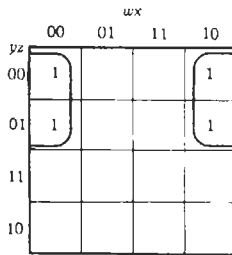
†This term was introduced by W. V. Quine. It is derived from the terminology of mathematical logic, but it has received widespread use in connection with logic design.



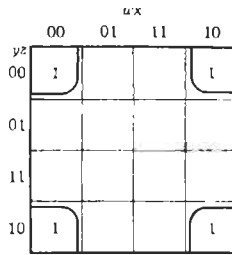
(A) $f = w'z$



(B) $f = y'z'$



(C) $f = x'y'$

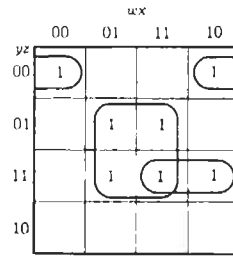


(D) $f = x'x'$

Fig. 17. Four-variable maps showing sets of cells corresponding to four fundamental products that can be combined.

essential prime implicant and must be included in the corresponding minimal sum.

In Fig. 19A, both prime implicants are essential and must be included in the minimal sum. A minimal sum does not always consist only of essential prime implicants. In Fig. 19B, only the essential prime implicants are shown. Cell 7 is not included in any of these, so another prime implicant that includes cell 7 must be



$f = \Sigma(0,5,7,8,11,13,15) = xz + x'y'z' + wyz$

Fig. 18. Map showing prime implicants.

present in the minimal sum. Fig. 19C shows the function of Fig. 19B after removal of the essential prime implicants. One of the two prime implicants shown must be included in the minimal sum, and the larger is chosen because the corresponding term contains fewer literals. The final minimal sum is $f = y'z + wz + w'yz' + xz$. There are some functions, such as that shown in Fig. 19D, that do not contain any essential prime implicants. For such functions, the minimum number of prime implicants required in the minimal sum can be determined by trial and error. The function of Fig. 19D has two minimal sums,

$$f = wx'y' + w'xy + x'yz'$$

and

$$f = wxz' + x'y'z + w'xy$$

The addition of d terms does not introduce any extra complexity into the procedure for determining minimal sums. Any d terms that are present are treated as 1 terms in forming the prime implicants, with the exception that no prime implicants containing only d terms are formed. The d terms are disregarded in choosing terms of the minimal sum. No prime implicants are included in order to ensure that each d term is contained in at least one prime implicant of the minimal sum.

The explanation of this procedure is that d terms are used to make the prime implicants as large as possible so as to include the maximum number of 1 cells and to contain as few literals as possible. No prime implicants need be included in the minimal sum because of the d terms, for it is not required that the function equal 1 for the d terms. An example of a function with d terms is given in Fig. 20.

It is often convenient to avoid determining all the prime implicants. This can sometimes be done by searching for 1 cells that are contained in only one prime implicant and thus determining the essential prime implicants. A 1 cell is selected, and the prime implicant or prime implicants that include the 1 cell are determined. If there is only one prime implicant, it is essential and must be included in the minimal sum. This procedure is continued until all the 1 cells are included in prime implicants of the minimal sum.

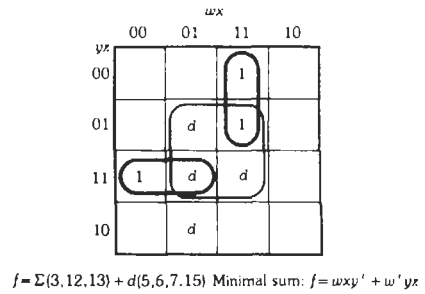
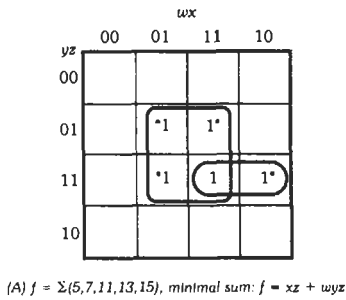


Fig. 20. Determination of minimal sum for a function with *d* terms. Prime implicants used in minimal sums are shown darkened.

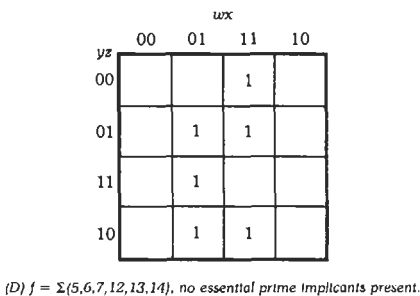
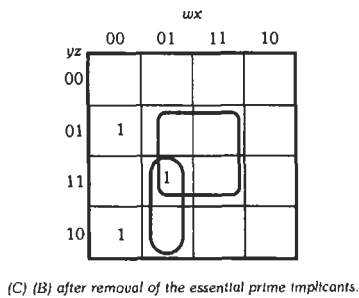
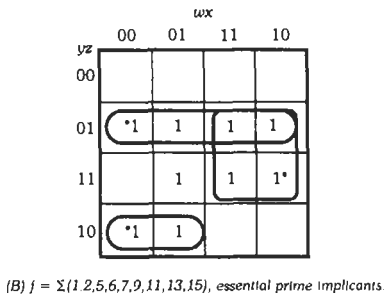


Fig. 19. Determination of minimal sums.

SEQUENTIAL CIRCUITS

The methods developed above are applicable to combinational circuits—circuits whose outputs are determined completely by their present inputs. Many digital circuits satisfy this restriction; however, there are also many circuits that do not. Circuits whose outputs depend not only on the present inputs but also on previous inputs are called *sequential circuits*.

The difference between a combinational circuit and a sequential circuit is analogous to the difference between

the two types of combination lock shown in Fig. 21. Lock A is open or closed depending only on the present setting of its dials; past settings are unimportant, just as, in a combinational circuit, past inputs are unimportant in determining the present circuit outputs. Whether lock B is open or not depends not only on which number the pointer is selecting but also on which numbers the pointer stopped at previously. Similarly, the output of a sequential circuit depends on previous as well as present inputs.

In order for the output of a sequential circuit to depend on past inputs, the circuit must have some mechanism to retain information about previous inputs. This mechanism is some type of memory element. The first memory element used in electronic logic circuits is now called a *latch*. A latch constructed of NOR gates is shown in Fig. 22. This circuit is called a *set-reset latch* or *SR latch* (S represents the set input and R represents the reset input).

The operation of the circuit is illustrated in the figure. When $S = 1$ and $R = 0$ as in Fig. 22A, the circuit is "set," and the output values are $Q = 1, \bar{Q} = 0$. The shading of the top NOR gate of this figure indicates that the output of that gate is "forced" by an external signal (S in this case). If S is then changed to 0, the circuit remains in the set condition, as shown in Fig. 22B. Changing S back to 1 returns the circuit to the situation of Fig. 22A. If, instead, R becomes 1, the circuit becomes "reset," and the conditions of Fig. 22C are present. The circuit remains in the reset state if R is returned to 0 as in Fig. 22D. Thus, when S and R are

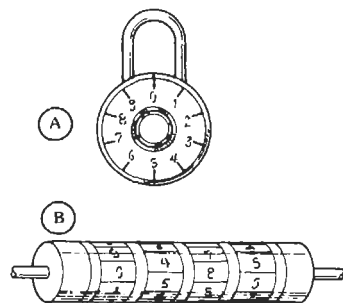


Fig. 21. Two types of combination lock.

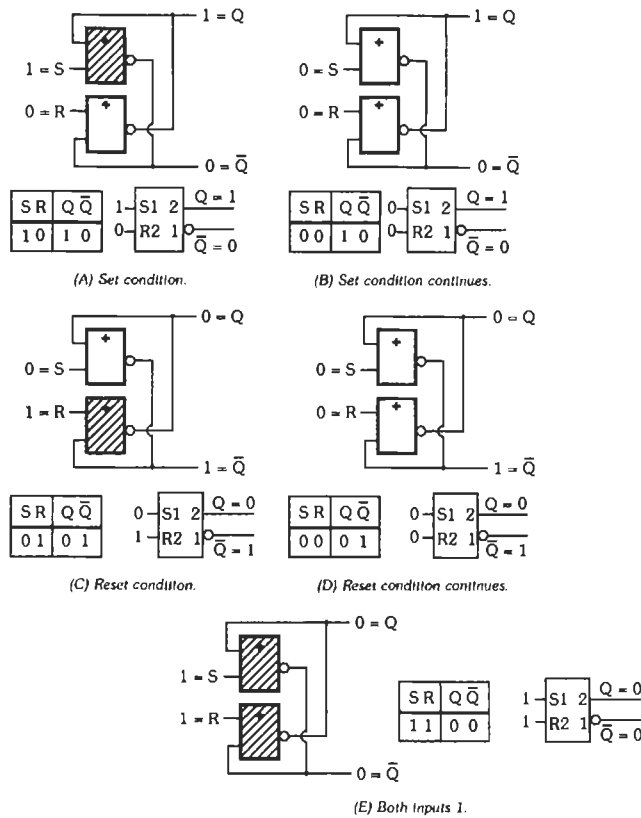


Fig. 22. Set-reset latch constructed of NOR gates.

both 0, the output indicates which of S or R was last equal to 1—the circuit “remembers” the last nonzero input condition. Waveforms illustrating the latch operation are shown in Fig. 23. Fig. 22 includes tables showing the circuit conditions, and the logic symbol for a set-reset latch with the appropriate signal values. A table typical of those used by manufacturers to summarize the circuit operation is shown in Table 7.

Definitions of the SR latch do not specify the circuit operation with both inputs at 1. This is because the circuit action with $S = R = 1$ is very dependent on the specific implementation. Usually, networks containing SR latches are designated so that the $S = R = 1$ input condition never occurs.

For a specific circuit design, the action of an SR latch with both inputs equal to 1 is, of course, fixed. The

result of placing 1s on both inputs to the NOR-gate SR latch is shown in Fig. 22E: both gates have their outputs forced to 0. A disturbing feature of this is that the outputs are no longer complementary. Despite the drawbacks, there are situations in which use of the $S = R = 1$ input state can be very useful.

An important characteristic of this circuit is that when an input value is changed, any effect on the output appears immediately after the new input appears. The new output is delayed only by the propagation times of the devices in the circuit. All latches have this property, but it is sometimes emphasized by calling them *transparent latches*.*

*Texas Instruments Inc. 1981 Supplement to *The TTL Data Book*, 2nd ed. 1981 p. 338.

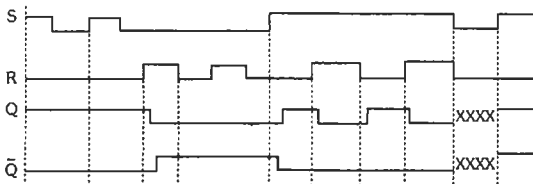


Fig. 23. Waveforms of SR latch.

TABLE 7. SR LATCH OPERATION

S	R	Q_n
0	0	Q_{n-1}
1	0	1
0	1	0
1	1	F

A word of warning is appropriate here. The terminology used in connection with bistable circuits such as the set-reset latch can be confusing since there is a lack of consistency and care in naming such circuits. It does appear that this situation is improving and will continue to get better as the new IEEE Standard—*Graphic Symbols for Logic Functions* (reference 3)—gains wider acceptance. The symbols used here follow the spirit of the new Standard, which is discussed in more detail in references 4 and 5.

The discussion in this chapter is an introduction to some of the basic issues involved in logic design. It is far from complete. For a more thorough coverage on this topic, the interested reader is referred to reference 5, from which much of the material of this chapter has been drawn.

REFERENCES

1. Shannon, C. E. "A Symbolic Analysis of Relay and Switching Circuits." *Trans. AIEE*, Vol. 57, 1938, pp. 713-723.
2. Birkhoff, G., and MacLane, S. *A Survey of Modern Algebra*. New York: The Macmillan Co., 1955.
3. ANSI/IEEE Std 91-1984, IEEE Publication SH08615, July 13, 1984.
4. Peatman, J. B. *Digital Hardware Design*. New York: McGraw-Hill Book Co., 1980.
5. McCluskey, E. J. *Logic Design Principles*. Englewood Cliffs, N.J.: Prentice-Hall, Inc., 1986.

44

Probability and Statistics

Surendra M. Gupta

Introduction 44-2

 Discrete Probability Distribution

 Continuous Probability Distribution

Mathematical Expectation 44-3

Important Theoretical Probability Distributions 44-4

Goodness of Fit 44-4

INTRODUCTION

Probability and statistics concepts have very useful applications in science and engineering. This chapter briefly introduces some of the important concepts of probability and statistics.

In probabilistic modeling one often refers to a *random experiment*, which represents an experiment that can be repeated time and again under similar circumstances, but which yields unpredictable results at each trial. For example, tossing a die is a random experiment where the result is a number from the sample space { 1, 2, 3, 4, 5, 6 }. A random experiment may consist of observing or measuring elements taken from a set that is known as a *population*. A real number associated with the result of a random experiment is called a *random variable* or *variate*. A variate may be *discrete* or *continuous*.

Discrete Probability Distribution

A discrete variate takes its value from a finite or denumerable set { x_1, x_2, \dots, x_n } described by its *probability function*, p_k , defined as the probability of obtaining x_k as a result of one trial. p_k has the following properties:

$$0 \leq p_k \leq 1$$

and

$$\sum_{\text{all } k} p_k = 1$$

A discrete variate can be described by its *cumulative probability function* as

$$P(x) = \sum_{x_k \leq x} p_k$$

A discrete variate of multiple dimensions is described by its *joint distribution*. For example, if (x_j, y_k) are the coordinates of a point in the (x, y) plane, then $p(x_j, y_k)$ is the joint probability distribution such that $x = x_j$ and $y = y_k$. In addition,

$$p_1(x_j) = \sum_{\text{all } k} p(x_j, y_k)$$

and

$$p_2(y_k) = \sum_{\text{all } j} p(x_j, y_k)$$

are, respectively, the *marginal probabilities* that $x = x_j$ independent of y and that $y = y_k$, independent of x . Also,

$$p(x_j | y_k) = p(x_j, y_k) / p_2(y_k), \quad p_2(y_k) > 0$$

and

$$p(y_k | x_j) = p(x_j, y_k) / p_1(x_j), \quad p_1(x_j) > 0$$

are, respectively, the *conditional probabilities* that $x = x_j$ given that $y = y_k$ has already occurred, and that $y = y_k$ given that $x = x_j$ has already occurred.

Finally, x_j and y_k are said to be *independent* if $p(x_j, y_k) = p_1(x_j)p_2(y_k)$ for all x_j and y_k .

Continuous Probability Distribution

A continuous variate takes its value from a non-denumerable set described by its *probability density function*, $p(x)$. The probability that one trial of the experiment gives a result between x and $x + dx$ is $p(x)dx$. The *cumulative distribution function*, $P(x)$, is defined as the probability that the continuous variate is less than or equal to x . It can be represented mathematically as

$$P(x) = \int_{-\infty}^x p(s)ds$$

Note that the following properties hold:

$$P(x_1) \geq P(x_2) \text{ if } x_1 \geq x_2$$

$$P(-\infty) = 0$$

$$P(+\infty) = \int_{-\infty}^{+\infty} p(s)ds = 1$$

$$P(x) \geq 0$$

$$p(x) = dP / dx$$

A continuous variate of multiple dimensions, or *multivariate*, is described by its *joint distribution*. For example, if (x, y) are the coordinates of a point in a plane, then $p(x, y)dxdy$ is the probability that the multivariate has its x -coordinate between x and $x + dx$ and its y -coordinate between y and $y + dy$. In addition,

$$p_1(x) = \int_{-\infty}^{+\infty} p(x, y)dy$$

and

$$p_2(y) = \int_{-\infty}^{+\infty} p(x, y)dx$$

are the *marginal probability density functions*, which means that $p_1(x)dx$ is the probability that the variate x lies between x and $x + dx$ independent of y , and $p_2(y)dy$ is the probability that the variate y lies between y and $y + dy$ independent of x . Also,

$$p(x | y) = p(x, y) / p_2(y), \quad p_2(y) > 0$$

and

$$p(y|x) = p(x, y) / p_1(x), \quad p_1(x) > 0$$

are, respectively, the *conditional probability functions* such that $p(x|y_0)dx$ is the probability that the variate x lies between x and $x + dx$ given that $y = y_0$ and $p(y|x_0)dy$ is the probability that the variate y lies between y and $y + dy$ given that $x = x_0$.

Finally, two variates x and y are said to be *independent* if $p(x, y) = p_1(x)p_2(y)$ for all x and y .

MATHEMATICAL EXPECTATION

Various forms of mathematical expectations are used to describe the main features of a probability distribution. The mean and median of the distribution establish the center. The mean represents the center of gravity of the probability density function while the median divides that function into two equal parts. The root mean square, variance, standard deviation, and the mean absolute deviation are measures of the spread about the mean. The moments of order three and four about the mean describe the asymmetry and peakedness, respectively, of the probability density function.

The *expected value of a function* of a random variable can be written as

$$E[g(x_k)] = \sum_{\text{all } k} g(x_k) p_k$$

for the function $g(x_k)$ of the discrete random variable x_k , and

$$E[g(x)] = \int_{-\infty}^{\infty} g(x) p(x) dx$$

for the function $g(x)$ of the continuous random variable x .

Using the definition of the expected value, the various properties of a distribution could be obtained as follows.

The *mean* of a distribution is

$$\mu = E[x_k] = \sum_{\text{all } k} x_k p_k$$

for the discrete case

and

$$\mu = E[x] = \int_{-\infty}^{\infty} x p(x) dx$$

for the continuous case

For the discrete case, the *median* is a value m such that the variate x_k has equal probability of being larger or smaller than m . For the continuous case, it has a value such that the variate x satisfies

$$\int_{-\infty}^m p(x) dx = \int_m^{\infty} p(x) dx = \frac{1}{2}$$

The *mode* is a value of x_k (or x) where the probability density p_k (or $p(x)$) is the largest. Note that there may be more than one mode.

The *root mean square* is

$$r = [E[x_k^2]]^{1/2} = \left[\sum_{\text{all } k} x_k^2 p_k \right]^{1/2}$$

for the discrete case

and

$$r = [E[x^2]]^{1/2} = \left[\int_{-\infty}^{\infty} x^2 p(x) dx \right]^{1/2}$$

for the continuous case

The *moment of order r about the origin* is

$$\nu_r = E[x_k^r] = \sum_{\text{all } k} x_k^r p_k$$

for the discrete case

and

$$\nu_r = E[x^r] = \int_{-\infty}^{\infty} x^r p(x) dx$$

for the continuous case

Note that $\nu_1 = \mu$.

The *moment of order r about the mean* is

$$\mu_r = E[(x_k - \mu)^r] = \sum_{\text{all } k} (x_k - \mu)^r p_k$$

for the discrete case

and

$$\mu_r = E[(x - \mu)^r] = \int_{-\infty}^{\infty} (x - \mu)^r p(x) dx$$

for the continuous case

The *variance* is

$$\sigma^2 = \mu_2 = E[(x_k - \mu)^2] = \sum_{\text{all } k} (x_k - \mu)^2 p_k$$

for the discrete case

and

$$\sigma^2 = \mu_2 = E[(x - \mu)^2] = \int_{-\infty}^{\infty} (x - \mu)^2 p(x) dx$$

for the continuous case

The *standard deviation* (or the *root mean square deviation from the mean*) is

$$\sigma = \{E[(x_k - \mu)^2]\}^{1/2} = \left[\sum_{\text{all } k} (x_k - \mu)^2 p_k \right]^{1/2}$$

for the discrete case

and

$$\sigma = \{E[(x - \mu)^2]\}^{1/2} = \left[\int_{-\infty}^{\infty} (x - \mu)^2 p(x) dx \right]^{1/2}$$

for the continuous case

The *mean absolute deviation* (or the *mean absolute error*) is

$$mae = E[|x_k - \mu|] = \sum_{\text{all } k} |x_k - \mu| p_k$$

for the discrete case

and

$$mae = E[|x - \mu|] = \int_{-\infty}^{\infty} |x - \mu| p(x) dx$$

for the continuous case

IMPORTANT THEORETICAL PROBABILITY DISTRIBUTIONS

A probability distribution describes the behavior of a random variable. Often, the observations generated from many different statistical experiments behave in similar ways. This means that the random variables generated by the different experiments may be more or less explained by the same probability distribution and hence could be represented by a single mathematical expression. As it turns out, one requires only a few of the standard probability distributions to describe most types of random variables encountered in practice.

Tables 1 and 2 present a summary of some common theoretical discrete and continuous probability distributions.

GOODNESS OF FIT

To determine whether the sample data fit a given theoretical distribution or not, a goodness of fit test is performed. One of the most widely used goodness of fit tests is the *chi-square test*. The chi-square test is based on comparing the probability density of the sample data with that of the theoretical distribution.

The first step is to divide the sample data into m disjoint sets. Let the number of sample points falling within the i th set be f_i . From the assumed theoretical distribution and the size of the sample, the expected number, g_i , of points in the i th set is computed. The deviation between this and the actual result is expressed by

$$D = \sum_{i=1}^m \frac{(f_i - g_i)^2}{g_i}$$

If the f_i are sufficiently large, D is approximately distributed according to the chi-square distribution with

TABLE 1. VARIOUS DISCRETE DISTRIBUTIONS AND THEIR PROPERTIES*

	Parameters	p_k	Mean	Variance
Bernoulli	$0 < p < 1$	$p^k q^{1-k}, k = 0, 1$	p	pq
Binomial	n $0 < p < 1$	$\binom{n}{k} p^k q^{n-k},$ $k = 0, 1, \dots, n$	np	npq
Geometric	$0 < p < 1$	$pq^k, k = 0, 1, \dots$	$\frac{q}{p}$	$\frac{q}{p^2}$
Hypergeometric	$N > 0$ $n, k \geq 0$	$\frac{\binom{r}{k} \binom{N-r}{n-k}}{\binom{N}{n}}$ $k = 0, 1, \dots, \text{Min}(r, n)$	$\frac{nr}{N}$	$\frac{nr(N-r)(N-n)}{N^2(N-1)}$
Negative Binomial	$0 < p < 1$ r is positive integer	$\binom{r+k-1}{k} p^r q^k,$ $k = 0, 1, \dots$	$\frac{rq}{p}$	$\frac{rq}{p^2}$
Poisson	$\alpha > 0$	$e^{-\alpha} \frac{\alpha^k}{k!}$ $k = 0, 1, \dots$	α	α

* $q = 1 - p$

TABLE 2. VARIOUS CONTINUOUS DISTRIBUTIONS AND THEIR PROPERTIES

	Parameters	$p(x)$	Mean	Variance
Chi-square	$n > 0$	$\frac{x^{n/2-1}}{2^{n/2}\Gamma(n/2)}e^{-x/2}$, $x > 0$	n	$2n$
Erlang-k	$k > 0, \mu > 0$	$\frac{\mu^k(\mu k x)^{k-1}}{(k-1)!}e^{-\mu k x}$, $x > 0$	$\frac{1}{\mu}$	$\frac{1}{k\mu^2}$
Exponential	$\alpha > 0$	$\alpha e^{-\alpha x}, x \geq 0$	$\frac{1}{\alpha}$	$\frac{1}{\alpha^2}$
Gamma	$\alpha > 0, \beta > 0$	$\frac{\alpha(\alpha x)^{\beta-1}}{\Gamma(\beta)}e^{-\alpha x}$, $x > 0$	$\frac{\beta}{\alpha}$	$\frac{\beta}{\alpha^2}$
Normal	$\sigma > 0$	$\frac{1}{\sigma\sqrt{2\pi}}e^{-\frac{1}{2}\frac{(x-\mu)^2}{\sigma^2}}$	μ	σ^2
Uniform	$a < b$	$\frac{1}{b-a}, a \leq x \leq b$	$\frac{a+b}{2}$	$\frac{(b-a)^2}{12}$

$m - 1 - r$ degrees of freedom, where r is the number of parameters that are estimated from the sample data for use in defining the theoretical distribution.

To illustrate the application of the chi-square test, consider the following example. During three successive one-hour periods, the number of telephone calls received at a station was 11, 15, and 23, while during two nonoverlapping two-hour periods, it was 40 and 37. How does this agree with a Poisson process?

Since the number of calls per hour is not specified, it is estimated from the sample as

$$\frac{11 + 15 + 23 + 40 + 37}{7} = 18$$

The deviation from the expected number is

$$\frac{7^2}{18} + \frac{3^2}{18} + \frac{5^2}{18} + \frac{4^2}{36} + \frac{1^2}{36} = 5.1$$

For $5 - 1 - 1 = 3$ degrees of freedom, this deviation exceeded about 15 percent of the time. The assumption of a Poisson process is therefore very good. It would have been significantly doubtful only if the deviation obtained was exceeded as rarely as 5 percent or less of the time, which corresponds to D larger than 7.8 in this application.

REFERENCES

Miller, I., Freund, J. E., and Johnson, R. A. *Probability and Statistics for Engineers*, 6th ed. New York: Prentice Hall, 1999.

Montgomery, D. C., and Runger, G. C. *Applied Statistics and Probability for Engineers*, 2nd ed. New York: John Wiley & Sons, 1998.

Ross, S. M. *Introduction to Probability and Statistics for Engineers and Scientists*, 2nd ed. San Diego: Academic Press, 1999.

45 Reliability and Life Testing

*Revised by
Douglas L. Marriott*

Definitions and Terminology	45-2
Reliability Definitions	45-5
Organization for Reliability	45-5
Elements of Reliability Assessment	45-6
Design Reviews	
Failure Mode Analysis	
Quantitative Reliability Assessment	
Failure Data Collection and Assessment	
Component Reliability	45-9
System Reliability	45-11
Sources of Reliability Data	45-13
Generic Data Sources	
Life Testing	
Probability and Statistical Inference	45-15
Confidence Limits	
Fitting a Distribution Using Chi-Squared Test	
Probability Paper	
Weibull Analysis	
Distribution-Free Tests of Goodness of Fit	
Bayesian Statistics	
Codes and Standards	45-25
The Use of Computers in Reliability	45-25

DEFINITIONS AND TERMINOLOGY

Availability: Probability of a system subject to repair operating satisfactorily on demand.

Average Life: The mean value for a normal distribution of lives. Generally applied to mechanical failures resulting from "wear-out."

Burn-In (also Initial Failure): Initially high failure rate encountered on first placing a component on test. Burn-in failures are usually associated with manufacturing defects and the debugging phase of early service.

Component: (Normally used interchangeably with the term "unit.") A component is defined as an article which is normally a combination of parts, subassemblies, or assemblies, and is a self-contained element of a complete operating equipment and performs a function necessary to the operation of that equipment. *Examples:* indicator unit, power unit, receiver, transmitter, rotating antenna, modulator unit, amplifier unit.

Confidence Level (Coefficient): The degree of desired trust or assurance in a given result. A confidence level is always associated with some assertion and measures the probability that a given assertion is true. For example, it could be the probability that a particular characteristic will fall within specified limits, i.e., the chance that the true value of P lies between $P = a$ and $P = b$.

Configuration Management: Management of and knowledge of where all specifications, procedures, and associated test results are located and assigned, so that it is possible to produce these controlled items and all reliability evaluations and predictions pertaining to the system.

Cumulative Distribution Function: If x is a random variable, then the cumulative distribution function of x is defined to be the function F such that for every real number t , $F(t)$ is the probability that a given outcome of x will not exceed t ; in symbols:

$$F(t) = \Pr(x \leq t)$$

Defect: Any deviation of a unit of product from specified requirements. A unit of product may contain more than one defect.

Degradation Failure: A failure that results from a gradual change in performance characteristics of an equipment or part with time.

Design Reviews:

(A) Preliminary design review: As soon as possible after a contract has been signed, a breadboard model should be built and its reliability estimated. Reliability engineering shall re-evaluate all parts and components and determine and recommend improvements in the design.

(B) Intermediate design reviews: While developing the system, conduct design reviews on a formal basis at

all suppliers as well as at the prime contractor. This program should be coordinated with the reliability growth program. Account must be taken of the contract requirements for reliability goals at scheduled points in the production schedule.

(C) Critical design review: When engineering believes the design is ready to be "frozen" and also when a satisfactory prototype has met the qualification and other reliability tests, a final design review shall be scheduled. This formal review takes into account all contract demands as modified by the most recent changes in the contract. If the product is adjudged to be satisfactory, the final design may be approved for production. The block system should be used and authorization should be given to production to make x units per the specifications and blueprints without any change.

Downtime: Time during which equipments are not capable of doing useful work because of malfunction. This does not include preventive-maintenance time. In other words, downtime is measured from the occurrence of a malfunction to the correction of that malfunction.

Equipment: Material or articles (such as sets) comprising an outfit. The term "equipment" sometimes is used instead of "set," i.e., one or more assemblies or a combination of items capable of independently performing a complete function. *Examples:* Radio receiver, digital computer, automobile. An equipment may contain several sets as components. An example would be two radio receivers assembled for dual-diversity reception. The combination would constitute the equipment.

Failure: A failure is a detected cessation of ability to perform a specified function or functions within previously established limits on the area of interest. It is beyond adjustment by the operator by means of controls normally accessible to him during the routine operation of the device. This requires that measurable limits be established to define satisfactory performance of the function.

Failure Mode Analysis: Research, development, and production engineers as well as the reliability engineers analyze the basic design and determine by simulation and logistics what possible failures might occur. Corrective measures for eliminating and preventing failures are built into the basic design. Standard forms for the failure mode analysis are made available, and the results must be given to the prime contractors and production engineers for evaluation.

Failure Modes and Effects Analysis (FMEA): Extension of failure mode analysis which also considers the consequences of the failure mode, such as "revealed fault," "fail-safe," etc.

Failure Rate and Hazard Rate: Failure rate is generally the rate at which failure occurs during an interval of time (given that it has not occurred before the start of

that interval) as a function of the total interval length. Hazard rate is an instantaneous failure rate and is defined as the limit of the failure rate as the time interval approaches 0. An example might be: A family takes an automobile trip of 120 miles and completes the trip in three hours. The average rate was 40 mph, although they drove at various rates of speed. The rate at any given instant could be determined by reading the speedometer at that instant. Therefore, the 40 mph is equivalent to the failure rate and the speed at any instant of time equals the hazard rate.

Inherent Reliability: The basic or generic failure rates of components have often been compiled by several companies as well as by some government agencies. A library covering failure rates should be part of a good reliability program. Such a listing gives concretely the reliability that can be guaranteed.

Lot Size: A specific quantity of similar material or collection of similar units from a common source; in inspection work, the quantity offered for inspection and acceptance at any one time. It may be a collection of raw material, parts, subassemblies inspected during production, or a consignment of finished products to be sent out for service.

Maintainability: The maintainability of an equipment in a specified maintenance environment is the probability that a failure will be repaired within a specified time after the failure occurs.

Mean Time Before Failure (MTBF): The total measured operating time of a population of equipments divided by the total number of failures of a repairable equipment is defined as the ratio of the total operating time to the total number of failures. The measured operating time of the equipments of the population that did not fail must be included. This measurement is normally made during that period of time between the early life and wear-out failures. In the case of a constant hazard rate, this ratio is the reciprocal of the failure rate.

The MTBF can be determined by dividing the product of the number of equipments tested, N , and the test time, t , by the number of failures, f , which occur during that time; i.e., MTBF, or often just m , is equal to Nt/f . The quantity m is the reciprocal of λ , i.e., $m = 1/\lambda$, and is related to the probability of survival by the exponential law $P_s = e^{-t/m}$. The figure of merit m (sometimes expressed as t) is convenient for use in determining if the reliability of an equipment is likely to be adequate for missions of specific lengths.

Mean Time to Failure (MTTF): The measured operating time of a single piece of equipment divided by the total number of failures of the equipment during the measured period of time. This measurement is normally made during that period of time between the early life and wear-out failures.

Mean Time to Repair (MTTR): The measured repair time divided by the total number of failures of the equipment.

Mission Success Rate: That percentage of the total missions uninterrupted by failure of the equipment. This figure of merit is more closely dependent on the reliability of the parts included in the system and on the design of the system than are either maintenance ratio or in-commission rate. However, this measure of reliability is valuable primarily to a using agency that has a regular schedule of missions. A mission success rate obtained by one agency is not typical of the equipment in general and will not necessarily apply for other agencies with different operating schedules.

Mode of Failure: The physical description of the manner in which a failure occurs and the operating condition of the equipment or part at the time of the failure.

Part Failure Rate: The rate at which a part fails to perform its intended function.

Probability: The limiting relative frequency in an infinite random series. If an event can occur in n ways and its failure in m ways, and if these $m + n$ ways are equally likely, then the mathematical probability that the event will occur in any one trial is the ratio $n/(n + m)$.

In other words, the probability of an event is the theoretical relative frequency with which it will occur, such relative frequency being the ratio of the number of times the event is observed under experimental conditions to the total of a great number of observations made under those conditions.

Common notation is:

$Pr(X)$ = Probability of event X

$Pr(X/Y)$ = Probability of event X
given Y has occurred

Probability Density Function (PDF): The relative frequency of a continuous random variable, obtained from the cumulative distribution function by differentiation.

$$f(t) = dF(t)/dt$$

Probability of Survival: A numerical expression of reliability with the accepted nomenclature of P_s and a range from 0 to 1.0 indicating the extremes of "impossibility" and "certainty."

In other words, the probability of a given equipment performing its intended function or the given use cycle is

$$R(t) = 1 - F(t)$$

Product Effectiveness: The entire reliability program must be tied in with the quality engineering programs for securing the most effective operation possible. Product effectiveness includes all the elements for

securing at minimum cost a product with maximum effectiveness. Programs for quality assurance contain schedules and procedures that encompass within them reliability, preventive maintenance, value engineering, human engineering, quality control, inspection, and tests that result in systems and products that will prove most effective when in operation.

Quality Assurance (QA): All those activities, including surveillance, inspection, control, and documentation, aimed at ensuring that the product will meet its performance specifications. Quality assurance generally has an independent role from production and reports directly to senior management.

Reliability Allocations: With an overall system reliability goal and where the confidence level is known, reliability values may be allocated to every component in the system by the use of available failure rates and weighting factors.

Reliability Demonstration: Critical tests must be programmed to provide sufficient valid data to determine the reliability of the system and all critical component parts. Provision should be made for processing all such data as expeditiously as possible to speed up all phases of the program with truly reliable materials and parts.

Reliability Goals: Requests for bids, specifications, and contracts requiring quality assurance and reliability generally describe completely the reliability goals. The reliability of the system with specified confidence levels is the principal goal. Values to be secured at specified points in time are listed on the growth curves. A reliability demonstration program is detailed to show that the reliability goals have been attained. How are such reliability goals expressed? The simplest statement covers only an expected value. The quality assurance paragraph in the specification will state simply: "The desired reliability is 99.7%." This is too simple, as the length of the mission or number of cycles of operation has not been detailed. The statement should be, "For 100 hours of operation the specified reliability is 99.7%."

Many contracts also introduce confidence levels. For example, after the initial design review, the reliability of the system must be 99% with a confidence level of 0.60. After qualifying, the system must have a reliability of 99.7% for a mission time of 10 hours with a confidence level of 0.99. Thus expressed, this goal has within it a final goal plus some information concerning the desired growth curve for reliability. In many programs, provision for the reliability demonstration program has not been made, or in making final arrangements for the finalized contract it is cancelled because of lack of funds. The reliability engineer should establish a very modest program to check the achievement of the reliability goals by means of an economic reliability demonstration program. Thus, it must be emphasized that the simplest possible reliability demonstration test should be made a part of the final program covered by

contract and funds. This provides vital evidence that the reliability goals have been achieved and that the customer reliability requirements have been met.

Reliability Growth Curves: Periodic reports, such as monthly or quarterly, should be prepared containing up-to-date predictions of the system's reliability. These should be presented graphically on the reliability growth curves for this system and its various components and parts.

Reliability Predictions: Many agencies and companies have compiled failure rates for parts, components, subassemblies, assemblies, and systems. These generic failure rates are used as basic data to predict a value for reliability.

Sample: One or more sample units selected at random from a quantity of product to represent that quantity of product for inspection purposes.

Sequential Sampling: Sampling inspection in which, after each unit is inspected, the decision is made to accept, to reject, or to inspect another unit. *Note:* Sequential sampling as defined here is sometimes called "unit sequential sampling" or "multiple sampling." Multiple sampling is preferred, to differentiate from sequential testing.

Significance Level: The level of confidence (i.e., the probability of being correct) at which a hypothesis, such as "goodness of fit," is to be accepted or rejected.

Specification Limits: The specification limit(s) is the requirement that a quality characteristic should meet. This requirement may be expressed as an upper or a lower specification limit (called a single specification limit) or both upper and lower specification limits (called a double specification limit).

System (General): A combination of parts, assemblies, and sets joined together to perform a specific operational function or functions. *Examples:* Piping system, refrigeration system, air conditioning system.

Test to Failure: Testing conducted on one or more items until a predetermined number of failures have been observed. Failures are induced by increasing electrical, mechanical, and/or environmental stress levels, usually in contrast to life tests in which failures occur after extended exposure to predetermined stress levels. A life test can be considered a test to failure using age as the stress.

Value Engineering: One feature of a good reliability program is a concurrent value engineering program. Many programs include incentive provisions. If an operation is usually expensive, provisions should be made for a series of improvements that fall in with the value engineering and incentive programs. Provisions for successive cost reduction programs as well as basic value engineering improvements should be added to each contract. It may be made a part of the incentive programs that are now usually included in military and government contracts for procurement.

Weapon System: A missile and all the necessary support equipment (either ground or airborne) necessary to launch and properly operate a missile.

Wear-Out Period: The wear-out period of an equipment is that period of equipment life, following the normal operating period, during which the equipment failure rate increases above the normal rate.

RELIABILITY DEFINITIONS

Reliability is the probability that a given product, system, or action will achieve its designated goal successfully under the specified environmental conditions and for a prescribed period of time or for the number of cycles of operation required for the mission or task.

Reliability involves three distinct concepts:

1. Attaining a specified level of performance
2. The probability of achieving that level
3. Maintaining that level for a specified time

ORGANIZATION FOR RELIABILITY

Reliability is one of the disciplines that form the assurance sciences. These disciplines govern the quality, safety, and dependability of a product throughout its life cycle. Whereas many of the aspects of assurance science, such as quality control, are aimed at ensuring that a product has been manufactured within acceptable limits, reliability is also related to aspects of design and development, in that its objective is to predict and assure an acceptable level of performance in service. For this reason, provision must be made in the structure of any manufacturing organization to integrate reliability activities with the product evolution process.

There is no single answer to the question of where the reliability function should exist in any organization. To some extent, the solution depends on the form of product and the organization itself. Comprehensive discussions of organization for product assurance and reliability are contained in References 1 and 2. Guidance is also available from professional bodies. For instance, the Institute of Electrical and Electronics Engineers and the American Society for Quality Control have defined a structure as shown in Fig. 1.

Today, reliability seems to fit most naturally into the product-assurance or product-effectiveness section of major space equipment manufacturers.

The most important organizational aspect is the fact that reliability and well-versed quality personnel use one tool which, when shared jointly, adds to the "experience-retention" capability; that tool is statistics. If the intent of any organization is to get the most out of the "fall of data," the organization of Fig. 1 is effective. Because of the problem of decision making, reliability is best placed at the staff level.

In a small organization, reliability is usually shared by engineering and by quality assurance. Some firms, to better satisfy the requirements of the Department of Defense, put reliability in the quality organization spectrum of operations, viz.,

- Quality and reliability assurance tasks
- Data reduction
- Planning
- Inspection, etc.
- Test plan
- Vendor survey
- Quality-control engineering

Reliability planning, to be effective, must fit into the total program. One of the best ways to gain acceptance of the reliability or maintainability effort is to establish

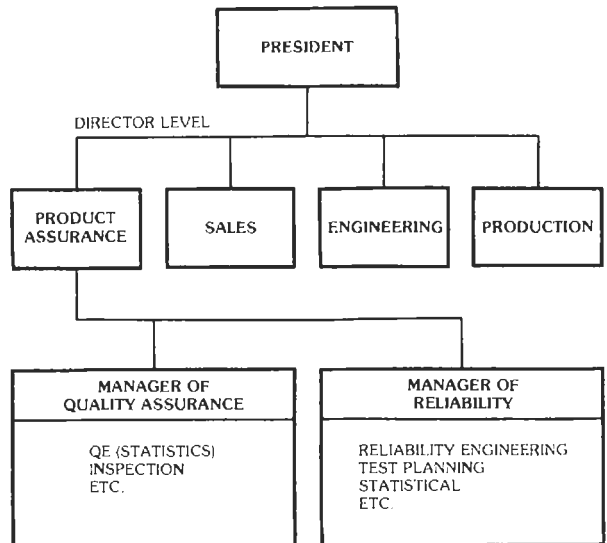


Fig. 1. Typical organization chart showing reliability chain of responsibility.

a proper sequence for its insertion into a master operating plan. For companies without a particular project, this plan usually is discussed with the engineering head and the production head. For corporations working on a program that requires reliability by specified intent, the customer invariably requires a reliability plan. The major items (or task delineators) in a reliability plan are noted in Chart 1.

In the plan itself, the time relationships of the tasks to reliability and overall project tasks are noted as milestones. Some of these milestones are then placed on more sophisticated PERT* or CPM† charts for management review of major effort.

The recognition of reliability as an integral part of product performance is more than simply a matter of early prediction of in-service behavior. It also identifies a number of major economic and safety implications. The realization that failures are inevitable adds an extra dimension to conventional design, consideration being given to the consequences of failures and actions which must be taken to mitigate them. One possible action, for instance, is to provide redundant or standby systems. Others are to give early attention to availability (the percentage of time a system is in a working state) and maintainability (the average time required to repair a system after failure).

Alternatively, potentially catastrophic failures can be prevented by designing in deliberate weak links, such as fuses or breakers that cause the system to fail prematurely in a safe manner ("fail-safe" design).

None of the above actions can be implemented unless product development and the reliability assessment program are closely coordinated. In the case of complex systems, for which high reliability must be achieved, the required performance in service may be impossible to obtain without conscious design for reliability.

ELEMENTS OF RELIABILITY ASSESSMENT

Reliability assessment, rather than occupying one place in product development, is involved in a more or

less continuous basis from the early conceptual stages through to the monitoring of service performance. As a result, reliability analysts may be called on to take part in the following activities:

- Design reviews
- Failure mode analysis
- Quantitative reliability assessment
- Failure data collection and assessment

Design Reviews

Several reviews of a system design are performed, usually by an interdisciplinary team of designers, production engineers, safety and reliability analysts, and others. The objectives are generally to ensure that the product is capable of meeting performance requirements, conforms with necessary standards, optimizes production facilities, and contains no unanticipated flaws. Guidelines for performing a design review have been set out in Reference 2. These emphasize careful documentation and the use of check lists prepared from past experience to ensure coverage of all relevant factors. The list of reliability terms at the beginning of this chapter gives a brief description of three common design reviews.

Failure Mode Analysis

An essential step in controlling reliability is a full investigation of all potential failure mechanisms and their causes, effects, and possible consequences. This is largely a qualitative exercise aimed at finding out what may happen, rather than how often, but it forms an essential foundation for all subsequent quantitative analysis. In fact, there are many instances where it is more important to have all the failure modes identified than to make accurate analyses of some while neglecting others altogether. This task is known by a variety of titles, one alternative being cause and effects analysis, but the most frequently used name is failure modes and effects analysis (FMEA). As in the case of design reviews, FMEA is helped considerably by a systematic approach using check lists and analysis sheets (See Reference 2 for examples.) In addition, the sequential structure of many forms of failure can be used to identify the failure itself, using logic diagrams or

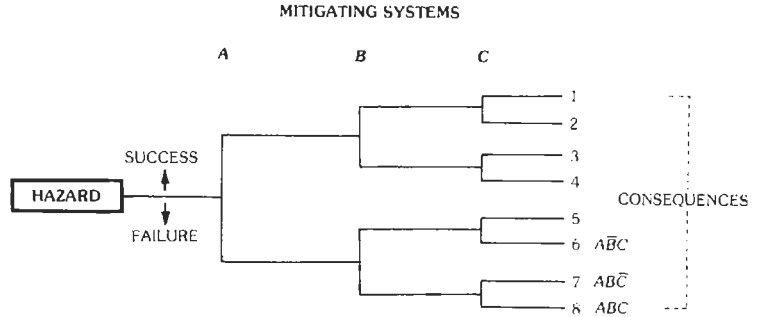
* Program Evaluation and Review Technique.

† Critical Path Method.

CHART 1. TASK DELINEATION

1	Program plan update	11	Change and configuration control
2	Education and manuals	12	Reports and project review
3	Design to specified reliability and maintainability	13	Corrective action control
4	Apportionment	14	Supplier control
5	Models and prediction	15	Manufacturing reliability and maintainability control
6	Cost-effectiveness analysis	16	Failure diagnosis
7	Failure modes and effects analysis	17	Data acquisition and reduction
8	Human factors	18	Verification
9	Stress/strength analysis	19	Summary
10	Design review	20	References

Fig. 2. A simple event tree.



decision trees. Two of the many variants available are shown in Figs. 2 and 3. Fig. 2 illustrates an event tree. This type of diagram is commonly used to map the end results of a sequence of safety device operations. Each branch point represents two possibilities, e.g., "device operates" or "device fails to operate." The event tree starts from an initiating event and progresses forward through subsequent events to a final state which may or may not be a significant failure state for the system. It is capable of pointing to undesirable event sequences which might not be identified by a less systematic approach. Unfortunately, event trees tend to generate many spurious branches and rapidly reach unmanageable proportions. Their main use is in mapping sequences where there is a clear chain of events, for instance, as in certain human operations (see Reference 3).

The most common form of diagram is the fault tree (Fig. 3). This diagram employs a top-down logic, expanding out from a known end result to identify events, or combinations of events, that lead to failure. Fault trees use the same conventions for logic gates that

are used in logic-circuit theory. They are not confined to only two branches at each branch point as is normally the case in event-tree construction. In theory, there is a problem with fault-tree analysis, in that the top event must first be identified by some other means. In practice, it is usually possible to identify failure events at a sufficiently high level to ensure exhaustive coverage of all possibilities. For instance, the top events for reliability analysis of a radio receiver would include:

- Fail to operate at all
- Fail to receive any signal
- Fading or hunting reception
- Drift
- Excessive interference and noise levels

More detailed failure mechanisms can then be identified within each category, on the next lower level of detail. Fig. 4 is a partially developed fault tree for the top event, "fail to operate." To illustrate the technique more fully, it has been assumed that the receiver is powered from a mains (power-line) supply with an automatic switching backup battery. This diagram is by no means complete.

Fault trees can grow rapidly beyond easy manipulation by hand, and, for formal analysis, it is advisable to resort to computer codes, of which several are currently available for fault tree analysis. (See Reference 4 for an extensive summary of currently available software.) On a less formal level, the fault tree is an extremely useful aid in that it imposes a disciplined approach to the problem of failure-mode identification and it provides a structure whereby details of potential failures can be recorded in a sequential manner, without getting lost in the process.

In addition to identifying individual events that lead to failure, a fault tree also contains the logic governing the necessary sequences of events that lead to failure. These sequences are known as "cut sets." In Fig. 5, any one of the combinations A and D, B and C, D and G, or J and B is a cut set. In more complicated problems, the cut sets may be very complex themselves, and not necessarily independent of one another. It can be shown that cut sets can be reduced to an equivalent set of minimal independent requirements for failure called "minimum cut sets." Two simple examples are shown in Fig. 5. In Fig. 5A, four events, A, B, C, and D

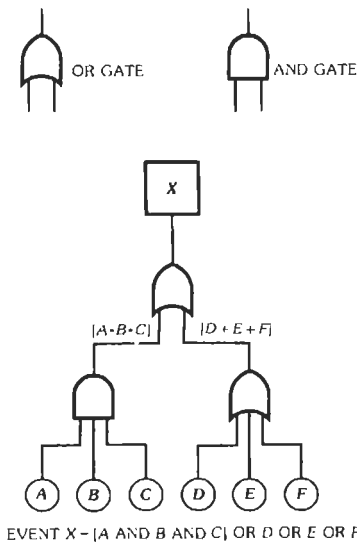
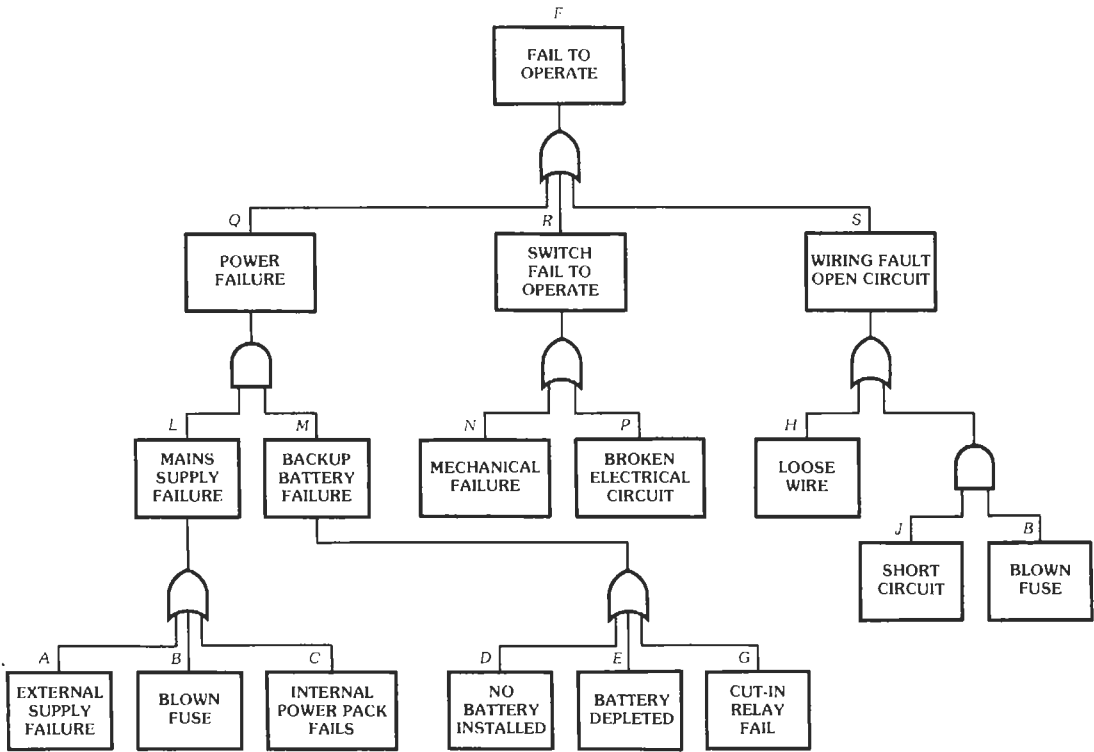


Fig. 3. A simple fault tree.



$$\begin{aligned}
 F &= Q + R + S \\
 &= L \cdot M + N + P + H + J \cdot B \\
 &= (A \cdot B + C) \cdot (D + E + G) + J \cdot B + N + P + H \\
 &= A \cdot (D + E + G) + B \cdot (D + E + G + J) + C \cdot (D + E + G) + (N + P + H)
 \end{aligned}$$

Fig. 4. Fault tree for radio-receiver failure mode "fail to operate on demand."

contribute to some failure event *X*. Combining these events by the rules of Boolean logic algebra, *X* occurs if *A* and *C* or *A* and *D* or *B* and *C* or *B* and *D* occur. These sequences are cut sets and, in this case, minimum cut sets. In Fig. 5B, a similar fault tree contains a common element *B* in both branches. In this case, the top event, *Y*, is only dependent on the occurrence of *A* and *C* or *B*. Here, *B* represents an important class of failures called "common mode" or "common cause" failures. It can be easily seen that if *A*, *B*, *C*, and *D* all have similar probabilities of occurrence, *Y* has a much higher probability of occurrence than *X*, because it only requires a single event instead of two simultaneous events. Failure to identify common modes or dependencies between different failure mechanisms is a frequent cause of excessively optimistic reliability predictions. It can be seen from the above simple example how fault tree analysis can be used to avoid this problem.

Quantitative Reliability Assessment

Quantitative reliability assessment is one of the most important tasks of the reliability analyst. Reliability quantification covers two dimensions—probability of

occurrence and magnitude of the consequence of failure. Assessment of failure consequence is probably more the task of design. Nevertheless, it cannot be neglected in reliability assessment, particularly when assigning reliability targets, or allocating reliability to subsystems. Intuitively, it seems sensible to assign a higher reliability target to a subsystem with potentially catastrophic failure consequences than to one that involves only a degradation in performance.

Consequences of failure can be extremely difficult to estimate. In the preliminary stages of development, a useful technique is called worst-case analysis. It takes the most pessimistic values of all parameters, regardless of the likelihood of this occurring. In electronic circuits, for instance, all components are assumed to be at the extreme limits of their tolerance bands, including any drift due to environmental factors. By this procedure, it is possible to identify undesirable modes of operation that might not be considered possible otherwise. The detection of potentially catastrophic modes of failure is probably the most useful outcome of this type of analysis. Regardless of their likelihood, failures of this type may require designing out of the system completely. It is only in those cases where the cost of failure can be reasonably balanced against the cost of

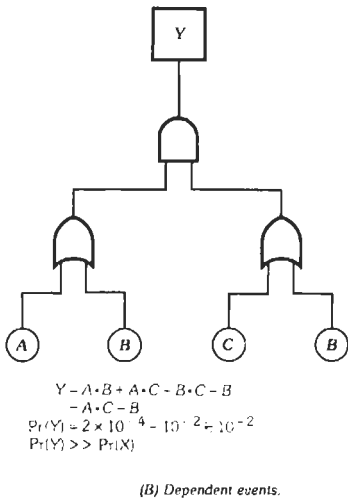
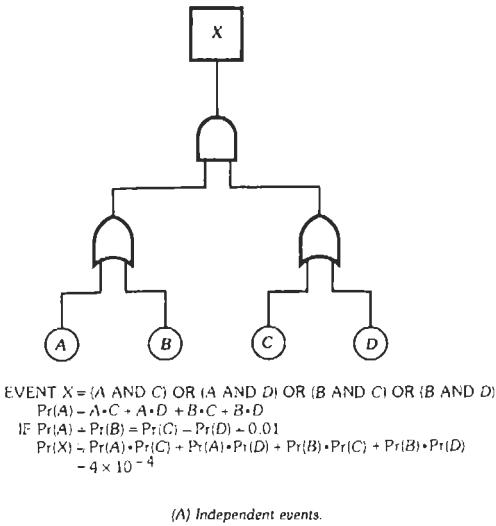


Fig. 5. Calculation of event probability.

improved design that the question of probability of occurrence becomes important.

The area of quantification most closely identified with the reliability analyst is the assessment of failure probability. This is a natural result of defining reliability in probabilistic terms. A wide range of mathematical techniques have been developed in the general fields of probability and statistics, supplemented with methods specifically directed toward reliability modeling. These methods will be reviewed in a special section later in this chapter.

Failure Data Collection and Assessment

The collection of failure data is part of a continuous feedback loop that is essential to reliability growth. The

process begins in the early stages of product development, where failures are likely to be common due to design faults and as yet uncontrolled production processes, and extends through to collection of data under normal service conditions. The reliability analyst should be closely involved in this process. In addition to receiving basic failure-rate data from this source, the reliability analyst also has an active role in evaluating failure reports, assessing whether the failures are random or represent a significant deviation from expected conditions, providing backup to quality control, and offering recommendations for remedial action.

For failure analysis to be successful, the data collection system must be carefully designed and controlled. Considerable thought has been given to this problem in recent years, and a number of reporting systems have been developed. A comprehensive review of data banks and their methods of operation is contained in Reference 4. Some of these are appropriate models for in-house reporting systems. This subject will be discussed further in a later section of this chapter on sources of data. One reporting system of particular relevance to integrated circuit reliability was developed by Lockheed Corporation on behalf of NASA (Reference 5). This report is both an excellent source of information and a model for setting up an in-house reporting and analysis system.

COMPONENT RELIABILITY

A population of components on test or in service under nominally similar conditions will decrease in time as a result of failures as shown in Fig. 6. The numbers failed and surviving after time t are $N_f(t)$ and $N_s(t)$, respectively. For a sufficiently large population number N , the relative frequencies represent the cumulative probability of failure and the probability of survival (i.e., the reliability), respectively:

$$\hat{F}(t) = N_f(t)/N$$

$$\hat{R}(t) = N_s(t)/N$$

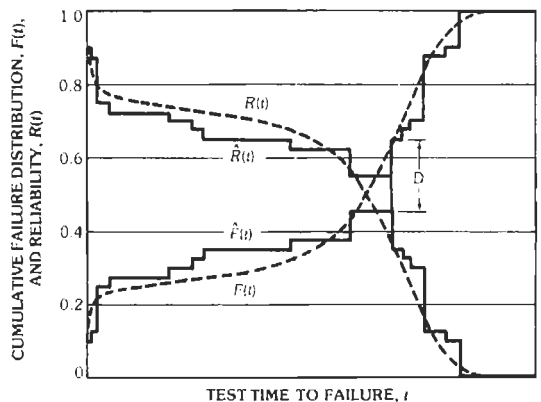


Fig. 6. Cumulative failure distribution and reliability plots for a test with 40 samples—experimental and theoretical curves shown.

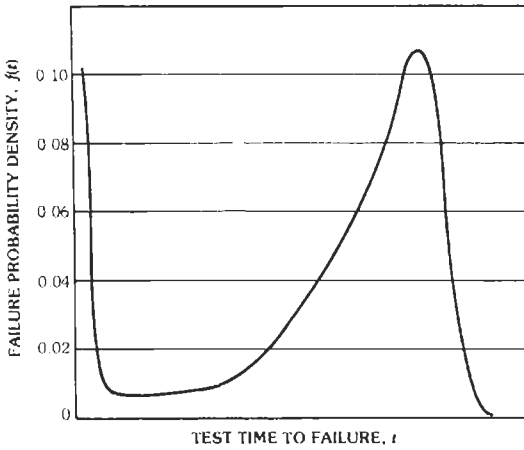


Fig. 7. Failure probability density derived from Fig. 6.

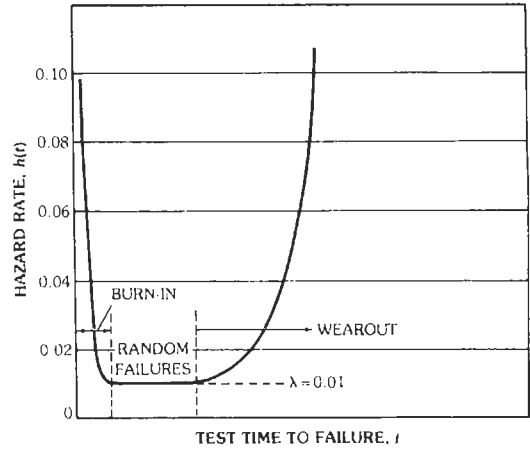


Fig. 8. Hazard rate for data given in Fig. 6.

As N tends to infinity, estimates $\hat{F}(t)$ and $\hat{R}(t)$ tend to the true failure and reliability functions, $F(t)$ and $R(t)$. Here, and subsequently in this chapter, a "hat" (^) will be used to represent an estimate of a random variable. These functions are shown also in Fig. 6. By differentiating $F(t)$, the probability density function of failures, $f(t)$, is obtained (Fig. 7). This represents the relative frequency of failures with respect to the original population. Three zones can be identified in this figure. They are an initial period of high failure density, corresponding to initial manufacturing defects, an intermediate period of relatively low failure rate, where the only cause is external and random, and a final peak of high density corresponding to wearout or material deterioration. These zones can be more readily distinguished by calculating the instantaneous failure rate, or hazard rate, $h(t)$, which is the failure rate expressed as a proportion of the population surviving at time t .

$$h(t) = f(t) / [1 - F(t)]$$

The hazard rate has a characteristic shape, illustrated in Fig. 8, known as the "bathtub curve." Here, initial, or burn-in, failures; random failures; and wearout failures are clearly identified by having decreasing, constant, and increasing hazard rates, respectively. This characteristic of the hazard rate is an important diagnostic tool, as will be demonstrated later in this chapter.

It can be shown (see Reference 4 for details) that

$$F(t) = 1 - \exp \left[- \int_0^t h(t) dt \right]$$

$$R(t) = \exp \left[- \int_0^t h(t) dt \right]$$

Knowing the hazard rate, e.g., from life tests, permits the reliability function, $R(t)$, to be calculated.

In electronic components, it is common practice to require a burn-in test period that, as far as service performance is concerned, effectively removes the initial stage of the hazard rate. It is also generally assumed that stress levels have been so chosen in design that there is little likelihood of the occurrence of wearout failures. Under these conditions, it is reasonable to use a restricted version of reliability theory that assumes a constant hazard rate and is relatively easy to manipulate mathematically. In fact, the majority of formal reliability theory is based on the constant-hazard-rate model (sometimes referred to loosely and erroneously in much of the literature as a "constant-failure-rate" theory). In this model,

$$h(t) = \lambda = \text{constant}$$

$$R(t) = 1 - F(t) = \exp(-\lambda t)$$

This is the so-called exponential distribution, a special case of the Poisson distribution.

In many cases, one is interested only in high component reliabilities, in which case there is a simple approximation. That is,

$$\text{if } \lambda t \ll 1, 1 - F(t) = R(t) = 1 - \lambda t$$

Where the simple models discussed above are not valid, other more complex distributions must be used. Table 1 summarizes some of the more commonly used distributions. In choosing a particular distribution, there is sometimes a degree of rationale, based on the physics of the problem. More often, however, the choice is empirical so as to best fit the available data. One distribution of predominant interest in reliability studies is the Weibull distribution (see Table 1). This distribution has some justification in that it is one of a class of extreme distributions, i.e., distributions of extreme values of groups of components. Its main

TABLE I. LIST OF COMMON FREQUENCY DISTRIBUTIONS

Name	Mass/Density Function	Mean	Variance
Binomial	$P_X(x) = \binom{n}{x} p^x (1-p)^{n-x}$ $x = 0, 1, 2, \dots$	np	$np(1-p)$
Poisson	$p_X(x) = (\theta^x / x!) e^{-\theta}$ $x = 0, 1, 2, \dots$	θ	θ
Normal	$f_X(x) = (1/\sqrt{2\pi\sigma^2}) \exp\{-(1/2\sigma^2)(x - \mu)^2\}$ $-\infty < x < \infty$	μ	σ^2
Lognormal	$f_X(x) = (1/\sqrt{2\pi\sigma x}) \exp\{-(1/2\sigma^2)(\ln x - \mu)^2\}$ $x > 0$	$\exp(\mu + \sigma^2/2)$	$\exp(2\mu + \sigma^2) [\exp(\sigma^2) - 1]$
Student-t	$f_X(t) = [\Gamma(m/2 + 1/2) / \sqrt{m\pi} \Gamma(m/2)] (1 + t^2/m)^{-(m+1)/2}$ $-\infty < t < \infty$	0	$m/(m-2)$
Gamma	$f_X(x) = \{\theta^n / \Gamma(n)\} x^{n-1} e^{-\theta x}$ $x > 0$	n/θ	n/θ^2
Exponential	$f_X(x) = \theta e^{-\theta x}$ $x > 0$	$1/\theta$	$1/\theta^2$
Weibull	$f_X(x) = \beta(x^{\beta-1} / \eta^\beta) \exp[-(x/\eta)^\beta]$ $x > 0$	$\eta \Gamma(\beta + 1) / \beta$	$\eta^2 \{\Gamma(\beta + 2) / \beta - [\Gamma(\beta + 1) / \beta]^2\}$
Chi-Square	$f_X(x) = [1/2^{n/2} \Gamma(n/2)] x^{n/2-1} e^{-x/2}$ $x > 0$	n	$2n$

attribute is its ability to represent any one of the three zones in the bathtub curve (Fig. 8) by choosing the shape parameter, β , to be less than, equal to, or greater than unity (see Fig. 9). This property of the Weibull distribution will be shown later to be a useful method of evaluating service data.

SYSTEM RELIABILITY

System reliability can be dealt with in two different ways, depending on the advancement of the project.

In the early stages of making a bid, the eventual configuration will not be known precisely, but a reliability estimate is required for the purpose of making the proposal. This estimate is inevitably crude and must draw on previous experience and generic data for its quantitative basis. The simplest approach, once the number and types of components are known approximately, is to assume the individual failure rates to be additive:

$$\lambda_s = \sum \lambda_i$$

The system reliability, $R_s(t)$, is obtained as

$$R_s(t) = \exp(-\lambda_s t)$$

This simple approach can also be used in reverse to assign component reliabilities so as to achieve a target system reliability that may have been set as a design requirement. Illustrating this simple method, if an assembly is made up of two components of equal weight and is to have an end reliability of 0.99, then each component should have a reliability of at least 0.995. The product $0.995 \times 0.995 = 0.990025$. If the assembly is made up of three units of equal weight, then each should have a reliability of 0.997, since $(0.997)^3 = 0.991026973$. The same procedure is used even if the

parts have widely different importance and hence extremely different reliability values.

Formalized variants of this approach are given in the AGREE report (Reference 6) and in a more recent version in Military Handbook MIL-HDBK-217B (Reference 7).

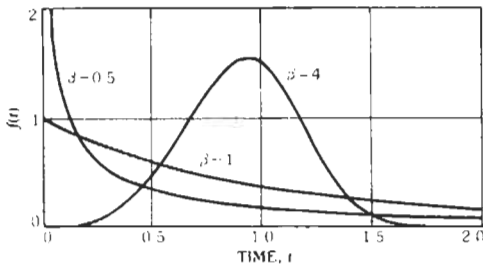
The AGREE report uses a more complex method, detailed on pages 192-195 of Reference 18. The AGREE report cautions against the use of its method for items of low importance. Its basis is the complexity of the units rather than their failure rates. Consideration is given to the number of hours the j th unit will be required to operate in T system hours and the total number of modules in the system. The importance factor of each unit must also be considered. All units do not operate the same length of time; hence, these times of operation must be considered, and many of these t_j values will be less than T hours. The AGREE report gives the allocated rate of the j th unit as

$$\tilde{\lambda}_j = n_j [-\log R^*(T)] / E_j t_j N$$

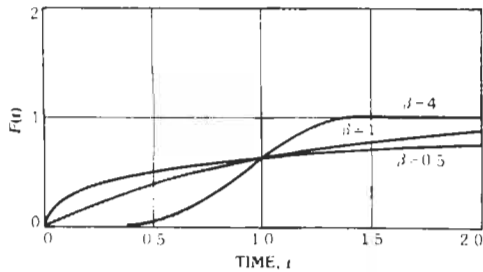
This is used in many instances to allocate reliabilities among the components.

MIL-HDBK-217B describes two methods of failure-rate prediction, "part stress analysis" and "parts count." The latter is the simpler of the two, and more appropriate to early design and bidding. For full explanations, reference is made to the handbook itself and to an expanded explanation in Reference 1. A parts count failure rate prediction lumps components into generic part types and assigns part quality levels to each type. The inputs are therefore

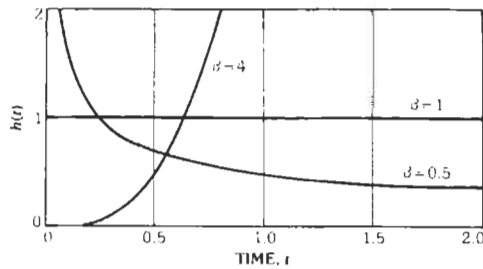
- Generic part types
- Quantities of parts in each type
- Part quality levels
- Equipment environment



(A) Probability density function.



(B) Probability distribution.



(C) Hazard function rate or instantaneous failure rate.

Fig. 9. The Weibull distribution, $F(t) = 1 - \exp[-(t/\eta)^\beta]$.

The total failure rate, λ_s , is then calculated from

$$\lambda_s = \sum_i (\lambda_{pc})_i$$

where $(\lambda_{pc})_i$ is the i th part category failure rate:

$$\lambda_{pc} = n \cdot \lambda_g \cdot \pi_Q$$

where,

- n = number of parts in category i ,
- λ_g = generic failure rate, provided in MIL-HDBK-217B, for a range of components, taking into account environmental factors such as "ground, benign," "space, flight," "naval,

unsheltered," up to the most severe, "missile, launch,"

π_Q = quality factor, specifying the standard of selection, e.g., "above spec.," "MIL-spec.," "below MIL-spec."

EXAMPLE 1. A preliminary design of an RC network for a ground mobile (G_M) environment gives the information shown in Table 2.

$$(\lambda_{pc})_1 = 32 \times 0.042 \times 10^{-6} \times 0.1/\text{hr} = 0.1344 \times 10^{-6}/\text{hr}$$

$$(\lambda_{pc})_2 = 5 \times 0.11 \times 10^{-6} \times 0.1/\text{hr} = 0.055 \times 10^{-6}/\text{hr}$$

$$(\lambda_{pc})_3 = 15 \times 0.044 \times 10^{-6} \times 0.1/\text{hr} = 0.066 \times 10^{-6}/\text{hr}$$

$$\lambda_s = 0.2554 \times 10^{-6}/\text{hr}$$

EXAMPLE 2. A system consisting of 1000 units as described in Example 1.

$$\lambda_{\text{system}} = 1000\lambda_s = 2.554 \times 10^{-4}/\text{hr}$$

The mean time between failures is

$$\text{MTBF} = 1/\lambda_{\text{system}} = 3915 \text{ hr}$$

The reliability at 1000 hours is

$$R(1000) = \exp(-2.554 \times 10^{-4} \times 10^3) = 0.775$$

The above approach assumes implicitly that all components in a system are connected in series and that a single failure causes system failure (weakest-link theory). This is conservative and unrealistic in large systems where backups are necessary if the system is to work at all. A more accurate assessment can be made once the system configuration has been finalized. In general, this may contain parallel, or redundant, elements in addition to series elements. A number of examples of redundant systems are given in Fig. 10. Note particularly that a series system is less reliable than its components, whereas a parallel system is more reliable. Series connections are generally forced on the designer by functional requirements, and redundancy is one method of retaining an acceptable level of reliability. Other techniques, such as standby and voting logic systems, are also possible (see Reference 4 for further details).

In practice, reliability of a system is not simply a function of the inherent failure rates of its components. The system will be subject to repair and scheduled maintenance. While periodic repair ensures that the reliability does not decay continuously to zero, the system will not be available on demand while it is being repaired. For complex systems with repair, therefore,

TABLE 2. DATA FOR EXAMPLE 1

Part Count	Generic Part Type*	Number of Parts	Generic Failure Rate, $\lambda_G/10^6 \text{ hr}^{-1}$	Quality Factor, π_Q
1	Film Resistor RNR 55182	32	0.042	0.1 (R)
2	Wire Resistor RWR 35007	5	0.110	0.1 (R)
3	Capacitor CKR 39014	15	0.044	0.1 (R)

* Data from MIL-HDBK-217B

it is more convenient to think in terms of "availability," A.

Product availability is defined as the probability that the system will operate satisfactorily at any point in time, where time includes not only operating life, but also active repair time and administrative and logistic time. An equation for availability is

$$A = \text{MTBF}/(\text{MTBF} + \text{MTTR})$$

where,

- A = availability,
- MTBF = mean time before failure,
- MTTR = mean time to repair.

The calculation of MTTR is related to repair hours, while the calculation of MTBF is related to component operating hours. Fig. 11 is a graph of the above equation.

It is evident from the figure that the effect of maintainability on availability increases as the ratio of MTBF to MTTR decreases. If an item has an inherently low MTBF, the MTTR must be very low to sustain a good level of availability.

In the design of any complex system, an optimum relationship should be established between reliability and maintainability, so that reliability is not increased beyond the point where very little availability gain is obtained because of lack of consideration of the effect of maximum maintainability.

To look at this another way, manufacturers of micro-electronic devices claim that MTBF is very high and repair time is nil or very low. Look again at Fig. 11 and the ratio of MTBF/MTTR in the region to the right of 100. Little is to be gained by designing a module that can be repaired extremely quickly, if to do the job special tools and costs are involved. In other words, the throwaway concept in this case is clearly justifiable. NOTE: There is some ambiguity in the use of the terms MTBF and MTTF. Some references use the same definitions as accepted in this chapter. Others, such as Reference 4, use MTTF as the inverse of failure rate, with MTBF to be given by

$$\text{MTBF} = \text{MTTF} + \text{MTTR}$$

Hence, by this convention, $A = \text{MTTF}/\text{MTBF}$

SOURCES OF RELIABILITY DATA

There are two sources of failure-rate data:

- Generic data banks
- Experimental determination by life testing

Generic Data Sources

Generic data are usually used for initial design and evaluation, whereas life testing is more appropriate to the latter stages of system development, where questions such as reliability growth and system substantiation are important.

Data banks and other sources of failure-rate data have grown considerably in recent years. A comprehensive review of existing data sources is given in Reference 4. Sources which are readily available are MIL-HDBK-217B (Reference 7), the results of the IEEE Project 500 on failure rates for electronic components (Reference 8), the Government/Industry Data Exchange Program (GIDEP) (Reference 9), and the Failure Data Handbook for Nuclear Facilities, available from NTIS (Reference 10). More detailed information can be obtained by subscribing to one of several data banks, such as the Nuclear Plant Reliability Data System (NPRDS) (Reference 11) or the System Reliability Service Data Bank (SYREL) (Reference 12), operated by the United Kingdom Atomic Energy Authority.

Life Testing

If it is possible to perform adequately, life testing of actual components is a more accurate method of determining failure rates, as long as testing conditions are chosen to simulate actual service conditions as realistically as possible.

Two basic methods are used in securing reliability data. In the first method, assume that 20 units are placed on test with the stipulation that each unit be operated 5000 hours, which requires over 200 days of continuous operation. During this period of operation, the time at which each of five failures occurs is recorded. If they occur at 4200, 4350, 4400, 4750, and 4900 hours, and if the test is stopped at the end of 5000 hours, the total time of operation is $100\ 000 - (800 + 650 + 600 + 250 + 100) = 97\ 600$ hours. If the

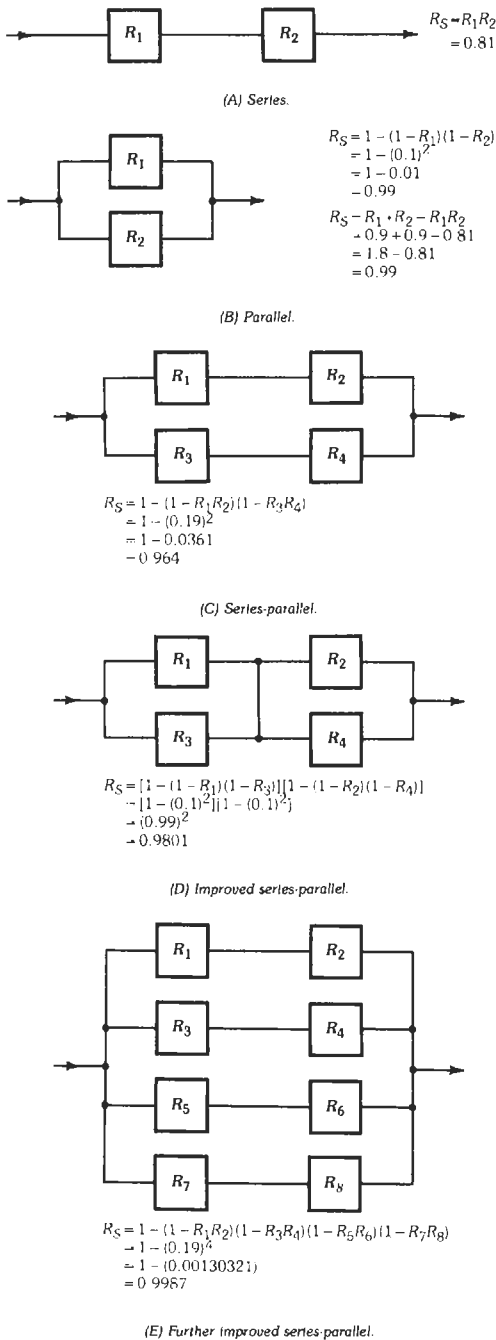


Fig. 10. Diagrams of redundant systems with relationships for computing system reliabilities. All parts are assumed to have equal reliability ($R_C = 0.90$).

technician making the test wishes to obtain the total operating time of 100 000 hours, then the remaining 15 units must operate $2400/15 = 160$ hours each more

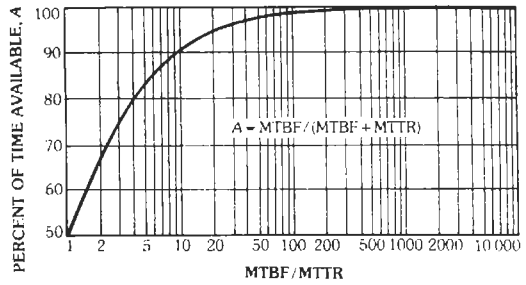


Fig. 11. Product availability.

than the original 5000 hours specified with no additional failures. Otherwise, the estimate of the failure rate should be $r = 5/97\ 600 = 5.12/100\ 000 = 5.12\%/1000$ hours. The value of $m = \text{MTBF}$ is $m = 97\ 600/5 = 19\ 520$ hours.

The second basic method for determining the failure rate is to operate all 20 units until they fail. This makes it difficult to plan a testing program, as some equipment may be tied up for two or even several times the period designated for the original test. However, failure rates found by the second method should be more representative of actual field failure rates. The data secured under these conditions for the 20 units might be as given in Chart 2.

This gives a failure rate much larger than the first estimate, i.e., $r = 20/153\ 200 = 13.05\%/1000$ hours. If the 5% value is near the actual central value, then the total operating time should be about 400 000 hours, or 20 000 hours (833 days) on the average for each unit.

It is desirable to establish confidence levels for reliability values. A distribution of probability values is used to give either a one-sided confidence value or a two-sided confidence value with both upper and lower limits. When determining failure rates, a one-sided level applies. A confidence level of 0.95 might be stipulated for a reliability of 0.95 or better (the goal stipulated in the contract). In the first set of data, the rate of 5%/1000 hours may be expressed as a reliability of 95% where the mission time is 1000 hours. If the mission time is 10 hours, then the corresponding reliability is 99.95%. These values are associated with expected values rather than having a range of values associated with a designated confidence-level value. Where a confidence limit is used, a limit is established that is associated with this limit. It will be expressed as a failure rate or an MTBF value.

Fixed time testing is covered in a useful publication of the Department of Defense (DOD), Handbook H108 (Reference 13). Both fixed and sequential testing are covered by a military standard, MIL-STD-781B (Reference 14). These references and their applications are well illustrated in Reference 1.

CHART 2. TIME IN HOURS TO FAILURE FOR RELIABILITY TESTS ($n = 20$ UNITS)

	6500		8000		10200		5500	
	4200		7000		4750		8900	
	9400		4350		4900		10500	
	4400		9100		8750		9150	
	7800		8100		9200		12500	
Sum =	32300	+	36550	+	37800	+	46550	= 153 200

PROBABILITY AND STATISTICAL INFERENCE

Confidence Limits

If 100 components are tested and two fail, the failure probability can be estimated approximately as

$$\hat{P}_F = 0.02$$

The corresponding reliability is therefore

$$\hat{R} = 1 - P_F = 0.98$$

However, these are only estimates based on a single test. If the test is repeated, a different estimate will be obtained, as shown in Table 3, which summarizes 20 such tests, each of 100 components. Empirically, it can be determined that the mean observed reliability is $\bar{R} = 0.975$, with a standard deviation of $\sigma_R = 0.01775$. By fitting a theoretical distribution such as a normal or Poisson distribution to this data, it is possible to determine the probability of the number of failures in some subsequent batch exceeding some value, say 6. The relative frequencies of different observed reliability estimates have been calculated in Table 4 using a normal distribution with sample mean and standard deviation, and using a Poisson distribution with the same observed overall failure rate, 0.025. Predictions are shown in Fig. 12. It appears that the normal distribution gives the better fit to the data.

This problem shows that the true reliability may be different from any observed estimate. However, it is possible to use the data summarized in Table 3 to calculate limits that will contain the true value a specified percentage of the time. These limits are known as *confidence limits*, and the percentage is called the *confidence level*.

In practice, the observed values of \hat{R} will be dispersed around the true reliability, R , but R is unknown. If it can be assumed that R is normally distributed with a true standard deviation equal to the estimate from the 20 tests, then it can be shown that the mean of n tests, \bar{R} , will be normally distributed with a standard deviation $\sigma_{\bar{R}} = \sigma_{\hat{R}}/\sqrt{n}$, where $n = 20$ in this case; furthermore \bar{R} tends to the true reliability as n increases. If we imagine many sets of n tests to be performed, the results will be distributed around the true (un-

known) reliability R , as shown in Fig. 13. The value $\bar{R} = 0.975$ obtained from Table 3 is also shown.

The probability that \bar{R} is greater than some amount L larger than R is

$$\Pr(\bar{R} > R + L) = \int_{R+L}^{\infty} n(R; \sigma_{\hat{R}}/\sqrt{n})dR$$

where $n(R; \sigma_{\hat{R}}/\sqrt{n})$ is the normal distribution with mean R and standard deviation $\sigma_{\hat{R}}/\sqrt{n}$.

This can be expressed more conveniently in terms of the standard normal variate, Z_X .

$$Z_X = (X - R)/(\sigma_{\hat{R}}/\sqrt{n}), Z_L = L/(\sigma_{\hat{R}}/\sqrt{n})$$

$$\Pr[(\bar{R} - R)/(\sigma_{\hat{R}}/\sqrt{n}) > Z_L] = \int_{Z_L}^{\infty} n(0;1)dZ$$

If as is usual \bar{R} is known and R is unknown, we need to determine the probability that R is less than \bar{R} by some amount L . From the previous equation

$$\Pr(R < \hat{R} - Z_L/\sigma_{\hat{R}}/\sqrt{n}) = \int_{Z_L}^{\infty} n(0;1)dZ = \alpha$$

The quantity α can be expressed in terms of the error function that is tabulated in any standard text for the standardized normal distribution, $n(0;1)$.

The quantity $(\bar{R} - L)$ is the lower confidence limit (LCL) for a confidence level of $(1 - \alpha)$. Table 5 gives a short list of α against standardized L , derived from tables of the error function.

EXAMPLE 2 (continued). Find the 97.5% LCL for the reliability, $R_{L,97.5}$.

$$\alpha = 0.025$$

Hence,

$$L/(\sigma_{\hat{R}}/\sqrt{n}) = 1.96 \text{ (from Table 5)}$$

$$\sigma_{\hat{R}}/\sqrt{n} = 0.01775/\sqrt{20}$$

$$L = 1.96 \times 0.01775/\sqrt{20}$$

$$= 7.77 \times 10^{-3}$$

$$R_{L,97.5} = 0.975 - 7.77 \times 10^{-3} = 0.967$$

TABLE 3. RESULTS OF 20 LIFE TEST GROUPS OF 100 COMPONENTS EACH

Test	Units on Test n_i	Failures d_i	Successes s_i	Reliability $R_i = s_i/n_i$
1	100	0	100	1.000
2	100	1	99	0.990
3	100	3	97	0.970
4	100	0	100	1.000
5	100	2	98	0.980
6	100	4	96	0.960
7	100	3	97	0.970
8	100	5	95	0.950
9	100	1	99	0.990
10	100	2	98	0.980
11	100	0	100	1.000
12	100	0	100	1.000
13	100	3	97	0.970
14	100	6	94	0.940
15	100	4	96	0.960
16	100	5	95	0.950
17	100	2	98	0.980
18	100	2	98	0.980
19	100	3	97	0.970
20	100	4	96	0.960
Total	2000	50	1950	(0.975)

$$m = 20, n = n_i = 100, \sum d_i = 50, \sum n_i = 2000.$$

$$R = (\sum n_i - \sum d_i) / \sum n_i = \sum s_i / \sum n_i$$

Average: $\bar{R} = 19.500/20 = 0.9750$

$$\sigma_R^2 = (\sum f_j R_j^2 / \sum f_j) - \bar{R}^2 = (19.0188/20) - (0.975)^2 = 0.950940 - 0.950625 = 0.000315.$$

Also $\sigma_R^2 = \sum f_j (\Delta R_j)^2 / \sum f_j$ (where $\Delta R_j = R_j - \bar{R}$) and $\sigma_R = 0.017748$; hence

$$\sigma_R^2 = \frac{1}{20}(0.001225 + 0.001250 + 0.000675 + 0.000100 + 0.000100 + 0.000450 + 0.002500)$$

$= \frac{1}{20}(0.006300) = 0.000315$, or using relative weights, $w_j f_j = 1.00$, $\sigma_R^2 = w_j f_j (\Delta R_j)^2 = 0.000315$, given in last column above. Hence $\sigma_R = (0.000315)^{1/2} = 0.017748$.

Similarly, an upper confidence level can be obtained by postulating \bar{R} to be to the left of R (see Fig. 14). Hence,

$$R_{U,97.5} = 0.975 + 7.77 \times 10^{-3} = 0.983$$

The above two quantities represent the one-sided LCL and UCL at the $(1 - \alpha)$ confidence level, respectively. Taken together, they constitute a two-sided confidence band at the $(1 - 2\alpha)$ confidence level. That is,

$$0.967 < R < 0.983 \text{ at } 95\% \text{ confidence level}$$

When the standard deviation of the population is unknown, it is no longer valid to use the above approach

based on a normal distribution. Instead, the student- t distribution must be used. The t distribution is flatter than the normal and gives wider confidence limits for the same data.

Confidence limits are calculated from the following relation:

$$UCL, LCL = R \pm t_\alpha S / \sqrt{n}$$

where,

t_α = confidence coefficient for level $(1 - \alpha)$ (See Table 6),

S = sample standard deviation,

n = number of components in sample,

r = number of degrees of freedom ($= n - 1$).

TABLE 4. CALCULATION TO FIT THEORETICAL DISTRIBUTIONS TO DATA OF TABLE 3
Derivation of Normal Law Theoretical Probabilities

Boundary Values	Deviations from \bar{R}		Normal Law Probabilities Corresponding to z	Reliability Values	Theoretical Frequency		
	Numerical	z Values σ_R Units			Individual		Cumulated No.
					Prob.	No.	
0.995	+0.020	+1.127	0.3701	1.000	0.1299	3	3
0.985	+0.010	+0.5634	0.2134	0.990	0.1567	3	6
0.975	0	0	0	0.980	0.2134	4	10
0.965	-0.010	-0.5634	0.2134	0.970	0.2134	4	14
0.955	-0.020	-1.127	0.3701	0.960	0.1567	3	17
0.945	-0.030	-1.690	0.4545	0.950	0.0844	2	19
				0.940	0.0455	1	20

Observed and Theoretical Individual and Cumulated Values for Failure Rate $r = 0.025$ Exponential and Normal Laws								
Observed Values			$m = 20$ Sets					
Reliability Values	Ind.	Cum.	Experimental, $rn = 2.5$				Normal Law	
			Probabilities		Number		Individual	Cumulated
			Ind.	Cum.	Ind.	Cum.		
1.00	4	4	0.082	0.082	2	2	3	3
0.99	2	6	0.205	0.287	4	6	3	6
0.98	4	10	0.257	0.544	5	11	4	10
0.97	4	14	0.214	0.758	4	15	4	14
0.96	3	17	0.134	0.892	3	18	3	17
0.95	2	19	0.067	0.959	1	19	2	19
0.94	1	20	0.041	1.000	1	20	1	20
Sum =	20				20		20	

EXAMPLE 3. Same problem as example 1, but with no assumption on standard deviation. The number of components in the sample is 20; therefore the number of degrees of freedom is 19. From Table 6, $t_{0.025} = 2.09$.

$$UCL = 0.975 + 2.09 \times 0.01775/\sqrt{20} = 0.9833$$

$$LCL = 0.975 - 2.09 \times 0.01775/\sqrt{20} = 0.9667$$

There is not much difference between this result and the previous one. This is because the sample size is relatively large. In general, the t distribution need only be used if the sample is smaller than 30.

In calculating confidence limits for other quantities, such as the MTBF, neither the normal nor t distribution is valid because the distribution of the random variable differs too much from a normal distribution (exponential in the case of MTBF). In the case of the exponential distribution, the appropriate relation for determining confidence limits is the chi-squared distribution.

When it is desired to specify that the true mean time between failures must exceed a given minimum value with a confidence level of $(1 - \alpha)$, the procedure for a one-sided confidence limit is applied. This provides a

tail area α and means there is a probability α that the m value actually observed by test will be smaller than the specified minimum and a probability of $1 - \alpha$ that it will be larger. Reference 15 denotes the one-sided confidence limit by the notation C_L to distinguish it from the two-sided lower limit L . Its value is given by

$$C_L = (2r/\chi^2_{\alpha;2r})\hat{m} = 2T/\chi^2_{\alpha;2r}$$

where,

$\chi^2_{\alpha;2r}$ is the value of χ^2 at the $(1 - \alpha)$ confidence level for $2r$ degrees of freedom.

tests are continued until the r th failure occurs with $r = 1, 2, \dots, d$,

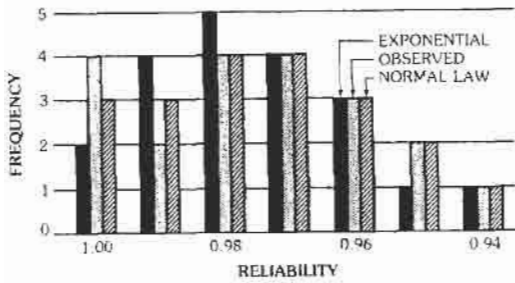
$T =$ accumulated test time $= \sum t_i$,

$m = T/r =$ an estimate of the mean time between failures,

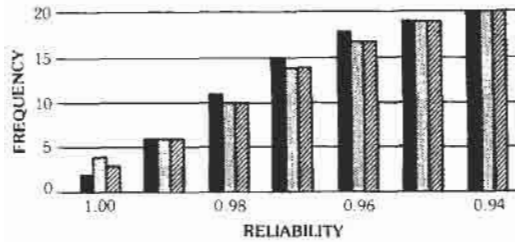
$1 - \alpha =$ confidence level prescribed.

Note that in this case $2r =$ degrees of freedom (d.f.).

However, a test can also be terminated at some preselected test time without a failure occurring exactly at that time. For such a case, Reference 16 has shown that for the accumulated hours of operating time $T = \sum t_i$, then



(A) Individual failures.



(B) Cumulated failures.

Fig. 12. Observed reliability distributions of individual and cumulated failures for 20 sets with $n = 100$ on test for 1000 hours each. The corresponding theoretical frequencies for exponential and normal law are also shown. Failure rate $r = 0.025$ and $m = 2.5$ for exponential; $\hat{R} = 0.975$ and $\sigma_R = 0.01775$ for normal law.

$$m \geq 2T/\chi^2_{\alpha, 2r+2}$$

where d.f. = $2r + 2$ and the case where $r = 0$ is covered. For $r = 0$, then

$$C_L = 2T/\chi^2_{\alpha, 2}$$

In the percent survival method, the accumulated operating time T is not measured, and only the straight

TABLE 5. SHORT LIST OF CONFIDENCE LEVELS VERSUS STANDARD NORMAL VARIATE FOR THE NORMAL DISTRIBUTION

Confidence Level, ($1 - \alpha$)	Standard Normal Variate, Z
0.900	1.282
0.925	1.440
0.950	1.645
0.960	1.751
0.965	1.812
0.970	1.881
0.975	1.960
0.980	2.054
0.985	2.170
0.990	2.326
0.995	2.580

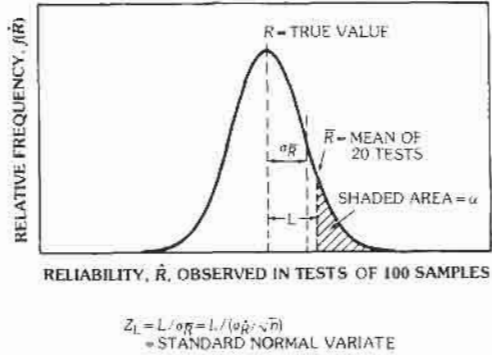


Fig. 13. Calculation of lower confidence limit when standard deviation of sample is known.

test duration time t_d is known, at which time r failures of n units on test are counted. In this method, confusion may exist between chance failures and failures due to actual wearout. The time to wearout must be known, and it is necessary to design and select parts from manufacturers that can be made so that their respective wearout time is many hours past the time of the mission. Again referring to both Reference 15 and to Epstein, for a one-sided confidence level of $1 - \alpha$, the lower-limit estimated reliability for t_d hours is

$$\hat{R}(t_d) = \frac{1}{1 + [(r + 1)/(n - r)]F_{\alpha, 2r+2, 2n-2r}}$$

where F is the upper α percentage point of the Fisher distribution (termed the F distribution) with the two corresponding degrees of freedom, $2r + 2$ and $2n - 2r$. For this estimate of reliability there is a probability of $1 - \alpha$ that the true reliability for t_d hours is equal to or larger than $\hat{R}(t_d)$. It must be noted that this reliability estimate is nonparametric and is valid for the exponential as well as the nonexponential case.

A general mathematical approach is used in many cases to determine the confidence levels for either one-sided or two-sided distributions for various density functions. Confidence levels and reliability values are

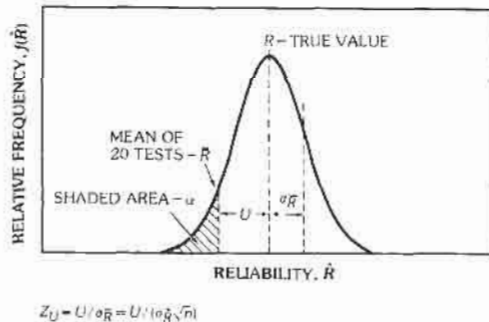


Fig. 14. Calculation of upper confidence limit when standard deviation of sample is known.

TABLE 6. SHORT LIST OF STUDENT-*t* DISTRIBUTION

<i>r</i>	<i>t</i> _{0.1}	<i>t</i> _{0.05}	<i>t</i> _{0.025}	<i>t</i> _{0.01}
1	3.078	6.314	12.706	31.82
2	1.886	2.920	4.303	6.965
4	1.533	2.132	2.776	3.747
9	1.383	1.833	2.262	2.821
19	1.328	1.729	2.093	2.539
29	1.311	1.699	2.045	2.462
∞	1.282	1.645	1.96	2.326

related by the two following general relations, where *P_b* = degree of belief, equivalent to the confidence level. One relation covers continuous distributions and makes use of the area under the density function secured by integration, while the second relation covers summations for integral values. These relations are

$$P_b = \int_0^{x^*} f(x)dx / \int_0^{x=\infty} f(x)dx$$

$$P_b = \sum_0^{x^*} F(x) / \sum_0^{x=\infty} F(x)$$

For the exponential density function, use of these relations gives

$$\begin{aligned}
 P_b &= \int_0^{t^*} \lambda \exp(-\lambda t) dt / \int_0^{\infty} \lambda \exp(-\lambda t) dt \\
 &= -\exp(-\lambda t)_0^{t^*} / -\exp(-\lambda t)_0^{\infty} \\
 &= [-\exp(-\lambda t^*) + 1] / (0 + 1) \\
 &= 1 - \exp(-\lambda t^*)
 \end{aligned}$$

The reliability *R(t)* for time *t* for *P_b* (the one-sided confidence level) is derived from the term λt^* , where no failures have been observed in time $T = \sum t_i$. From an exponential table determine $\lambda t^* = a$, corresponding to $1 - P_b$. Then $\lambda = a/t^* = a/T$, since *t** corresponds to the total time required for the test. The final reliability value is determined from

$$R(t) = \exp(-\lambda t) = \exp(-at/T)$$

If a test is terminated when the *r*th failure has occurred, the ratio $2r(\hat{m}/m)$ has a chi-square distribution with $2r$ degrees of freedom. The two-sided confidence interval at a confidence level of $(1 - \alpha)$ is

$$\hat{m}(2r/\chi^2_{\alpha/2;2r}) \leq m \leq \hat{m}(2r/\chi^2_{1-\alpha/2;2r})$$

Here \hat{m} represents the estimate of *m* derived from the samples tested and is the MTBF. The lower limit *L* is given by

$$L = (2r/\chi^2_{\alpha/2;2r})\hat{m} = 2T/\chi^2_{\alpha/2;2r}$$

while the upper confidence limit is given by

$$U = (2r/\chi^2_{1-\alpha/2;2r})\hat{m} = 2T/\chi^2_{1-\alpha/2;2r}$$

Herein $\hat{m} = T/r$ and can be derived from either a replacement or a nonreplacement test, while $T = \sum t_i$, the sum of the operating times accumulated by all the components during the test. When the test is terminated at time *t_d* without a failure occurring exactly at that time, then the degrees of freedom for the lower limit are changed from $2r$ to $2r + 2$. The upper and lower limits are given by

$$2T/\chi^2_{\alpha/2;2r+2} \leq m \leq 2T/\chi^2_{1-\alpha/2;2r}$$

From these limits giving lower and upper limits, *L* and *U*, in terms of mean time before failure, for any mission time *t*, then lower and upper limiting values for the reliability *R(t)* may be readily computed from

$$L::R_L:R_L(t) = \exp(-t/L)$$

$$U::R_U:R_U(t) = \exp(-t/U)$$

When the Gaussian (normal law) distribution applies or is used as a means of determining upper and lower limits for either *m* or *R(t)*, where \bar{m} , σ_m and \bar{R} , σ_R are known, either symmetric or nonsymmetric confidence limits may be determined from

$$L:\bar{m} - z_\alpha\sigma_m; \bar{R} - z_\alpha\sigma_R$$

$$U:\bar{m} + z_\beta\sigma_m; \bar{R} + z_\beta\sigma_R$$

where $\alpha + \beta = \gamma$ = probability for the specified confidence band.

As an aid to the calculation of confidence bands under certain stated conditions, several of the military specifications listed in Table 7 allow an easy calculation of these limits. One specification to note is MIL-R-22973.

EXAMPLE 4. Estimation of reliability for times different from test time.

In most cases, reliability values are associated with life usage or with the time of storage. A mission may require *t* hours to be accomplished. For example, it may require 10 hours to drive an automobile from Los Angeles to San Francisco, a distance of approximately 420 miles. What is *R(10)*, the reliability of accomplishing this mission in 10 hours at any time? In a prior example, $p = 2.5\%$ where it was assumed that each unit was tested 1000 hours. The failure rate then may be expressed as 2.5%/1000 hours. If the mission time is 10 hours, the reliability *R(t)*, assumed to be based on the exponential, is determined as follows. For $t = 10$ hours, $v = 2.5\%/1000$ hours = 0.000025/hr, and for this case

$$R(10) = \exp[-0.000025(10)]$$

$$= \exp(-0.00025) = 0.99975$$

TABLE 7. RELIABILITY SPECIFICATIONS

Specification	Title
MIL-A-8866	Airplane Strength and Rigidity Reliability Requirements, Repeated Loads and Fatigue
MIL-R-19610	General Specifications for Reliability of Production Electronic Equipment
MIL-R-22732	Reliability Requirements for Shipboard and Ground Electronic Equipment
MIL-R-22973	General Specification for Reliability Index Determination for Avionic Equipment Models
MIL-R-23094	General Specification for Reliability Assurance for Production Acceptance of Avionic Equipment
MIL-R-26484	Reliability Requirements for Development of Electronic Subsystems for Equipment
MIL-R-26667	General Specification for Reliability and Longevity Requirements, Electronic Equipment
MIL-R-27173	Reliability Requirements for Electronic Ground Checkout Equipment
M-REL-M-131-62	Reliability Engineering Program Provisions for Space System Contractors
NASA NPC 250-1	Reliability Program Provisions for Space System Contractors
NASA Circular No. 293	Integration of Reliability Requirements into NASA Procurements
LeRC-REL-1	Reliability Program Provisions for Research and Development Contracts
WR-41 (BUWEPS)	Naval Weapons Requirements, Reliability Evaluation
NAVSHIPS 900193	Reliability Stress Analysis for Electronic Equipment
NAVSHIPS 93820	Handbook for Prediction of Shipboard and Shore Electronic Equipment Reliability
NAVSHIPS 94501	Bureau of Ships Reliability Design Handbook
NAVWEPS 16-1-519	Handbook Preferred Circuits—Naval Aeronautical Electronic Equipment
PB 181080	Reliability Analysis Data for Systems and Components Design Engineers
PB 131678	Reliability Stress Analysis for Electronic Equipment, TR-1100
TR-80	Techniques for Reliability Measurement and Prediction Based on Field Failure Data
TR-98	A Summary of Reliability Prediction and Measurement Guidelines for Shipboard Electronic Equipment
AD-DCEA	Reliability Requirements for Production Ground Electronic Equipment
AD 114274	(ASTIA) Reliability Factors for Ground Electronic Equipment
AD 131152	(ASTIA) Air Force Ground Electronic Equipment-Reliability Improvement Program
AD 148556	(ASTIA) Philosophy and Guidelines—Prediction on Ground Electronic Equipment
Ad 148801	(ASTIA) Methods of Field Data Acquisition, Reduction and Analysis
AD 148977	(ASTIA) Prediction and Measurement of Air Force Ground Electronic Reliability
MIL-HDBK-217B	Reliability Prediction of Electronic Equipment
RADC 2623	Reliability Requirements for Ground Electronic Equipment
USAF BLTN 2629	Reliability Requirements for Ground Electronic Equipment
AR-705-25	Reliability Program for Material and Equipment
OP 400	General Instructions: Design, Manufacture and Inspection of Naval Ordnance Equipment
MIL-STD-105D	Sampling Procedures and Tables for Inspection by Attributes
MIL-STD-414	Sampling Procedures and Tables for Inspection by Variables of Percent Defective
MIL-STD-721	Definitions for Reliability Engineering
MIL-STD-756	Procedures for Prediction and Reporting Prediction of Reliability of Weapon Systems
MIL-STD-781B	Reliability Tests Exponential Distribution
MIL-STD-785	Requirements for Reliability Program (for Systems and Equipments)
DOD H-108	Sampling Procedure and Table for Life and Reliability Testing

This value of reliability is based on the expected value. For the exponential, the variance is equal to the expected value. Hence, since for this 10-hour mission $t = 0.00025$, $\sigma_t = (0.00025)^{1/2} = 0.01581$. For a 90% confidence level using the proper multiplying factor based on the normal law

$$t_{0.90} = 0.00025 + 1.282(0.01581) = 0.00025 + 0.02026842 = 0.02051842$$

The corresponding reliability for $t = 10$ hours is $R_{0.90}(t = 10) = \exp(-0.02052)$, $R(10)_{0.90} = 0.97968$. For a 95% confidence level based on the normal law:

$$(\lambda t)_{0.95} = 0.00025 + 1.645(0.01581) = 0.00025 + 0.0260075 = 0.0262575$$

For this expected value of λt with $t = 10$ hours, $R(10) = \exp(-0.02626) = 0.97408$. For a 99% confidence level based on the normal law:

$$(\lambda t)_{0.99} = 0.00025 + 2.326(0.01581) = 0.00025 + 0.03677406 = 0.03702406$$

Hence $R(t = 10)_{0.99} = \exp(-0.03702) = 0.96366$. For the six confidence levels often used, the reliability values for a one-tailed confidence level may be obtained from Table 8.

According to the data in Table 3, there existed only one set of 100 units = n out of 10 000 that had a reliability observed of 0.94. Associated with this value is a confidence level of $(10\ 000 - 100)/10\ 000 = 9900/10\ 000 = 0.99$ (Table 9). This reliability value is determined for an assumed operating period of 1000

TABLE 8. RELIABILITY VALUES FOR A MISSION OF $t = 10$ HOURS FOR $\lambda = 2.5\%/1000$ HOURS FOR 6 ONE-TAILED CONFIDENCE LEVELS FOR EXPONENTIAL: $R(\lambda t) = R[\lambda t + z(\lambda t)^{1/2}]$, z GIVEN IN NORMAL LAW (GAUSSIAN) TABLES FOR P_z TABULATED FOR CONFIDENCE LEVEL

Confidence Level	Normal Law z Values	Upper Limit for $(\lambda t)_z = 0.00025 + z(0.00025)^{1/2}$	$R(\lambda t)_z = \exp[-(\lambda t)_z]$
0.90	1.282	0.02052	0.97968
0.95	1.645	0.02626	0.97408
0.96	1.751	0.02794	0.97245
0.97	1.881	0.02999	0.97045
0.98	2.054	0.03272	0.96780
0.99	2.326	0.03702	0.96366

hours. Hence this gives $R(1000) = 0.94$ when the confidence level $P_C = 0.99$. Hence the value of λ is thus computed

$$R(1000) = 0.94 = \exp[-\lambda(1000)] = \exp(-0.0619)$$

Then $1000\lambda = 0.0619$, and $\lambda = 6.19\%/1000$ hours. For $t = 10$ hours, then

$$R(10) = \exp(-\lambda t) = \exp[-0.0619(10)/1000] = \exp(-0.000619) = 0.99938$$

Thus the actual data provide more optimistic estimates of the reliability based on field test results. Since the distribution as graphed appears to be almost rectangular in Fig. 12, the assumption of normality is pessimistic.

In these life tests, each failure must be carefully analyzed to determine whether it is a *chance* failure or a *wearout* failure. These results must be fed back to the design engineers to make certain that corrective measures for improving the life characteristics are taken and established as standard procedures.

Fitting a Distribution Using Chi-Squared Test

In a sample of n observed values of a random variable, n_1, n_2, \dots, n_k frequencies are observed in each of k intervals (Fig. 12). An assumed distribution

predicts e_1, e_2, \dots, e_k for the same intervals. The goodness-of-fit is tested by the criterion

$$\sum (n_i - e_i)^2 / e_i < \chi^2_{1-\alpha, r}$$

In this expression, $\chi^2_{1-\alpha, r}$ is the value of the χ^2 distribution at the cumulative probability level $(1 - \alpha)$ and with degrees of freedom $r = k - 1$ if the mean is known independently and $r = k - 2$ if the sample mean is used.

EXAMPLE 5. Using data from example 1, test the goodness of fit for normal and Poisson distributions.

From Table 4 for $k = 7$, obtain the data in Table 10. For the normal distribution,

$$\chi^2 = (1/3 + 1/3 + 0 + 0 + 0 + 0 + 0) = 0.667$$

For $\alpha = 0.05$ significance level and $r = k - 2 = 5$, standard tables of χ^2 , $\chi^2_{0.95, 5} = 11.1$, which is greater than χ^2 for the normal distribution. Hence the normal distribution is a valid model at the 5% level.

For the Poisson distribution,

$$\chi^2 = (4/2 + 4/4 + 1/5 + 0 + 0 + 1/1 + 0) = 4.2$$

Hence, the Poisson distribution is also valid at the 5% significance level, but it is a poorer fit to the data than the normal distribution.

TABLE 9. RELIABILITIES ASSOCIATED WITH ONE-TAILED CONFIDENCE LEVELS

Confidence Level	Multiplying Factor for Normal Law, z	$n = 10\ 000$ units	
		Observed Reliability	Normal Law Theoretical Reliability
0.90	1.282	95.07%	95.69%
0.95	1.645	94.50%	95.18%
0.96	1.751	94.38%	95.03%
0.97	1.881	94.25%	94.85%
0.98	2.054	94.12%	94.60%
0.99	2.326	94.00%	94.22%

TABLE 10. DATA FOR EXAMPLE 5

<i>i</i>	<i>n</i>	<i>e_i</i> (normal)	<i>e_i</i> (Poisson)
1	4	3	2
2	2	3	4
3	4	4	5
4	4	4	4
5	3	3	3
6	2	2	1
7	1	1	1

Probability Paper

An alternative method of fitting a distribution is to use probability plotting paper. This procedure attempts to produce the best fit to the cumulative distribution rather than the probability density function. A number of results obtained in a life test are tabulated in Table 11. The cumulative failure distribution, *F(t)*, can be approximated by

$$F(t_i) = i/n$$

where *i* is the number of failures at time *t_i*.

The principle in using probability paper is to transform the axes of the cumulative plot to give a straight line if the data conform with the assumed theoretical distribution. For instance, an exponential distribution gives a straight line if log [*F(t)*] is plotted against *t*, with a slope equal to the hazard rate, λ. In other cases, such as the normal distribution, special plotting papers must be prepared. These can be obtained for the more common distributions from suppliers of drawing materials.

Weibull Analysis

By far the most commonly used plotting paper in reliability studies is Weibull paper (see Fig. 15). A

TABLE 11. RESULTS OF COMPONENT LIFE TEST ON 15 UNITS

Test Number	Time to Failure (Hours)
1	8.35 × 10 ⁴
2	2.24 × 10 ⁴
3	4.84 × 10 ⁴
4	0.33 × 10 ⁴
5	9.90 × 10 ⁴
6	8.85 × 10 ⁴
7	0.83 × 10 ⁴
8	10.60 × 10 ⁴
9	6.55 × 10 ⁴
10	1.80 × 10 ⁴
11	4.30 × 10 ⁴
12	8.45 × 10 ³
13	9.27 × 10 ⁴
14	3.15 × 10 ⁴
15	6.02 × 10 ⁴

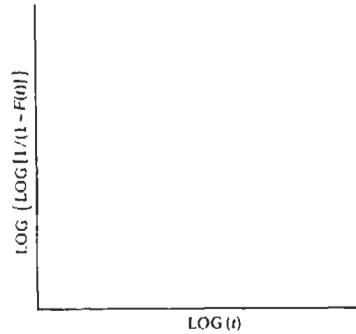


Fig. 15. Weibull plotting transformation.

Weibull distribution has the attribute that it can describe distributions with decreasing, constant, or increasing hazard rate, simply by appropriate choice of the shape parameter, β. Furthermore, the slope of the Weibull plot gives β directly.

EXAMPLE 6. Use of Weibull paper to analyze the data in Table 11. It is first necessary to place the results in increasing rank order *m*. Next the cumulative failure distribution *F(t)* is approximated by the median rank

$$\hat{F}_n(t) = m/(n + 1) \approx F(t)$$

This is slightly different from the approximation to *F(t)* quoted earlier, i.e., *F(t_i) = m/n*. The difference is that the median rank is a minimally biased estimate of *F(t)* and avoids some plotting problems near *m = 1* and *m = n* as well. Similarly, the reliability *R(t)* can be estimated by

$$\hat{R}_n(t) = (n - m)/(n + 1)$$

Table 12 shows the ordered list with the median ranks calculated. The results are plotted in Fig. 16. Note the following:

1. Fitting a distribution using probability paper requires relatively few points compared with the number needed to produce a histogram to approximate a probability density function. As few as five or six can be used to fit a cumulative distribution, whereas more than 20 are needed for a histogram.
2. The data in the Weibull plot clearly indicate three distinct zones, suggesting that three modes of failure are occurring.
3. From the estimates of the shape parameter, it can be deduced that the first two or three failures are burn-in (β = 0.6), the last five are clearly wearout (β = 3.6, which is incidentally very close to a normal distribution in shape), while the remainder are random failures with external cause (β = 1.08, approximately 1).

TABLE 12. CALCULATION OF MEDIAN RANKS FOR WEIBULL ANALYSIS OF DATA IN TABLE 11

Rank Order, <i>m</i>	Median Rank, <i>m</i> /(<i>n</i> + 1)	Time to Failure, (10 ⁴ hr)
1	0.0625	0.33
2	0.125	0.83
3	0.1875	1.80
4	0.25	2.24
5	0.3125	3.15
6	0.375	4.30
7	0.4375	4.84
8	0.50	6.02
9	0.5625	6.55
10	0.625	8.35
11	0.6875	8.45
12	0.75	8.85
13	0.8125	9.27
14	0.875	9.90
15	0.9375	10.60

The Weibull plot is therefore not just a means of fitting a distribution to data. It is also a useful diagnostic tool, indicating the number and nature of failure modes present.

Distribution-Free Tests of Goodness of Fit

As with all data-fitting procedures, probability-paper plots provide only an estimate of the true distribution. Confidence limits should be placed on any parameter estimates made. This is more difficult to do for several parameters simultaneously than for a single variable such as failure rate. Further details of this type of analysis are to be found in any standard text on statistics, e.g., Reference 17. One technique that will be illustrated here is a check to find out whether a set of data belongs to an already known and characterized population. This test is known as the Smirnov-Kolmogorov test. The procedure is illustrated in Fig. 6, which shows the actual *F*(*t*) of the population and a median-rank plot for a typical sample of data. It is required to test whether the data is a member of the population. The test requires that the maximum deviation between *F*(*t*) and *F_n*^α(*t*), *D_n*^α, must be less than some specified value, *D_n*^α which depends on the number of data points, *n*, and the desired significance level, (1 - α). The value for *D_n*^α is found in tables, e.g., in Reference 17. A short version of *D_n*^α values is given in Table 13.

EXAMPLE 7. Given a test *n* = 40, *D_n* = 0.2 (see Fig. 6), determine the appropriateness of the theoretical model at the 10% significance level.

From Table 13, *D₄₀*^{0.1} = 0.19 < *D_n*. Model rejected at 10% level.

NOTE. For more information on estimation and confidence limits, consult the references at the end of this chapter. Chapters 44 and 47 give additional statistical and mathematical information.

Bayesian Statistics

Previous discussion has been confined to the classical approach to estimation, which implicitly assumes no information other than the immediately acquired test data. Since 1964, increased interest has been expressed in the so-called Bayesian approach.

Bayes' theorem is an early result in probability theory which states that, if the occurrences of two events *A* and *B* are dependent, then

$$\Pr(A \text{ and } B) = \Pr(A/B) \cdot \Pr(B) = \Pr(B/A) \cdot \Pr(A)$$

This can be re-expressed as

$$\Pr(A/B) = \Pr(A) \cdot \Pr(B/A)/\Pr(B)$$

This equation can be used to infer a result using both test data and previous experience. For instance,

A = λ, the actual (unknown) failure rate

B = λ̂, the observed failure rate in a test

The question is posed, what is the actual failure rate given the observed rate in the test?

The focal point of the method is the so-called "prior" distribution, Pr(*A*).

TABLE 13. SHORT LIST OF CRITICAL VALUES OF *D_n* FOR THE SMIRNOV-KOLMOGOROV TEST

<i>n</i>	<i>D_n</i> ^{0.1}	<i>D_n</i> ^{0.05}	<i>D_n</i> ^{0.01}
5	0.51	0.56	0.67
10	0.37	0.41	0.49
20	0.26	0.29	0.36
30	0.22	0.24	0.29
40	0.19	0.21	0.25
50 and more	1.22/√ <i>n</i>	1.36/√ <i>n</i>	1.63/√ <i>n</i>

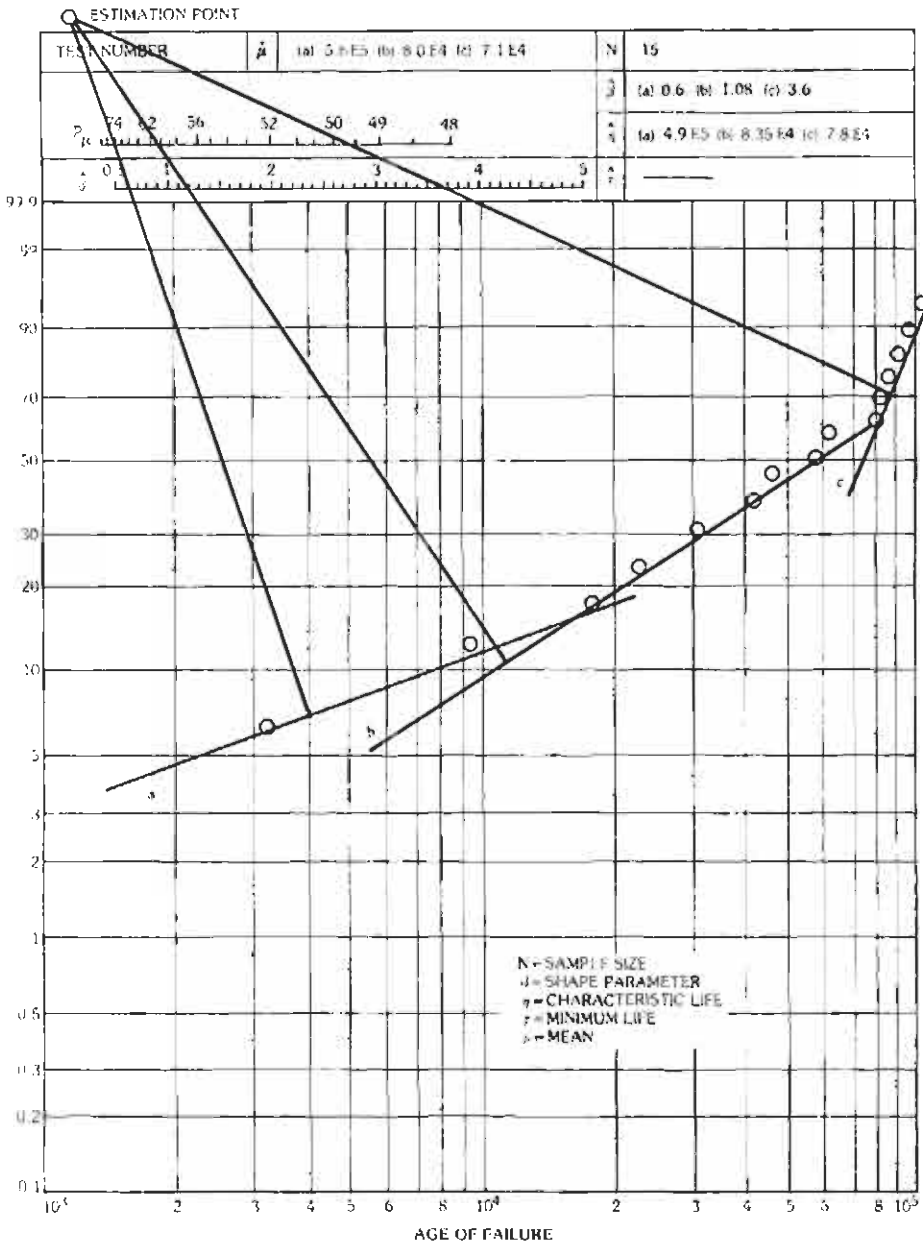


Fig. 16. Standardized Weibull paper with example.

would be considered a constant, but unknown. In Bayesian statistics, A is considered to be a random variable, and use is made of previous experience, or even subjective judgement, to choose an appropriate distribution for $\Pr(A)$. The Bayesian approach, therefore, has some appeal in engineering applications, where reasonable estimations based on previous experience are not only common but absolutely necessary in

some cases, if a decision is to be made at all. On the other hand, there is some controversy over the validity of the reasoning behind the Bayesian approach that should be explored before it is used, although most of the concern appears to be over points that do not affect the practical use of the approach. A readable introduction to Bayesian statistics is given in Reference 17. Applications to problems in reliability are given in

References 4 and 15. Generally, Bayesian estimates tend to be more optimistic than those derived from classical statistical methods.

CODES AND STANDARDS

Practically all military contracts contain clauses under quality assurance requiring reliability programs. This is becoming increasingly true in other fields as well, particularly civil aviation, nuclear power, and industries where equipment failure is potentially hazardous. Many contracts include provisions for a preliminary design review and also a critical design review before qualification tests. In addition, some contracts require a reliability demonstration test. Use is made in this latter case of all the engineering data that might be obtained in the qualification tests, if they precede the reliability demonstration test.

Table 7 presents a family tree of US Government documents establishing and supporting reliability requirements. New specifications are added frequently to build up the reliability factors and requirements. One area (not listed in the table) that is expanding steadily is special parts reliability specifications such as the MIL-R-38000 series. These specifications cover the acceptance and qualification testing of high-reliability parts.

THE USE OF COMPUTERS IN RELIABILITY

Most organizations today have extensive computing facilities. Invariably, such a facility will include software packages for performing statistical calculations. Even the most modest of desk-top computers has a package of this type, or one can be obtained as an option. The majority of standard calculations, such as curve fitting, evaluation of confidence limits, etc., can be done as a routine exercise, and cumbersome, error-prone tasks like referring to sets of tables for values of chi-square and the *t* distribution are now obsolete, except for educational demonstration.

Of particular interest to reliability analysts is an ongoing series in *Quality Progress*, the journal of the American Society for Quality Control, which offers computer programs for use on small desk-top models and hand-held programmable calculators, covering a range of reliability-related topics.

Several large computer codes are now generally available for various aspects of reliability analysis of complex systems. These are reviewed in Reference 4. Many of the codes described in Reference 4 are available from the authors of the reference. Some examples of the range of capabilities are

- KITT — Fault-tree analysis of complex systems (several versions).
- BACFIRE — Search routine for common-mode failures.
- FAMULS — Calculates cut sets for systems with

multiple control loops given signal flow graph representation of system.

- SCHE — Converts reliability block diagrams (see Fig. 10) into fault trees.
- HEUR — Reliability optimization under constraints of cost, weight, etc.

REFERENCES

1. Halpern, S. *The Assurance Sciences—An Introduction to Quality Control and Reliability*. Englewood Cliffs, N.J.: Prentice-Hall, Inc., 1978.
2. Marguglio, B. W. *Quality Systems in the Nuclear Industry and Other High Technology Industries*. ASTM/STP 616, 1977.
3. Swain, A. D., and Guttman, H. E. *Handbook of Human Reliability Analysis with Emphasis on Nuclear Power Plant Applications*. USNRC Report NUREG/CR-1278, October 1980.
4. Hensley, E. J., and Kumamoto, H. *Reliability Engineering and Risk Assessment*. Englewood Cliffs, N.J.: Prentice-Hall, Inc., 1981.
5. NASA. *Parts Materials, and Process Experience Summary*. Report CR114391, Feb. 1972.
6. Advisory Group on Reliability of Electronic Equipment (AGREE). Office of the Assistant Secretary of Defense (Research and Engineering). *Reliability of Military Electronic Equipment*. Washington, D.C.: Supt. of Documents, US Government Printing Office, 4 June 1957.
7. MIL-HDBK-217B, *Military Standardization Handbook, Reliability Prediction of Electronic Equipment*, September 20, 1974.
8. IEEE Std. 500-1977. *Guide to the Collection and Presentation of Electrical, Electronic and Sensing Component Reliability Data for Nuclear Power Generating Stations*. New York: IEEE.
9. Government/Industry Data Exchange Program (GIDEP), Fleet Missile Systems Analysis and Evaluation Group, Corona, California.
10. *Failure Data Handbook for Nuclear Facilities*, LNEC-Memo-69-7. Available from NTIS.
11. Nuclear Plant Reliability Data Systems (NPRDS), operated by Southwest Research Institute, San Antonio, Texas.
12. SYREL—System Reliability Service Data Bank, UKAEA, Culcheth, Warrington, United Kingdom.
13. H108, *Quality Control and Reliability Handbook, Sampling Procedures and Tables for Life and Reliability Testing*. Washington, DC: Office of the Assistant Secretary of Defense, April 19, 1960.
14. MIL-STD-781B, *Military Standard, Reliability Tests Exponential Distribution*. Washington, DC: Department of Defense, November 15, 1967.
15. Bazovsky, Igor. *Reliability Theory and Practice*. Englewood Cliffs, N.J.: Prentice-Hall, Inc., 1961.
16. Epstein, B. "Estimation From Life Test Data." *IRE Trans on Reliability and Quality Control*, Vol. RQC-9, April 1960.

17. Ang, A. H-S., and Tang, W. H. *Probability Concepts in Engineering Planning and Design*. New York: John Wiley and Sons, Inc., 1975.
18. Von Alven, William H., ed. *Reliability Engineering*. Englewood Cliffs, N.J.: Prentice-Hall, Inc., 1964 (23 contributors).

Other Useful References

- Abramowitz, Milton, and Stegun, Irene A. *Applied Mathematics Series AMS 55, Handbook of Mathematical Functions With Formulas, Graphs and Mathematical Tables* (third printing with corrections). Washington, DC: Supt. of Documents. US Government Printing Office, 1965.
- Ad Hoc Study Group on Parts Specification Management for Reliability, Office of the Director of Defense Research and Engineering and Office of the Assistant Secretary of Defense Supply and Logistics. *Parts Specification Management for Reliability*, Vols. 1 and 2, PSMR-1. Washington, DC: Supt. of Documents, US Government Printing Office, May 1960.
- Calabro, S. R. *Reliability Principles and Practices*. New York: McGraw-Hill Book Co., 1962.
- Staff of the Computation Laboratory. *Tables of the Error Function and of Its First Twenty Derivatives*. Cambridge, Mass.: Harvard University Press, 1952.
- Goldman, A. S., and Slattery, T. B. *Maintainability*. New York: John Wiley & Sons, Inc., 1964. (Contributions by S. Firstman, Rand Corp.; and J. Rigney, University of Southern California): Chapter on General Electric Co. "TEMPO."
- Gryna, Frank M. Jr., McAfee, Naomi J., Ryerson, Clifford M., and Zwerling, Stanley, eds. *Reliability Training Text*, 2nd ed. Sponsored by ASQC and IEEE, March 1960.
- Hald, A. *Statistical Tables and Formulas*. New York: John Wiley & Sons, Inc., 1952.
- Ireson, W. Grant, ed. *Reliability Handbook*. New York: McGraw-Hill Book Co., 1966. (19 contributors)
- Johnson, Norman L., and Leone, Fred C. *Statistics and Experimental Design in Engineering and the Physical Sciences*, Vols. 1 and 2. New York: John Wiley & Sons, Inc., 1964.
- Lambe, C. G. *Elements of Statistics*. London and New York: Longmans, Green and Co., 1952.
- Landers, Richard R. *Reliability and Product Assurance, A Manual for Engineering and Management*. Englewood Cliffs, N.J.: Prentice-Hall, Inc., 1963.
- Lloyd, David K., and Lipow, Myron. *Reliability: Management Methods and Mathematics*. Englewood Cliffs, N.J.: Prentice-Hall, Inc., 1962.
- Lowan, Arnold N., and Staff. *Applied Mathematics Series AMS 23, Tables of Normal Probability Function*, Washington, DC: Supt. of Documents. US Government Printing Office, 1953.
- Lowan, Arnold N., and Staff. *Applied Mathematics Series AMS 14, Tables of the Exponential Function, e^x* , 4th ed. Washington, DC: Supt. of Documents, US Government Printing Office, 1961.
- Molina, E. C. *Poisson's Exponential Binomial Limit, Table 1—Individual Terms; Table 2—Cumulated Terms*. New York: D. Van Nostrand Co., Inc., 1949.
- Zelen, Marvin, ed. "Statistical Theory of Reliability." Proceedings of an Advanced Seminar Conducted by the Mathematics Research, United States Army, at the University of Wisconsin, Madison, 8-10 May 1962. Madison, Wis.: The University of Wisconsin Press, 1963.

46 Cellular Telecommunications Systems

William C. Y. Lee

General Description	46-3
History	
Allocated Frequency Spectrum	
A Duplexed System	
Cellular System Elements	46-4
Cell Site (or Base Station)	
The Switch and BSC	
Data-Link Network	
Enhancers and Converters	
Mobile Subscribers' Units	
Operation of Cellular Systems	46-7
Mobile-Unit Initialization	
Mobile-Originated Call	
Network-Originated Call	
Call Termination	
Call Blocking	
Call Completion	
Call Drops	
Handoff Procedure	
Mobile Radio Environment	46-8
Uniqueness of the Mobile Radio Environment	
A Model of a Mobile-Radio Path Loss	
Characteristics of Multipath and Selective Fading	

Frequency Reuse and Design Aspect for FDMA and TDMA	46-17
Frequency Reuse Distance	
Cellular Design Aspect	
Frequency Reuse Factor, K	
Radio Capacity	46-19
Increase Capacity	
Diversity Schemes and Combining Techniques	
Switching Equipment and Traffic Models	46-21
Analog Switching Equipment	
Cellular Digital Switching Equipment	
Traffic Models	
Digital Cellular Systems	46-22
Digital AMPS (or North America Digital Cellular, NADC)	
GSM (Special Mobile Group) System	
Cellular CDMA System	
3G Systems	
4G System	
Other Personal Communication Systems	

GENERAL DESCRIPTION

History

The cellular telecommunications system was developed by AT&T Bell Laboratories. It was called the Advanced Mobile Phone Service (AMPS) System. The first AMPS system was deployed in Japan with a slight change in 1979. In 1983, the first cellular telecommunications system was operated in Chicago. Each city has two cellular system operators licensed by the FCC for achieving a duopoly competition policy. In 1984, one cellular system operator in each of the first thirty largest cities was in operation.

Since the cellular industry was growing very fast, the system capacity became a challenge and needed to be resolved. The answer was to go digital. In 1989, the first cellular digital standard (IS-54), called North American TDMA (Time Division Multiple Access), was issued. In 1991, the second digital cellular standard (IS-95), called CDMA (Code Division Multiple Access), was issued. In 1994, the standard IS-54 had a major modification and was renamed as IS-136. The capacity of IS-95 can be ten times that of AMPS, and the IS-136 is three times that of the AMPS. Both systems were commercialized in 1995 and both are mainly for voice.

In 1998, the Internet industry started to grow rapidly and transmitted data became important for wireless communication. The ITU (International Telecommunication Union) then issued an IMT-2000 standard also called the third-generation (3G) standard. There were three modes: WCDMA (Direct Sequence (DS) CDMA) developed by a standard group named 3GPP (Third Generation Partnership Project); CDMA 2000 (or multi-carrier (MC) CDMA developed by 3GPP2); and TDD (Time Division Duplexing). In TDD, there were UTRA-TDD developed by European community, and TD-SCDMA developed by Chinese standard. TDD system is used for a small-area system such as on a campus or in a shopping mall. There is a new large-area TDD system called LAS-TDD using a smart code technology. It will be described later in this chapter. Both CDMA 2000 (its 1X version, i.e., 1.25 MHz bandwidth), and WCDMA will be commercialized in 2003. The maximum data rate can be 2Mbps.

In the future, wireless communications and the Internet will be merged to become a wireless-Internet industry carrying both voice and data services.

Allocated Frequency Spectrum

In 1983, the FCC allocated a spectrum of 40 MHz to analog cellular systems. In 1986, another 10 MHz of spectrum was added. The total spectrum is 50 MHz. The same 50 MHz spectrum is operated in every city. Since there are two operators in each city, the 50 MHz

spectrum is divided into two 25 MHz bands denoted Band A and Band B. Each operator may operate only in its own band, either A or B, as licensed by the FCC.

In 1996, the PCS (personal communications service) spectrum was auctioned. There were six bands: A, B, C, D, E, F as shown in Table 1. Each market can have all of the six bands in operation. Each band can choose its own standard system among three: CDMA, TDMA, or DCS 1900 (a GSM version). A non-standard system called iDEN, invented by Motorola and operated by Nextel, also serves to the public at its SMR spectrum. The IMT-2000 system allocated spectrum issued by ITU is also shown in Table 1. The U.S. PCS bands for base transmit are overlaid on the IMT-2000 band for mobile transmit also shown in Table 1. Therefore, a new IMT-2000 spectrum is needed.

A Duplexed System

FDD System: A cellular system is a frequency division duplexed (FDD) system. In an analog system, a 25 MHz spectrum, 12.5 MHz is for base transmit and 12.5 MHz is for mobile transmit. In the North American system there are 832 frequency channels. Since this is a duplexed system, a frequency channel consists of a pair of two frequencies, forward link and reverse link. The forward link is from the cell to the mobile unit, and the reverse link is from the mobile unit to the cell site. The allocation of Band A and Band B is shown in Chart 1. In each band, Band A and Band B, there are 416 channels. Among them, 21 channels are setup channels (sometimes called control channels). The remaining 395 channels in each band are voice channels. Neither Band A nor Band B has a continuous spectrum, as shown in Chart 1.

In 1991, a European standard digital system called GSM was deployed in Germany. In 1994, the North American digital TDMA system was deployed in the United States and in 1995, the CDMA was deployed in Hong Kong and in the United States. The three systems' descriptions are shown in the digital cellular systems section. The two modes of third-generation (3G) systems, DS and MC, are also FDD-CDMA systems that we will mention later.

TDD system: The Time Division Duplex systems in 3G such as UTRA-TDD and TD-SCDMA are using a single frequency band for both transmitting and receiving. The current 3G TDD systems are developed for campus or small area communications. There is a new TDD system called the LAS-CDMA system, which uses a smart code, called LAS code, to isolate interference. The LAS-CDMA system is a large area synchronous (LAS) system. With its efficient spectrum by using smart code, the LAS-CDMA system could be a candidate for the future high capacity and high-speed data system. The system specification will be mentioned later.

TABLE I. DIAGRAM OF GLOBAL SPECTRUM ALLOCATION.

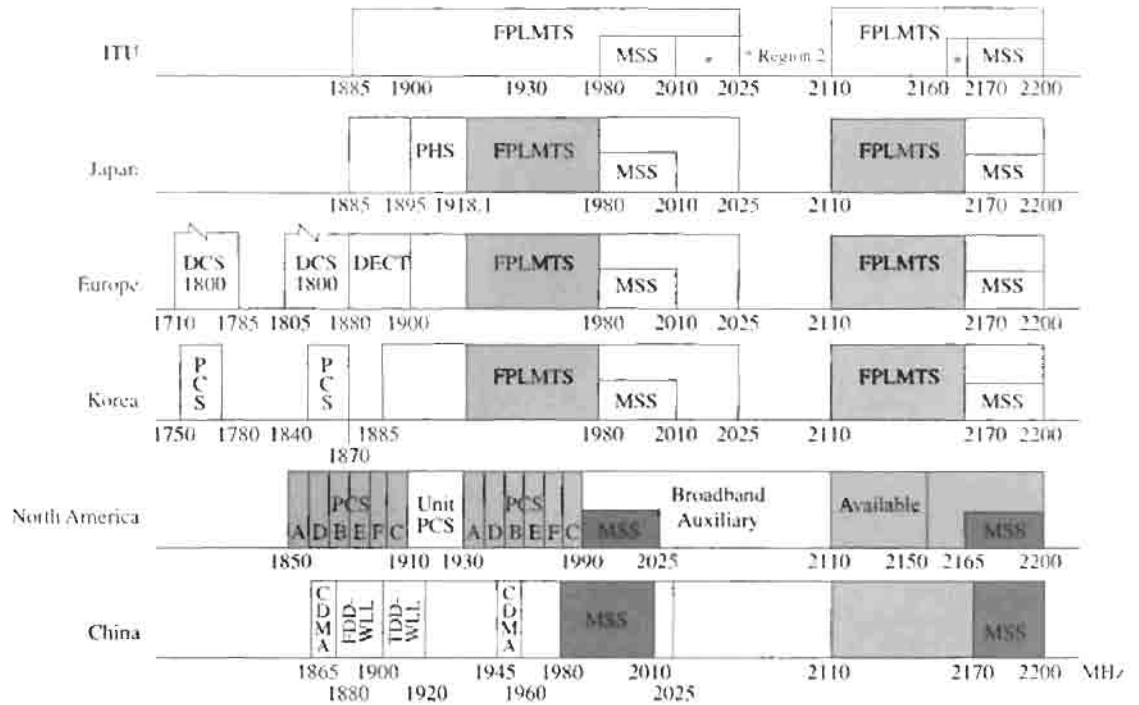


CHART I. ALLOCATION OF BANDS A AND B

	Numbering Scheme					
(Base Tx)→	869	870	880	890	891.5	894 MHz (Base Tx)
(Mobile Tx)→	824	828	835	845	846.5	849 MHz (Mobile Tx)
(No. of Channels)→	(New) 33(A)	333(A)	333(B)	(New) 50(A)	(New) 83(B)	
(Ch. Numbering)→	991	1 (1023)	333 334	666 667	716 717	799

CELLULAR SYSTEM ELEMENTS

The cellular system consists of a switch, a number of cell sites, a data-link network, enhancers and converters, and mobile subscribers' units. The general view of cellular telecommunications systems is shown in Fig. 1.

Cell Site (or Base Station)

The cell site consists of a cell-site antenna with a height of about 30 m (100 ft) to 45 m (150 ft), a controller, and a number of transceivers. The controller is used to handle the call process between the switch and the

mobile unit or handset via a setup channel. For an analog system, each cell site is assigned a different setup channel from the 21 setup channels. Also, each cell site is assigned one of three SATs (Supervisory Audio Tones). The three SATs are 5970 Hz, 6000 Hz, and 6030 Hz. The cell-site transceivers modulate the assigned SATs on all the forward voice channels from each cell site. Also, the cell-site transceivers are able to detect the SATs on the reverse voice channels sent back by the mobile unit. The choice of a proper cell site is dependent upon the planned signal coverage area, which is determined by the ERP (Effective Radiated Power), the antenna patterns, and the antenna height. In addition, a permit from the local unit of government is required.

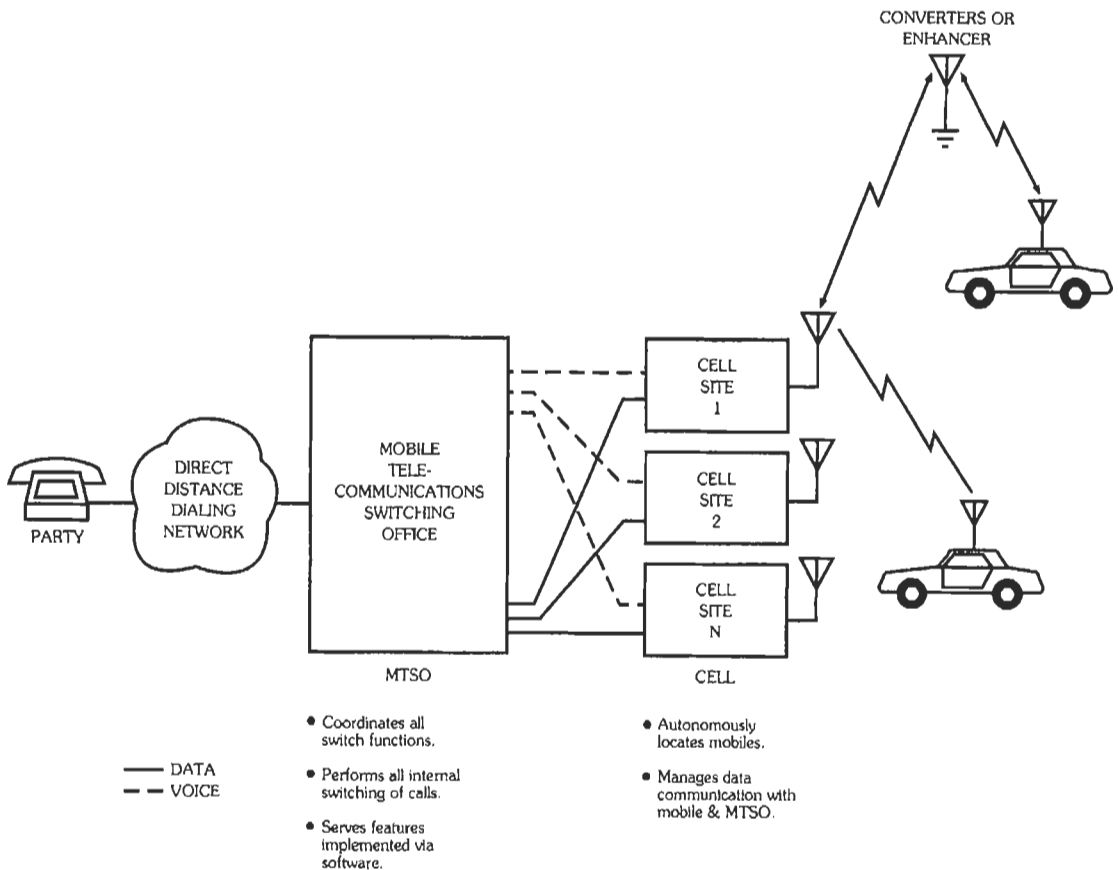


Fig. 1. A general view of cellular telecommunications systems.

In digital systems,

TDMA Systems (GSM or North American TDMA)—In these systems, there is a frequency correction channel (FCCH), a synch channel (SCH), a broadcast control channel (BCCH), a paging and access grant channel (PAGCH), a call broadcast channel (CBCH), and a traffic channel (TCH).

CDMA System—In this system, there is a pilot channel, a synch channel, a paging channel, and a traffic channel. In the pilot channel, an unmodulated, direct-sequence (DS) signal is transmitted continuously by each CDMA BS.

3G Systems—3G systems are wideband CDMA systems, which have 5MHz bandwidth. The European and Japanese standards bodies have developed WCDMA (DS) and the North American standards body developed CDMA 2000 (MC). These two are FDD systems. The TDD mode has two versions: TD-SCDMA, developed by China, and UTRA-TDD, developed by Europe. The differences among these systems are shown in Table 2.

The Switch and BSC

In analog systems, the Mobile Telephone Switching Office (MTSO) is a special-purpose switch that connects the call between mobile units and between the mobile unit and the landline telephone via the PSTN (Public Service Telephone Network). The function of the MTSO is to assign the voice channel to each cell, perform the handoffs (this operation will be described later) and the vertical service features, and monitor the calling information for billing.

In digital systems, the base station controller (BSC) was developed to handle all the control functions of the system, including handoffs. The transcoder and mobility management reside in the BSC. Therefore, in the advanced network, the switch and control are separated.

Data-Link Network

In analog systems, the data-link network carries the data between the cell sites and the MTSO. Each data link can carry multiple-channel data (10 kb/s data

TABLE 2. 3G HARMONIZATION AGREEMENT

Legend			
	FDD: Frequency Division Duplex		
	TDD: Time Division Duplex		
	CDM: Code Division Multiplexing		
	TDM: Time Division Multiplexing		
	IMT-2000 (3G)		
	FDD Direct Spread	FDD Multicarrier	TDD
Bandwidth	5 MHz	5 MHz/ 1.25 MHz	5 MHz/ 1.6 MHz
Chip Rate	3.84 Mcps	3.6864 Mcps/ 1.228 Mcps	3.84 Mcps/ 1.28 Mcps
Common Pilot	CDM	CDM	TDM
Dedicated Pilot	TDM	CDM	TDM
Synchronization	Asynchronous/ Synchronous	Synchronous as cdma2000	Synchronous

transmitted per channel). The data link can also carry PCM voice (64 kb/s per voice channel). There are generally three means of implementing the data links: 1) use T1-carrier wire lines; 2) use microwave radio links; and 3) use optical fibers.

In digital systems, the functional architecture and principal interfaces are shown in Fig. 2. The HLR (home location register) is a permanent database that contains subscriber information and is related to the subscriber's current location. VLR (visitor location register) is linking to one or more MSC's temporarily storing subscription data. OAM (operation, administration, and maintenance) is collecting the operational

data for troubleshooting and statistical analysis for performance. The architecture shown in Fig. 2 is mainly for voice communication using circuit switches which provide dedicated lines. The general packet radio system (GPRS) is called 2.5G and is used for packet transmission. Every packet has its sending address for delivery. The packet switch would not hold the line until the call termination. The packet transmission link is called the virtual connection or the connectionless transmission.

In 3G systems, there are two legacy networks, GSM MAP and ANSI 4, shown in Fig. 3. In order to have network-to-network interface (NNI) between two networks, the standards bodies (3GPP and 3GPP2) are developing new standards for hooks and extensions, as shown in Fig. 3. The hooks provide the specified functionality. The extensions provide the additional functions needed to meet the specified requirements. Also, in 3G system networks, two industry alliances were formed, 3GIP and MWiF (mobile wireless IP forum). These two forums are developing the wireless Internet network. The IP (Internet protocol) network can be the future 3G network, which is migrated from today's legacy network.

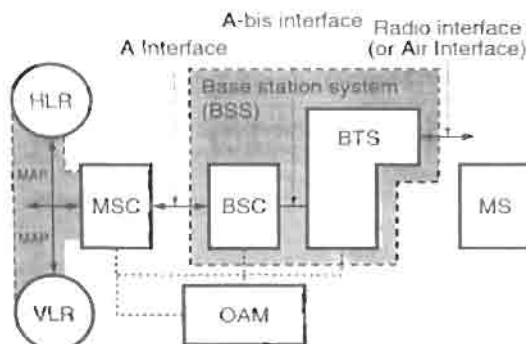


Fig. 2. Functional architecture and principal interfaces.

Enhancers and Converters

The enhancer is a repeater with amplification. It operates at the same cellular bandwidth. The signal range can be extended by the enhancer.

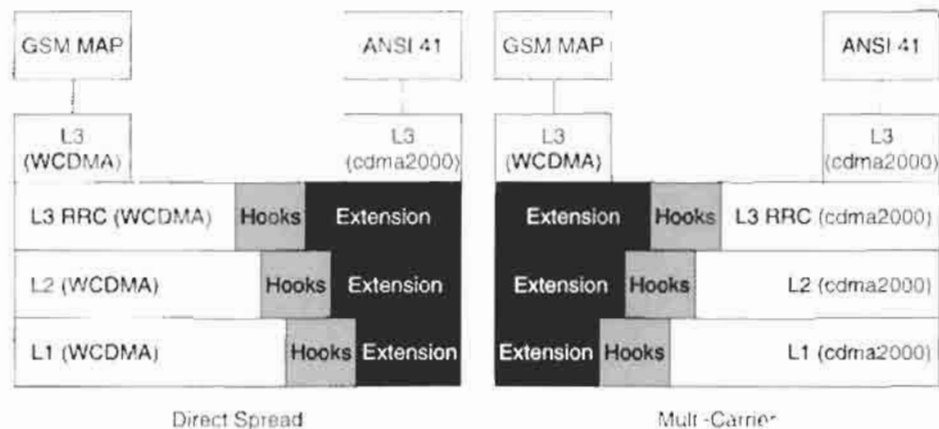


Fig. 3. Hooks and extensions.

If the converter changes the cellular frequency to a microwave frequency, it is called a microwave converter. If it converts the cellular frequency to an optical frequency, it is called an optical converter. If it down-converts to baseband frequency, it is called a baseband converter. The converted signal can be delivered to any destination without interference, then converted back to the cellular frequency at the destination.

Mobile Subscribers' Units

The mobile subscribers' units consist of three parts: the transceiver, the logic unit (handles signal processes), and the control unit (indicators). The mobile subscribers' units can be classified into two kinds, mobile units and handsets. Mobile units are mounted in vehicles; handsets are carried by the users.

OPERATION OF CELLULAR SYSTEMS

The cellular systems does not require the user to touch any adjustment buttons during a call. The operation can be described in seven different stages, plus a handoff procedure.

Mobile-Unit Initialization

In analog systems, when the mobile unit is turned on, the receiver scans the 21 setup channels, selects the strongest, and locks on for a certain time. Since each cell site is assigned a different setup channel, locking onto the strongest one means choosing the nearest cell site. This scheme is called a self-location scheme and is used at the idle stage. Every 60 seconds, the self-location procedure is repeated. This repetition is necessary due to the motion of the vehicle. A new nearest cell site can be determined after 60 seconds.

In digital systems,

For TDMA Systems—The mobile station (MS) guides the FCCH burst, then looks for an SCH burst on the same frequency to achieve synchronization. The MS then receives BCCH on several time slots and selects a proper cell, remaining for a period in the idle mode. The random-access channel (RACH) is the only common uplink channel that the MS chooses to access the calls.

For CDMA Systems—This allows the MS to acquire the timing of the forward CDMA channel, provides a phase reference for coherent demodulation, and provides a means for signal strength comparisons between base stations for determining when to hand off.

Mobile-Originated Call

In analog systems, when the user wants to make a call, a number is placed into an originating register in the mobile unit. When the "send" button is pushed, a request for service is sent on a preselected setup channel. The receiver has to check whether the setup channel is idle or not by examining the "busy" or "idle" bit in the forward setup channel. The cell site receives the request and in turn sends a request to the MTSO via a data link. The MTSO selects an appropriate voice channel from a group of allocated voice channels and notifies the cell site; then the cell site informs the mobile unit. At the same time, the MTSO also connects the wire-line telephone party through the telephone company central office.

In digital systems, the base station receives the signal from the RACH, an assigned traffic channel to the MS, informing MS through the PAGCH (paging ch.).

Network-Originated Call

In analog systems, a land-line party dials a mobile-unit number. The telephone company recognizes the

mobile number and forwards it to the MTSO. The MTSO then sends a paging message to certain cell sites depending on the mobile number. Paging is needed at the cell site because the MTSO does not know where the mobile unit is due to the self-location scheme being used by the mobile unit. Each cell site transmits the paging signal on its own assigned setup channel. The mobile unit recognizes its number on the strongest setup channel during the idle stage and responds to the particular cell site. The mobile unit also follows the message to tune to an assigned voice channel, at this time, the mobile unit alerts the user.

In digital systems, PAGCH is used to page the MS, and the MS will respond to the PAGCH message through RACH. As soon as the base station receives the response, it assigns a traffic channel to MS through PAGCH.

Call Termination

When the analog mobile user terminates a call, a special signal tone (10 kHz) is transmitted to the cell site to free the voice channel. The mobile unit resumes its initialization stage. In the digital mobile system, the functions are handled by different information bits.

Call Blocking

If during a mobile-originated call stage all the voice channels assigned to the nearest cell site are busy, then the analog mobile-originated call is repeated every 100 ms. After ten tries, a busy tone is turned on to inform the user. In digital systems, the operations can be changed dynamically depending on the demand.

Call Completion

After the mobile-originated call or network-originated call is completed, the voice channel is assigned to make a call connection. This process is called "call completion."

Call Drops

During a conversation, because of interference or weak-signal spots in certain areas, if the analog system cell site cannot receive an assigned SAT (generated at the cell site and sent back by the mobile unit) in 5 seconds, the cell site will turn off the transceiver. In digital systems, no SAT tone is needed. There are 256 digital color codes to assign to different cells and sectors. The call drops can be reduced.

Handoff Procedure

When a mobile unit moves out of the range of one cell and enters the range of another cell during a conversation on a voice channel, that voice channel has to change from one frequency channel to another. The system can make this handoff procedure without either

interrupting the call or alerting the user. The handoff process is a unique function in cellular-system operation. The signal threshold level setting for a handoff request cannot be too high, or there will be many unnecessary handoffs. However, if the threshold level is set too low, this will cause dropped calls.

In digital systems, TDMA is using hard handoff (break before make) but with a MAHO (mobile assisted handoff). The level of signal strengths of neighboring cells, measured by MS, will send back to the home base station for it to make the handoff decision.

CDMA is using soft handoff (make before break). In the handoff region, two cells (one home cell, one new cell) will simultaneously receive the MS signals. Also, many other potential cells will be considered (up to six cell fingers) by their pilot signal, to choose dynamically the new candidate cell in the soft handoff process.

MOBILE RADIO ENVIRONMENT

Uniqueness of the Mobile Radio Environment

Effect of the Natural Terrain Configuration—Due to the strong ground-reflected wave at the vehicle, both the direct wave and the reflected wave will be received. Since the ground-reflected wave has a 180° phase shift after reflection from the ground, the two waves tend to cancel and result in an excessive signal loss that is dependent on the natural terrain configuration.

Effect of Man-Made Structures—Because the mobile antenna is lower than most man-made structures, the multipath phenomenon prevails. Multipath waves arriving at the mobile unit cause signal fading called multipath fading or Rayleigh fading.

Effect of Man-Made Noise—Due to automobile ignition noise on the cellular frequency, the noise floor is much higher in the city than in a suburban area. Due to unintentional interference from signal sources such as government or military communications equipment or industrial machines, the noise floor can also be different in different areas.

Effect of Vehicle Motion—Different speeds of the vehicle cause different signal fading characteristics. When the speed is high, the fading is rapid; when the speed is low, the fading becomes slow. Also, the duration of the fading is short when the speed is high and is long when the speed is slow.

A Model of a Mobile-Radio Path Loss

A mobile-radio signal can be artificially separated into two parts as shown in Fig. 4. The first part is the radio-path propagation loss due to the terrain configura-

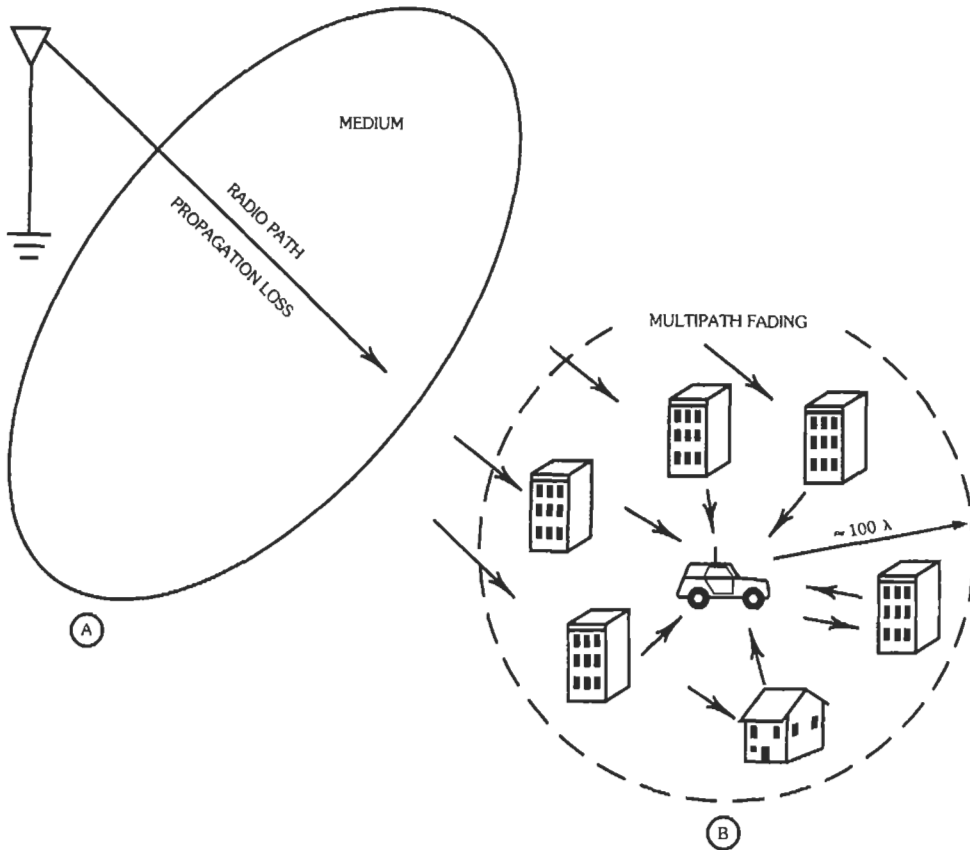


Fig. 4. Two parts of a mobile-radio environment: (A) propagation loss, (B) multipath fading.

ration, and the second part is the multipath fading due to the man-made structures.

Radio-Path Propagation Loss—The path-loss curve follows a log-normal distribution with a standard deviation of ± 8 dB. In theory:

$$\Delta P = 40 \log (d_1/d_0) \quad (\text{A } 40 \text{ dB/decade path loss}) \quad (\text{Eq. 1})$$

$$\Delta G = 20 \log (h_1'/h_1) \quad (\text{A base-station antenna-height gain of } 6 \text{ dB/octave}) \quad (\text{Eq. 2})$$

where ΔP is the difference in two received signal strengths based on the two different path lengths d_1 and d_0 , and ΔG is the difference in two received signal strengths based on the two different antenna heights h_1 and h_1' .

Local Mean Prediction from a Prediction Tool—The Lee model* can predict the local mean on each individual road provided that the geographical terrain

data is available for that location. This model combines the effects of man-made structures, which are different in different cities, and the natural terrain configuration. The received signal in decibels can be expressed as follows:

$$P_r = \text{received signal (due to the man-made structure + due to the natural terrain configuration) + correction factor}$$

The received predicted signal P_r in dBm can be obtained as follows for the indicated conditions.

For a nonobstructive path:

$$P_r = P_{ro} - \gamma \log (r/r_o) + 20 \log (h_c'/h_1) + \alpha \quad (\text{Eq. 3a})$$

For an obstructive path:

$$P_r = P_{ro} - \gamma \log (r/r_o) - L + \alpha \quad (\text{Eq. 3b})$$

For land-to-mobile path over water:

Use free-space formula.

For land-to-boat path on the water:

Use open-area formula.

* Reference 1.

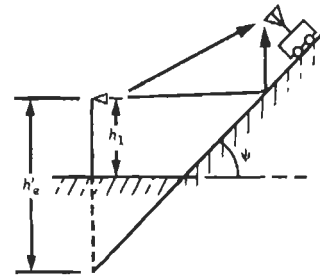
In equations 3a and 3b,

- r is the distance between the base and the mobile unit in miles or kilometers,
- h_e' is the effective antenna height shown in Fig. 5,
- L is the shadow loss shown in Fig. 6.
- P_{r0} is the received signal at a reference distance r_0 ;
- r_0 is usually equal to 1 mile (1.6 km),
- α is a correction factor which will be described later.

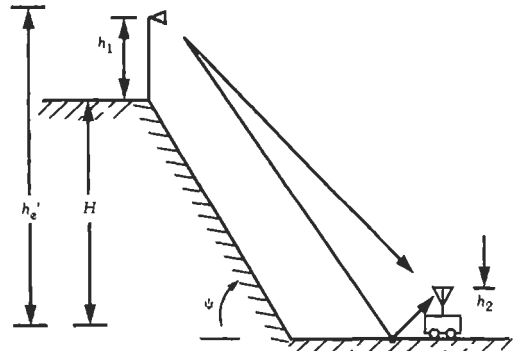
Typical values of P_{r0} and γ are listed in Table 3. If the area of interest is similar to an area listed in the table, the corresponding data from the table can be used for prediction purposes. The suburban data can be used generally for most areas. The values of P_{r0} and γ

TABLE 3. TYPICAL VALUES OF P_{r0} AND γ

Area	P_{r0} (dBm)		γ (dB/dec)
	1 mi	1 km	
Free Space	-45	-40	20
Open Area	-49	-40	43.5
Suburban	-61.7	-54	38.4
Philadelphia	-70	-62.5	36.8
Newark (Urban)	-64	-55	43.1
Tokyo	-84	-78	30.5
New York	-77	-67	48



(A) Mobile unit on hill.



(B) Base station on hill.

Fig. 5. Effective antenna heights, h_e' , for two cases.

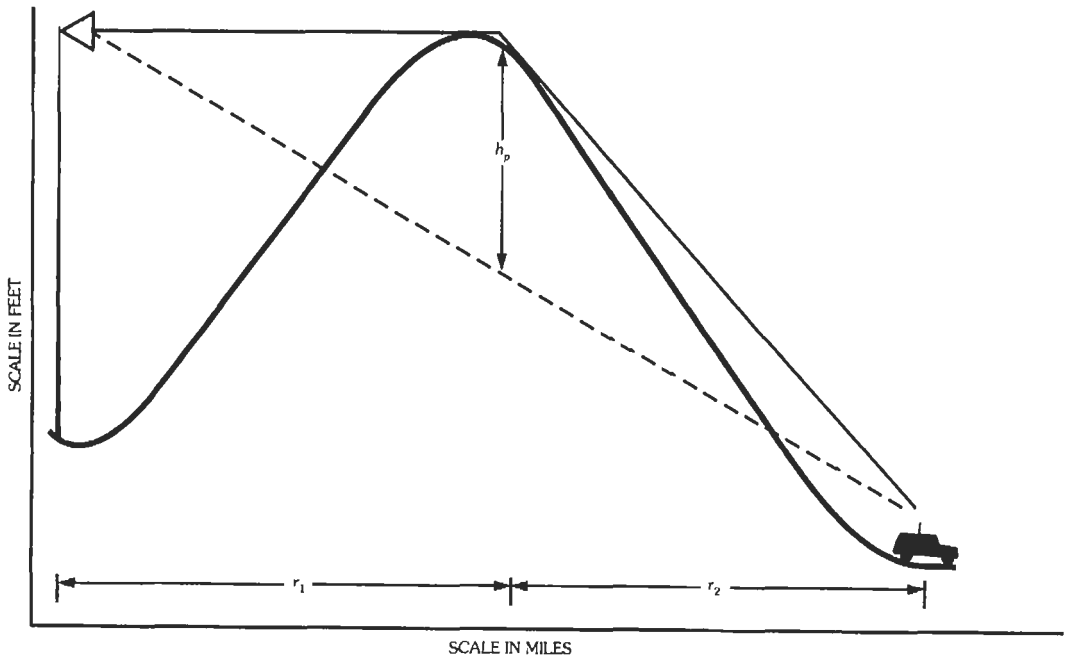


Fig. 6A. Shadow-loss prediction. (A) Prediction model.

also can be obtained very easily by taking the average values of five or more local mean data points at 1 mile and 10 miles. The averaging process attempts to make a real nonflat area into an imaginary flat area. Since the area becomes flat, the different average values of P_{ro} and γ shown for each city are caused solely by the man-made structures. Then the additional changes of the signal strength due to the variation of the terrain

configuration in each area can be added back to the first predicted signal obtained from a flat man-made area. For mobile units traveling along paths in a non-flat area with the terrain contour information provided, the additional effect due to the terrain variation can be predicted according to one of the four conditions listed above. Therefore, in equations 3a and 3b, the first and second parts are due to the man-made structures, and

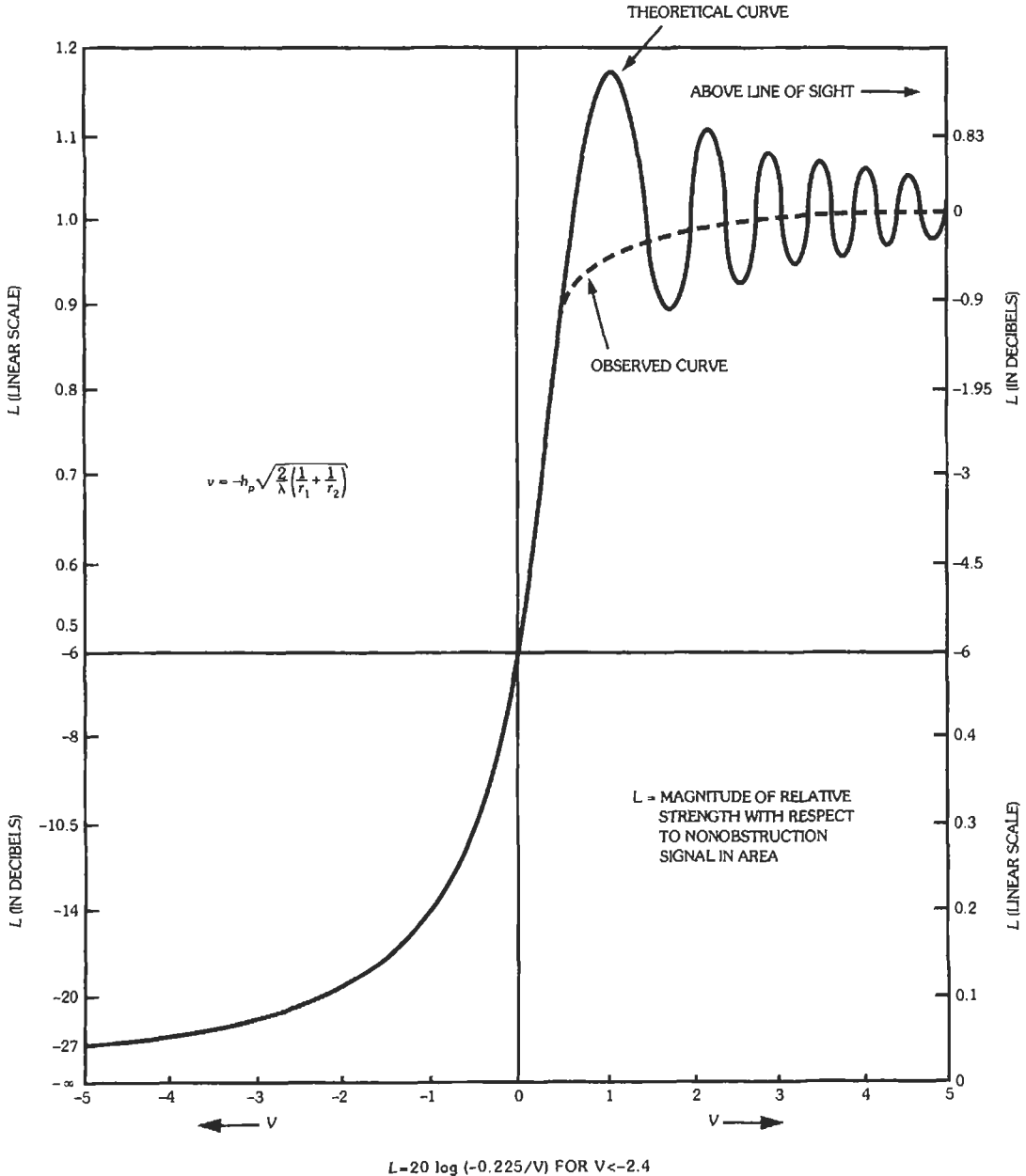


Fig 6B. Shadow-loss prediction. (B) Prediction curve.

the third part is due to the terrain-contour variation. This model can be run on a computer.

Standard Condition: This model can predict the power level in dBm that compares directly with experimental data. For achieving this capability, a standard condition has to be set, as follows:

- Frequency $f_0 = 900$ MHz
- Base-station antenna height (h_1) = 30.48 m (100 ft)
- Base-station power at the antenna = 10 watts (40 dBm)
- Base-station antenna gain $G_1 = 6$ dB above dipole gain (6 dBd)
- Mobile-unit antenna height = 3 m (10 ft)
- $r_o = 1.6$ km (or 1 mi)
- Mobile-unit antenna gain $G_m = 0$ dB above dipole gain (0 dBd)

Correction Factor: For a condition in an actual environment that is different from the standard condition, the following correction factor is used:

$$\alpha = \sum_{i=1}^5 \alpha_i$$

where

- $\alpha_1 = 20 \log (h_1'/30,48)$
- $\alpha_2 = 10 \log (h_2'/3)$
- $\alpha_3 = 10 \log (P_1/10)$ (P_1 is the transmit power)
- $\alpha_4 =$ Base-station antenna gain in dBd, $G_1 - 6$
- $\alpha_5 =$ Mobile-unit antenna gain in dBd, G_m

In the above equations, the heights are in meters, and the power is in watts.

General Formula for Mobile Radio Propagation Path Loss: Based on the suburban-area formula with correction factors inserted, the general formula can be expressed as

$$P_r = P_t - 134.4 - 38.4 \log r_1 + 20 \log h_1 + 10 \log h_2 + G_t + G_m \quad (\text{Eq. 4})$$

where P_r and P_t are in decibels above 1 mW, r_1 is in kilometers, h_1 and h_2 are in meters, and G_t and G_m are in decibels.

Characteristics of Multipath and Selective Fading

Cumulative Distribution Function (CDF)—The CDF of a multipath fading signal follows the Rayleigh statistics shown in Fig. 7 when no line-of-sight (LOS) path condition is present and follows Rician statistics if an LOS path condition is present. Both conditions are shown in Fig. 8. The normalized amplitude of an LOS path is defined as

$$\frac{a_o}{\sqrt{2}} = \frac{\text{amplitude of LOS path}}{\text{rms of the received signal}}$$

The Rayleigh fading condition occurs when $a_o = 0$, i.e., no LOS is present.

Level-Crossing Rate (LCR) and Average Duration of Fades (ADF)—These two statistical parameters are second-order statistics. They involve time, which is a function of the vehicle speed. When the vehicle travels fast, more level crossing and shorter average duration of fades can be seen in the received signal plot.

The LCR curve shown in Fig. 9 normalizes to n_R with the factor,

$$n_0 = \sqrt{2\pi} (V/\lambda)$$

where

- V is the vehicle speed in meters/second,
- λ is the wavelength in meters.

Therefore, after normalization the curve shown in Fig. 9 is independent of the frequency and the vehicle speed. For example, at a level of -10 dB, $n_R = 0.3$ from Fig. 9. Assume that a signal of 850 MHz is received at a mobile unit with a vehicle speed of 24 km/h (15 mi/h). Then $n_0 \approx 50$. The LCR \bar{n} equals

$$\bar{n}(-10 \text{ dB}) = 0.3 \times 50 = 15 \text{ crossings/second}$$

The average duration of fades is shown in Fig. 10. For example, at a level of -10 dB, $t_R = 0.33$. With the same assumption as was used in calculating LCR above, $t_0 = 1/n_0$, and ADF, \bar{t} is:

$$\bar{t}(-10 \text{ dB}) = 0.33/50 = 0.0066 \text{ s} = 6.6 \text{ ms}$$

Distribution of Duration of Fades—The distribution of duration of fades $F_r(U, R)$ is shown in Fig. 11, where $U = \tau/\bar{t}$. The distribution function $F_r(U, R)$ can be interpreted as the probability $P(t > \tau, R(t) < R)^2$ that $R(t) < R$ for an interval lasting longer than τ . The average duration of fades is \bar{t} , and R is the envelope with respect to its rms values. Furthermore, the vehicle speed V can be related to \bar{t} , which is a normalized factor for U . The probability $P(t) > \tau, R(t) < R$ is shown in Fig. 11 as a family of curves for different values of envelope levels. From the curves, we can find that the probability of occurrence of a fade whose duration is greater than twice that of an average duration of fades (ADF) at a level of 3 dB below its rms value, is 10%.

Selective Fading and Time Delay Spread—Selective fading is due to the time-dispersive medium. When an impulse (short pulse) is sent out from the base-station antenna, due to man-made structures the reflected waves travel in different paths and arrive at the mobile-unit antenna at different times. They form a signal $e(t)$. The delay spread (Δ) can be obtained from $e(t)$ as follows:

$$\Delta^2 = \int_0^\infty t^2 \cdot e(t) dt - \int_0^\infty t \cdot e(t) dt$$

In general, the delay spread is measured from the strongest received level to the range which is at its -20 dB

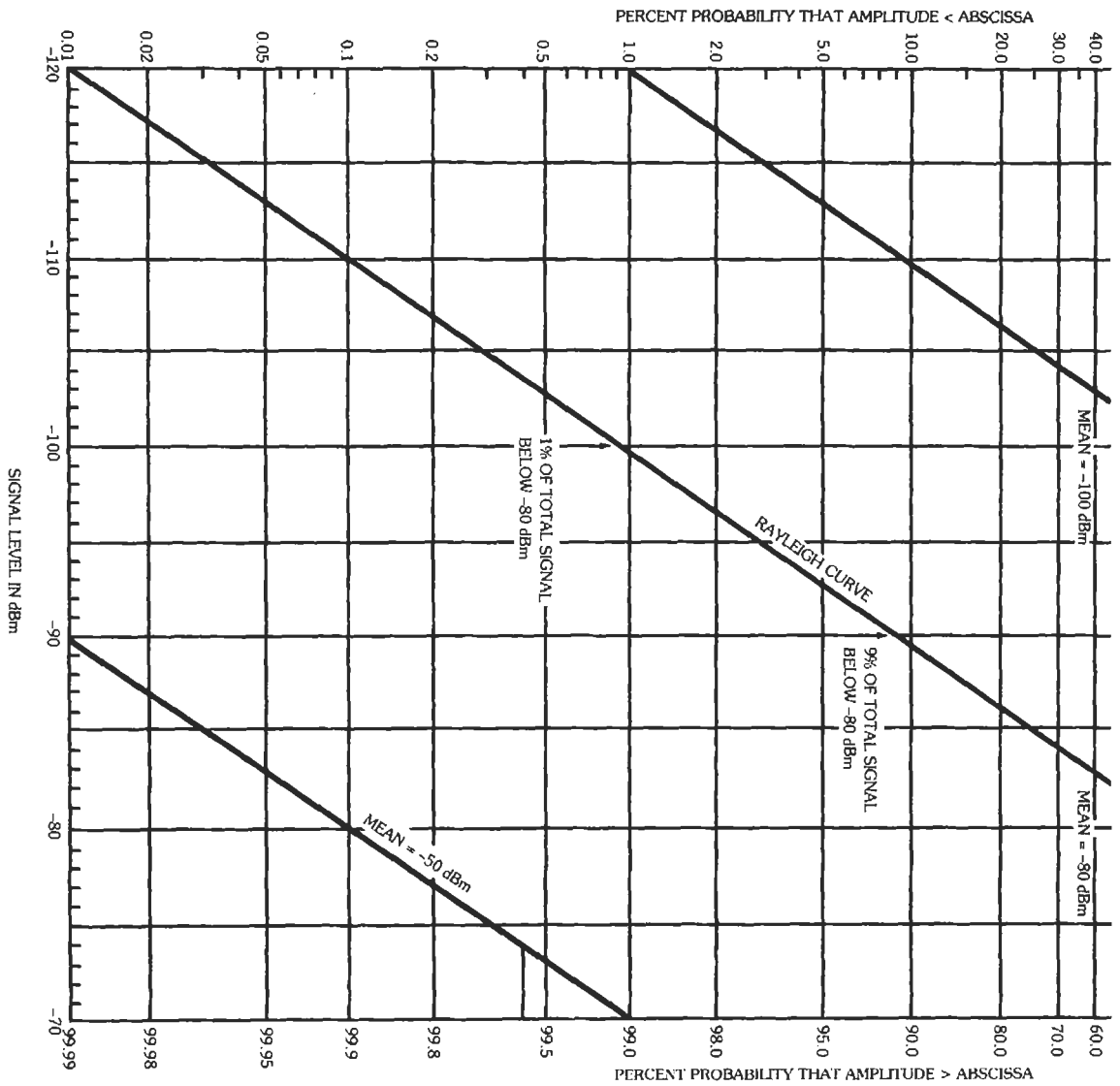


Fig. 7. Cumulative distribution function.

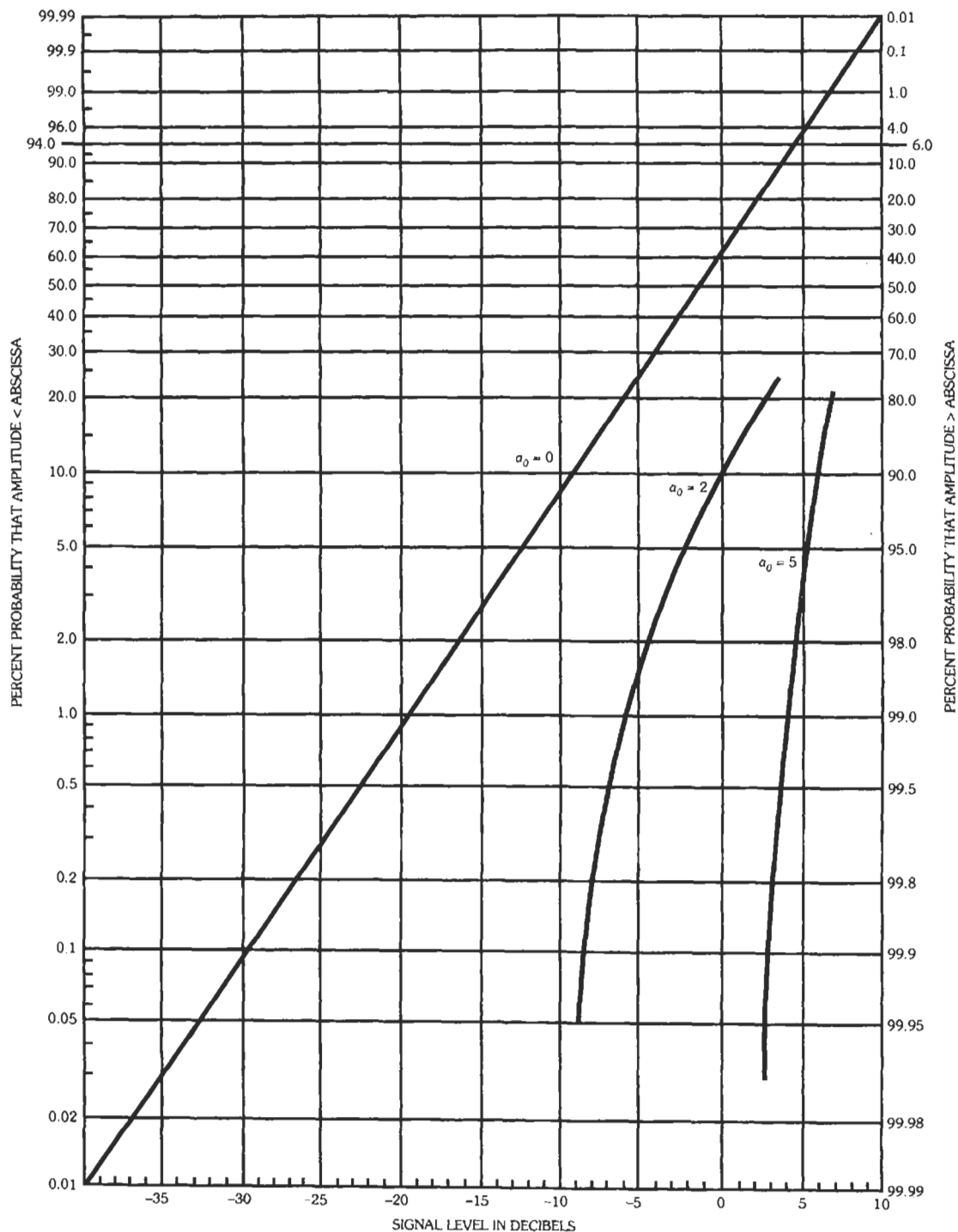


Fig. 8. CPD of Rician distribution. $\left(\frac{a_0}{\sqrt{2}} = \frac{\text{amplitude of direct path}}{\text{rms value}} \right)$

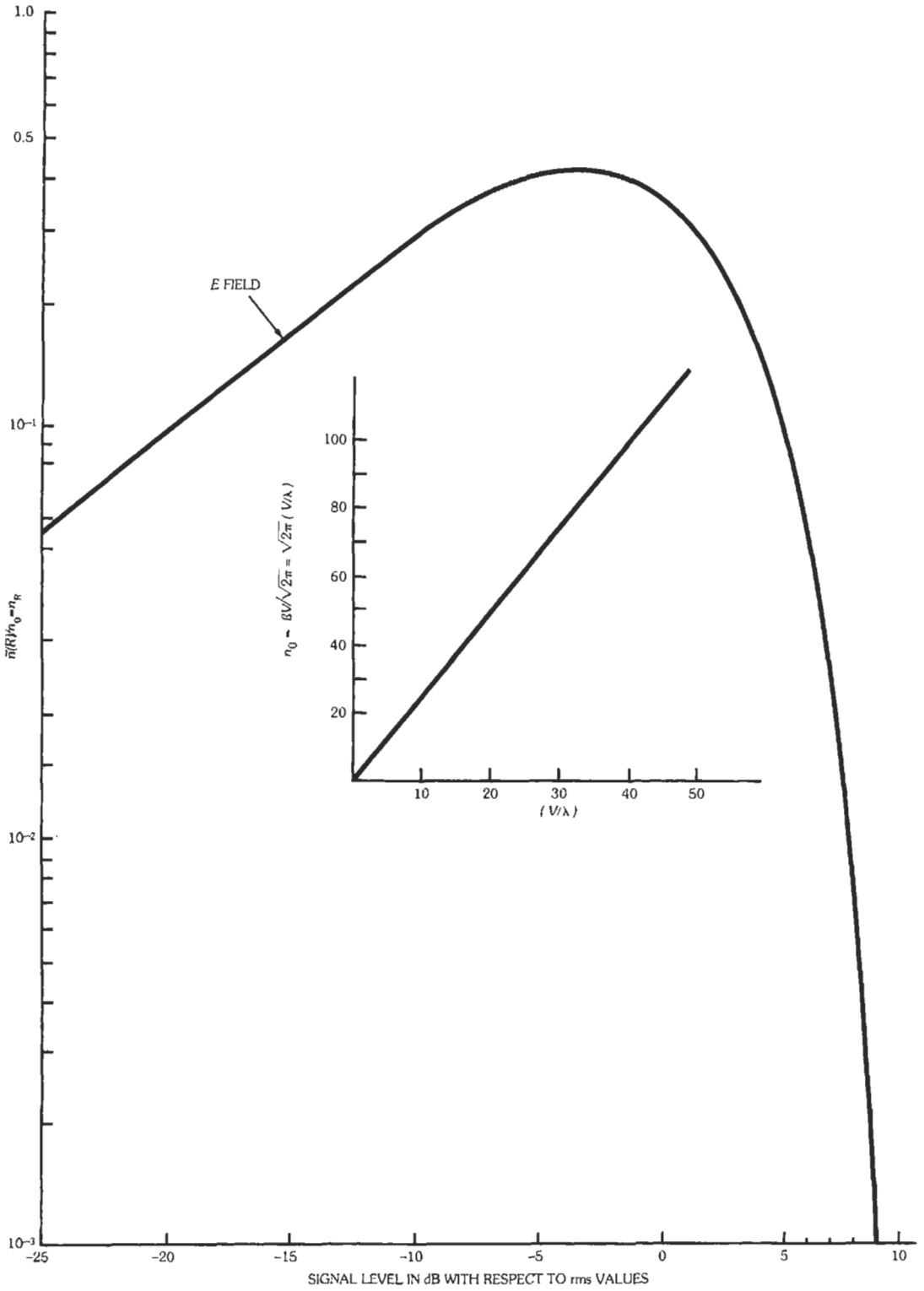


Fig. 9. Level-crossing rate.

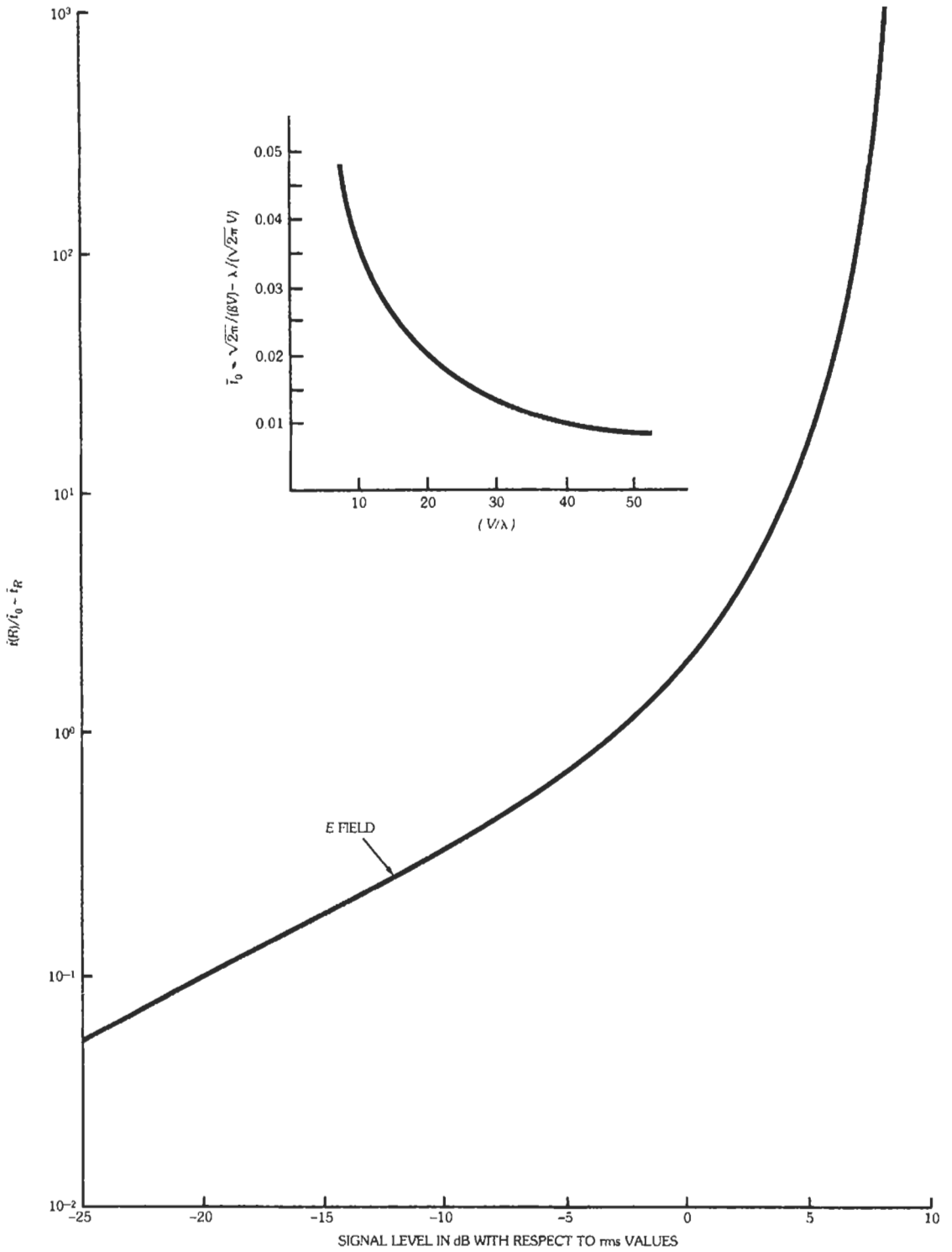


Fig. 10. Average duration of fades.

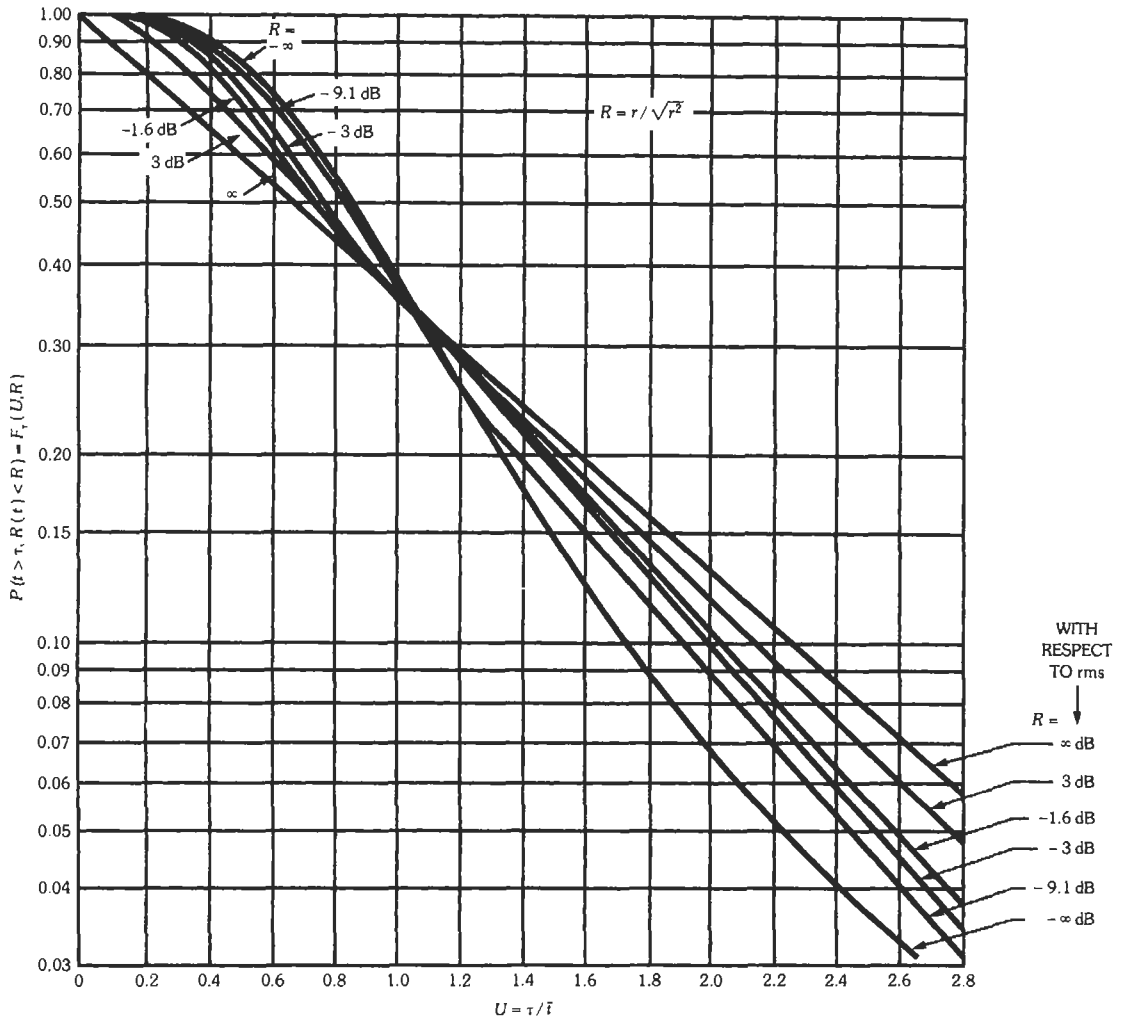


Fig. 11. The probability $F_r(U, R)$ that $R(t) < R$ for an interval lasting longer than τ . $R(t)$ is the envelope of a fading signal.

level. Mean time delay spreads ($\bar{\Delta}$) are shown in Table 4. This data does not include the possibility of reflections from high-rise buildings or nearby mountains, which would generate more time delay spread due to

the late arrival of the reflected waves. Sometimes a time delay caused by a mountain reflection can be up to 100 μs .

TABLE 4. MEAN TIME DELAY SPREADS

Type of Environment	Mean Delay Spread ($\bar{\Delta}$)
Open area	<0.2 μs
Suburban area	0.5 μs
Urban area	0.3 μs
Metropolitan area	0.35 μs
In building	<0.1 μs

FREQUENCY REUSE AND DESIGN ASPECT FOR FDMA AND TDMA

The frequency reuse concept is used in the cellular system to achieve spectrum efficiency and serve a high volume of calls. The frequency reuse concept can be used in the time domain and the space domain. The cellular system applies frequency reuse in the space domain. In many cities, the same frequency is being reused in a general geographical area more than 30 times.

Frequency Reuse Distance

When the same channel is to be reused in two cells, the two cells are called cochannel cells. The distance D is the separation of the two cochannel cells. The smaller the D , the greater the spectrum efficiency would be. The key parameter, a ratio of D/R , where R is the cell radius, is used to measure the spectrum efficiency.

Cellular Design Aspect

The cellular system is a high-capacity and spectrum-efficient system. In the AMPS system, it was found from subjective tests that 75% of the people

grade the system voice quality as “good” or “excellent” at the system’s required carrier-to-interference ratio (C/I), ≥ 18 dB (which is 63 in a linear scale). The level of $C/I = 18$ dB will be set at the boundary of the cell as shown in Fig. 12. The D/R ratio is obtained based on the six cochannel cells surrounding the cell of interest, and the mobile radio propagation rule of 40 dB/decade. Furthermore, assume for simplicity that all the D s are equal. Then

$$(C/I)_s = 63 = \frac{R^{-4}}{\sum_{i=1}^6 D_i^{-4}} \approx \frac{R^{-4}}{6D^{-4}} = \frac{(D/R)^4}{6} \quad (\text{Eq. 5})$$

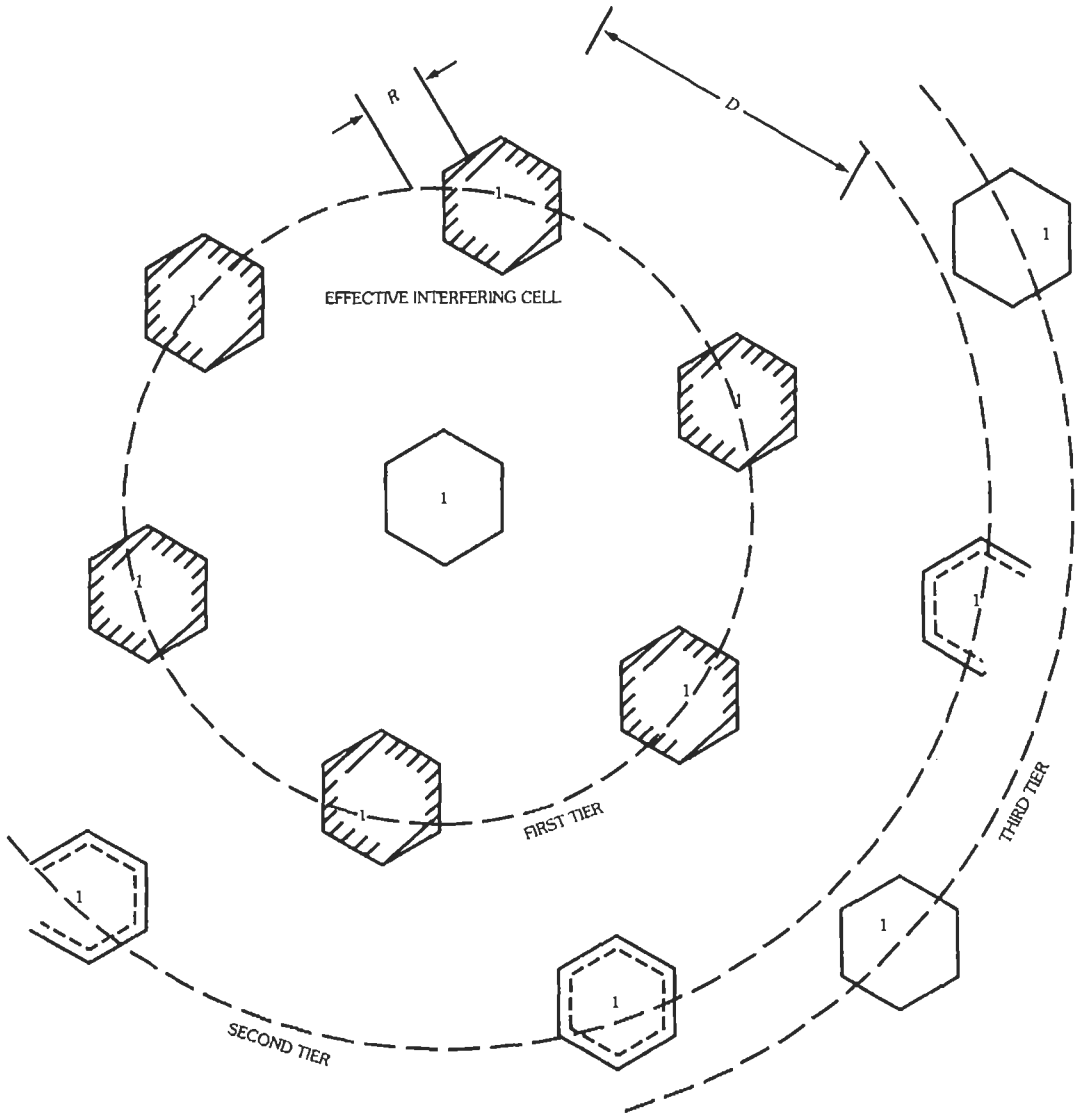


Fig. 12. Six effective interfering cells of cell 1.

or

$$(D/R)_s \approx 4.4$$

The computer simulation finds a required value of $(D/R)_s = 4.6$, which is slightly different from the value of 4.4 which was generated from Eq. 5. The $(D/R)_s$ ratio can be called Cochannel Interference Reduction Factor (CIRF), q ,

$$q = (D/R)_s \tag{Eq. 6}$$

In designing a cellular system, the key factor is q .

Frequency Reuse Factor, K

The factor q can be related to a finite set of cells in a frequency reuse pattern. K is a frequency reuse factor calculated in a hexagonal-shaped cellular system by

$$q = \sqrt{3K} \quad \text{or} \quad K = q^2/3 \tag{Eq. 7}$$

For $q = 4.6$, $K = 7$. This indicates seven frequency reuse patterns (Fig. 13). This means that the total voice channels, say 395 channels, should be divided by seven. Then the number of voice channels per cell is 57.

RADIO CAPACITY*

The radio capacity, m , is used to measure the spectrum efficiency.

$$m = B_t / \left[B_c \sqrt{\frac{2}{3}} (C/I)_s \right] \quad \text{(for FDMA and TDMA)} \tag{Eq. 8a}$$

$$= M/K \quad \begin{array}{l} \text{number of voice channels/cell} \\ \text{(for CDMA, } K \approx 1) \end{array} \tag{Eq. 8b}$$

where

B_t = total allocated spectrum,

B_c = channel bandwidth,

M = total number of voice channels,

K = frequency reuse factor,

$(C/I)_s$ = required carrier-to-interference ratio in linear values.

In an analog cellular system, B_c and $(C/I)_s$ are related. For example, to keep an accepted voice quality, $B_c = 30$ kHz and $(C/I)_s = 18$ dB. If B_c is reduced to 15 kHz, then $(C/I)_s$ has to be increased to 24 dB in order to maintain the same voice quality. Then the radio capacity (m) does not change. Equation 8a is used in analog, frequency division multiple access (FDMA), and time division multiple access (TDMA) cellular systems. Equation 8b is used in code division multiple access (CDMA) cellular systems.

Increase Capacity

Increase Radio Capacity—Try to increase capacity by designing a new digital system the voice quality of which can cope with a narrow bandwidth and a low required $(C/I)_s$. The radio capacity of Eq. 8 can be used to compare the spectrum efficiencies of any two different systems.

Increase System Capacity: Cell Splitting—The system capacity is measured by the total calling channels in a given geographical area. One physical frequency channel can become ten calling channels if the frequency is reused ten times. The system capacity can be increased without increasing radio capacity. Cell splitting uses this type of approach. When a cell radius is reduced to one-half the original radius, four small cells can be created. Each small cell has the same radio capacity as the original one. The system capacity then increases by four times. When a cell is small, the effect of interference on the performance becomes greater. Splitting cells will not achieve four times system capacity.

Microcells—A cell of radius less than 1 mile or 1 kilometer is called a microcell. A conventional microcell system can still increase system capacity based on cell splitting, but only by two or three times; it cannot increase the radio capacity. Also, a conventional microcell system may not be able to provide the fast handoffs required for small-cell operation. A newly invented microcell system† cannot only increase the system capacity but can also increase the radio capacity by reducing the frequency reuse factor K from 7 to 3. The key factor in increasing the radio capacity is to reduce the interference by simply aiming the signal at the mobile unit.

Diversity Schemes and Combining Techniques

General Description—In a mobile radio environment, the multipath fading and the time delay spread generate major problems for cellular communications. Diversity schemes provide the means to generate two or more uncorrelated fading signals at the mobile reception point. The information carried by the two different received signals is the same, but the multipath fading, at any given time, is different for the two signals. The cross correlation coefficient ρ , where $\rho \leq 1$, is used to measure their correlation. If $\rho = 1$, then the two fading signals are totally correlated. If $\rho = 0$, then the two fading signals are uncorrelated. The ideal reduction of signal fading at reception, after combining two fading signals, occurs when ρ approaches zero.

* Reference 2.

† Reference 3.

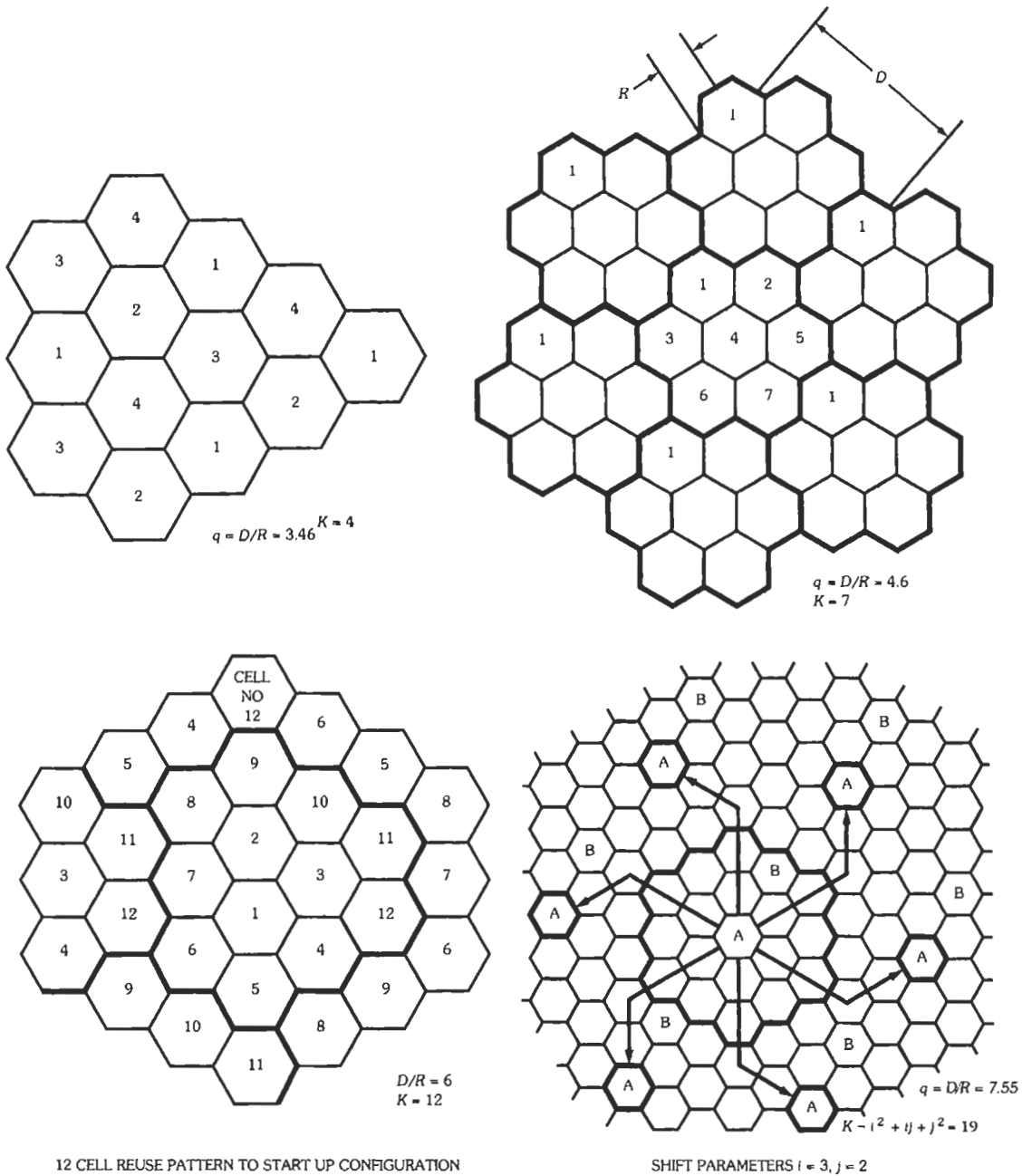


Fig. 13. *N*-cell reuse patterns.

Diversity Schemes—There are six common diversity schemes, as follows:

1. Space diversity (receives two fading signals from two separate antennas)
2. Field component diversity (receives E field and H field)
3. Polarization diversity (receives vertically and horizontally polarized waves, or left-hand and right-hand circularly polarized waves)

4. Time diversity (repeats the same signal more than once)
5. Frequency diversity (receives two fading signals on two different frequencies)
6. Angle diversity (receives two fading signals from two different directions)

Combining Techniques—There are three general techniques for combining two or more uncorrelated or partially correlated received fading signals. The signal fading is reduced after combination. The best method is the maximum-ratio combiner, which combines the signals to achieve the maximum signal-to-noise ratio. The next best technique is the equal-gain combiner. It combines the two fading signals in phase to reduce the signal fading. The third technique is the selective combiner. This technique always selects the strongest signal from among the received fading signals. The difference in performance between the maximum-ratio technique and the selective-combiner technique is 2 dB. The switched combined technique which switches the signal when it is below a threshold level is not applicable in cellular systems.

Antenna Separation Requirement—*At the Base Station:* The antenna separation for cellular systems is determined by the formula

$$h/d = 11$$

where

h is the antenna height,
 d is the spacing between two antennas.

If $h = 30$ m, then $d = 2.72$ m. The formula does not apply if the location of a mobile unit is in line with the two base-station antennas. Therefore, the orientation of two antennas mounted on a mast has to be carefully determined to obtain the maximum diversity gain in the cell.

At the Mobile Unit: A separation of a half wavelength between two mobile antennas is required at 850 MHz. Therefore, the separation between two antennas needs to be only 0.18 m (about 6 inches) at the cellular frequency of 850 MHz.

Diversity Scheme Reduces Time Delay Spread—The diversity scheme not only reduces the signal fading, it also reduces the time delay spread. Time delay spread does not affect an analog system, but it does affect a digital system. The time delay spread in a mobile radio environment always exists at the reception point whether the mobile unit is standing still or moving. Signal fading occurs only when the mobile unit is moving. In a digital system, an equalizer is normally used to reduce the intersymbol interference which is caused by the time delay spread. The strength of the equalizer equipment is based on the severity of the mobile radio environment. When the diversity scheme has been implemented, the equalizer design requirement can be relaxed or removed.

SWITCHING EQUIPMENT AND TRAFFIC MODELS

Analog Switching Equipment

The analog switch is usually a circuit switch which must hold a call throughout the duration of the call. The analog switching equipment consists of processor, memory, switching network, trunk circuitry, and miscellaneous service circuitry as shown in Fig. 14. The switching of a Class 5 telephone central office is switching from trunk frame to line frame. The difference in this cellular switching equipment is the switching from trunk frame to trunk frame via the switch network. One trunk side is connected to the mobile calls because the mobile unit does not have a fixed frequency channel associated with a mobile unit. The other trunk side is connected to PSTN.

Cellular Digital Switching Equipment

The digital switch is usually a message switch handling digitized messages. The digital switch can send the message or transmit the voice in digital form. Therefore, the digital signal format permits breaking a message into small pieces for faster transmission. Also, it can handle other calls while the switch alternates between the “on” and “off” modes periodically. Hence, the call-processing efficiency of digital switching is higher than that of analog switching. The other advantages of using a digital switch are small size, less power consumption, less human effort required to operate, and ease of maintenance. Digital switching equipment can be modular and is very flexible. Often, capacity can be increased simply by adding modules. Digital equipment can be either a centralized system (Fig. 15A) or a decentralized system (Fig. 15B).

Traffic Models

There are three telephone traffic models (see Chapter 39):

- Erlang B—Blocked calls clear condition specified blocking probability
- Erlang C—Blocked calls delay condition specified delay probability
- Poisson—Blocked calls held condition specified the held probability at a time period equal to an average holding time.

In the cellular industry, the system operators usually specify the blocking probability for designing the system and calculate the number of radios in each cell site to meet this requirement based on the Erlang B formula. Some operators are more conservative and use the Erlang C formula to calculate the number of radios. However, they have to substitute the value of the specified blocking probability into the delay proba-

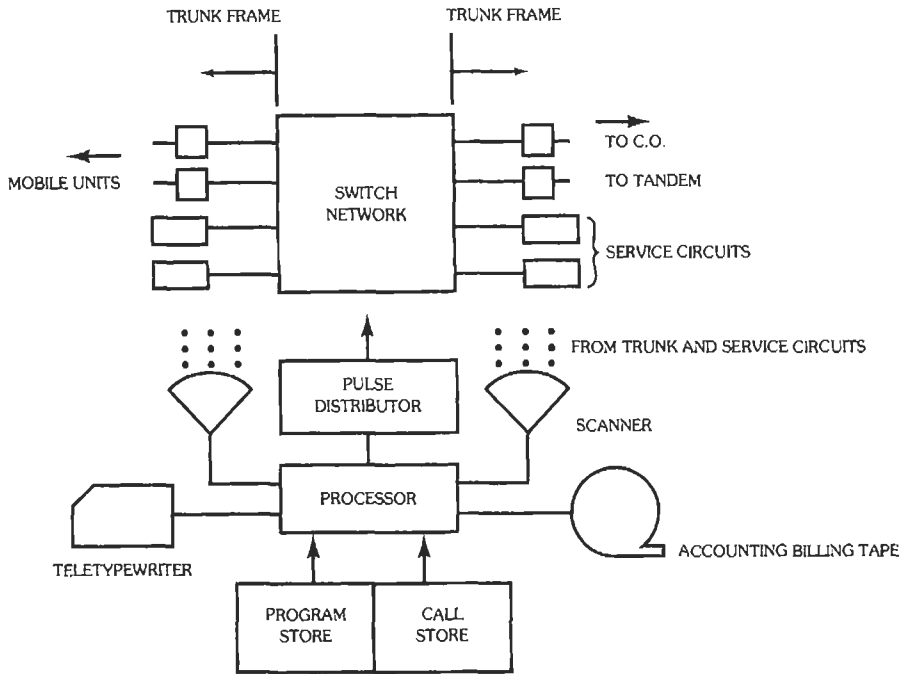


Fig. 14. A modified analog switching system for cellular mobile systems.

bility in the Erlang C formula. Also, by substituting the specified blocking probability into the held probability of the Poisson formula, another different number of radios would be required. Among the three formulas, any one can be used within an error in general less than 15%. In the future, a cellular traffic model will be generated based on the past traffic data in cellular telecommunication systems.

EXAMPLE: A cell site has 57 voice channels. Assume a blocking probability of 2%. From the Erlang B table, the offered load $A = 46.8$ erlangs. Assume the average talking time per call is 100 s; then the number of calls per busy hour, Q_1 , is

$$Q_1 = (46.8 \times 3600)/100 = 1685 \text{ calls/busy hour}$$

and the number of calls per radio, Q_2 , becomes

$$Q_2 = 1685/57 = 29.5 \text{ calls/radio}$$

DIGITAL CELLULAR SYSTEMS

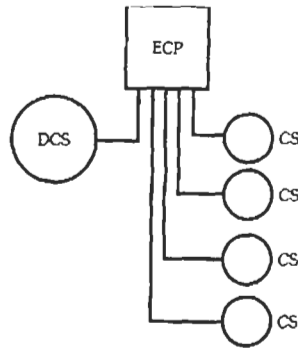
Digital AMPS (or North America Digital Cellular, NADC)

- Voted cellular digital standard at TIA
- Allocated the same spectrum as analog cellular system
- Requires a dual-mode (analog/digital) subscriber's unit

- Handoff requirements— analog to digital, digital to digital, digital to analog
- Number of channels—416 channels for Band A and 416 channels for Band B
- 21 analog setup channels are used for digital system
- Channel bandwidth—30 kHz
- Multiple access method—TDMA (3 time slots per channel)
- Modulation— $\frac{\pi}{4}$ -DQPSK
- Detection—differential detection
- Speed coding—VSELP (vector sum excited LPC)
- Speech frame—20 ms/frame
- Speech code rate—8 kbps
- Channel coding—interleaving convolutional code
- Equalizer—equalize the time delay spread up to 60 μ s
- Duration of slots—6.6 ms

GSM (Special Mobile Group) System

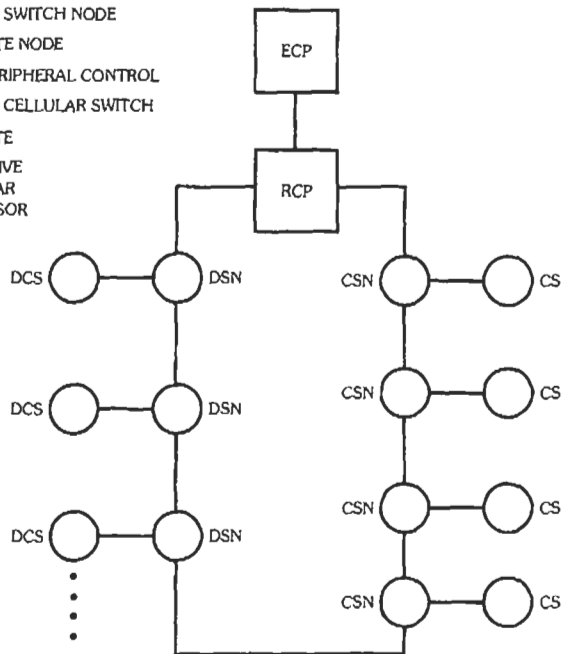
- GSM is a European digital cellular standard.
- GSM is allocated a spectrum exclusively.
- Frequency
 - 935–960 MHz base-to-mobile
 - 890–915 MHz mobile-to-base
- Number of channels—124 radio carriers



CS: CELL SITE
 ECP: EXECUTIVE CELLULAR PROCESSOR
 DCS: DIGITAL CELLULAR SWITCH

(A) Centralized system.

DSN: DIGITAL SWITCH NODE
 CSN: CELL SITE NODE
 RPC: RING PERIPHERAL CONTROL
 DCS: DIGITAL CELLULAR SWITCH
 CS: CELL SITE
 ECP: EXECUTIVE CELLULAR PROCESSOR



(B) Decentralized system.

Fig. 15. Cellular switching equipment.

- Channel bandwidth—200 kHz
- Multiple access method—TDMA (8 time slots per channel)
- Modulation—GMSK
- Detection—coherent detection
- Speech coding—RPE-LPC (regular pulse excited-LPC)
- Speech frame—20 ms/frame

- Speech coding rate—13 kbps
- Channel coding—convolutional code
- Equalizer—equalize the time delay spread up to 16 μ s
- Transmission rate—270 kbps
- Duration of slots—0.557 ms (the frame of 8 slots is 4.615 ms)

Cellular CDMA System

- This is a code division multiple access (CDMA) system.
- CDMA is in the process of being standardized by TIA.
- Bandwidth—1.23 MHz
- Modulation—BPSK
- $E_b/I_0 = 7$ dB
- Speech coding rate—8 kbps (variable rate for nonvoice conditions)
- Power control increment—0.5 dB
- Channel coding—convolutional code
- Detection—Forward channel uses coherent detection; reverse channel uses noncoherent detection.
- Transmit power
 - cell site—1.25 W
 - mobile—300 mW
- Range—13 miles

3G Systems

Three modes are FDD direct spread, FDD multicarrier, and TDD as listed in Table 2.

In FDD modes, several high-speed data channel systems have been introduced as shown in Table 5. Depending on different modulations, the peak data rate can be from a range of 2.4 to 5.2 Mbps.

4G System

A TDD system with a set of smart code (LAS code)[†] can reduce the multiuser access interference and the multi-path interference along a single-spectrum band transmission. It brings a high hope that using a code technology to change the system environment from an interference environment to a noise environment is a winner for the 4G system.

Other Personal Communication Systems

PCN (Personal Communication Network)

- Three licenses in UK
- Frequency—50–100 MHz band in 1.8 GHz range
- System—modified GSM

DECT (Digital European Cordless Telephone)

- Frequency—1880–1900 MHz
- Multiple access—TDMA/TDD

- Transmission rate—38.8 kbps/slot channel
- Modulation—GFSK
- Channel bandwidth—1.728 MHz/channel
- TDD (Time division duplex)
 - 12 slots (base-to-portable)
 - 12 slots (portable-to-base)
- Frame length—10 ms/24 time slots
- Number of channels—5
- Transmit power—10 mW

CT-2 (Cordless Telephone—2)

This system acts as a “portable payphone booth”

- Applied country—four licenses issued in UK
- Transmission signal format—TDD (time division duplexing)
- Number of channels—40
- Channel bandwidth—100 kHz/TDD channel
- No call delivery capability
- No handoffs

REFERENCES

1. Lee, W. C. Y. “Lee’s Model,” *IEEE Transactions on Vehicular Technology*, Feb. 1988, pp. 69–71. Also in Chapter 4 of *Mobile Cellular Telecommunications Systems*. McGraw-Hill, 1989.
2. Lee, W. C. Y. “Spectrum Efficiency in Cellular.” *IEEE Transactions on Vehicular Technology*, May 1989, pp. 69–75.
3. Lee, W. C. Y. “A New Microcell Architecture.” *IEEE Communications Magazine*, Nov. 1991, pp. 19–23.
4. Lee, W. C. Y. *Lee’s Essentials of Wireless Communications*. McGraw-Hill, 2001, p. 289

BIBLIOGRAPHY

1. Bell Laboratories. *High Capacity Mobile Telephone System Technical Report*. Dec. 1971, submitted to the FCC.
2. Blecher, F. H. “Advanced Mobile Phone Service.” *IEEE Transactions on Vehicular Technology*, Vol. VT-29, May 1980, pp. 238–244.
3. *Bell System Technical Journal*, special issue on “Advanced Mobile Phone Service (AMPS),” Vol. 58, Jan. 1979.
4. EIA Interim Standard, *Cellular System Mobile Station—Land Station Compatibility Specification*, IS-3-D EIA, March 1987.
5. TIA Interim Standard, IS-54, *Digital Cellular System Mobile Station—Land Station Compatibility Specification*, Jan., 1990.
6. TIA Interim Standard, IS-55, *Digital Cellular System Mobile Station Specification*. Jan. 1991.
7. TIA Interim Standard, IS-56, *Digital Cellular System Land Station Specification*. Jan. 1991.
8. GSM—European Digital Cellular Standard.
9. Qualcomm Presentation at CTIA, “Cellular CDMA,” July 21, 1991.

TABLE 5. SEVERAL HIGH-SPEED DATA CHANNEL SYSTEMS

	HDR*	1XTREME	E-1XRTT	1XPLUS
Proposed by	Qualcomm	Motorola, Nokia	Nortel	Ericsson
Peak FL Data Rate	2.4 Mbps	5.2 Mbps	3.7 Mbps	2.4 Mbps
Peak RL Data Rate	307.2 kbps	614.4 kbps	460.8 kbps	307.2 kbp
Average Throughput	~600 kbps	~400 kbps	~600 kbps	~600 kbps
FL Access Method	TDM (fat pipe)	CDM (shared pipe)	TDM (fat pipe)	TDM
FL Modulation	QPSK, 8-PSK, 16-QAM	QPSK, 8-PSK, 16-QAM, 64-QAM	FL: QPSK, 8-PSK, 16-QAM, 32-QAM	Same as HDR
RL Modulation	BPSK	BPSK	BPSK	BPSK
FL Power Control	None, use full power available	Closed loop @ 800 Hz	None, use full power available	Same as HDR
RL Power Control	Open/closed loop @ 800 Hz	Open/closed loop @ 800 Hz	Open/closed loop @ 800 Hz	Same as HDR
Voice Service Support	No	Yes	Yes if concurrent with data	Yes
Transmit Diversity	None specified	Space Time Block Coding	Open-loop and closed-loop array technologies	None specified
Net. Arch - RLP/MAC	Decentralized	Not discussed	Decentralized	Decentralized
Net. Arch - Mob. Mgmt.	Decentralized	Not discussed	Centralized	Centralized
Mobility	Mobile	Nomadic	Mobile	Mobile

* The system had a public trial in San Diego in November 1999.

10. Lee, W. C. Y. "Overview of Cellular CDMA." *IEEE Transactions on Vehicular Technology*. Special issue on digital cellular, May 1992.
11. Lee, W. C. Y. *Mobile Communications Design Fundamentals*. 1st ed. Indianapolis: Howard W. Sams & Co., 1986. 2nd ed. New York: John Wiley & Sons, Inc., 1993.
12. Lee, W. C. Y. *Mobile Cellular Telecommunications*, 2nd ed. McGraw-Hill, 1995.
13. Lee, W. C. Y. *Mobile Communications Engineering*, 2nd ed. McGraw-Hill, 1998.
14. Lee, W. C. Y. *Lee's Essentials of Wireless Communications*. McGraw-Hill, 2001.

47

Mathematical Equations

Mensuration Equations 47-4

Areas and Lengths Associated With Plane Figures
Surface Areas and Volumes of Solid Figures

Algebraic and Trigonometric Equations (Including Complex
Quantities) 47-6

Quadratic Equation
Solution of Cubic Equations
Solution of Quartic Equations
Complex Quantities
Properties of e
Properties of Logarithms
Sums
Combinations and Permutations
Bernoulli Numbers
Trigonometric Identities
Approximations for Small Angles
Inequalities

Plane Trigonometry 47-11

Right Triangles
Oblique Triangles

Spherical Trigonometry 47-12

Right Spherical Triangles ($\gamma = 90^\circ$)
Oblique Triangles

Hyperbolic Functions 47-13

Hyperbolic Trigonometry 47-14

Conformal Model
Projective Model
Right Hyperbolic Triangles ($\gamma = 90^\circ$)
Oblique Hyperbolic Triangles

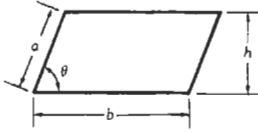
Plane Analytic Geometry	47-16
Straight Line	
Transformation of Rectangular Coordinates	
Circle	
Parabola	
Ellipse	
Hyperbola	
Solid Analytic Geometry	47-19
Coordinates	
Direction Cosines and Numbers	
Distance Between Two Points	
Equations of a Plane	
Equations of a Straight Line	
Ellipsoid	
Prolate Spheroid	
Oblate Spheroid	
Paraboloid of Revolution	
Hyperboloid of Revolution	
Differential Calculus	47-21
List of Derivatives	
Table of Integrals	47-21
Indefinite Integrals	
Logarithmic Integrands	
Exponential Integrands	
Trigonometric Integrands	
Definite Integrals	
Series	47-29
Taylor's Series for a Single Variable	
Taylor's Series for Two Variables	
Miscellaneous Series	
Matrix Algebra	47-31
Notation	
Operations	
Determinant	
Linear Transformations	
Eigenvectors and Eigenvalues	
Further Definitions and Properties	
Hermitian Forms	

Vector-Analysis Equations	47-35
Rectangular Coordinates	
Gradient, Divergence, Curl, and Laplacian in Coordinate Systems	
Other Than Rectangular	
Space Curves	
Laplace Transform	47-38
Table of Laplace Transforms	47-38
General Equations	
Miscellaneous Functions	
Inverse Transforms	
Selected Functions	47-39
Exponential Integrals	
Cosine and Sine Integrals	
Gamma Function	
Psi and Polygamma Functions	
Error Function	
Fresnel Integrals	
Elliptic Integrals	
Bessel Functions	
Orthogonal Polynomials	
Numerical Analysis	47-43
Algorithms for Solving $F(x) = 0$	
Algorithm for Solving $F(x, y) = G(x, y) = 0$	
Interpolation Polynomial	
Interpolation at Equidistant Points	
Integration	
Differentiation	
Error in Arithmetic Operations	

MENSURATION EQUATIONS

Areas and Lengths Associated With Plane Figures

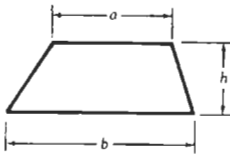
Parallelogram:



$$\text{Area} = bh$$

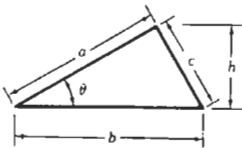
$$= ab \sin \theta$$

Trapezoid:



$$\text{Area} = \frac{1}{2}h(a + b)$$

Triangle:



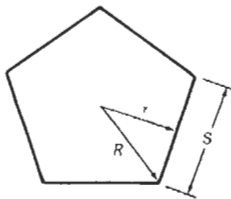
$$\text{Area} = \frac{1}{2}bh$$

$$= \frac{1}{2}ab \sin \theta$$

$$= [s(s - a)(s - b)(s - c)]^{1/2}$$

$$\text{where } s = \frac{1}{2}(a + b + c)$$

Regular Polygon:



$$\text{Area} = \frac{1}{2}nrS$$

$$= nr^2 \tan (180^\circ/n)$$

$$= \frac{1}{4}nS^2 \cot (180^\circ/n)$$

$$= \frac{1}{2}nR^2 \sin (360^\circ/n)$$

where,

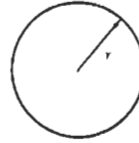
n = number of sides,

S = length of one side,

R = long radius,

r = short radius = $R \cos(180^\circ/n) = \frac{1}{2}S \cot(180^\circ/n)$.

Circle:



$$\text{Area} = \pi r^2$$

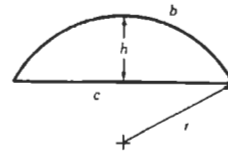
$$\text{Circumference} = 2\pi r$$

where,

r = radius,

$\pi = 3.1416$.

Segment of Circle:



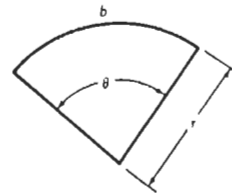
$$\text{Area} = \frac{1}{2}[br - c(r - h)]$$

where,

b = length of arc,

c = length of chord = $[4(2hr - h^2)]^{1/2}$.

Sector of Circle:



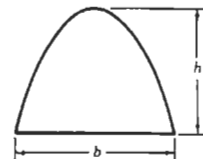
$$\text{Area} = br/2 = \pi r^2(\theta/360^\circ)$$

where,

$b = (\pi r \theta / 180^\circ)$,

θ is in degrees.

Parabola:

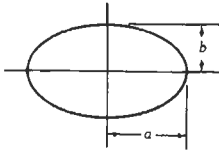


$$\text{Area} = \frac{2}{3}bh$$

$$\text{Arc length} = [4h^2 + (b^2/4)]^{1/2}$$

$$+ (b^2/8h) \ln \{4h + 2[4h^2 + (b^2/4)]^{1/2}/b\}$$

Ellipse:

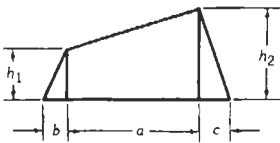


Area = πab

Circumference = $4aE(k)$

where $E(k)$ is a complete elliptic integral with $k = (a^2 - b^2)^{1/2}/a$ of the second kind, $a > b$.

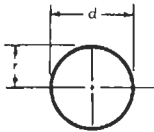
Trapezium:



Area = $\frac{1}{2}[a(h_1 + h_2) + bh_1 + ch_2]$

Surface Areas and Volumes of Solid Figures

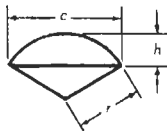
Sphere:



Surface = $4\pi r^2 = 12.5664r^2 = \pi d^2$

Volume = $(4\pi r^3/3) = 4.1888r^3$

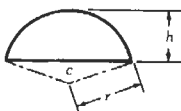
Sector of Sphere:



Total surface = $(\pi r/2)(4h + c)$

Volume = $(2\pi r^2 h/3) = 2.0944r^2 h$
 $= (2\pi r^2/3)[r - (r^2 - \frac{1}{4}c^2)^{1/2}]$
 $c = [4(2hr - h^2)]^{1/2}$

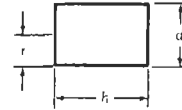
Segment of Sphere:



Spherical surface = $2\pi rh = \frac{1}{4}\pi(c^2 + 4h^2)$

Volume = $\pi h^2(r - \frac{1}{3}h)$
 $= \pi h^2[(c^2 + 4h^2)/8h - \frac{1}{3}h]$
 $c = [4(2hr - h^2)]^{1/2}$

Cylinder:



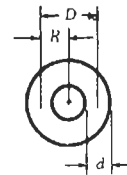
Cylindrical surface = $\pi dh = 3.1416dh$

Total surface = $2\pi r(r + h)$

Volume = $\pi r^2 h = 0.7854d^2 h$
 $= c^2 h/4\pi = 0.0796c^2 h$

$c =$ circumference

Torus or Ring of Circular Cross Section:



Surface = $4\pi^2 Rr = 39.4784Rr = 9.8696Dd$

Volume = $2\pi^2 Rr^2 = 19.74Rr^2 = 2.467Dd^2$

where,

$D = 2R =$ diameter to centers of cross-section of torus,

$r = d/2$.

Pyramid:



Volume = $Ah/3$

When base is a regular polygon

Volume = $\frac{1}{3}h\{nr^2[\tan(360^\circ/2n)]\}$
 $= \frac{1}{3}h\{\frac{1}{4}(ns^2)[\cot(360^\circ/2n)]\}$

where,

$A =$ area of base,

$n =$ number of sides,

$r =$ short radius of base.

See "Regular Polygon" in subsection on plane figures.

Pyramidal frustum:

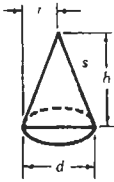


$$\text{Volume} = \frac{1}{3}h[a + A + (aA)^{1/2}]$$

where,

A = area of base,
 a = area of top.

Cone With Circular Base:

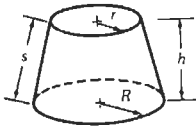


$$\text{Conical area} = \pi r s = \pi r^2 + h^2)^{1/2}$$

$$\text{Volume} = \pi r^2 h / 3 = 1.047 r^2 h = 0.2618 d^2 h$$

where s = slant height.

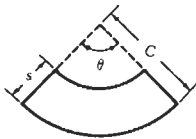
Conic Frustum:



$$\text{Volume} = (\pi h / 3)(R^2 + Rr + r^2)$$

$$\text{Area of conic surface} = \pi s(R + r)$$

Development of Conic Surface:

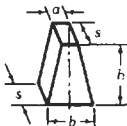


$$C = sR / (R - r)$$

$$\theta = 360R / C$$

where θ is in degrees.

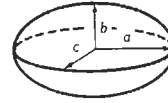
Wedge Frustum:



$$\text{Volume} = \frac{1}{2}hs(a + b)$$

where h = height between parallel bases.

Ellipsoid:



$$\text{Volume} = (4\pi abc / 3)$$

$$\text{Surface} = 2\pi\{c^2 + [b/(a^2 - c^2)^{1/2}][c^2 F(\phi, k)$$

$$+ (a^2 - c^2)E(\phi, k)] \text{ if } a > b > c$$

$$= 2\pi a\{a + [c^2/(a^2 - c^2)^{1/2}]$$

$$\times \ln[a + (a^2 - c^2)^{1/2}/c]$$

$$\text{if } a = b > c \text{ (oblate ellipsoid)}$$

$$= 2\pi c\{c + a^2(a^2 - c^2)^{-1/2} \arcsin[(a^2 - c^2)^{1/2}/a]$$

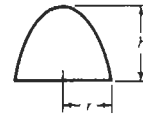
$$\text{if } a > b = c \text{ (prolate ellipsoid)}$$

$F(\phi, k)$ and $E(\phi, k)$ are incomplete elliptic integrals of the first and second kinds, respectively.

$$\phi = \arcsin [(a^2 - c^2)^{1/2}/a]$$

$$k = (a/b)[(b^2 - c^2)/(a^2 - c^2)]^{1/2}$$

Paraboloid:



$$\text{Volume} = (\pi r^2 h / 2)$$

$$\text{Curved surface} = (\pi r / 6h^2)[(r^2 + 4h^2)^{3/2} - r^3]$$

ALGEBRAIC AND TRIGONOMETRIC EQUATIONS (INCLUDING COMPLEX QUANTITIES)

Quadratic Equation

If $ax^2 + bx + c = 0$, then

$$x = [-b \pm (b^2 - 4ac)^{1/2}] / 2a$$

$$= 2c[-b \mp (b^2 - 4ac)^{1/2}]$$

Solution of Cubic Equations*

Given $z^3 + a_2z^2 + a_1z + a_0 = 0$, let

$$q = \frac{1}{3}a_1 - \frac{1}{9}a_2^2; \quad r = \frac{1}{6}(a_1a_2 - 3a_0) - \frac{1}{27}a_2^3$$

* Abramovitz, M. and Stegun, I. A. *Handbook of Mathematical Functions*. Washington, D. C.: National Bureau of Standards; p. 17.

If $q^3 + r^2 > 0$, one real root and a pair of complex conjugate roots

$q^3 + r^2 = 0$, all roots real and at least two are equal

$q^3 + r^2 < 0$, all roots real (irreducible case)

Let

$$s_1 = [r + (q^3 + r^2)^{1/2}]^{1/3}$$

$$s_2 = [r - (q^3 + r^2)^{1/2}]^{1/3}$$

then

$$z_1 = (s_1 + s_2) - (a_2/3)$$

$$z_2 = -\frac{1}{2}(s_1 + s_2) - (a_2/3) + (j\sqrt{3}/2)(s_1 - s_2)$$

$$z_3 = -\frac{1}{2}(s_1 + s_2) - (a_2/3) - (j\sqrt{3}/2)(s_1 - s_2)$$

If z_1, z_2, z_3 are the roots of the cubic equation

$$z_1 + z_2 + z_3 = -a_2$$

$$z_1z_2 + z_1z_3 + z_2z_3 = a_1$$

$$z_1z_2z_3 = -a_0$$

Solution of Quartic Equations*

Given $z^4 + a_3z^3 + a_2z^2 + a_1z + a_0 = 0$, find the real root u_1 of the cubic equation

$$u^3 - a_2u^2 + (a_1a_3 - 4a_0)u - (a_1^2 + a_0a_3^2 - 4a_0a_2) = 0$$

and determine the four roots of the quartic as solutions of the two quadratic equations

$$v^2 + \{(a_3/2) \mp [(a_3^2/4) + u_1 - a_2]^{1/2}\}v + (u_1/2) \mp [(u_1/2)^2 - a_0]^{1/2} = 0$$

If all roots of the cubic equation are real, use the value of u_1 that gives real coefficients in the quadratic equation.

Complex Quantities

In the following equations, all quantities are real except $j = (-1)^{1/2}$.

$$(A + jB) + (C + jD) = (A + C) + j(B + D)$$

$$(A + jB)(C + jD) = (AC - BD) + j(BC + AD)$$

$$\frac{A + jB}{C + jD} = \frac{AC + BD}{C^2 + D^2} + j \frac{BC - AD}{C^2 + D^2}$$

$$\frac{1}{A + jB} = \frac{A}{A^2 + B^2} - j \frac{B}{A^2 + B^2}$$

Polar Form:

$$A + jB = \rho(\cos\theta + j \sin\theta) = \rho e^{j\theta}$$

De Moivre's Equation:

$$(A + jB)^v = \rho^v(\cos v\theta + j \sin v\theta)$$

where,

$$\rho = (A^2 + B^2)^{1/2} > 0,$$

$$\cos\theta = A/\rho,$$

$$\sin\theta = B/\rho.$$

For nonintegral v , this quantity is many-valued.

Complex Conjugate:

$$A + jB = (A + jB)^* = A - jB$$

Analytic Function: Let $f(z)$ be a function of the complex variable $z = x + jy$. Function $f(z)$ is analytic at a point $z = z_0$, if

$$\lim_{\Delta z \rightarrow 0} [f(z_0 + \Delta z) - f(z_0)]/\Delta z$$

exists independent of the manner in which Δz approaches zero. Analyticity is equivalent to the existence of a Taylor series about the point in question. In addition, $f(z)$ is analytic if and only if the Cauchy-Riemann equations hold:

$$\partial u/\partial x = \partial v/\partial y$$

$$\partial v/\partial x = -\partial u/\partial y$$

where,

$$u = \text{Re}f(z),$$

$$v = \text{Im}f(z).$$

If $f(z)$ is analytic for all z , it is called an entire function.

Properties of e

$$e = \lim_{n \rightarrow \infty} (1 + n^{-1})^n = \sum_{k=0}^{\infty} (k!)^{-1} = 2.71828$$

$$e^{\pm jx} = \cos x \pm j \sin x = \exp(\pm jx)$$

Properties of Logarithms

If $\log_a x = N$, then $a^N = x$

$$\log_a x = \log_a b \log_b x$$

$$\log_a xy = \log_a x + \log_a y$$

$$\log_a x/y = \log_a x - \log_a y$$

* Abramovitz, M. and Stegun, I. A. *Handbook of Mathematical Functions*. Washington, D. C.: National Bureau of Standards.

$$\log_a x^y = y \log_a x$$

$$\log_a b = 1/\log_b a$$

$$a^{\log_a x} = x$$

$$\log_a 1 = 0$$

$$\log_a a = 1$$

$$\log_e x = \ln x = \log_e 10 \log_{10} x = 2.30259 \log_{10} x$$

$$\log_{10} x = \log_{10} e \log_e x = 0.43429 \log_e x$$

Sums

In this subsection, the following symbols will be used:

$\Gamma(\alpha)$ = Gamma function of α

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$

in which n and k are positive integers, $n \geq k$, and

$$\binom{n}{0} = 1$$

Arithmetic Progression:

$$\sum_{k=0}^{n-1} (a + kd) = a + (a + d) + (a + 2d) + \dots + [a + (n - 1)d]$$

$$= \frac{1}{2}n[2a + (n - 1)d]$$

Geometric Progression:

$$\sum_{k=0}^{n-1} ar^k = a + ar + ar^2 + \dots + ar^{n-1}$$

$$= [a(r^n - 1)/(r - 1)], \quad \text{for } r \neq 1$$

$$= na, \quad \text{for } r = 1$$

Sums of Powers of Integers:

$$\sum_{k=1}^n k^2 = 1^2 + 2^2 + 3^2 + \dots + n^2$$

$$= [n(n + 1)(2n + 1)/6]$$

$$\sum_{k=1}^n k^3 = 1^3 + 2^3 + 3^3 + \dots + n^3$$

$$= [n^2(n + 1)^2/4]$$

$$\sum_{k=1}^n k^4 = 1^4 + 2^4 + 3^4; \dots + n^4$$

$$= \frac{1}{30} n(n + 1)(2n + 1)(3n^2 + 3n - 1)$$

$$\sum_{k=1}^n k^r = 1^r + 2^r + 3^r + \dots + n^r$$

$$= \frac{n^{r+1}}{r + 1} + \frac{1}{2}n^r + \sum_{k=1}^{[r/2]} (2k)^{-1} \binom{r}{2k - 1} B_{2k} n^{r-2k+1}$$

where,

r is a positive integer,
 $[r/2]$ is the largest integer less than or equal to $r/2$,
 B_{2k} is the $2k$ th Bernoulli number.

These sums are tabulated for $n = 1, 2, 3, \dots, 100$ and for $r = 1, 2, 3, \dots, 10$.*

Sums of integral powers of odd integers may be obtained from the above. For example,

$$\sum_{k=0}^n (2k + 1)^2 = 1^2 + 3^2 + 5^2 + \dots + (2n + 1)^2$$

$$= 1^2 + 2^2 + 3^2 + \dots + (2n + 1)^2 - 2^2 - 4^2 - 6^2 - \dots - (2n)^2$$

$$= \sum_{k=1}^{2n+1} k^2 - 2^2 \sum_{k=1}^n k^2$$

$$= \frac{1}{3}(n + 1)(2n + 1)(2n + 3)$$

Sums of Powers of Reciprocals of Integers:

$$\zeta(z) = \sum_{k=1}^{\infty} (1/k^z) = (1/1^z) + (1/2^z) + (1/3^z) + \dots$$

For $\text{Re}z > 1$, $\zeta(z)$ is the Riemann zeta function of z

$$\sum_{k=1}^{\infty} (1/k^2) = (1/1^2) + (1/2^2) + (1/3^2) + \dots$$

$$= \pi^2/6 = 1.64493$$

$$\sum_{k=1}^{\infty} (1/k^3) = (1/1^3) + (1/2^3) + (1/3^3) + \dots$$

$$= 1.20206$$

$$\sum_{k=1}^{\infty} (1/k^4) = (1/1^4) + (1/2^4) + (1/3^4) + \dots$$

$$= \pi^4/90 = 1.08232$$

$$\sum_{k=1}^{\infty} (1/k^{2r}) = (1/1^{2r}) + (1/2^{2r}) + (1/3^{2r}) + \dots$$

$$= [2^{2r-1} \pi^{2r} |B_{2r}| / (2r)!],$$

r integral

* Abramovitz, M. and Stegun, I. A. *Handbook of Mathematical Functions*. Washington, D. C.: National Bureau of Standards.

where B_{2r} is the $2r$ th Bernoulli number. Values of $\zeta(z)$ are tabulated for $z = 2, 3, 4, \dots, 42$, to 20 decimal places in the *Handbook of Mathematical Functions*. The truncated sums

$$\sum_{k=1}^n (1/k^z)$$

are related to the polygamma functions.

Finite Sums of Binomial Coefficients:

$$\begin{aligned} \sum_{k=0}^n \binom{m+k}{m} &= \binom{m}{m} + \binom{m+1}{m} + \dots \\ &+ \binom{m+n}{m} \\ &= \binom{m+n+1}{m+1} \end{aligned}$$

$$\begin{aligned} \sum_{k=0}^n \binom{r}{k} \binom{s}{n-k} &= \binom{r}{0} \binom{s}{n} \\ &+ \binom{r}{1} \binom{s}{n-1} + \dots + \binom{r}{n} \binom{s}{0} \\ &= \binom{r+s}{n} \quad (r \geq n \text{ and } s \geq n) \end{aligned}$$

Binomial Theorem:
Nonnegative integral exponent:

$$\begin{aligned} (a+b)^n &= a^n + na^{n-1}b + \frac{1}{2}[n(n-1)]a^{n-2}b^2 \\ &+ \dots + b^n = \sum_{k=0}^n \binom{n}{k} a^{n-k} b^k \end{aligned}$$

For other n , integral negative and nonintegral, the series will be infinite. For convergence it is assumed $|a| > |b|$. (If $|a| < |b|$, interchange a and b .)

Negative integral exponent:

$$\begin{aligned} (a+b)^{-n} &= a^{-n} - na^{-n-1}b \\ &+ \frac{1}{2}[n(n+1)]a^{-n-2}b^2 - \dots \\ &= \sum_{k=0}^{\infty} (-1)^k \binom{n+k-1}{k} a^{-n-k} b^k \end{aligned}$$

Nonintegral exponent:

$$\begin{aligned} (a+b)^\alpha &= a^\alpha + \alpha a^{\alpha-1}b + \frac{1}{2}\alpha(\alpha-1)a^{\alpha-2}b^2 + \dots \\ &= \sum_{k=0}^{\infty} \frac{\Gamma(\alpha+1)}{k!\Gamma(\alpha-k+1)} a^{\alpha-k} b^k \end{aligned}$$

Multinomial Series:

$$\begin{aligned} (x_1 + x_2 + \dots + x_r)^n &= \sum_{n_1=0}^n \sum_{n_2=0}^{n-n_1} \sum_{n_3=0}^{n-n_1-n_2} \dots \sum_{n_{r-1}=0}^{n-n_1-n_2-\dots-n_{r-2}} n! \prod_{k=1}^r \frac{x_k^{n_k}}{n_k!} \end{aligned}$$

where the interpretation $n_r = n - n_1 - n_2 - \dots - n_{r-1}$ is to be used in the final product.

Combinations and Permutations

A combination is a selection from a number of things in which the order of the selected objects is disregarded, whereas a permutation is a selection in which the order is taken into consideration. For example, if from the letters a, b , and c a group of two is selected, then ab, bc, ac are the combinations and ab, ba, bc, cb, ac, ca are the permutations.

The number of different combinations of n (dissimilar) things taken r at a time is

$$\binom{n}{r} = C_r^n = \frac{n!}{r!(n-r)!}$$

The number of different permutations of n (dissimilar) things taken r at a time is

$$P_r^n = n!(n-r)! = n \times (n-1) \times \dots \times (n-r+1)$$

Bernoulli Numbers

Definition:

$$B_n = (d^n/dx^n)[x/(e^x - 1)]|_{x=0}$$

Values for Small n:

$$B_0 = 1, B_1 = -\frac{1}{2}, B_2 = \frac{1}{6}, B_4 = -\frac{1}{30},$$

$$B_6 = \frac{1}{42}, B_8 = -\frac{1}{30}, B_{10} = \frac{5}{66},$$

$$B_{2n+1} = 0, \text{ for all integral } n > 0$$

Identities:

$$B_{2n} = (-1)^{n-1} \frac{2 \cdot (2n)!}{(2\pi)^{2n}} \sum_{k=1}^{\infty} \frac{1}{k^{2n}}, n \geq 1$$

$$B_n = \sum_{k=0}^n \binom{n}{k} B_k, n > 1$$

The latter identity may be used for recursive calculation of the B_n 's.

Trigonometric Identities

$$\begin{aligned} 1 &= \sin^2 A + \cos^2 A = \sin A \csc A = \tan A \cot A \\ &= \cos A \sec A \end{aligned}$$

$$\sin A = \cos A / \cot A = 1 / \csc A = \cos A \tan A$$

$$= \pm(1 - \cos^2 A)^{1/2}$$

$$\cos A = \sin A / \tan A = 1 / \sec A = \sin A \cot A$$

$$= \pm(1 - \sin^2 A)^{1/2}$$

$$\tan A = \sin A / \cos A = 1 / \cot A = \sin A \sec A$$

$$\sin A = (e^{jA} - e^{-jA})/2j$$

$$\cos A = (e^{jA} + e^{-jA})/2$$

$$\sin(A \pm B) = \sin A \cos B \pm \cos A \sin B$$

$$\cos(A \pm B) = \cos A \cos B \mp \sin A \sin B$$

$$\tan(A \pm B) = (\tan A \pm \tan B)/(1 \mp \tan A \tan B)$$

$$= (\tan A \cot B \pm 1)/(\cot B \mp \tan A)$$

$$\cot(A \pm B) = (\cot A \cot B \mp 1)/(\cot B \pm \cot A)$$

$$= (\cot A \mp \tan B)/(1 \pm \cot A \tan B)$$

$$\sin 2A = 2 \sin A \cos A$$

$$\cos 2A = \cos^2 A - \sin^2 A$$

$$\tan 2A = (2 \tan A)/(1 - \tan^2 A)$$

$$\sin 3A = 3 \sin A - 4 \sin^3 A$$

$$\cos 3A = -3 \cos A + 4 \cos^3 A$$

$$\tan 3A = (3 \tan A - \tan^3 A)/(1 - 3 \tan^2 A)$$

$$\cos nA = \operatorname{Re}(\cos A + j \sin A)^n$$

$$\sin nA = \operatorname{Im}(\cos A + j \sin A)^n$$

$$\sin \frac{1}{2}A = \pm[(1 - \cos A)/2]^{1/2}$$

$$\cos \frac{1}{2}A = \pm[(1 + \cos A)/2]^{1/2}$$

$$\tan \frac{1}{2}A = \sin A/(1 + \cos A) = (1 - \cos A)/\sin A$$

$$\sin A \pm \sin B = 2 \sin \frac{1}{2}(A \pm B) \cos \frac{1}{2}(A \mp B)$$

$$\cos A + \cos B = 2 \cos \frac{1}{2}(A + B) \cos \frac{1}{2}(A - B)$$

$$\cos B - \cos A = 2 \sin \frac{1}{2}(A + B) \sin \frac{1}{2}(A - B)$$

$$\tan A \pm \tan B = [\sin(A \pm B)/\cos A \cos B]$$

$$\cot A \pm \cot B = [\sin(B \pm A)/\sin A \sin B]$$

$$\sin^2 A - \sin^2 B = \sin(A + B) \sin(A - B)$$

$$\cos^2 A - \sin^2 B = \cos(A + B) \cos(A - B)$$

$$\tan \frac{1}{2}(A \pm B) = (\sin A \pm \sin B)/(\cos A + \cos B)$$

$$\cot \frac{1}{2}(A \mp B) = (\sin A \pm \sin B)/(\cos B - \cos A)$$

$$\cos^2 A = \frac{1}{2}(\cos 2A + 1)$$

$$\cos^3 A = \frac{1}{4}(\cos 3A + 3 \cos A)$$

$$\cos^4 A = \frac{1}{8}(\cos 4A + 4 \cos 2A + 3)$$

$$\sin^2 A = \frac{1}{2}(-\cos 2A + 1)$$

$$\sin^3 A = \frac{1}{4}(-\sin 3A + 3 \sin A)$$

$$\sin^4 A = \frac{1}{8}(\cos 4A - 4 \cos 2A + 3)$$

$$\sin A \cos B = \frac{1}{2}[\sin(A + B) + \sin(A - B)]$$

$$\cos A \cos B = \frac{1}{2}[\cos(A + B) + \cos(A - B)]$$

$$\sin A \sin B = \frac{1}{2}[\cos(A - B) - \cos(A + B)]$$

$$\sin A + m \sin B = \rho \sin C$$

with

$$\rho^2 = 1 + m^2 + 2m \cos(B - A)$$

and

$$\tan(C - A) = [m \sin(B - A)]/[1 + m \cos(B - A)]$$

$$\sum_i A_i \exp(j\theta_i) = \rho e^{j\psi}$$

with

$$\tan \psi = \left(\sum_i A_i \sin \theta_i / \sum_i A_i \cos \theta_i \right)$$

and

$$\rho = \left[\sum_i A_i^2 + \sum_{i < j} \sum_j A_i A_j \cos(\theta_i - \theta_j) \right]^{1/2}$$

In the previous notation

$$\sum_i A_i \cos \theta_i = \rho \cos \psi$$

$$\sum_i A_i \sin \theta_i = \rho \sin \psi$$

apply.

$$\sin x + \sin 2x + \sin 3x + \cdots + \sin mx$$

$$= [\sin \frac{1}{2}mx \sin \frac{1}{2}(m + 1)x / \sin \frac{1}{2}x]$$

$$\cos x + \cos 2x + \cos 3x + \cdots + \cos mx$$

$$= [\sin \frac{1}{2}mx \cos \frac{1}{2}(m + 1)x / \sin \frac{1}{2}x]$$

$$\begin{aligned} \sin x + \sin 3x + \sin 5x + \dots + \sin(2m - 1)x &= (\sin^2 mx / \sin x) \\ \cos x + \cos 3x + \cos 5x + \dots + \cos(2m - 1)x &= (\sin 2mx / 2 \sin x) \\ \frac{1}{2} + \cos x + \cos 2x + \dots + \cos mx &= [\sin(m + \frac{1}{2})x / 2 \sin \frac{1}{2}x] \end{aligned}$$

Angle (degrees)	Sine	Cosine	Tangent
0	0	1	0
30	$\frac{1}{2}$	$\frac{1}{2}\sqrt{3}$	$\frac{1}{3}\sqrt{3}$
45	$\frac{1}{2}\sqrt{2}$	$\frac{1}{2}\sqrt{2}$	1
60	$\frac{1}{2}\sqrt{3}$	$\frac{1}{2}$	$\sqrt{3}$
90	1	0	$\pm\infty$
180	0	-1	0
270	-1	0	$\pm\infty$
360	0	1	0
0-90	+	+	+
90-180	+	-	-
180-270	-	-	+
270-360	-	+	-

versine: $\text{vers } \theta = 1 - \cos \theta$

haversine: $\text{hav } \theta = \frac{1}{2}(1 - \cos \theta) = \sin^2 \frac{1}{2}\theta$

Approximations for Small Angles

$$\left. \begin{aligned} \sin \theta &= (\theta - \theta^3/6 \dots) \\ \tan \theta &= (\theta + \theta^3/3 \dots) \\ \cos \theta &= (1 - \theta^2/2 \dots) \end{aligned} \right\} \theta \text{ in radians}$$

$\sin \theta = \theta$

with less than 1-percent error up to $\theta = 0.24$ radian = 14.0°
 with less than 10-percent error up to $\theta = 0.78$ radian = 44.5°

$\tan \theta = \theta$

with less than 1-percent error up to $\theta = 0.17$ radian = 10.0°
 with less than 10-percent error up to $\theta = 0.54$ radian = 31.0°

Inequalities

$$\begin{aligned} \sin x &\leq x \leq \tan x, & \text{for } 0 \leq x < \pi/2 \\ \sin x &\geq (2/\pi)x, & \text{for } 0 \leq x \leq \pi/2 \\ \cos x &< \sin x/x \leq 1, & \text{for } 0 < x \leq \pi \end{aligned}$$

where x is in radians.

PLANE TRIGONOMETRY

Right Triangles

Refer to Fig. 1.

$$\begin{aligned} C &= 90^\circ \\ B &= 90^\circ - A \\ \sin A &= \cos B = a/c \\ \tan A &= a/b \\ c^2 &= a^2 + b^2 \\ \text{area} &= \frac{1}{2}ab \end{aligned}$$

Oblique Triangles

Refer to Fig. 2.

Sum of Angles:

$$A + B + C = 180^\circ \quad (\text{Eq. 1})$$

Law of Cosines:

$$a^2 = b^2 + c^2 - 2bc \cos A \quad (\text{Eq. 2})$$

Law of Sines:

$$a/\sin A = b/\sin B = c/\sin C \quad (\text{Eq. 3})$$

Law of Tangents:

$$\frac{a - b}{a + b} = \frac{\tan \frac{1}{2}(A - B)}{\tan \frac{1}{2}(A + B)} \quad (\text{Eq. 4})$$

Half-Angle Equation (Fig. 3):

$$\tan \frac{1}{2}A = r/(p - a) \quad (\text{Eq. 5})$$

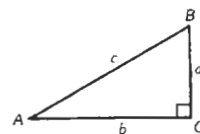


Fig. 1. Right triangle.

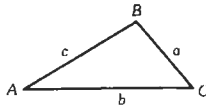


Fig. 2. Oblique triangle.

where,

$$2p = a + b + c$$

$$r = [(p - a)(p - b)(p - c)/p]^{1/2}$$

Solving an Oblique Triangle:

Given	Use Eq.	To Obtain
aBC	(1) (3)	A bc
Abc	(1) (4)	B + C hence B - C B,C
abc	(5) or (2)	ABC
abA ambiguous case	(3) and (1)	BCc

SPHERICAL TRIGONOMETRY

Spherical triangles are bounded by the arcs of great circles. These are circles formed by the intersection of a sphere with planes passing through the center of the sphere. In the following equations, α, β, γ are the angles, and a, b, c are the corresponding opposite sides, respectively. The sides are measured by the arcs subtended by the arcs; for example, a side extending from the Equator to the North Pole is a 90° side.

Right Spherical Triangles ($\gamma = 90^\circ$)

Refer to Fig. 4.

$$\text{cosec } \gamma = \text{cosec } \alpha \text{ cosec } \beta = \text{cot } \alpha \text{ cot } \beta$$

$$\text{cos } \alpha = \text{sin } \beta \text{ cos } \alpha = \text{tan } b \text{ cot } c$$

$$\text{cos } \beta = \text{sin } \alpha \text{ cos } \beta = \text{tan } a \text{ cot } c$$

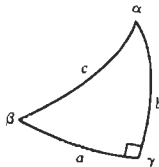


Fig. 4. Right spherical triangle.

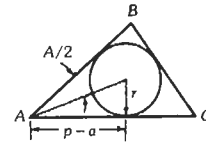


Fig. 3. Half angle.

$$\text{sin } a = \text{sin } c \text{ sin } \alpha = \text{tan } b \text{ cot } \beta$$

$$\text{sin } b = \text{sin } c \text{ sin } \beta = \text{tan } a \text{ cot } \alpha \quad (\text{Eq. 6})$$

Oblique Triangles

Refer to Fig. 5.

Law of Cosines for Sides:

$$\text{cos } a = \text{cos } b \text{ cos } c + \text{sin } b \text{ sin } c \text{ cos } \alpha$$

$$\text{cos } b = \text{cos } c \text{ cos } a + \text{sin } c \text{ sin } a \text{ cos } \beta$$

$$\text{cos } c = \text{cos } a \text{ cos } b + \text{sin } a \text{ sin } b \text{ cos } \gamma \quad (\text{Eq. 7a})$$

Law of Cosines for Angles:

$$\text{cos } \alpha = -\text{cos } \beta \text{ cos } \gamma + \text{sin } \beta \text{ sin } \gamma \text{ cos } a$$

$$\text{cos } \beta = -\text{cos } \gamma \text{ cos } \alpha + \text{sin } \gamma \text{ sin } \alpha \text{ cos } b$$

$$\text{cos } \gamma = -\text{cos } \alpha \text{ cos } \beta + \text{sin } \alpha \text{ sin } \beta \text{ cos } c \quad (\text{Eq. 7b})$$

Law of Sines:

$$\text{sin } a / \text{sin } \alpha = \text{sin } b / \text{sin } \beta = \text{sin } c / \text{sin } \gamma \quad (\text{Eq. 8})$$

Napier's Analogies:

$$\frac{\text{sin } \frac{1}{2}(\alpha - \beta)}{\text{sin } \frac{1}{2}(\alpha + \beta)} = \frac{\text{tan } \frac{1}{2}(a - b)}{\text{tan } \frac{1}{2}c} \quad (\text{Eq. 9a})$$

$$\frac{\text{cos } \frac{1}{2}(\alpha - \beta)}{\text{cos } \frac{1}{2}(\alpha + \beta)} = \frac{\text{tan } \frac{1}{2}(a + b)}{\text{tan } \frac{1}{2}c} \quad (\text{Eq. 9b})$$

$$\frac{\text{sin } \frac{1}{2}(a - b)}{\text{sin } \frac{1}{2}(a + b)} = \frac{\text{tan } \frac{1}{2}(\alpha - \beta)}{\text{cot } \frac{1}{2}\gamma} \quad (\text{Eq. 9c})$$

$$\frac{\text{cos } \frac{1}{2}(a - b)}{\text{cos } \frac{1}{2}(a + b)} = \frac{\text{tan } \frac{1}{2}(\alpha + \beta)}{\text{cot } \frac{1}{2}\gamma} \quad (\text{Eq. 9d})$$

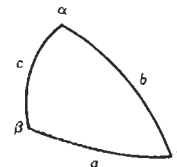


Fig. 5. Oblique spherical triangle.

Half-Angle Equations:

$$\begin{aligned} \tan\frac{1}{2}\alpha &= \tan r/\sin(p - a) \\ \tan\frac{1}{2}\beta &= \tan r/\sin(p - b) \\ \tan\frac{1}{2}\gamma &= \tan r/\sin(p - c) \end{aligned} \quad (\text{Eq. 10a})$$

where $2p = a + b + c$ and

$$\tan^2 r = \frac{\sin(p - a) \sin(p - b) \sin(p - c)}{\sin p}$$

$$\begin{aligned} \sin^2 \frac{1}{2}\alpha &= [\sin(p - b) \sin(p - c)/\sin b \sin c] \\ \cos^2 \frac{1}{2}\alpha &= [\sin p \sin(p - a)/\sin b \sin c] \\ \tan^2 \frac{1}{2}\alpha &= [\sin(p - b) \sin(p - c)/\sin p \sin(p - a)] \end{aligned} \quad (\text{Eq. 10b})$$

and equations obtained by cyclical permutation for β and γ .

Half-Side Equations:

$$\begin{aligned} \tan\frac{1}{2}a &= \tan R \sin(\alpha - E) \\ \tan\frac{1}{2}b &= \tan R \sin(\beta - E) \\ \tan\frac{1}{2}c &= \tan R \sin(\gamma - E) \end{aligned} \quad (\text{Eq. 11a})$$

where $2E = \alpha + \beta + \gamma - \pi$ is the spherical excess and

$$\tan^2 R = \frac{\sin E}{\sin(\alpha - E) \sin(\beta - E) \sin(\gamma - E)}$$

$$\begin{aligned} \sin^2 \frac{1}{2}a &= -[\sin E \sin(E - \alpha)/\sin \beta \sin \gamma] \\ \cos^2 \frac{1}{2}a &= [\sin(E - \beta) \sin(E - \gamma)/\sin \beta \sin \gamma] \\ \tan^2 \frac{1}{2}a &= -[\sin E \sin(E - \alpha)/\sin(E - \beta) \sin(E - \gamma)] \end{aligned} \quad (\text{Eq. 11b})$$

and equations obtained by cyclical permutation for b and c .

Area: On a sphere of radius one, the area of a triangle is equal to the spherical excess

$$2E = \alpha + \beta + \gamma - \pi$$

L'Huilier's Theorem:

$$\begin{aligned} \tan^2 \frac{1}{2}E &= \tan\frac{1}{2}p \tan\frac{1}{2}(p - a) \\ &\quad \times \tan\frac{1}{2}(p - b) \tan\frac{1}{2}(p - c) \end{aligned} \quad (\text{Eq. 12})$$

Solving an Oblique Triangle*:

Given	Use Eq.	To Obtain
abc	(10)	$\alpha\beta\gamma$
$\alpha\beta\gamma$	(11)	abc
$ab\gamma$	(9)	$\alpha \pm \beta$, hence α, β , then c
$\alpha\beta c$	(9)	$a \pm b$, hence a, b , then γ
$ab\alpha$ ambiguous case	(8) (9)	β $c\gamma$
$\alpha\beta a$ ambiguous case	(8) (9)	b $c\gamma$

HYPERBOLIC FUNCTIONS*

$$\begin{aligned} \sinh x &= (e^x - e^{-x})/2 \\ \cosh x &= (e^x + e^{-x})/2 \\ \tanh x &= \sinh x/\cosh x \\ &= [1 - \exp(-2x)]/[1 + \exp(-2x)] \\ &= 1/\coth x \end{aligned}$$

$$\operatorname{sech} x = 1/\cosh x$$

$$\operatorname{csch} x = 1/\sinh x$$

$$\sinh(-x) = -\sinh x$$

$$\cosh(-x) = \cosh x$$

$$\tanh(-x) = -\tanh x$$

$$\coth(-x) = -\coth x$$

$$\sinh jx = j \sin x$$

$$\cosh jx = \cos x$$

$$\tanh jx = j \tan x$$

$$\coth jx = -j \cot x$$

$$\cosh^2 x - \sinh^2 x = 1$$

$$1 - \tanh^2 x = 1/\cosh^2 x$$

$$\coth^2 x - 1 = 1/\sinh^2 x$$

$$\sinh 2x = 2 \sinh x \cosh x$$

$$\cosh 2x = \cosh^2 x + \sinh^2 x$$

$$\sinh(x \pm jy) = \sinh x \cos y \pm j \cosh x \sin y$$

$$\cosh(x \pm jy) = \cosh x \cos y \pm j \sinh x \sin y$$

* See also great-circle calculations in Chapter 33.

* Tables of hyperbolic functions appear in Chapter 48.

$$\tanh(x \pm y) = (\tanh x \pm \tanh y)/(1 \pm \tanh x \tanh y)$$

If $y = \operatorname{gd}x$ (gudermannian of x) is defined by

$$x = \log_e \tan\left(\frac{1}{4}\pi + \frac{1}{2}y\right)$$

then

$$\sinh x = \operatorname{tany}$$

$$\cosh x = \operatorname{secy}$$

$$\tanh x = \operatorname{siny}$$

$$\tanh(x/2) = \tan(y/2)$$

HYPERBOLIC TRIGONOMETRY

Hyperbolic (or pseudospherical) trigonometry applies to triangles drawn in the hyperbolic type of non-Euclidean space. Reflection charts, used in transmission-line theory and waveguide analysis, are models of this hyperbolic space.*

Conformal Model

The space is limited to the inside of a unit circle, Γ . Geodesics (or "straight lines" for the model) are arcs of circle orthogonal to Γ as shown in Fig. 6. The hyperbolic distance between two points A and B is defined by

$$[AB] = \log_e[(BI/BJ):(AI/AJ)]$$

where I and J are the intersections with Γ of geodesic AB . The distance $[AB]$ is expressed in nepers. For engineering purposes, a unit corresponding to the decibel and equal to $1/8.686$ neper is sometimes used.

As this model is conformal, the angle between two lines is the ordinary angle between the tangents at their common point.

* G. A. Deschamps, *Hyperbolic Protractor for Microwave Impedance Measurements and Other Purposes*. New York: International Telephone and Telegraph Corp., 1953.

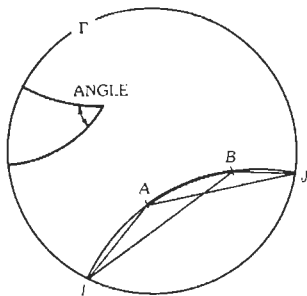


Fig. 6. Conformal model.

Projective Model

The space is again composed of the points inside a circle, Γ . Geodesics are straight-line segments limited to the inside of Γ (IJ in Fig. 7).

The hyperbolic distance AB is defined by

$$\langle AB \rangle = \frac{1}{2} \log_e [(BI/BJ):(AI/AJ)]$$

and can be measured directly by means of a hyperbolic protractor. The angles for this model do not appear in true size except when at the center of Γ . An angle such as BAC , when it is considered in reference to the projective model, will be called an *elliptic* angle. It can be evaluated, as shown in Fig. 8, by projecting B and C through the hyperbolic midpoint of OA onto B' and C' on circle Γ , then measuring $B'OC'$ as in Euclidean geometry.

The two models in Fig. 9 drawn inside the same circle Γ can be set into a distance-preserving correspondence by the transformation: $\mathcal{B}(M) = M'$ defined by

$$[OM] = \langle OM' \rangle$$

or in terms of ordinary distances

$$OM' = 2OM/(1 + OM^2)$$

The hyperbolic distance to the center O being denoted by u

$$OM = \tanh(u/2)$$

and

$$OM' = \tanh u$$

The points on Γ are at an infinite distance from any point inside Γ .

In the following equations, the sides are expressed in nepers, the angles in radians. The three points A, B, C are assumed to be inside circle Γ .

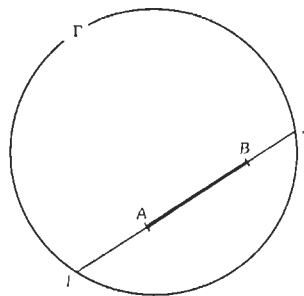


Fig. 7. Projective model.

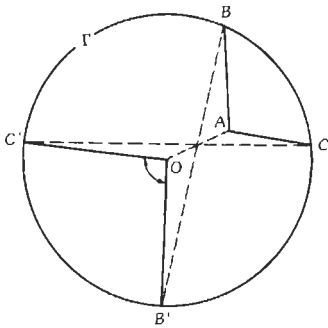


Fig. 8. Construction of angle on projective model.

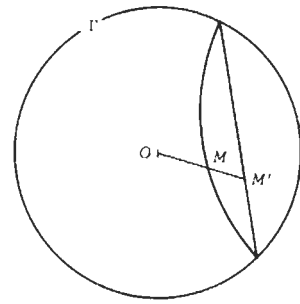


Fig. 9. Correspondence between the two models.

Right Hyperbolic Triangles
($\gamma = 90^\circ$)

Refer to Figs. 10 and 11.

$$\begin{aligned} \cosh c &= \cosh a \cosh b \\ &= \cot \alpha \cot \beta \\ \cos \alpha &= \sin \beta \cosh a \\ &= \tanh b \coth c \\ \cos \beta &= \sin \alpha \cosh b \\ &= \tanh a \coth c \end{aligned}$$

When B is at infinity, i.e., on Γ

$$\begin{aligned} \cos A &= \tanh b \\ \cot A &= \sinh b \\ \csc A &= \cosh b \\ \tan \frac{1}{2} A &= \exp b \end{aligned}$$

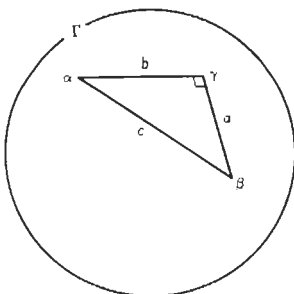


Fig. 10. Projective representation of right hyperbolic triangle.

or

$$(\pi/2) - A = \text{gdb}$$

(See definition of gd in section "Hyperbolic Functions.")

CB and AB are "parallel." A is also called angle of parallelism and is noted by

$$\begin{aligned} A &= \Pi(b) \\ &= \pi/2 - \text{gdb} \end{aligned}$$

Oblique Hyperbolic Triangles

Law of Cosines:

$$\cosh a = \cosh b \cosh c - \sinh b \sinh c \cos \alpha$$

and permutations (Eq. 13a)

$$\cos \alpha = -\cos \beta \cos \gamma + \sin \beta \sin \gamma \cosh a$$

and permutations (Eq. 13b)

Law of Sines (Fig. 12):

$$\sinh a / \sin \alpha = \sinh b / \sin \beta = \sinh c / \sin \gamma \quad (\text{Eq. 14})$$

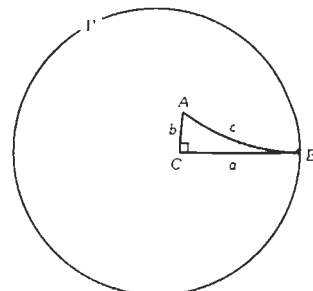


Fig. 11. Conformal representation of right hyperbolic triangle with B at infinity.

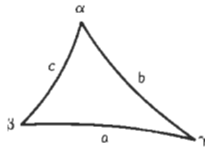


Fig. 12. Oblique hyperbolic triangle.

Napier's Analogies:

$$\frac{\sin \frac{1}{2}(\alpha - \beta)}{\sin \frac{1}{2}(\alpha + \beta)} = \frac{\tanh \frac{1}{2}(a - b)}{\tanh \frac{1}{2}c} \quad (\text{Eq. 15a})$$

$$\frac{\cos \frac{1}{2}(\alpha - \beta)}{\cos \frac{1}{2}(\alpha + \beta)} = \frac{\tan \frac{1}{2}(a + b)}{\tanh \frac{1}{2}c} \quad (\text{Eq. 15b})$$

$$\frac{\sinh \frac{1}{2}(a - b)}{\sinh \frac{1}{2}(a + b)} = \frac{\tan \frac{1}{2}(\alpha - \beta)}{\cot \frac{1}{2}\gamma} \quad (\text{Eq. 15c})$$

$$\frac{\cosh \frac{1}{2}(a - b)}{\cosh \frac{1}{2}(a + b)} = \frac{\tan \frac{1}{2}(\alpha + \beta)}{\cot \frac{1}{2}\gamma} \quad (\text{Eq. 15d})$$

Half-Angle Equations:

$$\tan \frac{1}{2}\alpha = \frac{\sinh p}{\sinh(p - a)}$$

and permutations where $2p = a + b + c$ and

$$\tanh^2 r = \frac{\sinh(p - a) \sinh(p - b) \sinh(p - c)}{\sinh p} \quad (\text{Eq. 16a})$$

$$\sin^2 \frac{1}{2}\alpha = \frac{\sinh(p - b) \sinh(p - c)}{\sinh b \sinh c}$$

$$\cos^2 \frac{1}{2}\alpha = \frac{\sinh p \sinh(p - a)}{\sinh b \sinh c}$$

$$\tan^2 \frac{1}{2}\alpha = \frac{\sinh(p - b) \sinh(p - c)}{\sinh p \sinh(p - a)} \quad (\text{Eq. 16b})$$

Half-Side Equations:

$$\coth \frac{1}{2}a = \frac{\coth R}{\sin(\Delta + \alpha)}$$

and permutations where $2\Delta = \pi - \alpha - \beta - \gamma$ is the hyperbolic defect and

$$\tanh^2 R = \frac{\sin \Delta}{\sin(\Delta + \alpha) \sin(\Delta + \beta) \sin(\Delta + \gamma)} \quad (\text{Eq. 17a})$$

$$\sinh^2 \frac{1}{2}a = \frac{\sin \Delta \sin(\Delta + \alpha)}{\sin \beta \sin \gamma}$$

$$\cosh^2 \frac{1}{2}a = \frac{\sin(\Delta + \beta) \sin(\Delta + \gamma)}{\sin \beta \sin \gamma}$$

$$\tanh^2 \frac{1}{2}a = \frac{\sin \Delta \sin(\Delta + \alpha)}{\sin(\Delta + \beta) \sin(\Delta + \gamma)} \quad (\text{Eq. 17b})$$

Area: The hyperbolic area of a triangle is equal to the hyperbolic defect.

$$2\Delta = \pi - (\alpha + \beta + \gamma) \quad (\text{Eq. 18})$$

Solving an Oblique Hyperbolic Triangle: Solution of an oblique hyperbolic triangle is analogous to that for an oblique spherical triangle, as follows.

Given	Use Eq.	To Obtain
abc	(16)	$\alpha\beta\gamma$
$\alpha\beta\gamma$	(17)	abc
$ab\gamma$	(15)	$\alpha \pm \beta$, hence α, β , then c
$\alpha\beta c$	(15)	$a \pm b$, hence a, b , then γ
$ab\alpha$ ambiguous case	(14) (15)	β $c\gamma$
$\alpha\beta a$ ambiguous case	(14) (15)	b $c\gamma$

PLANE ANALYTIC GEOMETRY

In the following, x and y are coordinates of a variable point in a rectangular-coordinate system.

Straight Line

General Equation:

$$Ax + By + C = 0$$

$A, B,$ and C are constants.

Slope-Intercept Form (Fig. 13):

$$y = sx + b$$

$b = y$ -intercept

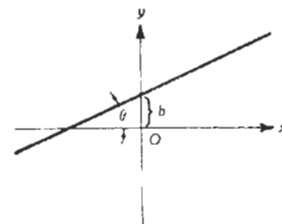


Fig. 13. Slope-intercept.

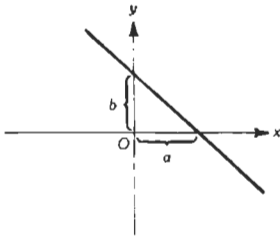


Fig. 14. Intercept-intercept.

$$s = \tan\theta$$

$$= \text{slope}$$

Intercept-Intercept Form (Fig. 14):

$$(x/a) + (y/b) = 1$$

a = x-intercept

b = y-intercept

Point-Slope Form (Fig. 15):

$$y - y_1 = s(x - x_1)$$

$$s = \tan\theta$$

(x_1, y_1) = coordinates of known point on line.

Point-Point Form:

$$(y - y_1)/(y_1 - y_2) = (x - x_1)/(x_1 - x_2)$$

where (x_1, y_1) and (x_2, y_2) are coordinates of two different points on the line.

Normal Form:

$$\frac{A}{\pm(A^2 + B^2)^{1/2}}x + \frac{B}{\pm(A^2 + B^2)^{1/2}}y + \frac{C}{\pm(A^2 + B^2)^{1/2}} = 0$$

The sign of the radical is chosen so that

$$\frac{C}{\pm(A^2 + B^2)^{1/2}} < 0$$

Distance From Point (x_1, y_1) to a Line: Substitute coordinates of the point in the normal form of the line. Thus

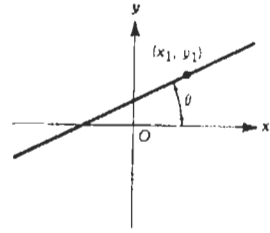


Fig. 15. Point-slope.

$$\text{distance} = \frac{A}{\pm(A^2 + B^2)^{1/2}}x_1 + \frac{B}{\pm(A^2 + B^2)^{1/2}}y_1 + \frac{C}{\pm(A^2 + B^2)^{1/2}}$$

Angle Between Two Lines:

$$\tan\phi = (s_1 - s_2)/(1 + s_1s_2)$$

where,

ϕ = angle between the lines,

s_1 = slope of one line,

s_2 = slope of other line.

When the lines are mutually perpendicular, $\tan\phi = \pm\infty$, whence

$$s_1 = -1/s_2$$

Transformation of Rectangular Coordinates

Translation:

$$x_1 = h + x_2$$

$$y_1 = k + y_2$$

$$x_2 = x_1 - h$$

$$y_2 = y_1 - k$$

(h, k) = coordinates of new origin referred to old origin

Rotation (Fig. 16):

$$x_1 = x_2 \cos\theta - y_2 \sin\theta$$

$$y_1 = x_2 \sin\theta + y_2 \cos\theta$$

$$x_2 = x_1 \cos\theta + y_1 \sin\theta$$

$$y_2 = -x_1 \sin\theta + y_1 \cos\theta$$

(x_1, y_1) = "old" coordinates

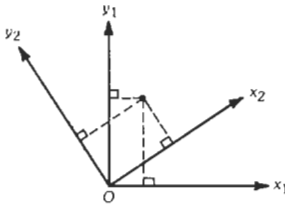


Fig. 16. Rotation of axes.

(x_2, y_2) = "new" coordinates

θ = counterclockwise angle of rotation of axes

Circle

The equation of a circle of radius r with center at (m, n) is

$$(x - m)^2 + (y - n)^2 = r^2$$

Tangent Line to a Circle: At (x_1, y_1) is

$$y - y_1 = -[(x_1 - m)/(y_1 - n)](x - x_1)$$

Normal Line to a Circle: At (x_1, y_1) is

$$y - y_1 = [(y_1 - n)/(x_1 - m)](x - x_1)$$

Parabola

Fig. 17 shows an x -parabola centered at the origin open to the right.

Focus: F

Directrix: D

Vertex: O

Latus rectum: AA'

e = eccentricity = 1

$MP/FP = 1$ for any point P on the parabola

x-Parabola:

$$(y - k)^2 = \pm 2p(x - h)$$

where (h, k) are the coordinates of the vertex, and the sign used is plus or minus when the parabola is open to the right or to the left, respectively. The semilatus rectum is p .

y-Parabola:

$$(x - h)^2 = \pm 2p(y - k)$$

where (h, k) are the coordinates of the vertex. Use plus sign if parabola is open above, and minus sign if open below.

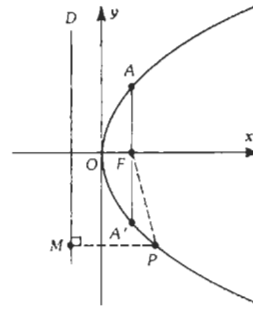


Fig. 17. Parabola.

Tangent Lines to a Parabola:

(x_1, y_1) = point of tangency

For x -parabola

$$y - y_1 = \pm [p/(y_1 - k)](x - x_1)$$

Use plus sign if parabola is open to the right, minus sign if open to the left.

For y -parabola

$$y - y_1 = \pm [(x_1 - h)/p](x - x_1)$$

Use plus sign if parabola is open above, minus sign if open below.

Normal Lines to a Parabola:

(x_1, y_1) = point of contact

For x -parabola

$$y - y_1 = \mp [(y_1 - k)/p](x - x_1)$$

Use minus sign if parabola is open to the right, plus sign if open to the left.

For y -parabola

$$y - y_1 = \mp [p/(x_1 - h)](x - x_1)$$

Use minus sign if parabola is open above, plus sign if open below.

Ellipse

Fig. 18 shows an ellipse centered at the origin. If the ellipse is centered at (h, k) instead, the equations that follow must be modified by replacing x, x_1, y, y_1 by $x - h, x_1 - h, y - k, y_1 - k$, respectively.

Foci: F, F'

Directrices: D, D'

e = eccentricity < 1

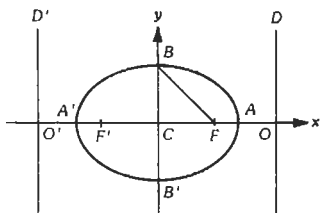


Fig. 18. Ellipse.

$$2a = A'A = \text{major axis}$$

$$2b = BB' = \text{minor axis}$$

$$2c = FF' = \text{focal distance}$$

Then

$$OC = ae$$

$$BF = a$$

$$FC = ae$$

$$1 - e^2 = b^2/a^2$$

Equation of Ellipse:

$$(x^2/a^2) + (y^2/b^2) = 1$$

Sum of the Focal Radii:

$$\text{To any point on ellipse} = 2a$$

Equation of Tangent Line to Ellipse:

$$(x_1, y_1) = \text{point of tangency}$$

$$(xx_1/a^2) + (yy_1/b^2) = 1$$

Equation of Normal Line to an Ellipse:

$$y - y_1 = (a^2y_1/b^2x_1)(x - x_1)$$

Hyperbola

Fig. 19 shows an *x*-hyperbola centered at the origin. If the hyperbola is centered at (h, k) instead, the equations that follow must be modified by replacing x, x_1, y, y_1 by $x - h, x_1 - h, y - k, y_1 - k$, respectively.

Foci: F, F'

Directrices: D, D'

$$e = \text{eccentricity} > 1$$

$$2a = \text{transverse axis} = A'A$$

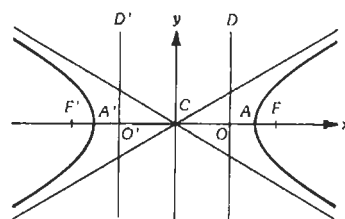


Fig. 19. Hyperbola.

$$CO = ae$$

$$CF = ae$$

Equation of *x*-Hyperbola:

$$(x^2/a^2) - (y^2/b^2) = 1$$

where,

$$b^2 = a^2(e^2 - 1)$$

Equation of Conjugate (*y*-) Hyperbola:

$$(y^2/b^2) - (x^2/a^2) = 1$$

Tangent Line to *x*-Hyperbola:

$$(x_1, y_1) = \text{point of tangency}$$

$$a^2y_1y - b^2x_1x = -a^2b^2$$

Normal Line to *x*-Hyperbola:

$$y - y_1 = -(a^2y_1/b^2x_1)(x - x_1)$$

Asymptotes to Hyperbola:

$$y = \pm(b/a)x$$

SOLID ANALYTIC GEOMETRY

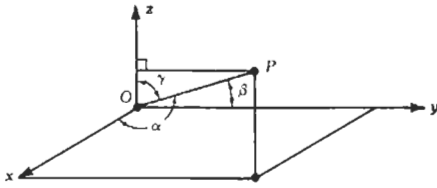
In the following, $x, y,$ and z are the coordinates of a variable point in space in a right-handed rectangular-coordinate system. See Fig. 20.

Coordinates

The coordinates of a point are given by a triplet, e.g., (x_0, y_0, z_0) .

Direction Cosines and Numbers

In Fig. 20, $\alpha, \beta,$ and γ are the angles line OP makes with the $x, y,$ and z axes, respectively. $\text{Cos}\alpha, \text{Cos}\beta,$ and $\text{Cos}\gamma$ are the direction cosines of line OP . Numbers

Fig. 20. Coordinates of point P .

proportional to the direction cosines are called direction numbers.

Distance Between Two Points

$$d = [(x_1 - x_0)^2 + (y_1 - y_0)^2 + (z_1 - z_0)^2]^{1/2}$$

Equations of a Plane

General Form:

$$Ax + By + Cz + D = 0$$

Three-Point Form:

$$\begin{vmatrix} x - x_0 & y - y_0 & z - z_0 \\ x_1 - x_0 & y_1 - y_0 & z_1 - z_0 \\ x_2 - x_0 & y_2 - y_0 & z_2 - z_0 \end{vmatrix} = 0$$

where (x_0, y_0, z_0) , (x_1, y_1, z_1) , and (x_2, y_2, z_2) are three points on the plane.

Intercept Form:

$$(x/a) + (y/b) + (z/c) = 1$$

where a , b , and c are the intercepts on the x , y , and z axes, respectively.

Point-Direction Form:

$$A(x - x_1) + B(y - y_1) + C(z - z_1) = 0$$

where,

A , B , and C are the direction numbers of a normal to the plane,

(x_1, y_1, z_1) is a point on the plane.

Normal Form:

$$x \cos \alpha + y \cos \beta + z \cos \gamma - p = 0$$

where,

$\cos \alpha$, $\cos \beta$, and $\cos \gamma$ are the direction cosines of a normal to the plane,

p is the distance of the plane to the origin.

Equations of a Straight Line

Point-Direction Form:

$$(x - x_1)/A = (y - y_1)/B = (z - z_1)/C$$

or

$$x = At + x_1$$

$$y = Bt + y_1$$

$$z = Ct + z_1 \text{ (parametric)}$$

where,

A , B , and C are direction numbers of the line, (x_1, y_1, z_1) is a point on the line.

Two-Point Form:

$$\begin{aligned} (x - x_1)/(x_1 - x_0) &= (y - y_1)/(y_1 - y_0) \\ &= (z - z_1)/(z_1 - z_0) \end{aligned}$$

or

$$x = x_0 + (x_1 - x_0)t$$

$$y = y_0 + (y_1 - y_0)t$$

$$z = z_0 + (z_1 - z_0)t \text{ (parametric)}$$

where (x_0, y_0, z_0) and (x_1, y_1, z_1) are two points on the line.

Ellipsoid

$$(x^2/a^2) + (y^2/b^2) + (z^2/c^2) = 1$$

where a , b , c are the semiaxes of the ellipsoid or the intercepts on the x , y , and z axes, respectively.

Prolate Spheroid

$$a^2(y^2 + z^2) + b^2x^2 = a^2b^2$$

where $a > b$, and x -axis = axis of revolution.

Oblate Spheroid

$$b^2(x^2 + z^2) + a^2y^2 = a^2b^2$$

where $a > b$, and y -axis = axis of revolution.

Paraboloid of Revolution

$$y^2 + z^2 = 2px$$

x -axis = axis of revolution

Hyperboloid of Revolution

Revolving an x -hyperbola about the x -axis results in the hyperboloid of two sheets

$$a^2(y^2 + z^2) - b^2x^2 = -a^2b^2$$

Revolving an x -hyperbola about the y -axis results in the hyperboloid of one sheet

$$b^2(x^2 + z^2) - a^2y^2 = a^2b^2$$

DIFFERENTIAL CALCULUS

List of Derivatives

In the following, u, v, w are differentiable functions of x , and c is a constant.

General Equations:

$$dc/dx = 0$$

$$dx/dx = 1$$

$$(d/dx)(u + v - w) = (du/dx) + (dv/dx) - (dw/dx)$$

$$(d/dx)(cv) = c(dv/dx)$$

$$(d/dx)(uv) = u(dv/dx) + v(du/dx)$$

$$(d/dx)(v^c) = cv^{c-1}(dv/dx)$$

$$\frac{d}{dx} \left(\frac{u}{v} \right) = \frac{v(du/dx) - u(dv/dx)}{v^2}$$

$$dy/dx = (dy/dv)(dv/dx), \quad \text{if } y = y(v)$$

$$dy/dx = (dx/dy)^{-1}, \quad \text{if } dx/dy \neq 0$$

Transcendental Functions:

$$(d/dx) \ln v = v^{-1}(dv/dx)$$

$$(d/dx) (c^v) = c^v \ln c (dv/dx)$$

$$(d/dx) (e^v) = e^v (dv/dx)$$

$$(d/dx) (u^v) = vu^{v-1}(du/dx) + (\ln u)u^v(dv/dx)$$

$$(d/dx) (\sin v) = \cos v (dv/dx)$$

$$(d/dx) (\cos v) = -\sin v (dv/dx)$$

$$(d/dx) (\tan v) = \sec^2 v (dv/dx)$$

$$(d/dx) (\cot v) = -\csc^2 v (dv/dx)$$

$$(d/dx) (\sec v) = \sec v \tan v (dv/dx)$$

$$(d/dx) (\csc v) = \csc v \cot v (dv/dx)$$

$$(d/dx) (\arcsin v) = (1 - v^2)^{-1/2} (dv/dx)$$

$$(d/dx) (\arccos v) = -(1 - v^2)^{-1/2} (dv/dx)$$

$$(d/dx) (\arctan v) = (1 + v^2)^{-1} (dv/dx)$$

$$(d/dx) (\operatorname{arccot} v) = -(1 + v^2)^{-1} (dv/dx)$$

$$(d/dx) (\operatorname{arcsec} v) = [v(v^2 - 1)^{1/2}]^{-1} (dv/dx)$$

$$(d/dx) (\operatorname{arccsc} v) = -[v(v^2 - 1)^{1/2}]^{-1} (dv/dx)$$

$$(d/dx) (\sinh v) = \cosh v (dv/dx)$$

$$(d/dx) (\cosh v) = \sinh v (dv/dx)$$

$$(d/dx) (\tanh v) = \operatorname{sech}^2 v (dv/dx)$$

TABLE OF INTEGRALS

Indefinite Integrals

General Equations:

$$\int af(x)dx = a \int f(x)dx$$

$$\int [f(x) + g(x)]dx = \int f(x)dx + \int g(x)dx$$

$$(d/dx) \int f(x)dx = f(x)$$

$$(d/dx) \int_u^v f(y, x)dy$$

$$= f(v, x) \cdot (dv/dx) - f(u, x) \cdot (du/dx) + \int_u^v \frac{\partial f(y, x)}{\partial x} dy$$

$$\int f'(x) g(x)dx = f(x) g(x) - \int f(x) g'(x)dx$$

$$\int f(x)dx = \int f[h(y)]h'(y)dy$$

$$\int_{x_1}^{x_2} f(x)dx = \int_{-x_1}^{-x_2} f(-x)dx$$

$$\int f(\sin x, \cos x)dx = \int f\left(\frac{2z}{1+z^2}, \frac{1-z^2}{1+z^2}\right) \frac{dz}{1+z^2}$$

Elementary Forms:

$$\int x^m dx = x^{m+1}/(m+1), m \neq -1$$

$$\int (dx/x) = \ln |x|$$

$$\int e^x dx = e^x$$

$$\int \ln x dx = x \ln x - x$$

$$\int \sin x dx = -\cos x$$

$$\int \cos x dx = \sin x$$

$$\int \tan x dx = -\ln |\cos x|$$

$$\int \cot x dx = \ln |\sin x|$$

$$\int \csc x dx = \ln \left| \tan \frac{1}{2} x \right|$$

$$\int \sec x dx = \ln |\sec x + \tan x|$$

$$\int \sinh x dx = \cosh x$$

$$\int \cosh x dx = \sinh x$$

$$\int \tanh x dx = \ln |\cosh x|$$

$$\int \coth x dx = \ln |\sinh x|$$

$$\int \operatorname{sech} x dx = 2 \tan^{-1} e^x$$

$$\int \operatorname{csch} x dx = \ln \left| \tanh \frac{1}{2} x \right|$$

$$\int dx/(1+x^2) = \tan^{-1} x$$

$$\int dx/(1-x^2) = \frac{1}{2} \ln |(1+x)/(1-x)|$$

$$\int dx/(x^2-1) = \frac{1}{2} \ln |(x-1)/(x+1)|$$

$$\int dx/(1-x^2)^{1/2} = \sin^{-1} x$$

$$\int dx/(x^2 \pm 1)^{1/2} = \ln |x + (x^2 \pm 1)^{1/2}|$$

$$\int (x^2 \pm 1)^{1/2} dx = \frac{1}{2} [x(x^2 \pm 1)^{1/2} \pm \ln |x + (x^2 \pm 1)^{1/2}|]$$

$$\int (1-x^2)^{1/2} dx = \frac{1}{2} [x(1-x^2)^{1/2} + \sin^{-1} x]$$

Forms Containing $ax + b$, $a \neq 0$, $b \neq 0$:

$$\int (ax + b)^n dx = [(ax + b)^{n+1}/a(n+1)], n \neq -1$$

$$\int (ax + b)^{-1} dx = (1/a) \ln |ax + b|$$

$$\int x dx/(ax + b) = a^{-2} [ax + b - b \ln (ax + b)]$$

$$\int x dx/(ax + b)^2 = a^{-2} [b/(ax + b) + \ln (ax + b)]$$

$$\int dx/[x(ax + b)] = b^{-1} \ln [x/(ax + b)]$$

$$\int dx/[x(ax + b)^2] = [b(ax + b)]^{-1} + b^{-2} \ln [x/(ax + b)]$$

$$\int dx/[x^2(ax + b)] = -(bx)^{-1} + (a/b^2) \ln [(ax + b)/x]$$

$$\int dx/[x^2(ax + b)^2] = -(2ax + b)/[b^2 x(ax + b)] + (2a/b^3) \ln [(ax + b)/x]$$

$$\int x^m(ax + b)^n dx = [x^{m+1}(ax + b)^n / (m + n + 1)] + [bn / (m + n + 1)] \int x^m(ax + b)^{n-1} dx$$

$$\int \frac{x^m dx}{(ax + b)^n} = [b(1 - n)]^{-1} \left[\frac{-x^{m+1}}{(ax + b)^{n-1}} + (m - n + 2) \int \frac{x^m}{(ax + b)^{n-1}} dx \right], n \neq 1$$

$$\int \frac{x^m dx}{ax + b} = \frac{x^m}{am} + (b/a) \int \frac{x^{m-1}}{ax + b} dx, m \neq 0$$

Forms Containing $(ax + b)^{1/2}$, $a \neq 0$, $b \neq 0$:

$$\int x(ax + b)^{1/2} dx = \frac{2(3ax - 2b)[(ax + b)^3]^{1/2}}{15a^2}$$

$$\int x^m(ax + b)^{1/2} dx = [2/a(2m + 3)] \left\{ x^m[(ax + b)^3]^{1/2} - mb \int x^{m-1}(ax + b)^{1/2} dx \right\}$$

$$\int \frac{(ax + b)^{1/2}}{x} dx = 2(ax + b)^{1/2} + (b)^{1/2} \ln \left| \frac{(ax + b)^{1/2} - (b)^{1/2}}{(ax + b)^{1/2} + (b)^{1/2}} \right|, b > 0$$

$$= 2(ax + b)^{1/2} - 2(-b)^{1/2} \tan^{-1}[-(ax + b)/b]^{1/2}, b < 0$$

$$\int \frac{(ax + b)^{1/2}}{x^m} dx = -[(m - 1)b]^{-1} \frac{[(ax + b)^3]^{1/2}}{x^{m-1}} + \frac{1}{2}[(2m - 5)a] \int \frac{(ax + b)^{1/2} dx}{x^{m-1}}, m \neq 1$$

$$\int \frac{xdx}{(ax + b)^{1/2}} = [2(ax - 2b)/3a^2](ax + b)^{1/2}$$

$$\int \frac{x^m dx}{(ax + b)^{1/2}} = \frac{2x^m(ax + b)^{1/2}}{(2m + 1)a} - \frac{2bm}{(2m + 1)a} \int \frac{x^{m-1} dx}{(ax + b)^{1/2}}$$

$$\int \frac{dx}{x(ax + b)^{1/2}} = (b)^{-1/2} \ln \left| \frac{(ax + b)^{1/2} - (b)^{1/2}}{(ax + b)^{1/2} + (b)^{1/2}} \right|, b > 0$$

$$= [2/(-b)^{1/2}] \tan^{-1}[(ax + b)/-b]^{1/2}, b < 0$$

$$\int \frac{dx}{x^n(ax + b)^{1/2}} = -\frac{(ax + b)^{1/2}}{(n - 1)bx^{n-1}} - \frac{(2n - 3)a}{(2n - 2)b} \int \frac{dx}{x^{n-1}(ax + b)^{1/2}}, n \neq 1$$

Forms Containing $R = ax^2 + bx + c$, $a \neq 0$, $x > 0$: Let $q = 4ac - b^2$.

$$\int \frac{dx}{R} = [2/(q)^{1/2}] \tan^{-1}[(2ax + b)/(q)^{1/2}], q > 0$$

$$= (-q)^{-1/2} \ln \frac{2ax + b - (-q)^{1/2}}{2ax + b + (-q)^{1/2}}, q < 0$$

(If $q = 0$, R is a perfect square.)

$$\int \frac{dx}{R^n} = \frac{2ax + b}{(n - 1)qR^{n-1}} + \frac{2(2n - 3)a}{q(n - 1)} \int \frac{dx}{R^{n-1}}$$

$$\int \frac{x^m}{R^n} dx = \frac{x^{m-1}}{(2n-m-1)aR^{n-1}} - \frac{n-m}{2n-m-1} (b/a) \int \frac{x^{m-1} dx}{R^n} + \frac{m-1}{2n-m-1} (c/a) \int \frac{x^{m-2} dx}{R^n}$$

$$\int \frac{dx}{x^m R^n} = -[(m-1)cx^{m-1}R^{n-1}]^{-1} - \frac{m+n-2}{m-1} (b/c) \int \frac{dx}{x^{m-1}R^n} - \frac{m+2n-3}{m-1} (a/c) \int \frac{dx}{x^{m-2}R^n}$$

Forms Containing $(R)^{1/2} = (ax^2 + bx + c)^{1/2}$, $a \neq 0$: Let $q = 4ac - b^2$.

$$\int \frac{dx}{(R)^{1/2}} = (a)^{-1/2} \ln \left[(R)^{1/2} + x(a)^{1/2} + \frac{b}{2(a)^{1/2}} \right], a > 0$$

$$= \frac{-1}{(-a)^{1/2}} \sin^{-1} \left[\frac{2ax + b}{(-q)^{1/2}} \right], a < 0$$

$$\int (R)^{1/2} dx = \frac{(2ax + b)(R)^{1/2}}{4a} + \frac{q}{8a} \int \frac{dx}{(R)^{1/2}}$$

$$\int \frac{dx}{R^n (R)^{1/2}} = \frac{2(2ax + b)(R)^{1/2}}{(2n-1)qR^n} + \frac{8a(n-1)}{(2n-1)q} \int \frac{dx}{R^{n-1} (R)^{1/2}}$$

$$\int R^n (R)^{1/2} dx = \frac{(2ax + b)R^n (R)^{1/2}}{4(n+1)a} + \frac{(2n+1)q}{8(n+1)a} \int \frac{R^n dx}{(R)^{1/2}}$$

$$\int \frac{xdx}{(R)^{1/2}} = \frac{(R)^{1/2}}{a} - \frac{b}{2a} \int \frac{dx}{(R)^{1/2}}$$

$$\int x(R)^{1/2} dx = \frac{R(R)^{1/2}}{3a} - \frac{b}{2a} \int (R)^{1/2} dx$$

$$\int \frac{x^m dx}{R^n (R)^{1/2}} = a^{-1} \int \frac{x^{m-2} dx}{R^{n-1} (R)^{1/2}} - (b/a) \int \frac{x^{m-1} dx}{R^n (R)^{1/2}} - (c/a) \int \frac{x^{m-2} dx}{R^n (R)^{1/2}}$$

$$\int \frac{x^m R^n}{(R)^{1/2}} dx = \frac{x^{m-1} R^n (R)^{1/2}}{(2n+m)a} - \frac{(2n+2m+1)b}{2a(2n+m)} \int \frac{x^{m-1} R^n dx}{(R)^{1/2}} - \frac{(m-1)c}{(2n+m)a} \int \frac{x^{m-2} R^n dx}{(R)^{1/2}}$$

$$\int \frac{dx}{x^m R^n (R)^{1/2}} = -\frac{(R)^{1/2}}{(m-1)cx^{m-1}R^n} - \frac{(2n+2m-3)b}{2c(m-1)} \int \frac{dx}{x^{m-1}R^n (R)^{1/2}} - \frac{(2n+m-2)a}{(m-1)c} \int \frac{dx}{x^{m-2}R^n (R)^{1/2}}$$

$$\int \frac{R^n dx}{x^m (R)^{1/2}} = -\frac{R^{n-1} (R)^{1/2}}{(m-1)x^{m-1}} + \frac{(2n-1)b}{2(m-1)} \int \frac{R^{n-1} dx}{x^{m-1} (R)^{1/2}} + \frac{(2n-1)a}{m-1} \int \frac{R^{n-1} dx}{x^{m-2} (R)^{1/2}}$$

Logarithmic Integrands

$$\int \ln ax dx = x(\ln ax - 1)$$

$$\int \log_b x dx = \log_b e (\ln x - 1)x = [(\ln x - 1)x / \ln b]$$

$$\int (\ln x)^n dx = x(\ln x)^n - n \int (\ln x)^{n-1} dx$$

$$\int x^m \ln x dx = x^{m+1}[(m+1)^{-1} \ln x - (m+1)^{-2}], m \neq -1$$

$$\int \frac{\ln x}{x} dx = \frac{1}{2}(\ln x)^2$$

$$\int x^m (\ln x)^n dx = [x^{m+1}(\ln x)^n / (m+1)] - [n/(m+1)] \int x^m (\ln x)^{n-1} dx, m \neq -1$$

$$\int [(\ln x)^n / x] dx = (\ln x)^{n+1} / (n+1), n \neq -1$$

$$\int dx / (x \ln x) = \ln |\ln x|$$

$$\int \frac{x^m dx}{(\ln x)^n} = -\frac{x^{m+1}}{(n-1)(\ln x)^{n-1}} + \frac{m+1}{n-1} \int \frac{x^m dx}{(\ln x)^{n-1}}, n \neq 1$$

$$\int (x^m / \ln x) dx = \text{Ei}[(n+1) \ln x], m \neq -1$$

Exponential Integrands

$$\int a^{bx} dx = a^{bx} / (b \ln a)$$

$$\int x^m e^x dx = x^m e^x - m \int x^{m-1} e^x dx$$

$$\int dx / (a + be^{mx}) = (x/a) - (am)^{-1} \ln |a + be^{mx}|, a \text{ and } m \neq 0$$

$$\int \frac{dx}{(a + be^{mx})^{1/2}} = [m(a)^{1/2}]^{-1} \ln \frac{(a + be^{mx})^{1/2} - (a)^{1/2}}{(a + be^{mx})^{1/2} + (a)^{1/2}}, a > 0$$

$$= \{2/[m(-a)^{1/2}]\} \text{arc tan}[(a + be^{mx})^{1/2} / (-a)^{1/2}], a < 0$$

Trigonometric Integrands

$$\int \sin^2 x dx = \frac{1}{2}(x - \sin x \cos x)$$

$$\int \sin^n x dx = -[(\sin^{n-1} x \cos x) / n] + (n-1) / n \int \sin^{n-2} x dx$$

$$\int \cos^2 x dx = \frac{1}{2}(x + \sin x \cos x)$$

$$\int \cos^n x dx = [(\cos^{n-1} x \sin x) / n] + (n-1) / n \int \cos^{n-2} x dx$$

$$\int \sin x \cos^m x dx = -\cos^{m+1} x / (m + 1)$$

$$\int \sin^m x \cos x dx = \sin^{m+1} x / (m + 1)$$

$$\begin{aligned} \int \sin^n x \cos^m x &= \frac{\cos^{m-1} x \sin^{n+1} x}{m + n} + \frac{m - 1}{m + n} \int \cos^{m-2} x \sin^n x dx \\ &= -\frac{\sin^{n-1} x \cos^{m+1} x}{m + n} + \frac{n - 1}{m + n} \int \cos^m x \sin^{n-2} x dx \end{aligned}$$

$$\int \frac{\sin^n x dx}{\cos^m x} = (m - 1)^{-1} \left[\frac{\sin^{n-1} x}{\cos^{m-1} x} - (n - 1) \int \frac{\sin^{n-2} x dx}{\cos^{m-2} x} \right]$$

$$\int \frac{\cos^m x}{\sin^n x} dx = -(n - 1)^{-1} \left[\frac{\cos^{m-1} x}{\sin^{n-1} x} - (m - 1) \int \frac{\cos^{m-2} x dx}{\sin^{n-2} x} \right]$$

$$\begin{aligned} \int \frac{dx}{\sin^m x \cos^n x} &= [(n - 1) \sin^{m-1} x \cos^{n-1} x]^{-1} + \frac{m + n - 2}{n - 1} \int \frac{dx}{\sin^m x \cos^{n-2} x} \\ &= -[(m - 1) \sin^{m-1} x \cos^{n-1} x]^{-1} + \frac{m + n - 2}{m - 1} \int \frac{dx}{\sin^{m-2} x \cos^n x} \end{aligned}$$

$$\int \tan^n x dx = [\tan^{n-1} x / (n - 1)] - \int \tan^{n-2} x dx$$

$$\int \cot^n x dx = -[\cot^{n-1} x / (n - 1)] - \int \cot^{n-2} x dx$$

$$\int \sec^2 x dx = \tan x$$

$$\int \sec^n x dx = [\sin x / (n - 1) \cos^{n-1} x] + [(n - 2) / (n - 1)] \int \sec^{n-2} x dx, n \neq 1$$

$$\begin{aligned} \int \frac{dx}{a + b \cos x + c \sin x} &= \frac{2}{(a^2 - b^2 - c^2)^{1/2}} \tan^{-1} \frac{(a - b) \tan(x/2) + c}{(a^2 - b^2 - c^2)^{1/2}}, a^2 > b^2 + c^2 \\ &= (b^2 + c^2 - a^2)^{-1/2} \ln \left| \frac{(a - b) \tan(x/2) + c - (b^2 + c^2 - a^2)^{1/2}}{(a - b) \tan(x/2) + c + (b^2 + c^2 - a^2)^{1/2}} \right|, a^2 < b^2 + c^2, a \neq b \\ &= c^{-1} \ln |a + c \tan(x/2)|, a = b \\ &= \frac{-2}{c + (a - b) \tan(x/2)}, a^2 = b^2 + c^2 \end{aligned}$$

$$\int x^n \sin ax dx = -\sum_{k=0}^n k! \binom{n}{k} \frac{x^{n-k}}{a^{k+1}} \cos[ax + (k\pi/2)], n \text{ nonnegative integer}$$

$$\int x^n \cos ax dx = \sum_{k=0}^n k! \binom{n}{k} \frac{x^{n-k}}{a^{k+1}} \sin[ax + (k\pi/2)], n \text{ nonnegative integer}$$

Inverse Trigonometric Integrals, x and $a > 0$:

$$\int \sin^{-1}(x/a)dx = x \sin^{-1}(x/a) + (a^2 - x^2)^{1/2}$$

$$\int \cos^{-1}(x/a)dx = x \cos^{-1}(x/a) - (a^2 - x^2)^{1/2}$$

$$\int \tan^{-1}(x/a)dx = x \tan^{-1}(x/a) - \frac{1}{2}a \ln(a^2 + x^2)$$

$$\int \cot^{-1}(x/a)dx = x \cot^{-1}(x/a) + \frac{1}{2}a \ln(a^2 + x^2)$$

$$\int \sec^{-1}(x/a)dx = x \sec^{-1}(x/a) - a \ln[x + (x^2 - a^2)^{1/2}]$$

$$\int \csc^{-1}(x/a)dx = x \csc^{-1}(x/a) + a \ln[x + (x^2 - a^2)^{1/2}]$$

$$\int x \sin^{-1} x dx = \frac{1}{4}[(2x^2 - 1)\sin^{-1} x + x(1 - x^2)^{1/2}]$$

$$\int x \cos^{-1} x dx = \frac{1}{4}[(2x^2 - 1)\cos^{-1} x - x(1 - x^2)^{1/2}]$$

$$\int x^n \sin^{-1} x dx = [x^{n+1} \sin^{-1} x / (n + 1)] - (n + 1)^{-1} \int [x^{n+1} / (1 - x^2)^{1/2}] dx$$

$$\int x^n \cos^{-1} x dx = [x^{n+1} \cos^{-1} x / (n + 1)] + (n + 1)^{-1} \int [x^{n+1} / (1 - x^2)^{1/2}] dx$$

Miscellaneous Integrals:

$$\int \frac{dx}{ax^3 + b} = \frac{p}{3b} \left[\frac{1}{2} \ln \frac{(x+p)^2}{x^2 - px + p^2} + \sqrt{3} \arctan \frac{x\sqrt{3}}{2p - x} \right], p = (b/a)^{1/3}$$

$$\int \frac{dx}{x^4 + a^4} = (4a^3\sqrt{2})^{-1} \left[\ln \left(\frac{x^2 + ax\sqrt{2} + a^2}{x^2 - ax\sqrt{2} + a^2} \right) + 2 \tan^{-1} \left(\frac{ax\sqrt{2}}{a^2 - x^2} \right) \right]$$

$$\int \frac{dx}{x^4 - a^4} = (1/4a^3) \left[\ln \left(\frac{x - a}{x + a} \right) - 2 \tan^{-1}(x/a) \right]$$

Definite Integrals

$$\int_0^1 x^u (1 - x)^v dx = [\Gamma(u + 1)\Gamma(v + 1)/\Gamma(u + v + 2)], u > -1, v > -1$$

$$= [u!v!/(u + v + 1)!], \text{ if } u, v \text{ nonnegative integers}$$

$$\int_0^\infty x^u dx / (1 + x^v) = (\pi/v) \csc[(u + 1)\pi/v], 0 < u + 1 < v$$

$$\int_0^1 \frac{x^{2n+1} dx}{(1-x^2)^{1/2}} = \frac{(2n)(2n-2)\cdots 6\cdot 4\cdot 2}{(2n+1)(2n-1)\cdots 5\cdot 3\cdot 1}, n \text{ positive integer}$$

$$\int_0^1 \frac{x^{2n} dx}{(1-x^2)^{1/2}} = \frac{(2n-1)(2n-3)\cdots 5\cdot 3\cdot 1}{(2n)(2n-2)\cdots 6\cdot 4\cdot 2} \cdot \frac{1}{2}\pi, n \text{ positive integer}$$

$$\int_0^\infty x^b e^{-ax} dx = \Gamma(b+1)/a^{b+1}, b > -1, a > 0$$

$$= b!/a^{b+1}, b \text{ nonnegative integer}$$

$$\int_0^\infty e^{-ax^2} dx = \frac{1}{2}(\pi/a)^{1/2}, a > 0$$

$$\int_0^\infty x^b e^{-ax^2} dx = \frac{\Gamma[(b+1)/2]}{2a^{(b+1)/2}}, b > -1$$

$$\int_0^\infty \exp[-(ax^2 + bx + c)] dx = \frac{1}{2}(\pi/a)^{1/2} \exp[(b^2 - 4ac)/4a] [1 - \operatorname{erf}(\frac{1}{2}ba^{1/2})]$$

$$\int_0^\infty \exp\{-[x^2 + (a^2/x^2)]\} dx = \frac{1}{2} \exp(-2|a|) (\pi)^{1/2}$$

$$\int_0^{\pi/2} \sin^n x dx = \int_0^{\pi/2} \cos^n x dx = \frac{1}{2}(\pi)^{1/2} \frac{\Gamma[\frac{1}{2}(n+1)]}{\Gamma(\frac{1}{2}n+1)}, n > -1$$

$$\int_0^\infty (\sin mx/x) dx = \frac{1}{2}\pi, m > 0$$

$$= 0, m = 0$$

$$= -\frac{1}{2}\pi, m < 0$$

$$\int_0^\infty (\sin x \cdot \cos mx/x) dx = 0, |m| > 1$$

$$= \frac{1}{4}\pi, m = \pm 1$$

$$= \frac{1}{2}\pi, -1 < m < 1$$

$$\int_0^\infty (\sin x/x)^2 dx = \frac{1}{2}\pi$$

$$\int_0^\infty \cos(x^2) dx = \int_0^\infty \sin(x^2) dx = \frac{1}{2}(\frac{1}{2}\pi)^{1/2}$$

$$\int_0^\infty [\cos mx/(1+x^2)] dx = \frac{1}{2}\pi e^{-|m|}$$

$$\int_0^\infty (\cos x/x^{1/2}) dx = \int_0^\infty (\sin x/x^{1/2}) dx = (\frac{1}{2}\pi)^{1/2}$$

$$\int_0^\infty \exp(-a^2x^2) \cos bx dx = [(\pi)^{1/2}/2 |a|] \exp(-b^2/4a^2)$$

$$\int_0^1 [\ln x/(1-x)] dx = -\frac{1}{6}\pi^2$$

$$\int_0^1 [\ln x/(1+x)] dx = -\frac{1}{12}\pi^2$$

$$\int_0^1 [\ln x/(1-x^2)] dx = -\frac{1}{8}\pi^2$$

$$\int_0^1 [\ln x/(1-x^2)^{1/2}] dx = -\frac{1}{2}\pi \ln 2$$

$$\int_0^1 (\ln x)^n dx = (-1)^n n!$$

$$\int_0^1 x^m (\ln x^{-1})^n dx = [\Gamma(n+1)/(m+1)^{n+1}], m > -1, n > -1$$

$$\int_0^{\pi/2} \ln \sin x dx = -\frac{1}{2}\pi \ln 2$$

$$\int_0^{\pi/2} \ln \cos x dx = -\frac{1}{2}\pi \ln 2$$

$$\int_0^{\pi} x \ln \sin x dx = -\frac{1}{2}\pi^2 \ln 2$$

$$\int_0^{\pi} \ln(a + b \cos x) dx = \pi \ln \frac{1}{2} [a + (a^2 - b^2)^{1/2}], a \geq b$$

SERIES

Taylor's Series for a Single Variable

$$\begin{aligned} f(z) &= f(a) + f'(a)(z-a) + \frac{1}{2}[f''(a)](z-a)^2 + \dots + [f^{(n)}(a)/n!](z-a)^n + R_n \\ &= \sum_{k=0}^n (k!)^{-1} f^{(k)}(a)(z-a)^k + R_n \end{aligned}$$

where the remainder is bounded by $|R_n| \leq [M/(n+1)!] |z-a|^{n+1}$ in which

$$M = \max_{0 \leq \theta \leq 1} |f^{(n+1)}[a + \theta(z-a)]|$$

If $a, z,$ and f are real, then there exists a real $\theta, 0 < \theta < 1,$ such that the remainder is

$$R_n = \{f^{(n+1)}[a + \theta(z-a)]/(n+1)!\} (z-a)^{n+1}$$

When $a = 0,$ this series is often called **Maclaurin's series.**

Taylor's Series for Two Variables

$$\begin{aligned}
 f(x, y) &= f(a, b) + \left[\frac{\partial f(x, y)}{\partial x} \right]_{x=a, y=b} (x - a) + \left[\frac{\partial f(x, y)}{\partial y} \right]_{x=a, y=b} (y - b) \\
 &+ \frac{1}{2} \left[\frac{\partial^2 f(x, y)}{\partial x^2} \right]_{x=a, y=b} (x - a)^2 + 2 \left[\frac{\partial^2 f(x, y)}{\partial x \partial y} \right]_{x=a, y=b} (x - a)(y - b) + \left[\frac{\partial^2 f(x, y)}{\partial y^2} \right]_{x=a, y=b} (y - b)^2 + \dots \\
 &= \sum_{k=0}^n (k!)^{-1} \{ [(x - a)(\partial/\partial \xi) + (y - b)(\partial/\partial \eta)]^k f(\xi, \eta) \}_{\xi=a, \eta=b} + R_n
 \end{aligned}$$

where the remainder is bounded by $|R_n| \leq [M/(n + 1)!](|x - a| + |y - b|)^{n+1}$ in which

$$M = \max_{\substack{k=0, 1, 2, \dots, n+1 \\ 0 \leq \theta \leq 1}} |\partial^{n+1} f(\xi, \eta) / \partial \xi^k \partial \eta^{n+1-k}|, \quad \xi = a + \theta(x - a), \quad \eta = b + \theta(y - b)$$

If $f(x, y)$ is a real function of two variables, there exists a number θ , $0 < \theta < 1$ such that the remainder is

$$R_n = [(n + 1)!]^{-1} \{ [(x - a)(\partial/\partial \xi) + (y - b)(\partial/\partial \eta)]^{n+1} f(\xi, \eta) \}, \quad \xi = a + \theta(x - a), \quad \eta = b + \theta(y - b)$$

Miscellaneous Series

$$\ln(1 + x) = x - \frac{1}{2}x^2 + \frac{1}{3}x^3 - \frac{1}{4}x^4 + \dots = -\sum_{k=1}^{\infty} (-1)^k \frac{x^k}{k}, \quad |x| < 1$$

$$e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots = \sum_{k=0}^{\infty} \frac{x^k}{k!}, \quad |x| < \infty$$

$$\sin x = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \dots = \sum_{k=0}^{\infty} (-1)^k \frac{x^{2k+1}}{(2k+1)!}, \quad |x| < \infty$$

$$\cos x = 1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \frac{x^6}{6!} + \dots = \sum_{k=0}^{\infty} (-1)^k \frac{x^{2k}}{(2k)!}, \quad |x| < \infty$$

$$\sinh x = x + \frac{x^3}{3!} + \frac{x^5}{5!} + \frac{x^7}{7!} + \dots = \sum_{k=0}^{\infty} \frac{x^{2k+1}}{(2k+1)!}, \quad |x| < \infty$$

$$\cosh x = 1 + \frac{x^2}{2!} + \frac{x^4}{4!} + \frac{x^6}{6!} + \dots = \sum_{k=0}^{\infty} \frac{x^{2k}}{(2k)!}, \quad |x| < \infty$$

$$\tan x = x + \frac{x^3}{3} + \frac{2x^5}{15} + \frac{17x^7}{315} + \frac{62x^9}{2835} + \dots, \quad |x| < \frac{1}{2}\pi$$

$$\cot x = \frac{1}{x} - \frac{x}{3} - \frac{x^3}{45} - \frac{2x^5}{945} - \frac{x^7}{4725} - \dots, \quad |x| < \pi$$

$$\arcsin x = x + \frac{1}{2} \frac{x^3}{3} + \frac{1 \cdot 3}{2 \cdot 4} \frac{x^5}{5} + \frac{1 \cdot 3 \cdot 5}{2 \cdot 4 \cdot 6} \frac{x^7}{7} + \dots, \quad |x| < 1$$

$$\operatorname{arctan} x = x - \frac{1}{3}x^3 + \frac{1}{5}x^5 - \frac{1}{7}x^7 + \dots, \quad |x| < 1$$

$$\operatorname{arsinh} x = x - \frac{1}{2} \frac{x^3}{3} + \frac{1 \cdot 3}{2 \cdot 4} \frac{x^5}{5} - \frac{1 \cdot 3 \cdot 5}{2 \cdot 4 \cdot 6} \frac{x^7}{7} + \dots, \quad |x| < 1$$

$$\operatorname{artanh} x = x + \frac{1}{3}x^3 + \frac{1}{5}x^5 + \frac{1}{7}x^7 + \dots, \quad |x| < 1$$

MATRIX ALGEBRA

Notation

A matrix of order $m \times n$ is a rectangular array of numbers, real or complex, consisting of m rows and n columns:

$$A = \begin{bmatrix} a_{11} & a_{12} & a_{13} & \cdots & a_{1n} \\ a_{21} & a_{22} & a_{23} & \cdots & a_{2n} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ a_{m1} & a_{m2} & a_{m3} & \cdots & a_{mn} \end{bmatrix} = [a_{ij}]$$

A row (column) vector is a $1 \times n$ ($n \times 1$) matrix. An $n \times n$ matrix is called a square matrix. A matrix with all entries equal to zero is called a zero matrix; it is denoted by $\mathbf{0}$. A square zero matrix with the elements on the main diagonal, that is, the diagonal extending from the upper left corner to the lower right corner, replaced by ones is called an identity matrix. It is denoted by \mathbf{I}_n or \mathbf{I} . A square matrix with all entries zero above or below the main diagonal is called triangular. For example,

$$(1,2,5), \begin{bmatrix} 2 \\ 3 \end{bmatrix}, \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}, \begin{bmatrix} 1 & 4 & 3 & 4 \\ 0 & 3 & -1 & 5 \\ 0 & 0 & 4 & -2 \\ 0 & 0 & 0 & 2 \end{bmatrix}, \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

are row, column, zero, triangular, and identity matrices, respectively.

Operations

Addition and Subtraction: If \mathbf{A} and \mathbf{B} are matrices of the same order with elements a_{ij} and b_{ij} , respectively, the matrix

$$C = A \pm B$$

has elements

$$c_{ij} = a_{ij} \pm b_{ij}, \quad i = 1, 2, \dots, m; j = 1, 2, \dots, n$$

Multiplication by a Number: If k is a number, real or complex, the matrix

$$C = kA$$

has elements

$$c_{ij} = ka_{ij}, \quad i = 1, 2, \dots, m; j = 1, 2, \dots, n$$

Multiplication of Two Matrices: Let \mathbf{A} and \mathbf{B} be two matrices of orders $m \times n$ and $n \times p$, respectively. The matrix

$$C = AB$$

will have elements

$$c_{ij} = \sum_{k=1}^n a_{ik} b_{kj}, \quad i = 1, 2, \dots, m; j = 1, 2, \dots, p$$

Matrix \mathbf{C} has order $m \times p$. Note that the product \mathbf{BA} is defined only when $m = p$. In general $\mathbf{AB} \neq \mathbf{BA}$ even when $m = n = p$, so it is necessary to distinguish between premultiplication and postmultiplication.

For a square matrix \mathbf{A} of order n , powers are defined

$$A^0 = I, A^1 = A, A^2 = A \cdot A, A^3 = A \cdot A^2, \text{ etc.}$$

A polynomial function of \mathbf{A} is a square matrix of order n given by

$$P(A) = a_n A^n + a_{n-1} A^{n-1} + \cdots + a_1 A^1 + a_0 A^0$$

where the a_i are real or complex numbers.

Division of Two Matrices: Not defined.

Determinant

Definition: The determinant of a square matrix \mathbf{A} of order n is usually defined

$$|A| = \sum \pm a_{1i_1} a_{2j_2} \cdots a_{nr}$$

where the second subscripts i, j, \dots, r form a permutation (rearrangement) of the integers $1, 2, \dots, n$. The sum is taken over all permutations with a plus (minus) sign if the permutation is even (odd). A permutation is called even (odd) if an even (odd) number of inversions is necessary to attain the natural or ascending order. For example, $4132 \rightarrow 1432 \rightarrow 1342 \rightarrow 1324 \rightarrow 1234$; therefore, 4132 is an even permutation.

A square matrix is said to be singular if its determinant is zero and nonsingular otherwise.

Laplace's Development: By an expansion known as Laplace's development, the determinant of a matrix \mathbf{A} of order n can be expressed in terms of determinants of matrices of order $n - 1$.

$$|A| = \sum_{j=1}^n (-1)^{i+j} a_{ij} M_{ij} \quad (\text{expansion by the } i\text{th row})$$

$$= \sum_{i=1}^n (-1)^{i+j} a_{ij} M_{ij} \quad (\text{expansion by the } j\text{th column})$$

The notation M_{ij} represents the determinant of the matrix formed by deleting the i th row and j th column of \mathbf{A} . Expanding by any row or column will lead to the

same value of $|A|$. The M_{ij} are in turn evaluated in terms of determinants of order $n - 2$. This process is continued until, say, second-order determinants are obtained.

$$\begin{vmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{vmatrix} = b_{11}b_{22} - b_{12}b_{21}$$

A first-order determinant has a value equal to its only entry. M_{ij} is called the minor and $(-1)^{i+j} \times M_{ij} = A_{ij}$ is called the cofactor of the element a_{ij} .

Laplace's development is valuable for a literal expansion. For numerical evaluation of determinants of large order, say greater than four, the Gauss algorithm described in the next subsection requires less effort.

Linear Transformations

The linear transformation or set of equations

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n &= y_1 \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n &= y_2 \\ &\vdots \\ a_{m1}x_1 + a_{m2}x_2 + \dots + a_{mn}x_n &= y_m \end{aligned}$$

may be compactly written in matrix form

$$AX = Y$$

where $A = [a_{ij}]$ is a matrix of order $m \times n$, and X and Y are column vectors.

Inverse Matrix: If $m = n$ and $|A| \neq 0$, then there exists an inverse matrix denoted by A^{-1} such that $AA^{-1} = A^{-1}A = I$. The inverse transformation expressing the x 's in terms of the y 's may then be compactly written

$$X = A^{-1}Y$$

This inverse transformation may be effected using *Cramer's rule*

$$x_i = (1/|A|) \sum_{k=1}^n A_{ki} y_k, \quad i = 1, 2, \dots, n$$

where A_{ki} denotes the cofactor associated with the element a_{ki} in the original matrix A . Cramer's rule provides a useful literal expansion of the solution. However, numerical evaluation of the determinants involved, say Laplace's development, requires of the order of $n!$ operations.

The following *Gauss algorithm* requires only of the order of n^3 operations and is therefore preferred for numerical evaluation when n is large: Renumber the x_k 's if necessary to make $a_{11} \neq 0$. Normalize the first equation by dividing it by a_{11} . If $a_{21} = 0$, leave the second equation intact. If $a_{21} \neq 0$, eliminate x_1 by subtracting the normalized first equation multiplied by a_{21} from the second equation. Similarly, eliminate x_1 from the remaining $n - 2$ equations. The result is

$$\begin{aligned} x_1 + (a_{12}/a_{11})x_2 + (a_{13}/a_{11})x_3 + \dots + (a_{1n}/a_{11})x_n &= (y_1/a_{11}) \\ [a_{22} - a_{21}(a_{12}/a_{11})]x_2 + [a_{23} - a_{21}(a_{13}/a_{11})]x_3 + \dots + [a_{2n} - a_{21}(a_{1n}/a_{11})]x_n &= y_2 - a_{21}(y_1/a_{11}) \\ &\vdots \\ [a_{n2} - a_{n1}(a_{12}/a_{11})]x_2 + [a_{n3} - a_{n1}(a_{13}/a_{11})]x_3 + \dots + [a_{nn} - a_{n1}(a_{1n}/a_{11})]x_n &= y_n - a_{n1}(y_1/a_{11}) \end{aligned}$$

The entire process is now repeated with the first equation omitted. There then results a set of the form

$$\begin{aligned} x_1 + b_{12}x_2 + b_{13}x_3 + \dots + b_{1n}x_n &= c_{11}y_1 \\ x_2 + b_{23}x_3 + \dots + b_{2n}x_n &= c_{21}y_1 + c_{22}y_2 \\ b_{33}x_3 + \dots + b_{3n}x_n &= c_{31}y_1 + c_{32}y_2 + c_{33}y_3 \\ &\vdots \\ b_{n3}x_3 + \dots + b_{nn}x_n &= c_{n1}y_1 + c_{n2}y_2 + c_{nn}y_n \end{aligned}$$

Again the process is repeated with the first two equations omitted. Continuing in this manner yields a triangular form

$$\begin{aligned} x_1 + b_{12}x_2 + b_{13}x_3 + b_{14}x_4 + \dots + b_{1n}x_n &= c_{11}y_1 \\ x_2 + b_{23}x_3 + b_{24}x_4 + \dots + b_{2n}x_n &= c_{21}y_1 + c_{22}y_2 \\ x_3 + d_{34}x_4 + \dots + d_{3n}x_n &= e_{31}y_1 + e_{32}y_2 + e_{33}y_3 \\ x_4 + \dots + d_{4n}x_n &= e_{41}y_1 + e_{42}y_2 + e_{43}y_3 + e_{44}y_4 \\ &\vdots \\ x_n &= e_{n1}y_1 + e_{n2}y_2 + e_{n3}y_3 + \dots + e_{nn}y_n \end{aligned}$$

Note that the last equation gives the value of x_n . This may be substituted in the next to the last equation to obtain x_{n-1} and so on. If the y 's are literal as shown above, the process will yield the inverse transformation

$$\mathbf{X} = \mathbf{A}^{-1}\mathbf{Y}$$

where \mathbf{A}^{-1} is the inverse of the original matrix. If the y 's are numerical, labor is saved by combining the values on the right side of each equation at each step of the algorithm.

Since the determinant of a triangular matrix is equal to the product of the elements on the main diagonal

$$|\mathbf{A}| = a_{11}a_{22}'a_{33}' \cdots a_{nn}'$$

where a_{kk}' is the quantity the k th equation is divided by in the above Gauss algorithm. This is useful for evaluating determinants of large order since it requires only of the order of n^3 operations.

A matrix \mathbf{A} may be viewed as consisting of column or row vectors. The largest number of linearly independent column vectors (which is the same as the largest number of linearly independent row vectors) is called the *rank* of the matrix, $\rho(\mathbf{A})$. A set of vectors \mathbf{V}_i is linearly independent if

$$\sum_i a_i \mathbf{V}_i = \mathbf{0}$$

implies that $a_i = 0$ for $i = 1, 2, \dots$

The rank is equal to the order of the largest nonvanishing determinant of the submatrix by deleting rows and columns of the original matrix. Consider the matrices

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}$$

and

$$\mathbf{B} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} & y_1 \\ a_{21} & a_{22} & \cdots & a_{2n} & y_2 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} & y_m \end{bmatrix}$$

The equations

$$\mathbf{AX} = \mathbf{Y}$$

have a solution if and only if

$$\rho(\mathbf{A}) = \rho(\mathbf{B})$$

in which case the equations are said to be consistent.

If $\rho(\mathbf{A}) < n = m$, that is, $|\mathbf{A}| = 0$, and if the equations are consistent, then the Gauss algorithm will terminate before n steps. That is, the coefficients of all the x_k 's will be zero in the remaining $n - \rho(\mathbf{A})$ equations. Therefore, among the x_k 's there will be certain ones, $n - \rho(\mathbf{A})$ in number, which may be assigned arbitrary values. Similarly, if $m \neq n$, the Gauss algorithm will yield an equivalent set of $\rho(\mathbf{A})$ equations that has the same solution as the original set. Again $n - \rho(\mathbf{A})$ (possibly zero) of the x_k 's may be assigned arbitrary values.

It is not necessary to know beforehand whether the equations are inconsistent. If they are inconsistent, the algorithm will yield an "equation" in which the coefficients of the x_k 's on the left side are zero but there is a nonzero combination of the y_k 's on the right side. Since the right side of the "equation" may contain accumulated round-off errors, an analysis of the error propagation in the Gauss algorithm may be necessary to determine whether a small right side is caused by inconsistent equations or by round-off errors.

Eigenvectors and Eigenvalues

An eigenvector of the square matrix \mathbf{A} of order n is nonzero vector \mathbf{X} such that

$$\mathbf{AX} = \lambda\mathbf{X}$$

The scalar λ is called an eigenvalue of \mathbf{A} , and \mathbf{X} is called an eigenvector corresponding to or associated with λ . The eigenvalues may be determined from the characteristic equation

$$|\mathbf{A} - \lambda\mathbf{I}| = 0$$

The corresponding eigenvectors \mathbf{X} may then be found by solving

$$(\mathbf{A} - \lambda\mathbf{I})\mathbf{X} = \mathbf{0}$$

The solution may be obtained by the Gauss algorithm. If the eigenvalues λ are distinct, an explicit solution may be obtained by taking a nontrivial row of cofactors from $\mathbf{A} - \lambda_i\mathbf{I}$. This is possible since the rank of $\mathbf{A} - \lambda_i\mathbf{I}$ is $n - 1$ and, therefore, there exists a nonvanishing subdeterminant of order $n - 1$.

Note that eigenvectors are determined only to within a multiplicative constant.

Further Definitions and Properties

The matrix whose elements are a_{ji}^* is called the conjugate transpose of $\mathbf{A} = [a_{ij}]$; it is denoted by \mathbf{A}^\dagger . The conjugate transpose of a product is

$$(\mathbf{AB})^\dagger = \mathbf{B}^\dagger\mathbf{A}^\dagger$$

If $\mathbf{A} = \mathbf{A}^\dagger$ ($\mathbf{A} = -\mathbf{A}^\dagger$), \mathbf{A} is said to be Hermitian

(skew-Hermitian). If $A^\dagger A = I$, A is called unitary. If $A^\dagger A = AA^\dagger$, A is called normal. Diagonal, unitary, Hermitian, and skew-Hermitian matrices are special cases of normal matrices. If the matrix A is real (that is, all entries are real), the terms Hermitian, skew-Hermitian, and unitary are usually replaced by symmetric, skew-symmetric, and unitary, respectively. The eigenvalues of Hermitian, skew-Hermitian, and unitary matrices are real, pure imaginary, and of unit absolute value, respectively.

The inner product of vectors x and y is

$$X^\dagger Y = \sum_{i=1}^n x_i^* y_i$$

where x_i^* is the complex conjugate of x_i . (A square matrix of order one is considered here as a scalar.)

The length of a vector is given by

$$\|X\| = (X^\dagger X)^{1/2} = \left(\sum_{i=1}^n |x_i|^2 \right)^{1/2}$$

Two vectors are said to be orthogonal if

$$x^\dagger Y = 0$$

Orthogonality corresponds to perpendicularity in two and three dimensions. Inner products obey the following inequalities

$$|X^\dagger Y| \leq \|X\| \|Y\| \text{ (Schwarz)}$$

$$\|X + Y\| \leq \|X\| + \|Y\| \text{ (Triangle)}$$

Two matrices A and B are called similar if there exists a nonsingular matrix C such that

$$B = CAC^{-1}$$

Of particular interest is the case in which matrix B is diagonal

$$B = \begin{bmatrix} \lambda_1 & 0 & 0 & \cdots & 0 \\ 0 & \lambda_2 & 0 & \cdots & 0 \\ 0 & 0 & \lambda_3 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & \lambda_n \end{bmatrix}$$

The diagonal elements λ_i are the eigenvalues of the matrix A . Not all matrices are similar to diagonal ones. However, if a matrix of order n has n linearly independent eigenvectors, as is the case with normal matrices or when all its eigenvalues are distinct, it can be diagonalized. This may be done as follows: Let

$$A \begin{bmatrix} x_{1i} \\ x_{2i} \\ \vdots \\ x_{ni} \end{bmatrix} = \lambda_i \begin{bmatrix} x_{1i} \\ x_{2i} \\ \vdots \\ x_{ni} \end{bmatrix}$$

Then

$$A \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nn} \end{bmatrix}$$

$$= \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nn} \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 & 0 & \cdots & 0 \\ 0 & \lambda_2 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & \lambda_n \end{bmatrix}$$

The matrix formed from the eigenvectors is nonsingular since the eigenvectors are linearly independent. Hence

$$C^{-1} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nn} \end{bmatrix}$$

If a matrix A is reducible to diagonal form, polynomial functions of A are readily calculated

$$f(A) = f(C^{-1}BC)$$

$$= C^{-1} \begin{bmatrix} f(\lambda_1) & 0 & 0 & \cdots & 0 \\ 0 & f(\lambda_2) & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & f(\lambda_n) \end{bmatrix} C$$

Matrices can always be reduced to the Jordan canonical form, in which the eigenvalues are on the main diagonal and there are ones in certain places just above the main diagonal, and zeros elsewhere. In this form, operations such as polynomial functions are simplified.

Hermitian Forms

A Hermitian form is a polynomial

$$X^\dagger AX = \sum_{j=1}^n \sum_{i=1}^n a_{ij} x_i^* x_j$$

where A is a Hermitian matrix. Hermitian forms are real-valued and satisfy the inequality

$$\lambda_n \|X\|^2 \leq X^\dagger AX \leq \lambda_1 \|X\|^2$$

in which λ_1 and λ_n are the largest and smallest eigenvalues of A . By a suitable change of variable

$$Y = TX$$

any Hermitian form $X^\dagger AX$ may be reduced.

$$\mathbf{X}^\dagger \mathbf{A} \mathbf{X} = \mathbf{Y}^\dagger \mathbf{T} \mathbf{A} \mathbf{T} \mathbf{Y} = \sum_{i=1}^n \lambda_i |y_i|^2$$

by choosing \mathbf{T} to be that unitary matrix which diagonalizes \mathbf{A} .

The maximum value of $\mathbf{X}^\dagger \mathbf{A} \mathbf{X}$ for all unit vectors \mathbf{X} , that is, $\|\mathbf{X}\| = 1$, is the largest eigenvalue of \mathbf{A} , say λ_1 . A vector yielding this largest value will be a corresponding eigenvector, say \mathbf{X}_1 . The maximum of $\mathbf{X}^\dagger \mathbf{A} \mathbf{X}$ overall unit vectors orthogonal to \mathbf{X}_1 will be another eigenvalue, say λ_2 . A vector yielding λ_2 will be a corresponding eigenvector, say \mathbf{X}_2 . The process is repeated considering unit vectors orthogonal to both \mathbf{X}_1 and \mathbf{X}_2 . In this way, the eigenvalues of any Hermitian matrix may be found.

VECTOR-ANALYSIS EQUATIONS

Rectangular Coordinates

(In the following, vectors are indicated in bold-faced type.)

Notation:

$$\mathbf{a} = a\hat{\mathbf{a}}$$

a = magnitude of \mathbf{a}

$\hat{\mathbf{a}}$ = unit vector in direction of \mathbf{a}

Associative Law: For addition

$$\mathbf{a} + (\mathbf{b} + \mathbf{c}) = (\mathbf{a} + \mathbf{b}) + \mathbf{c} = \mathbf{a} + \mathbf{b} + \mathbf{c}$$

Commutative Law: For addition

$$\mathbf{a} + \mathbf{b} = \mathbf{b} + \mathbf{a}$$

Scalar or "Dot" Product (Fig. 21):

$$\begin{aligned} \mathbf{a} \cdot \mathbf{b} &= \mathbf{b} \cdot \mathbf{a} \\ &= ab \cos \theta \end{aligned}$$

where θ = angle included by \mathbf{a} and \mathbf{b} .

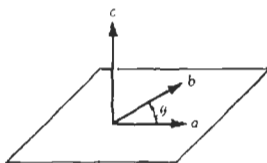


Fig. 21. "Dot" product.

Vector or "Cross" Product:

$$\begin{aligned} \mathbf{a} \times \mathbf{b} &= -\mathbf{b} \times \mathbf{a} \\ &= ab \sin \theta \hat{\mathbf{c}} \end{aligned}$$

where,

θ = smallest angle swept in rotating \mathbf{a} into \mathbf{b} ,

$\hat{\mathbf{c}}$ = unit vector perpendicular to plane of \mathbf{a} and \mathbf{b} , and directed in the sense of travel of a right-hand screw rotating from \mathbf{a} to \mathbf{b} through the angle θ .

Distributive Law for Scalar Multiplication:

$$\mathbf{a} \cdot (\mathbf{b} + \mathbf{c}) = \mathbf{a} \cdot \mathbf{b} + \mathbf{a} \cdot \mathbf{c}$$

Distributive Law for Vector Multiplication:

$$\mathbf{a} \times (\mathbf{b} + \mathbf{c}) = \mathbf{a} \times \mathbf{b} + \mathbf{a} \times \mathbf{c}$$

Scalar Triple Product:

$$\mathbf{a} \cdot (\mathbf{b} \times \mathbf{c}) = (\mathbf{a} \times \mathbf{b}) \cdot \mathbf{c} = \mathbf{c} \cdot (\mathbf{a} \times \mathbf{b}) = \mathbf{b} \cdot (\mathbf{c} \times \mathbf{a})$$

Vector Triple Product:

$$\mathbf{a} \times (\mathbf{b} \times \mathbf{c}) = (\mathbf{a} \cdot \mathbf{c})\mathbf{b} - (\mathbf{a} \cdot \mathbf{b})\mathbf{c}$$

$$(\mathbf{a} \times \mathbf{b}) \cdot (\mathbf{c} \times \mathbf{d}) = (\mathbf{a} \cdot \mathbf{c})(\mathbf{b} \cdot \mathbf{d}) - (\mathbf{a} \cdot \mathbf{d})(\mathbf{b} \cdot \mathbf{c})$$

$$(\mathbf{a} \times \mathbf{b}) \times (\mathbf{c} \times \mathbf{d}) = (\mathbf{a} \times \mathbf{b} \cdot \mathbf{d})\mathbf{c} - (\mathbf{a} \times \mathbf{b} \cdot \mathbf{c})\mathbf{d}$$

Del Operator:

$$\nabla \equiv \mathbf{i}(\partial/\partial x) + \mathbf{j}(\partial/\partial y) + \mathbf{k}(\partial/\partial z)$$

where $\mathbf{i}, \mathbf{j}, \mathbf{k}$ are unit vectors in the directions of the x, y, z coordinate axes, respectively.

Gradient:

$$\text{grad } \phi = \nabla \phi$$

$$= \mathbf{i}(\partial \phi / \partial x) + \mathbf{j}(\partial \phi / \partial y) + \mathbf{k}(\partial \phi / \partial z),$$

in Cartesian coordinates

$$\text{grad}(\phi + \psi) = \text{grad } \phi + \text{grad } \psi$$

$$\text{grad}(\phi\psi) = \phi \text{grad } \psi + \psi \text{grad } \phi$$

Divergence:

$$\text{div } \mathbf{a} = \nabla \cdot \mathbf{a} = (\partial a_x / \partial x) + (\partial a_y / \partial y) + (\partial a_z / \partial z),$$

in Cartesian coordinates

where a_x, a_y, a_z are components of \mathbf{a} in the directions of the x, y, z coordinate axes, respectively.

$$\text{div}(\mathbf{a} + \mathbf{b}) = \text{div } \mathbf{a} + \text{div } \mathbf{b}$$

$$\text{div}(\phi \mathbf{a}) = \phi \text{div } \mathbf{a} + \mathbf{a} \cdot \text{grad } \phi$$

Curl:

$$\text{curl } \mathbf{a} = \nabla \times \mathbf{a}$$

$$= \mathbf{i} \left(\frac{\partial a_z}{\partial y} - \frac{\partial a_y}{\partial z} \right) + \mathbf{j} \left(\frac{\partial a_x}{\partial z} - \frac{\partial a_z}{\partial x} \right) + \mathbf{k} \left(\frac{\partial a_y}{\partial x} - \frac{\partial a_x}{\partial y} \right)$$

$$= \begin{vmatrix} \mathbf{i} & \mathbf{j} & \mathbf{k} \\ \partial/\partial x & \partial/\partial y & \partial/\partial z \\ a_x & a_y & a_z \end{vmatrix}, \quad \text{in Cartesian coordinates}$$

$\text{curl}(\mathbf{a} + \mathbf{b}) = \text{curl} \mathbf{a} + \text{curl} \mathbf{b}$

$\text{curl}(\phi \mathbf{a}) = \text{grad} \phi \times \mathbf{a} + \phi \text{curl} \mathbf{a}$

$\text{curl grad} \phi = 0$

$\text{div curl} \mathbf{a} = 0$

$\text{div}(\mathbf{a} \times \mathbf{b}) = \mathbf{b} \cdot \text{curl} \mathbf{a} - \mathbf{a} \cdot \text{curl} \mathbf{b}$

$\text{Laplacian} \equiv \nabla^2 = \nabla \cdot \nabla$

$\nabla^2 \phi = (\partial^2 \phi / \partial x^2) + (\partial^2 \phi / \partial y^2) + (\partial^2 \phi / \partial z^2),$

in Cartesian coordinates

$\text{curl curl} \mathbf{a} = \text{grad div} \mathbf{a} - (\mathbf{i} \nabla^2 a_x + \mathbf{j} \nabla^2 a_y + \mathbf{k} \nabla^2 a_z)$

$= \nabla(\nabla \cdot \mathbf{a}) - \nabla^2 \mathbf{a}$

Directional Derivative: Derivative of ϕ in the direction of \mathbf{s}

$d\phi/ds = \hat{\mathbf{s}} \cdot \nabla \phi$

Integral Relations: In the following equations, τ is a volume bounded by a closed surface, S . The unit vector \mathbf{n} is normal to the surface and is directed outward. The symbol dS represents an element of surface area. If the surface is represented by $z = f(x, y)$ then

$dS = [1 + (\partial f/\partial x)^2 + (\partial f/\partial y)^2]^{1/2} dx dy$

$\int_{\tau} \nabla \phi d\tau = \int_S \phi \mathbf{n} dS$

$\int_{\tau} \nabla \cdot \mathbf{a} d\tau = \int_S \mathbf{a} \cdot \mathbf{n} dS$ (Gauss' theorem)

$\int_{\tau} \nabla \times \mathbf{a} d\tau = \int_S \mathbf{n} \times \mathbf{a} dS$

$\int_{\tau} (\psi \nabla^2 \phi - \phi \nabla^2 \psi) d\tau, \int_S [\psi(\partial \phi / \partial n) - \phi(\partial \psi / \partial n)] dS$

where $\partial/\partial n$ is the derivative in the direction of \mathbf{n} (Green's theorem).

In the two following equations, S is an open surface bounded by a contour C , with distance along C represented by s .

$\int_S \mathbf{n} \times \nabla \phi dS = \int_C \phi ds$

$\int_S (\nabla \times \mathbf{a}) \cdot \mathbf{n} dS = \int_C \mathbf{a} \cdot ds$ (Stokes' theorem)

where,

$\mathbf{s} = s\hat{\mathbf{s}},$

$\hat{\mathbf{s}}$ is the unit tangent vector along C .

Gradient, Divergence, Curl, and Laplacian in Coordinate Systems Other Than Rectangular

Cylindrical Coordinates: (ρ, ϕ, z) , unit vectors $\hat{\rho}, \hat{\phi}, \mathbf{k}$, respectively

$\text{grad} \psi = \nabla \psi = (\partial \psi / \partial \rho) \hat{\rho} + \rho^{-1} (\partial \psi / \partial \phi) \hat{\phi} + (\partial \psi / \partial z) \mathbf{k}$

Let $\mathbf{a} = a_{\rho} \hat{\rho} + a_{\phi} \hat{\phi} + a_z \mathbf{k}$. Then

$\text{div} \mathbf{a} = \nabla \cdot \mathbf{a} = \rho^{-1} (\partial / \partial \rho) (\rho a_{\rho}) + \rho^{-1} (\partial a_{\phi} / \partial \phi) + (\partial a_z / \partial z)$

$\text{curl} \mathbf{a} = \nabla \times \mathbf{a} = [\rho^{-1} (\partial a_z / \partial \phi) - (\partial a_{\phi} / \partial z)] \hat{\rho} + [(\partial a_{\rho} / \partial z) - (\partial a_z / \partial \rho)] \hat{\phi} + [\rho^{-1} (\partial / \partial \rho) (\rho a_{\phi}) - \rho^{-1} (\partial a_{\rho} / \partial \phi)] \mathbf{k}$

$\nabla^2 \psi = \rho^{-1} (\partial / \partial \rho) [\rho (\partial \psi / \partial \rho)] + \rho^{-2} (\partial^2 \psi / \partial \phi^2) + (\partial^2 \psi / \partial z^2)$

Spherical Coordinates: (r, θ, ϕ) , unit vectors $\hat{r}, \hat{\theta}, \hat{\phi}$

$r =$ distance to origin

θ = polar angle

ϕ = azimuthal angle

$$\text{grad}\psi = \nabla\psi = (\partial\psi/\partial r)\hat{\mathbf{r}} + r^{-1}(\partial\psi/\partial\theta)\hat{\boldsymbol{\theta}} + (r\sin\theta)^{-1}(\partial\psi/\partial\phi)\hat{\boldsymbol{\phi}}$$

Let $\mathbf{a} = a_r\hat{\mathbf{r}} + a_\theta\hat{\boldsymbol{\theta}} + a_\phi\hat{\boldsymbol{\phi}}$. Then

$$\text{div}\mathbf{a} = \nabla \cdot \mathbf{a} = r^{-2}(\partial/\partial r)(r^2a_r) + (r\sin\theta)^{-1}(\partial/\partial\theta)(a_\theta\sin\theta) + (r\sin\theta)^{-1}(\partial a_\phi/\partial\phi)$$

$$\begin{aligned} \text{curl}\mathbf{a} = \nabla \times \mathbf{a} &= (r\sin\theta)^{-1}[(\partial/\partial\theta)(a_\phi\sin\theta) - (\partial a_\theta/\partial\phi)]\hat{\mathbf{r}} \\ &\quad + r^{-1}[(\sin\theta)^{-1}(\partial a_r/\partial\phi) - (\partial/\partial r)(ra_\phi)]\hat{\boldsymbol{\theta}} + r^{-1}[(\partial/\partial r)(ra_\theta) - (\partial a_r/\partial\theta)]\hat{\boldsymbol{\phi}} \end{aligned}$$

$$\nabla^2\psi = r^{-2}(\partial/\partial r)[r^2(\partial\psi/\partial r)] + (r^2\sin\theta)^{-1}(\partial/\partial\theta)[\sin\theta(\partial\psi/\partial\theta)] + (r^2\sin^2\theta)^{-1}(\partial^2\psi/\partial\phi^2)$$

Orthogonal Curvilinear Coordinates:

Coordinates: u_1, u_2, u_3

Metric coefficients: h_1, h_2, h_3 ($ds^2 = h_1^2 du_1^2 + h_2^2 du_2^2 + h_3^2 du_3^2$)

Unit vectors: $\mathbf{i}_1, \mathbf{i}_2, \mathbf{i}_3$ ($d\mathbf{s} = \mathbf{i}_1 h_1 du_1 + \mathbf{i}_2 h_2 du_2 + \mathbf{i}_3 h_3 du_3$)

$$\text{grad}\psi = \nabla\psi = h_1^{-1}(\partial\psi/\partial u_1)\mathbf{i}_1 + h_2^{-1}(\partial\psi/\partial u_2)\mathbf{i}_2 + h_3^{-1}(\partial\psi/\partial u_3)\mathbf{i}_3$$

$$\text{div}\mathbf{a} = \nabla \cdot \mathbf{a} = (h_1 h_2 h_3)^{-1}[(\partial/\partial u_1)(h_2 h_3 a_1) + (\partial/\partial u_2)(h_3 h_1 a_2) + (\partial/\partial u_3)(h_1 h_2 a_3)]$$

$$\begin{aligned} \text{curl}\mathbf{a} = \nabla \times \mathbf{a} &= (h_2 h_3)^{-1}[(\partial/\partial u_2)(h_3 a_3) - (\partial/\partial u_3)(h_2 a_2)]\mathbf{i}_1 + (h_3 h_1)^{-1}[(\partial/\partial u_3)(h_1 a_1) - (\partial/\partial u_1)(h_3 a_3)]\mathbf{i}_2 \\ &\quad + (h_1 h_2)^{-1}[(\partial/\partial u_1)(h_2 a_2) - (\partial/\partial u_2)(h_1 a_1)]\mathbf{i}_3 \end{aligned}$$

$$= (h_1 h_2 h_3)^{-1} \begin{vmatrix} h_1 \mathbf{i}_1 & h_2 \mathbf{i}_2 & h_3 \mathbf{i}_3 \\ \partial/\partial u_1 & \partial/\partial u_2 & \partial/\partial u_3 \\ h_1 a_1 & h_2 a_2 & h_3 a_3 \end{vmatrix}$$

$$\nabla^2\psi = (h_1 h_2 h_3)^{-1} \left[\frac{\partial}{\partial u_1} \left(\frac{h_2 h_3}{h_1} \frac{\partial \psi}{\partial u_1} \right) + \frac{\partial}{\partial u_2} \left(\frac{h_3 h_1}{h_2} \frac{\partial \psi}{\partial u_2} \right) + \frac{\partial}{\partial u_3} \left(\frac{h_1 h_2}{h_3} \frac{\partial \psi}{\partial u_3} \right) \right]$$

Space Curves

A curve may be represented vectorially as $\mathbf{r} = \mathbf{r}(s)$. See Fig. 22.

A unit tangent \mathbf{t} is then given by

$$\mathbf{t} = d\mathbf{r}/ds$$

The principal normal \mathbf{n} is given by

$$\mathbf{n} = (1/k)(d\mathbf{t}/ds)$$

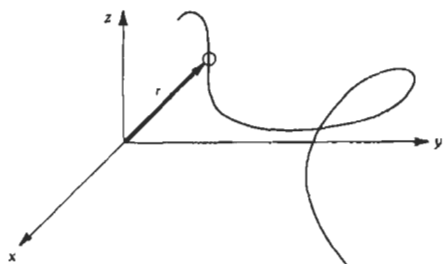


Fig. 22. Space curve.

where k is the curvature. The radius of curvature $R = 1/k$. For a plane curve $y=f(x)$, the curvature may be computed from

$$k = |y''|/[1 + (y')^2]^{3/2}$$

The binormal is defined by

$$\mathbf{b} = \mathbf{t} \times \mathbf{n}$$

These vectors satisfy Frenet's equations

$$d\mathbf{n}/ds = -k\mathbf{t} + \tau\mathbf{b}$$

$$d\mathbf{b}/ds = -\tau\mathbf{n}$$

where τ is the torsion. The torsion is zero everywhere if and only if the curve lies in a plane.

LAPLACE TRANSFORM

The Laplace transform of a function $f(t)$ is defined by the expression

$$F(p) = \int_0^{\infty} f(t)e^{-pt} dt$$

If this integral converges for some $p = p_0$, real or

complex, then it will converge for all p such that $\text{Re}p > \text{Re}p_0$.

The inverse transform may be found by

$$f(t) = (j2\pi)^{-1} \int_{c-j\infty}^{c+j\infty} F(z)e^{tz} dz, t > 0$$

where there are no singularities to the right of the path of integration.

TABLE OF LAPLACE TRANSFORMS

General Equations

	Function	Transform*
Shifting theorem	$f(t - a), f(t) = 0, t < 0$	$e^{-ap}F(p), a > 0$
Convolution	$\int_0^t f_1(\lambda)f_2(t - \lambda)d\lambda$	$F_1(p)F_2(p)$
Linearity	$a_1f_1(t) + a_2f_2(t), (a_1, a_2 \text{ const})$	$a_1F_1(p) + a_2F_2(p)$
Derivative	$df(t)/dt$	$-f(0) + pF(p)$
Integral	$\int f(t)dt$	$p^{-1} \left[\int f(t)dt \right]_{t=0} + [f(p)/p], \text{Re}p > 0$
Periodic function	$f(t) = f(t + r)$	$\int_0^r f(\lambda)e^{-p\lambda} d\lambda / (1 - e^{-pr}), r > 0$
	$f(t) = -f(t + r)$	$\int_0^r f(\lambda)e^{-\lambda} d\lambda / (1 + e^{-pr}), r > 0$
	$f(at), a > 0$	$F(p/a)/a$
	$e^{at}f(t)$	$F(p - a), \text{Re}p > \text{Re}a$
	$t^n f(t)$	$(-1)^n [d^n F(p)/dp^n]$
Final-value theorem	$f(\infty)$	$\lim_{p \rightarrow 0} pF(p)$
Initial-value theorem	$f(0 +)$	$\lim_{p \rightarrow \infty} pF(p)$

* $F(p)$ denotes the Laplace transform of $f(t)$.

Miscellaneous Functions*

	Function	Transform
Step	$u(t - a) = 0, 0 \leq t < a$ $= 1, t \geq a$	e^{-ap}/p
Impulse	$\delta(t)$	1
	$t^a, \text{Re}a > -1$	$\Gamma(a + 1)/p^{a+1}$
	e^{at}	$1/(p - a), \text{Re}p > \text{Re}a$
	$t^a e^{bt}, \text{Re}a > -1$	$\Gamma(a + 1)/(p - b)^{a+1}, \text{Re}p > \text{Re}b$
	$\cos at$ $\sin at$	$p/(p^2 + a^2)$ $a/(p^2 + a^2)$ } $\text{Re}p > \text{Im}a $

Function	Transform
coshat	$p/(p^2 - a^2)$
sinhat	$1/(p^2 - a^2)$
lnr	$-(\gamma + \ln p)/p$, γ is Euler's constant = 0.57722
$1/(t + a)$, $a > 0$	$e^{-ap}E_1(ap)$
e^{-at}	$\frac{1}{2}(\pi/a)^{1/2}e^{p^2/4a} \operatorname{erfc}[p/2(a)^{1/2}]$
Bessel function $J_\nu(at)$, $\operatorname{Re} \nu > -1$	$r^{-1}[(r - p)/a]^\nu$, $r = (p^2 + a^2)^{1/2}$, $\operatorname{Re} p > \operatorname{Re} a $
Bessel function $I_\nu(at)$, $\operatorname{Re} \nu > -1$	$R^{-1}[(R - p)/a]^\nu$, $R = (p^2 - a^2)^{1/2}$, $\operatorname{Re} p > \operatorname{Re} a $

* For an extensive listing, refer to A. Erdélyi, ed., *Tables of Integral Transforms*, Vol. 1, Bateman Manuscript Project, New York: McGraw-Hill Book Co., 1954.

Inverse Transforms*

Transform	Function
1	$\delta(t)$
$1/(p + a)$	e^{-at}
$1/(p + a)^\nu$, $\operatorname{Re} \nu > 0$	$t^{\nu-1}e^{-at}/\Gamma(\nu)$
$1/[(p + a)(p + b)]$	$(e^{-at} - e^{-bt})/(b - a)$
$p/[(p + a)(p + b)]$	$(ae^{-at} - be^{-bt})/(a - b)$
$1/(p^2 + a^2)$	$a^{-1} \operatorname{sinc} at$
$1/(p^2 - a^2)$	$a^{-1} \operatorname{sinhat} at$
$p/(p^2 + a^2)$	cosat
$p/(p^2 - a^2)$	coshat
$1/(p^2 + a^2)^{1/2}$	$J_0(at)$
e^{-ap}/p	$u(t - a)$
e^{-ap}/p^ν , $\operatorname{Re} \nu > 0$	$(t - a)^{\nu-1}u(t - a)/\Gamma(\nu)$
$(1/p)e^{-ap}$	$J_0[2(at)^{1/2}]$
$(1/p^\nu)e^{-ap}$	$(t/a)^{\nu-1/2}J_{\nu-1}[2(at)^{1/2}]$
$(1/p^\nu)e^{ap}$	$(t/a)^{\nu-1/2}I_{\nu-1}[2(at)^{1/2}]$
$(1/p) \ln p$	$-\gamma - \ln t$, $\gamma = 0.57722$

* Refer to A. Erdélyi, ed., *Tables of Integral Transforms*, Vol. 1, Bateman Manuscript Project, New York: McGraw-Hill Book Co., 1954.

Series Expansions:

$$E_1(z) = -\gamma - \ln z - \sum_{n=1}^{\infty} [(-1)^n z^n / (n \cdot n!)]$$

$|\operatorname{arg} z| < \pi$

$$Ei(x) = \gamma + \ln x + \sum_{n=1}^{\infty} x^n / (n \cdot n!)$$

in which γ is Euler's constant, $\gamma = 0.57722$.

Asymptotic Expansions:

$$E_1(z) \sim (e^{-z}/z)[1 - z^{-1} + (2!/z^2) - (3!/z^3) + \dots]$$

$|\operatorname{arg} z| < \pi$

$$Ei(x) \sim (e^x/x)[1 + x^{-1} + (2!/x^2) + (3!/x^3) + \dots]$$

SELECTED FUNCTIONS

Exponential Integrals

Definitions:

$$E_1(z) = \int_z^{\infty} (e^{-t}/t)dt, \quad |\operatorname{arg} z| < \pi$$

in which the path of integration does not cross the negative t -axis and also excludes the origin. For $\operatorname{arg} z = \pi$ the following function is used:

$$E_i(x) = -pv \int_{-x}^{\infty} (e^{-t}/t)dt, \quad x > 0$$

where pv stands for the Cauchy principal value.

$$E_1(-x \pm j0) = E_i(x) \mp j\pi$$

Cosine and Sine Integrals

Definitions:

$$Si(z) = \int_0^z (\operatorname{sint}/t)dt$$

$$Ci(z) = \gamma + \ln z + \int_0^z [(cost - 1)/t]dt, \quad |\operatorname{arg} z| < \pi$$

where γ is Euler's constant, $\gamma = 0.57722$.

Series Expansions:

$$Si(z) = \sum_{n=0}^{\infty} \frac{(-1)^n z^{2n+1}}{(2n+1)(2n+1)!}$$

$$Ci(z) = \gamma + \ln z + \sum_{n=1}^{\infty} \frac{(-1)^n z^{2n}}{2n(2n)!}$$

Asymptotic Expansions:

$$Si(z) \sim \frac{1}{2}\pi - \left(1 - \frac{2!}{z^2} + \frac{4!}{z^4} - \frac{6!}{z^6} + \dots\right) \frac{\cos z}{z}$$

$$- \left(1 - \frac{3!}{z^2} + \frac{5!}{z^4} - \frac{7!}{z^6} + \dots\right) \frac{\sin z}{z^2}$$

$$Ci(z) \sim \left(1 - \frac{2!}{z^2} + \frac{4!}{z^4} - \frac{6!}{z^6} + \dots\right) \frac{\sin z}{z}$$

$$- \left(1 - \frac{3!}{z^2} + \frac{5!}{z^4} - \frac{7!}{z^6} + \dots\right) \frac{\cos z}{z^2}$$

Gamma Function

Definition:

$$\Gamma(z) = \int_0^\infty t^{z-1} e^{-t} dt$$

$$\Gamma(-z) = -\pi / [\Gamma(z+1) \sin \pi z], \quad \text{Re } z > 0$$

The function $\Gamma(z)$ is an analytic function everywhere except at the negative integers.

Identities:

$$z\Gamma(z) = \Gamma(z+1)$$

$$\Gamma\left(\frac{1}{2} + z\right)\Gamma\left(\frac{1}{2} - z\right) = \pi \sec \pi z$$

Special Values:

$$\Gamma(z) = (z-1)! \quad z = 1, 2, 3, \dots$$

$$\Gamma\left(z + \frac{1}{2}\right) = [(2z)! / 2^{2z} z!](\pi)^{1/2}$$

n	1	2	3	4	5	6	7	8	9	10
$\Gamma(n) = (n-1)!$	1	1	2	6	24	120	720	5040	40 320	362 880

Asymptotic Expansion (Stirling's Equation):

$$\Gamma(z) \sim e^{-z} z^{z-1/2} (2\pi)^{1/2} [1 + (1/12z) + (1/288z^2) + \dots]$$

Note $z! = z\Gamma(z) \sim e^{-z} z^{z+1/2} (2\pi)^{1/2} [1 + (1/12z) + (1/288z^2) + \dots]$

Psi and Polygamma Functions

Definitions:

$$\psi^{(n)}(z) = \frac{d^{n+1}}{dz^{n+1}} [\ln \Gamma(z)] \quad n = 0, 1, 2, \dots$$

$\psi^{(0)}(z) = \psi(z)$ is known as the psi function, and $\psi^{(n)}(z)$ as the polygamma function of order n , $n = 1, 2, \dots$

Special Values:

$$\psi^{(n)}(1) = (-1)^{n+1} n! \zeta(n+1), \quad n = 1, 2, \dots$$

$$\psi(1) = -\gamma$$

$$\psi^{(n)}\left(\frac{1}{2}\right) = (-1)^{n+1} n! (2^{n+1} - 1) \zeta(n+1),$$

$n = 1, 2, \dots$

$$\psi\left(\frac{1}{2}\right) = -\gamma - 2 \ln 2$$

in which $\zeta(n+1)$ is the Riemann zeta function of $n+1$.

Series Expansions:

$$\psi^{(n)}(1+z) = \sum_{k=0}^\infty \frac{(-1)^{n+k+1} (n+k)! \zeta(n+k+1) z^k}{k!},$$

$|z| < 1, n = 1, 2, \dots$

$$\psi(1+z) = -\gamma - \sum_{k=1}^\infty (-1)^k \zeta(k+1) z^k, \quad |z| < 1$$

Asymptotic Expansions:

$$\psi^{(n)}(z) \sim \{(-1)^n (n-1)! / z^n\} [1 + (n/2z) + \dots],$$

as $z \rightarrow \infty, |\arg z| < \pi, n = 1, 2, \dots$

$$\psi(z) \sim \ln z - (1/2z) + \dots$$

Psi and Polygamma Functions

Error Function

Definitions:

Definitions:

$$\text{erfc } z = (2/\pi^{1/2}) \int_0^z \exp(-t^2) dt$$

$$\text{erfc } z = (2/\pi^{1/2}) \int_z^\infty \exp(-t^2) dt$$

$$= 1 - \text{erf } z$$

The path of integration for large t must remain within $|\arg t| < \pi/4$ in the latter integral.

Series Expansion:

$$\operatorname{erfc} z = (2/\pi^{1/2}) \sum_{n=0}^{\infty} [(-1)^n z^{2n+1}/n!(2n+1)], |z| < \infty$$

Asymptotic Expansion:

$$\begin{aligned} \operatorname{erfc} z \sim 1 - [\exp(-z^2)/z\pi^{1/2}] \\ \times [1 - (1/2z^2) + (3/4z^4) - \dots], \\ \text{as } z \rightarrow \infty, |\arg z| < 3\pi/4 \end{aligned}$$

Inequality:

$$\begin{aligned} \frac{(2/\pi^{1/2}) \exp(-x^2)}{x + (x^2 + 2)^{1/2}} < \operatorname{erfc} x \\ \leq \frac{(2/\pi^{1/2}) \exp(-x^2)}{x + [x^2 + (4/\pi)]^{1/2}}, x \geq 0 \end{aligned}$$

Derivatives:

$$\begin{aligned} (d^{(n+1)}/dz^{(n+1)})\operatorname{erfc} z \\ = (-1)^n (2/\pi^{1/2}) H_n(z) e^{-z^2}, n = 0, 1, 2, \dots \end{aligned}$$

where $H_n(z)$ is the Hermite polynomial of order n .

Relation to Gaussian Distribution:

$$\begin{aligned} \operatorname{Prob}(x \leq x) = [\sigma(2\pi)^{1/2}]^{-1} \int_{-\infty}^x \exp\left(-\frac{(t-\mu)^2}{2\sigma^2}\right) dt \\ = \frac{1}{2} \left[1 + \operatorname{erf}\left(\frac{x-\mu}{\sigma\sqrt{2}}\right) \right] \end{aligned}$$

Fresnel Integrals

Definitions:

$$C(z) = \int_0^z \cos[(\pi/2)t^2] dt$$

$$S(z) = \int_0^z \sin[(\pi/2)t^2] dt$$

Series Expansions:

$$C(z) = \sum_{n=0}^{\infty} \frac{(-1)^n (\pi/2)^{2n}}{(2n)!(4n+1)} z^{4n+1}$$

$$S(z) = \sum_{n=0}^{\infty} \frac{(-1)^n (\pi/2)^{2n+1}}{(2n+1)!(4n+3)} z^{4n+3}, |z| < \infty$$

Asymptotic Expansions:

$$\begin{aligned} C(z) \sim \frac{1}{2} + \left(1 - \frac{1 \cdot 3}{(\pi z^2)^2} + \frac{1 \cdot 3 \cdot 5 \cdot 7}{(\pi z^2)^4} - \dots \right) \frac{\sin \frac{1}{2} \pi z^2}{\pi z} \\ - \left[(1/\pi z^2) - \frac{1 \cdot 3 \cdot 5}{(\pi z^2)^3} \right. \\ \left. + \frac{1 \cdot 3 \cdot 5 \cdot 7 \cdot 9}{(\pi z^2)^5} - \dots \right] \frac{\cos \frac{1}{2} \pi z^2}{\pi z} \\ S(z) \sim \frac{1}{2} - \left[1 - \frac{1 \cdot 3}{(\pi z^2)^2} + \frac{1 \cdot 3 \cdot 5 \cdot 7}{(\pi z^2)^4} - \dots \right] \frac{\cos \frac{1}{2} \pi z^2}{\pi z} \\ - \left[(1/\pi z^2) - \frac{1 \cdot 3 \cdot 5}{(\pi z^2)^3} \right. \\ \left. + \frac{1 \cdot 3 \cdot 5 \cdot 7 \cdot 9}{(\pi z^2)^5} - \dots \right] \frac{\sin \frac{1}{2} \pi z^2}{\pi z} \end{aligned}$$

Elliptic Integrals

First Kind:

$$F(\phi, k) = \int_0^\phi \frac{d\theta}{(1 - k^2 \sin^2 \theta)^{1/2}}$$

Second Kind:

$$E(\phi, k) = \int_0^\phi (1 - k^2 \sin^2 \theta)^{1/2} d\theta$$

If $\phi = \pi/2$, the elliptic integrals are said to be complete, and they are denoted by K or $K(k)$ and E or $E(k)$.

Bessel Functions*

Refer to Fig. 23.

Definitions: Bessel functions are solutions to Bessel's differential equation

$$z^2(d^2w/dz^2) + z(dw/dz) + (z^2 - \nu^2)w = 0$$

They are divided into

First kind

$$J_\nu(z) = (z/2)^\nu \sum_{k=0}^{\infty} \frac{(-1)^k (z/2)^{2k}}{k! \Gamma(k + \nu + 1)}$$

Second kind

$$Y_\nu(z) = \frac{J_\nu(z) \cos \nu \pi - J_{-\nu}(z)}{\sin \nu \pi}, \nu \text{ not an integer}$$

* For an extensive treatment of Bessel functions, see G. N. Watson, *A Treatise on the Theory of Bessel Functions*, New York: Cambridge University Press, 1943.

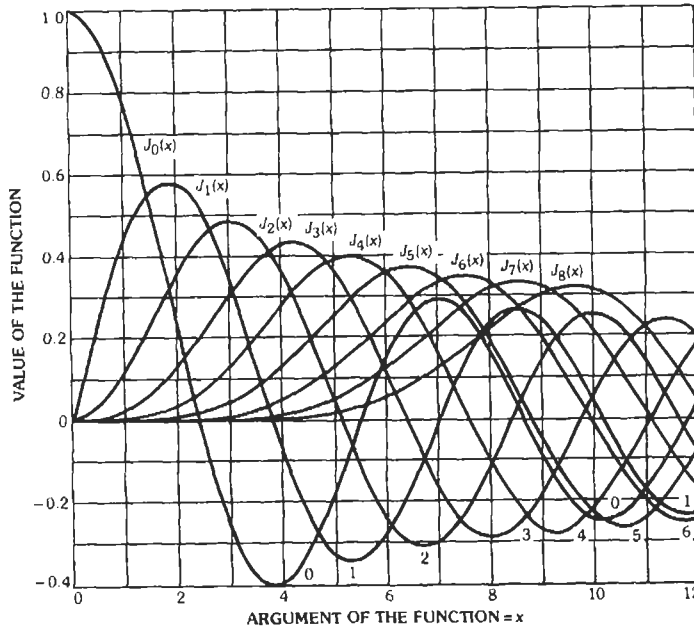


Fig. 23. Bessel functions for the first eight orders.

$$Y_n(z) = \lim_{\nu \rightarrow n} Y_\nu(z), \text{ } n \text{ integral}$$

$$\text{Re } \nu > -\frac{1}{2}$$

Third kind

$$H_\nu^{(1)}(z) = J_\nu(z) + jY_\nu(z)$$

$$H_\nu^{(1)}(z) = -(j/\pi)e^{-j\pi\nu/2} \int_{-\infty}^{\infty} \exp(jz \cosh t - \nu t) dt,$$

$$0 < \arg z < \pi$$

$$H_\nu^{(2)}(z) = J_\nu(z) - jY_\nu(z)$$

$$H_\nu^{(2)}(z) = (j/\pi)e^{j\pi\nu/2} \int_{-\infty}^{\infty} \exp(-jz \cosh t - \nu t) dt,$$

$$-\pi < \arg z < 0$$

The second and third kinds are sometimes called Neumann and Hankel functions, respectively.

The modified Bessel functions are solutions to Bessel's differential equation with z replaced by jz .

$$I_\nu(z) = \exp(-j\pi\nu/2) J_\nu [z \exp(j\pi/2)],$$

$$-\pi < \arg z \leq \pi/2$$

$$= \exp(j3\pi\nu/2) J_\nu [z \exp(-j3\pi/2)],$$

$$\pi/2 < \arg z \leq \pi$$

$$K_\nu(z) = (j\pi/2) \exp(j\pi\nu/2) H_\nu^{(1)} [z \exp(j\pi/2)],$$

$$-\pi < \arg z \leq \pi/2$$

$$= -(j\pi/2) \exp(-j\pi\nu/2)$$

$$\times H_\nu^{(2)} [z \exp(-j\pi/2)], \text{ } -\pi/2 < \arg z \leq \pi$$

Asymptotic Expressions: For $|z| \rightarrow \infty$

$$J_\nu(z) \sim [(2/\pi z)]^{1/2} \cos(z - \nu\pi/2 - \pi/4), \text{ } |\arg z| < \pi$$

$$Y_\nu(z) \sim [(2/\pi z)]^{1/2} \sin(z - \nu\pi/2 - \pi/4), \text{ } |\arg z| < \pi$$

$$H_\nu^{(1)}(z) \sim [(2/\pi z)]^{1/2} \exp[j(z - \nu\pi/2 - \pi/4)],$$

$$-\pi < \arg z < 2\pi$$

$$H_\nu^{(2)}(z) \sim [(2/\pi z)]^{1/2} \exp[-j(z - \nu\pi/2 - \pi/4)],$$

$$-2\pi < \arg z < \pi$$

Recurrence Relations:

$$C_{\nu-1}(z) + C_{\nu+1}(z) = (2\nu/z)C_\nu(z)$$

$$C_{\nu-1}(z) - C_{\nu+1}(z) = 2C_\nu'(z)$$

where C denotes $J, Y, H^{(1)}$ or $H^{(2)}$

Integral Representations:

$$J_\nu(z) = \frac{2(z/2)^\nu}{\Gamma(\nu + \frac{1}{2})\pi^{1/2}} \int_0^{\pi/2} \cos(z \sin \phi) (\cos \phi)^{2\nu} d\phi,$$

Series Containing Bessel Functions:

$$\exp(-ju \sin x) = \sum_{n=-\infty}^{\infty} J_n(u) \exp(-jnx)$$

$$\cos(u \sin x) = J_0(u) + 2 \sum_{n=1}^{\infty} J_{2n}(u) \cos 2nx$$

$$\sin(u \sin x) = 2 \sum_{n=1}^{\infty} J_{2n-1}(u) \sin(2n - 1)x$$

$$\cos(u \cos x) = J_0(u) + 2 \sum_{n=1}^{\infty} (-1)^n J_{2n}(u) \cos 2nx$$

$$\sin(u \cos x) = 2 \sum_{n=1}^{\infty} (-1)^{n+1} J_{2n-1}(u) \times \cos(2n - 1)x$$

$$1 = \sum_{n=-\infty}^{\infty} J_n^2(x)$$

Orthogonal Polynomials

Any set of polynomials $\{f_n(x)\}$ with the property

$$\int_a^b w(x)f_n(x)f_m(x)dx = 0, \text{ for } m \neq n$$

$$= h_n, \text{ for } m = n$$

is called a set of orthogonal polynomials on the interval (a, b) with respect to the weight function $w(x)$. These functions occur in the Gauss quadrature equations among other places. Chebishev polynomials are involved in the theory of the Chebishev filter; Hermite polynomials arise in the refinements of the central limit

theorem, the so-called Edgeworth series, etc. The important properties are summarized in Table 1.

NUMERICAL ANALYSIS

Algorithms for Solving $F(x) = 0$

Bisection Method and Regula Falsi (Rule of False Position): First determine x_1 and x_2 such that $F(x_1)F(x_2) < 0$, i.e., x_1 and x_2 are points at which the function has opposite signs.

Bisection Method: Calculate

$$x_3 = (x_1 + x_2)/2$$

Regula Falsi: Calculate

$$x_3 = [x_1F(x_2) - x_2F(x_1)]/[F(x_2) - F(x_1)]$$

To obtain the next approximation, take x_3 and x_i , $i = 1$ or 2 , such that $F(x_3)F(x_i) < 0$, and repeat the procedure.

Newton-Raphson: Take some initial value x_1 and calculate successively

$$x_{n+1} = x_n - [F(x_n)/F'(x_n)], n = 1, 2, 3 \dots$$

This method may not converge. When it converges, the rate of convergence is generally faster than the bisection method or the regula falsi.

Algorithm for Solving $F(x, y) = G(x, y) = 0$

The following is an extension of the Newton-Raphson method described above. Take some initial values x_1 and y_1 and calculate successively

TABLE I. PROPERTIES OF ORTHOGONAL POLYNOMIALS

$f_n(x)$	Name	a	b	$w(x)$	h_n	Explicit Expression*
$T_n(x)$	Chebichev	-1	+1	$(1 - x^2)^{-1/2}$	$\frac{1}{2}\pi, n \neq 0; \pi, n = 0$	$\frac{1}{2}n! \sum_{m=0}^{[n/2]} (-1)^m \frac{(n-m-1)!}{m!(n-2m)!} (2x)^{n-2m}$
$H_n(x)$	Hermite	$-\infty$	∞	e^{-x^2}	$2^n n! \pi^{1/2}$	$n! \sum_{m=0}^{[n/2]} (-1)^m \frac{(2x)^{n-2m}}{m!(n-2m)!}$
$L_n(x)$	Laguerre	0	∞	e^{-x}	1	$\sum_{m=0}^n (-1)^m \binom{n}{m} \frac{x^m}{m!}$
$P_n(x)$	Legendre	-1	1	1	$2/(2n + 1)$	$2^{-n} \sum_{m=0}^{[n/2]} (-1)^m \binom{n}{m} \binom{2n-2m}{n} x^{n-2m}$

* $[n/2]$ denotes the largest integer less than or equal to $n/2$.

$$x_{n+1} = x_n + \left[\left(\frac{\partial F}{\partial y} G - F \frac{\partial G}{\partial y} \right) / \left(\frac{\partial F}{\partial x} \frac{\partial G}{\partial y} - \frac{\partial F}{\partial y} \frac{\partial G}{\partial x} \right) \right] \Bigg|_{(x,y)=(x_n,y_n)}$$

$$y_{n+1} = y_n + \left[\left(F \frac{\partial G}{\partial x} - \frac{\partial F}{\partial x} G \right) / \left(\frac{\partial F}{\partial x} \frac{\partial G}{\partial y} - \frac{\partial F}{\partial y} \frac{\partial G}{\partial x} \right) \right] \Bigg|_{(x,y)=(x_n,y_n)}$$

for $n = 1, 2, 3, \dots$

By using Taylor's series to first derivatives and Cramer's rule, this algorithm may be further extended to the case of m simultaneous equations in m variables.

Interpolation Polynomial

The polynomial of lowest degree that passes through n points $(x_i, y_i), i = 1, 2, \dots, n$, is given by

$$P(x) = \sum_{i=1}^n \left(y_i \prod_{k=1, k \neq i}^n \frac{x - x_k}{x_i - x_k} \right)$$

where $x_i \neq x_k$ for $i \neq k$.

Interpolation at Equidistant Points

Let

$$f_i = f(x_0 + ih)$$

$$g_{ij} = g(x + ih, y + jk)$$

Then $f(x_0 + ph)$ may be approximated by

$$(1 - p)f_0 + pf_1, \text{ given two points}$$

$$[p(p - 1)/2]f_{-1} + (1 - p^2)f_0 + [p(p + 1)/2]f_1, \text{ given three points}$$

$$[-p(p - 1)(p - 2)/6]f_{-1} + [(p^2 - 1)(p - 2)/2]f_0 - [p(p + 1)(p + 2)/2]f_1 + [p(p^2 - 1)/6]f_2, \text{ given four points}$$

$g(x_0 + ph, y_0 + qk)$ may be approximated by

$$(1 - p - q)g_{00} + pg_{10} + qg_{01}, \text{ given three points}$$

$$(1 - p)(1 - q)g_{00} + p(1 - q)g_{10} + q(1 - p)g_{01} + pqg_{11}, \text{ given four points}$$

given four points

Integration

Refer to Fig. 24. Let

$$f_i = f(x_0 + ih)$$

$$\int_{x_0}^{x_0 + mh} f(x)dx \approx h \left(\frac{1}{2}f_0 + f_1 + f_2 + \dots + f_{n-1} + \frac{1}{2}f_n \right) \text{ (Trapezoidal Rule)}$$

$$\int_{x_0}^{x_0 + 2nh} f(x)dx \approx \frac{1}{3}h \left(f_0 + 4f_1 + 2f_2 + 4f_3 + 2f_4 + \dots + 2f_{2n-2} + 4f_{2n-1} + f_{2n} \right) \text{ (Simpson's Rule)}$$

Differentiation

If a function $f(x)$ is known at two points a and b , then

$$f'(a) \approx [f(b) - f(a)]/(b - a)$$

If $f(x)$ is given at a discrete set of points, one may differentiate the interpolation polynomial stated above. The equation simplifies when the derivative is calculated at interpolation points that are equidistant, for example

$$f_{-1}' \approx (1/2h)(-3f_{-1} + 4f_0 - f_1)$$

$$f_0' \approx (1/2h)(-f_{-1} + f_1)$$

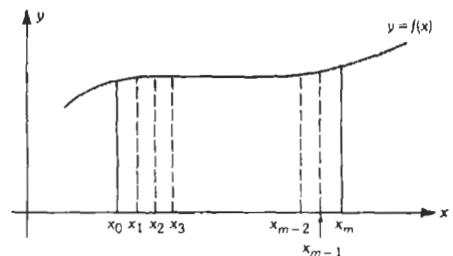
$$f_1' \approx (1/2h)(f_{-1} - 4f_0 + 3f_1)$$

where,

$$f_i = f(x_0 + ih), i = -1, 0, 1$$

Error in Arithmetic Operations

Let the quantities A and B be known to within errors of a and b , respectively, where a and b are small relative to A and B . Then



For Simpson's rule, m must be even.

Fig. 24. Integration.

Operation	Maximum Error (to first order in a and b)
$A \pm B$	$a + b$
$A \cdot B$	$ AB (a/A + b/B)$
A/B	$ A/B (a/A + b/B)$

The magnitude, q , of small random errors often has a normal distribution

$$(h/\pi^{1/2}) \exp(-h^2 q^2)$$

where h is called the index of precision. The rms of q , denoted by σ , is, in terms of h

$$\sigma = 1/h\sqrt{2}$$

Let $f(x_1, x_2, \dots, x_n)$ be a function of n variables x_i , which have normally distributed independent errors q_i with indices of precision h_i . Then $f(x_1, x_2, \dots, x_n)$ will have an error Q , which is also normally distributed with index of precision H . For small errors

$$H \approx \{[h_1^{-1}(\partial f/\partial q_1)]^2 + [h_2^{-1}(\partial f/\partial q_2)]^2 + \dots + [h_n^{-1}(\partial f/\partial q_n)]^2\}^{-1/2}$$

48 Mathematical Tables

Hyperbolic Sines	48-2
Hyperbolic Cosines	48-3
Hyperbolic Tangents	48-4
Multiples of 0.4343	48-4
Multiples of 2.3026	48-5
Logarithms to Base 2 and Powers of 2	48-5
Random Digits	48-6
Exponentials	48-7
Normal or Gaussian Distribution	48-8
Bessel Functions	48-9

HYPERBOLIC SINES

TABLE 1. HYPERBOLIC SINES [$\sinh x = \frac{1}{2}(e^x - e^{-x})$]

x	0	1	2	3	4	5	6	7	8	9	avg. diff.
0.0	0.0000	0.0100	0.0200	0.300	0.0400	0.0500	0.0600	0.0701	0.0801	0.0901	100
.1	0.1002	0.1102	0.1203	0.1304	0.1405	0.1506	0.1607	0.1708	0.1810	0.1911	101
.2	0.2013	0.2115	0.2218	0.2320	0.2423	0.2526	0.2629	0.2733	0.2837	0.2941	103
.3	0.3045	0.3150	0.3255	0.3360	0.3466	0.3572	0.3678	0.3785	0.3892	0.4000	106
.4	0.4108	0.4216	0.4325	0.4434	0.4543	0.4653	0.4764	0.4875	0.4986	0.5098	110
0.5	0.5211	0.5324	0.5438	0.5552	0.5666	0.5782	0.5897	0.6014	0.6131	0.6248	116
.6	0.6367	0.6485	0.6605	0.6725	0.6846	0.6967	0.7090	0.7213	0.7336	0.7461	122
.7	0.7586	0.7712	0.7838	0.7966	0.8094	0.8223	0.8353	0.8484	0.8615	0.8748	130
.8	0.8881	0.9015	0.9150	0.9286	0.9423	0.9561	0.9700	0.9840	0.9981	1.012	138
.9	1.027	1.041	1.055	1.070	1.085	1.099	1.114	1.129	1.145	1.160	15
1.0	1.175	1.191	1.206	1.222	1.238	1.254	1.270	1.286	1.303	1.319	16
.1	1.336	1.352	1.369	1.386	1.403	1.421	1.438	1.456	1.474	1.491	17
.2	1.509	1.528	1.546	1.564	1.583	1.602	1.621	1.640	1.659	1.679	19
.3	1.698	1.718	1.738	1.758	1.779	1.799	1.820	1.841	1.862	1.883	21
.4	1.904	1.926	1.948	1.970	1.992	2.014	2.037	2.060	2.083	2.106	22
1.5	2.129	2.153	2.177	2.201	2.225	2.250	2.274	2.299	2.324	2.350	25
.6	2.376	2.401	2.428	2.454	2.481	2.507	2.535	2.562	2.590	2.617	27
.7	2.646	2.674	2.703	2.732	2.761	2.790	2.820	2.850	2.881	2.911	30
.8	2.942	2.973	3.005	3.037	3.069	3.101	3.134	3.167	3.200	3.234	33
.9	3.268	3.303	3.337	3.372	3.408	3.443	3.479	3.516	3.552	3.589	36
2.0	3.627	3.665	3.703	3.741	3.780	3.820	3.859	3.899	3.940	3.981	39
.1	4.022	4.064	4.106	4.148	4.191	4.234	4.278	4.322	4.367	4.412	44
.2	4.457	4.503	4.549	4.596	4.643	4.691	4.739	4.788	4.837	4.887	48
.3	4.937	4.988	5.039	5.090	5.142	5.195	5.248	5.302	5.356	5.411	53
.4	5.466	5.522	5.578	5.635	5.693	5.751	5.810	5.869	5.929	5.989	58
2.5	6.050	6.112	6.174	6.237	6.300	6.365	6.429	6.495	6.561	6.627	64
.6	6.695	6.763	6.831	6.901	6.971	7.042	7.113	7.185	7.258	7.332	71
.7	7.406	7.481	7.557	7.634	7.711	7.789	7.868	7.948	8.028	8.110	79
.8	8.192	8.275	8.359	8.443	8.529	8.615	8.702	8.790	8.879	8.969	87
.9	9.060	9.151	9.244	9.337	9.431	9.527	9.623	9.720	9.819	9.918	96
3.0	10.02	10.12	10.22	10.32	10.43	10.53	10.64	10.75	10.86	10.97	11
.1	11.08	11.19	11.30	11.42	11.53	11.65	11.76	11.88	12.00	12.12	12
.2	12.25	12.37	12.49	12.62	12.75	12.88	13.01	13.14	13.27	13.40	13
.3	13.54	13.67	13.81	13.95	14.09	14.23	14.38	14.52	14.67	14.82	14
.4	14.97	15.12	15.27	15.42	15.58	15.73	15.89	16.05	16.21	16.38	16
3.5	16.54	16.71	16.88	17.05	17.22	17.39	17.57	17.74	17.92	18.10	17
.6	18.29	18.47	18.66	18.84	19.03	19.22	19.42	19.61	19.81	20.01	19
.7	20.21	20.41	20.62	20.83	21.04	21.25	21.46	21.68	21.90	22.12	21
.8	22.34	22.56	22.79	23.02	23.25	23.49	23.72	23.96	24.20	24.45	24
.9	24.69	24.94	25.19	25.44	25.70	25.96	26.22	26.48	26.75	27.02	26
4.0	27.29	27.56	27.84	28.12	28.40	28.69	28.98	29.27	29.56	29.86	29
.1	30.16	30.47	30.77	31.08	31.39	31.71	32.03	32.35	32.68	33.00	32
.2	33.34	33.67	34.01	34.35	34.70	35.05	35.40	35.75	36.11	36.48	35
.3	36.84	37.21	37.59	37.97	38.35	38.73	39.12	39.52	39.91	40.31	39
.4	40.72	41.13	41.54	41.96	42.38	42.81	43.24	43.67	44.11	44.56	43
4.5	45.00	45.46	45.91	46.37	46.84	47.31	47.79	48.27	48.75	49.24	47
.6	49.74	50.24	50.74	51.25	51.77	52.29	52.81	53.34	53.88	54.42	52
.7	54.97	55.52	56.08	56.64	57.21	57.79	58.37	58.96	59.55	60.15	58
.8	60.75	61.36	61.98	62.60	63.23	63.87	64.51	65.16	65.81	66.47	64
.9	67.14	67.82	68.50	69.19	69.88	70.58	71.29	72.01	72.73	73.46	71
5.0	74.20										

If $x > 5$, $\sinh x = \frac{1}{2}(e^x)$ and $\log_{10} \sinh x = (0.4343)x + 0.6990 - 1$, correct to four significant figures.

HYPERBOLIC COSINES

TABLE 2. HYPERBOLIC COSINES [$\cosh x = \frac{1}{2}(e^x + e^{-x})$]

x	0	1	2	3	4	5	6	7	8	9	avg diff
0.0	1.000	1.000	1.000	1.000	1.001	1.001	1.002	1.002	1.003	1.004	1
.1	1.005	1.006	1.007	1.008	1.010	1.011	1.013	1.014	1.016	1.018	2
.2	1.020	1.022	1.024	1.027	1.029	1.031	1.034	1.037	1.039	1.042	3
.3	1.045	1.048	1.052	1.055	1.058	1.062	1.066	1.069	1.073	1.077	4
.4	1.081	1.085	1.090	1.094	1.098	1.103	1.108	1.112	1.117	1.122	5
0.5	1.128	1.133	1.138	1.144	1.149	1.155	1.161	1.167	1.173	1.179	6
.6	1.185	1.192	1.198	1.205	1.212	1.219	1.226	1.233	1.240	1.248	7
.7	1.255	1.263	1.271	1.278	1.287	1.295	1.303	1.311	1.320	1.329	8
.8	1.337	1.346	1.355	1.365	1.374	1.384	1.393	1.403	1.413	1.423	10
.9	1.433	1.443	1.454	1.465	1.475	1.486	1.497	1.509	1.520	1.531	11
1.0	1.543	1.555	1.567	1.579	1.591	1.604	1.616	1.629	1.642	1.655	13
.1	1.669	1.682	1.696	1.709	1.723	1.737	1.752	1.766	1.781	1.796	14
.2	1.811	1.826	1.841	1.857	1.872	1.888	1.905	1.921	1.937	1.954	16
.3	1.971	1.988	2.005	2.023	2.040	2.058	2.076	2.095	2.113	2.132	18
.4	1.151	2.170	2.189	2.209	2.229	2.249	2.269	2.290	2.310	2.331	20
1.5	2.352	2.374	2.395	2.417	2.439	2.462	2.484	2.507	2.530	2.554	23
.6	2.577	2.601	2.625	2.650	2.675	2.700	2.725	2.750	2.776	2.802	25
.7	2.828	2.855	2.882	2.909	2.936	2.964	2.992	3.021	3.049	3.078	28
.8	3.107	3.137	3.167	3.197	3.228	3.259	3.290	3.321	3.353	3.385	31
.9	3.418	3.451	3.484	3.517	3.551	3.585	3.620	3.655	3.690	3.726	34
2.0	3.762	3.799	3.835	3.873	3.910	3.948	3.987	4.026	4.065	4.104	38
.1	4.144	4.185	4.226	4.267	4.309	4.351	4.393	4.436	4.480	4.524	42
.2	4.568	4.613	4.658	4.704	4.750	4.797	4.844	4.891	4.939	4.988	47
.3	5.037	5.087	5.137	5.188	5.239	5.290	5.343	5.395	5.449	5.503	52
.4	5.557	5.612	5.667	5.723	5.780	5.837	5.895	5.954	6.013	6.072	58
2.5	6.132	6.193	6.255	6.317	6.379	6.443	6.507	6.571	6.636	6.702	64
.6	6.769	6.836	6.904	6.973	7.042	7.112	7.183	7.255	7.327	7.400	70
.7	7.473	7.548	7.623	7.699	7.776	7.853	7.932	8.011	8.091	8.171	78
.8	8.253	8.335	8.418	8.502	8.587	8.673	8.759	8.847	8.935	9.024	86
.9	9.115	9.206	9.298	9.391	9.484	9.579	9.675	9.772	9.869	9.968	95
3.0	10.07	10.17	10.27	10.37	10.48	10.58	10.69	10.79	10.90	11.01	11
.1	11.12	11.23	11.35	11.46	11.57	11.69	11.81	11.92	12.04	12.16	12
.2	12.29	12.41	12.53	12.66	12.79	12.91	13.04	13.17	13.31	13.44	13
.3	13.57	13.71	13.85	13.99	14.13	14.27	14.41	14.56	14.70	14.85	14
.4	15.00	15.15	15.30	15.45	15.61	15.77	15.92	16.08	16.25	16.41	16
3.5	16.57	16.74	16.91	17.08	17.25	17.42	17.60	17.77	17.95	18.13	17
.6	18.31	18.50	18.68	18.87	19.06	19.25	19.44	19.64	19.84	20.03	19
.7	20.24	20.44	20.64	20.85	21.06	21.27	21.49	21.70	21.92	22.14	21
.8	22.36	22.59	22.81	23.04	23.27	23.51	23.74	23.98	24.22	24.47	23
.9	24.71	24.96	25.21	25.46	25.72	25.98	26.24	26.50	26.77	27.04	26
4.0	27.31	27.58	27.86	28.14	28.42	28.71	29.00	29.29	29.58	29.88	29
.1	30.18	30.48	30.79	31.10	31.41	31.72	32.04	32.37	32.69	33.02	32
.2	33.35	33.69	34.02	34.37	34.71	35.06	35.41	35.77	36.13	36.49	35
.3	36.86	37.23	37.60	37.98	38.36	38.75	39.13	39.53	39.93	40.33	39
.4	40.73	41.14	41.55	41.97	42.39	42.82	43.25	43.68	44.12	44.57	43
4.5	45.01	45.47	45.92	46.38	46.85	47.32	47.80	48.28	48.76	49.25	47
.6	49.75	50.25	50.75	51.26	51.78	52.30	52.82	53.35	53.89	54.43	52
.7	54.98	55.53	56.09	56.65	57.22	57.80	58.38	58.96	59.56	60.15	58
.8	60.76	61.37	61.99	62.61	63.24	63.87	64.52	65.16	65.82	66.48	64
.9	67.15	67.82	68.50	69.19	69.89	70.59	71.30	72.02	72.74	73.47	71
5.0	74.21										

If $x > 5$, $\cosh x = \frac{1}{2}(e^x)$, and $\log_{10} \cosh x = (0.4343)x + 0.6990 - 1$, correct to four significant figures.

HYPERBOLIC TANGENTS

TABLE 3. HYPERBOLIC TANGENTS [$\tanh x = (e^x - e^{-x})/(e^x + e^{-x}) = \sinh x/\cosh x$]

x	0	1	2	3	4	5	6	7	8	9	avg diff
0.0	.0000	.0100	.0200	.0300	.0400	.0500	.0599	.0699	.0798	.0898	100
.1	.0997	.1096	.1194	.1293	.1391	.1489	.1587	.1684	.1781	.1878	98
.2	.1974	.2070	.2165	.2260	.2355	.2449	.2543	.2636	.2729	.2821	94
.3	.2913	.3004	.3095	.3185	.3275	.3364	.3452	.3540	.3627	.3714	89
.4	.3800	.3885	.3969	.4053	.4136	.4219	.4301	.4382	.4462	.4542	82
0.5	.4621	.4700	.4777	.4854	.4930	.5005	.5080	.5154	.5227	.5299	75
.6	.5370	.5441	.5511	.5581	.5649	.5717	.5784	.5850	.5915	.5980	67
.7	.6044	.6107	.6169	.6231	.6291	.6352	.6411	.6469	.6527	.6584	60
.8	.6640	.6696	.6751	.6805	.6858	.6911	.6963	.7014	.7064	.7114	52
.9	.7163	.7211	.7259	.7306	.7352	.7398	.7443	.7487	.7531	.7574	45
1.0	.7616	.7658	.7699	.7739	.7779	.7818	.7857	.7895	.7932	.7969	39
.1	.8005	.8041	.8076	.8110	.8144	.8178	.8210	.8243	.8275	.8306	33
.2	.8337	.8367	.8397	.8426	.8455	.8483	.8511	.8538	.8565	.8591	28
.3	.8617	.8643	.8668	.8693	.8717	.8741	.8764	.8787	.8810	.8832	24
.4	.8854	.8875	.8896	.8917	.8937	.8957	.8977	.8996	.9015	.9033	20
1.5	.9052	.9069	.9087	.9104	.9121	.9138	.9154	.9170	.9186	.9202	17
.6	.9217	.9232	.9246	.9261	.9275	.9289	.9302	.9316	.9329	.9342	14
.7	.9354	.9367	.9379	.9391	.9402	.9414	.9425	.9436	.9447	.9458	11
.8	.9468	.9478	.9488	.9498	.9508	.9518	.9527	.9536	.9545	.9554	9
.9	.9562	.9571	.9579	.9587	.9595	.9603	.9611	.9619	.9626	.9633	8
2.0	.9640	.9647	.9654	.9661	.9668	.9674	.9680	.9687	.9693	.9699	6
.1	.9705	.9710	.9716	.9722	.9727	.9732	.9738	.9743	.9748	.9753	5
.2	.9757	.9762	.9767	.9771	.9776	.9780	.9785	.9789	.9793	.9797	4
.3	.9801	.9805	.9809	.9812	.9816	.9820	.9823	.9827	.9830	.9834	4
.4	.9837	.9840	.9843	.9846	.9849	.9852	.9855	.9858	.9861	.9863	3
2.5	.9866	.9869	.9871	.9874	.9876	.9879	.9881	.9884	.9886	.9888	2
.6	.9890	.9892	.9895	.9897	.9899	.9901	.9903	.9905	.9906	.9908	2
.7	.9910	.9912	.9914	.9915	.9917	.9919	.9920	.9922	.9923	.9925	2
.8	.9926	.9928	.9929	.9931	.9932	.9933	.9935	.9936	.9937	.9938	1
.9	.9940	.9941	.9942	.9943	.9944	.9945	.9946	.9947	.9949	.9950	1
3.0	.9951	.9959	.9967	.9973	.9978	.9982	.9985	.9988	.9990	.9992	4
4.0	.9993	.9995	.9996	.9996	.9997	.9998	.9998	.9998	.9999	.9999	1
5.0	.9999										

If $x > 5$, $\tanh x = 1.0000$ to four decimal places.

MULTIPLES OF 0.4343

TABLE 4. MULTIPLES OF 0.4343 ($0.43429448 = \log_{10} e$)

x	0	1	2	3	4	5	6	7	8	9
0.0	0.0000	0.0434	0.0869	0.1303	0.1737	0.2171	0.2606	0.3040	0.3474	0.3909
1.0	0.4343	0.4777	0.5212	0.5646	0.6080	0.6514	0.6949	0.7383	0.7817	0.8252
2.0	0.8686	0.9120	0.9554	0.9989	1.0423	1.0857	1.1292	1.1726	1.2160	1.2595
3.0	1.3029	1.3463	1.3897	1.4332	1.4766	1.5200	1.5635	1.6069	1.6503	1.6937
4.0	1.7372	1.7806	1.8240	1.8675	1.9109	1.9543	1.9978	2.0412	2.0846	2.1280
5.0	2.1715	2.2149	2.2583	2.3018	2.3452	2.3886	2.4320	2.4755	2.5189	2.5623
6.0	2.6058	2.6492	2.6926	2.7361	2.7795	2.8229	2.8663	2.9098	2.9532	2.9966
7.0	3.0401	3.0835	3.1269	3.1703	3.2138	3.2572	3.3006	3.3441	3.3875	3.4309
8.0	3.4744	3.5178	3.5612	3.6046	3.6481	3.6915	3.7349	3.7784	3.8218	3.8652
9.0	3.9087	3.9521	3.9955	4.0389	4.0824	4.1258	4.1692	4.2127	4.2561	4.2995

MULTIPLES OF 2.3026

TABLE 5. MULTIPLES OF 2.3026 (2.3025851 = 1/0.4343 = log_e 10)

x	0	1	2	3	4	5	6	7	8	9
0.0	0.0000	0.2303	0.4605	0.6908	0.9210	1.1513	1.3816	1.6118	1.8421	2.0723
1.0	2.3026	2.5328	2.7631	2.9934	3.2236	3.4539	3.6841	3.9144	4.1447	4.3749
2.0	4.6052	4.8354	5.0657	5.2959	5.5262	5.7565	5.9867	6.2170	6.4472	6.6775
3.0	6.9078	7.1380	7.3683	7.5985	7.8288	8.0590	8.2893	8.5196	8.7498	8.9801
4.0	9.2103	9.4406	9.6709	9.9011	10.131	10.362	10.592	10.822	11.052	11.283
5.0	11.513	11.743	11.973	12.204	12.434	12.664	12.894	13.125	13.355	13.585
6.0	13.816	14.046	14.276	14.506	14.737	14.967	15.197	15.427	15.658	15.888
7.0	16.118	16.348	16.579	16.809	17.039	17.269	17.500	17.730	17.960	18.190
8.0	18.421	18.651	18.881	19.111	19.342	19.572	19.802	20.032	20.263	20.493
9.0	20.723	20.954	21.184	21.414	21.644	21.875	22.105	22.335	22.565	22.796

LOGARITHMS TO BASE 2 AND POWERS OF 2

TABLE 6.

x	log ₂ x
0.1	-3.32193
0.2	-2.32193
0.3	-1.73697
0.4	-1.32193
0.5	-1.00000
0.6	-0.73697
0.7	-0.51457
0.8	-0.32193
0.9	-0.15200
1.0	0.00000
1.1	0.13750
1.2	0.26303
1.3	0.37851
1.4	0.48543
1.5	0.58496
1.6	0.67807
1.7	0.76553
1.8	0.84800
1.9	0.92600
2.0	1.00000
10	3.32193
100	6.64386
1000	9.96578
2 ^y	y

log₂ x = log₂ 10 log₁₀ x = log₂ e log_e x

2^y = e^{y log_e 2} = 10^{y log₁₀ 2}

log₂ 10 = 3.32193 = 1/log₁₀ 2

log₁₀ 2 = 0.30103 = 1/log₂ 10

log₂ e = 1.44269 = 1/log_e 2

log_e 2 = 0.69315 = 1/log₂ e

TABLE 7.

y	2 ^y
0.1	1.072
0.2	1.149
0.3	1.231
0.4	1.320
0.5	1.414
0.6	1.515
0.7	1.625
0.8	1.741
0.9	1.866
1	2
2	4
3	8
4	16
5	32
6	64
7	128
8	256
9	512
10	1 024
11	2 048
12	4 096
13	8 192
14	16 384
15	32 768
16	65 536
17	131 072
18	262 144
19	524 288
20	1 048 576
21	2 097 152
22	4 194 304
23	8 388 608
24	16 777 216
25	33 554 432
26	67 108 864
27	134 217 728
28	268 435 456
29	536 870 912
30	1 073 741 824
31	2 147 483 648
32	4 294 967 296
log ₂ x	x

TABLE 8. RANDOM DIGITS

49 31 97 45 80	57 47 01 46 00	57 16 83 04 58	23 89 20 78 25	18 53 20 38 74	66 22 07 90 50	29 22 37 05 41	67 11 58 45 84
88 78 67 69 63	12 12 72 50 14	71 88 66 53 34	38 01 30 93 79	22 93 62 20 58	49 17 11 10 27	22 68 18 01 10	31 59 50 92 46
84 86 69 52 02	43 98 37 26 55	40 41 85 95 04	52 38 30 72 32	66 39 77 65 10	81 15 00 07 04	74 58 09 03 54	43 74 42 21 78
11 84 92 64 82	20 46 19 94 50	28 83 37 66 61	47 27 79 29 35	89 73 02 32 72	65 42 03 50 91	69 09 37 13 64	08 10 79 69 52
54 96 61 75 94	57 39 37 32 67	37 88 36 21 24	62 19 94 95 42	81 82 17 53 23	96 06 89 17 24	40 45 69 12 34	58 09 06 53 42
10 95 93 33 49	80 71 99 67 51	44 88 23 35 92	66 23 41 38 21	94 37 78 25 54	53 58 61 14 32	72 92 76 73 49	83 96 25 89 12
22 78 40 77 83	35 90 30 00 91	19 08 21 38 73	07 18 42 15 66	68 48 54 99 91	53 16 51 98 65	61 86 93 30 93	81 12 90 64 81
86 03 76 17 91	33 81 56 39 68	45 31 62 92 83	89 31 85 58 06	07 33 00 71 84	86 78 86 45 77	40 04 81 65 20	07 63 81 07 97
80 03 76 50 89	85 91 97 43 91	22 78 85 54 33	31 18 87 48 82	10 99 31 49 30	35 07 23 64 29	68 77 39 76 69	28 65 68 99 38
72 75 18 43 59	15 76 91 36 15	08 29 38 61 93	05 02 62 12 55	20 80 11 51 78	64 45 38 33 57	09 77 43 07 51	49 74 01 13 85
79 24 13 53 47	66 85 17 92 47	46 13 93 66 89	82 58 71 35 86	93 36 91 30 44	69 68 67 81 62	66 37 80 29 19	34 01 25 00 80
43 59 33 95 55	97 34 55 84 94	26 56 69 53 23	32 99 38 99 88	19 36 05 50 49	94 95 17 63 41	84 01 93 06 90	25 65 67 29 96
29 52 26 27 13	33 70 11 71 86	06 76 55 71 41	48 61 71 82 82	47 79 88 98 90	06 89 36 54 83	17 70 12 12 92	14 88 01 53 86
88 83 64 72 90	67 27 47 83 62	35 38 49 03 80	12 31 78 97 02	69 22 33 20 07	03 51 36 11 49	32 54 69 20 72	62 52 22 15 04
65 90 56 62 53	91 48 23 06 89	49 33 37 84 82	36 19 91 13 55	34 51 15 07 21	84 85 03 41 59	97 13 86 19 19	97 78 92 85 75
44 79 86 93 71	07 86 59 17 56	45 59 51 40 44	56 80 69 91 26	54 03 15 93 29	58 96 35 22 20	35 29 22 79 24	55 46 74 30 36
35 51 09 91 39	32 03 12 79 25	79 81 91 50 54	76 17 41 22 06	66 72 28 55 15	04 72 39 24 11	02 73 70 81 68	30 04 36 34 50
50 12 59 32 23	64 20 94 97 14	11 97 16 22 34	74 85 74 64 01	71 05 90 74 96	38 40 41 81 26	28 26 13 78 44	12 54 31 43 98
25 17 39 00 38	63 87 14 04 18	11 45 28 93 18	53 08 42 19 93	45 47 88 60 66	31 13 53 32 43	80 57 33 06 06	48 64 45 30 08
68 45 99 00 94	44 99 59 37 18	38 74 68 12 71	96 26 09 81 37	97 24 69 11 21	89 43 72 03 93	77 15 38 85 52	26 84 31 28 44
22 98 22 59 36	96 41 73 48 45	85 14 95 75 04	15 05 93 68 49	84 98 36 83 12	25 51 95 61 58	86 30 00 76 89	14 00 67 77 53
48 24 36 29 93	47 13 28 52 48	35 22 97 28 37	36 75 27 16 55	35 55 40 29 35	72 88 96 87 72	19 85 03 96 50	65 22 21 55 63
93 51 41 49 15	67 96 08 22 03	40 11 72 43 46	32 18 98 70 74	04 36 81 76 32	50 96 27 19 08	94 46 46 64 32	62 24 31 36 74
69 70 79 83 03	93 06 91 62 16	60 87 59 75 45	68 65 29 21 60	81 31 16 04 79	69 98 53 09 52	23 92 14 97 30	21 71 89 23 14
87 46 79 17 94	70 81 41 27 43	03 76 93 25 51	74 80 14 16 92	03 82 38 98 87	55 82 87 44 52	72 77 52 37 16	42 85 37 47 93
81 00 68 14 98	59 37 53 05 02	94 07 79 22 09	31 50 66 96 06	80 42 26 54 37	38 79 75 62 61	27 81 64 67 04	82 73 50 33 39
15 45 88 14 81	50 18 74 33 75	94 37 60 06 66	94 14 52 23 99	61 30 74 94 68	43 34 44 37 00	20 20 77 70 88	17 16 72 45 31
33 46 91 25 10	23 09 54 80 16	42 35 41 13 47	90 92 00 38 64	83 87 38 25 57	10 00 28 00 93	59 28 30 44 94	60 72 52 14 31
67 19 80 71 76	65 99 61 83 17	81 14 94 32 91	10 81 74 43 48	38 11 01 68 55	28 92 29 37 58	88 73 13 63 16	51 38 35 76 19
58 03 79 22 61	85 50 45 56 90	10 63 17 82 38	00 15 74 62 59	43 89 29 11 89	87 22 65 69 35	84 76 26 79 36	75 00 00 17 95
93 68 30 96 64	53 92 74 98 85	20 75 49 23 55	57 95 51 09 40	14 95 42 22 99	40 15 65 26 85	29 22 33 83 83	30 31 57 09 99
32 74 80 21 21	11 97 29 69 14	28 06 56 95 64	06 83 55 68 45	01 71 19 84 39	09 44 63 39 37	49 09 54 02 38	81 69 71 24 74
49 21 19 29 63	38 62 56 53 12	62 17 57 33 53	84 97 21 77 26	62 32 85 53 28	45 73 89 39 40	27 46 62 69 27	53 34 51 13 79
63 36 56 42 24	69 47 55 75 12	11 04 45 04 83	68 82 19 74 26	73 00 46 21 09	81 90 77 10 77	57 46 37 00 45	65 12 34 90 70
63 57 62 63 73	44 61 04 37 48	00 33 16 34 22	99 62 27 67 57	34 21 88 94 45	05 60 95 23 36	50 55 89 22 42	52 73 28 15 02
41 07 84 70 36	65 52 46 84 66	67 15 72 64 19	37 97 81 65 11	99 15 90 19 68	45 88 68 68 75	28 41 39 59 18	44 15 64 69 59
70 84 68 95 58	64 17 31 53 81	87 71 35 08 41	46 27 02 65 08	92 85 82 99 49	15 81 79 33 72	56 65 74 31 93	58 13 05 42 73
68 80 06 44 92	20 16 23 27 07	10 28 18 25 25	74 15 58 67 49	27 39 69 74 77	65 55 47 16 01	13 12 16 88 67	95 76 35 96 67
44 97 78 95 25	51 26 96 37 47	91 36 77 40 33	67 02 06 90 92	37 10 34 53 09	30 12 94 33 80	96 99 68 93 56	22 78 46 01 84
79 35 46 38 47	24 39 55 36 79	40 56 03 69 14	69 17 63 19 18	57 34 79 70 12	48 42 82 06 06	60 74 22 22 26	89 99 32 45 97

EXPONENTIALS

TABLE 9. EXPONENTIALS (e^n AND e^{-n})

n	e^n diff	n	e^n diff	n	e^n (^)	n	e^{-n} diff	n	e^{-n}	n	$e^{-n}(^*)$			
0.00	1.000	10	0.50	1.649	16	1.0	2.718	0.00	1.000	-10	0.50	.607	1.0	.368
.01	1.010	10	.51	1.665	17	.1	3.004	.01	0.990	-10	.51	.600	.1	.333
.02	1.020	10	.52	1.682	17	.2	3.320	.02	.980	-10	.52	.595	.2	.301
.03	1.030	10	.53	1.699	17	.3	3.669	.03	.970	-10	.53	.589	.3	.273
.04	1.041	10	.54	1.716	17	.4	4.055	.04	.961	-10	.54	.583	.4	.247
0.05	1.051	11	0.55	1.733	18	1.5	4.482	0.05	.951	-9	0.55	.577	1.5	.223
.06	1.062	11	.56	1.751	17	.6	4.953	.06	.942	-10	.56	.571	.6	.202
.07	1.073	10	.57	1.768	18	.7	5.474	.07	.932	-9	.57	.566	.7	.183
.08	1.083	11	.58	1.786	18	.8	6.050	.08	.923	-9	.58	.560	.8	.165
.09	1.094	11	.59	1.804	18	.9	6.686	.09	.914	-9	.59	.554	.9	.150
0.10	1.105	11	0.60	1.822	18	2.0	7.389	0.10	.905	-9	0.60	.549	2.0	.135
.11	1.116	11	.61	1.840	19	.1	8.166	.11	.896	-9	.61	.543	.1	.122
.12	1.127	12	.62	1.859	19	.2	9.025	.12	.887	-9	.62	.538	.2	.111
.13	1.139	11	.63	1.878	18	.3	9.974	.13	.878	-9	.63	.533	.3	.100
.14	1.150	12	.64	1.896	20	.4	11.02	.14	.869	-8	.64	.527	.4	.0907
0.15	1.162	12	0.65	1.916	19	2.5	12.18	0.15	.861	-9	0.65	.522	2.5	.0821
.16	1.174	11	.66	1.935	19	.6	13.46	.16	.852	-8	.66	.517	.6	.0743
.17	1.185	12	.67	1.954	20	.7	14.88	.17	.844	-9	.67	.512	.7	.0672
.18	1.197	12	.68	1.974	20	.8	16.44	.18	.835	-8	.68	.507	.8	.0608
.19	1.209	12	.69	1.994	20	.9	18.17	.19	.827	-8	.69	.502	.9	.0550
0.20	1.221	13	0.70	2.014	20	3.0	20.09	0.20	.819	-8	0.70	.497	3.0	.0498
.21	1.234	12	.71	2.034	20	.1	22.20	.21	.811	-8	.71	.492	.1	.0450
.22	1.246	13	.72	2.054	21	.2	24.53	.22	.803	-8	.72	.487	.2	.0408
.23	1.259	12	.73	2.075	21	.3	27.11	.23	.795	-8	.73	.482	.3	.0369
.24	1.271	13	.74	2.096	21	.4	29.96	.24	.787	-8	.74	.477	.4	.0334
0.25	1.284	13	0.75	2.117	21	3.5	33.12	0.25	.779	-8	0.75	.472	3.5	.0302
.26	1.297	13	.76	2.138	22	.6	36.60	.26	.771	-8	.76	.468	.6	.0273
.27	1.310	13	.77	2.160	21	.7	40.45	.27	.763	-7	.77	.463	.7	.0247
.28	1.323	13	.78	2.181	22	.8	44.70	.28	.756	-8	.78	.458	.8	.0224
.29	1.336	14	.79	2.203	23	.9	49.40	.29	.748	-7	.79	.454	.9	.0202
0.30	1.350	13	0.80	2.226	22	4.0	54.60	0.30	.741	-8	0.80	.449	4.0	.0183
.31	1.363	14	.81	2.248	22	.1	60.34	.31	.733	-7	.81	.445	.1	.0166
.32	1.377	14	.82	2.270	23	.2	66.69	.32	.726	-7	.82	.440	.2	.0150
.33	1.391	14	.83	2.293	23	.3	73.70	.33	.719	-7	.83	.436	.3	.0136
.34	1.405	14	.84	2.316	24	.4	81.45	.34	.712	-7	.84	.432	.4	.0123
0.35	1.419	14	0.85	2.340	23	4.5	90.02	0.35	.705	-7	0.85	.427	4.5	.0111
.36	1.433	15	.86	2.363	24	.5	98.4	.36	.698	-7	.86	.423	.5	.00674
.37	1.448	14	.87	2.387	24	6.0	109.7	.37	.691	-7	.87	.419	6.0	.00248
.38	1.462	15	.88	2.411	24	7.0	123.92	.38	.684	-7	.88	.415	7.0	.000912
.39	1.477	15	.89	2.435	25	8.0	2981.	.39	.677	-7	.89	.411	8.0	.000335
0.40	1.492	15	0.90	2.460	24	9.0	8103.	0.40	.670	-6	0.90	.407	9.0	.000123
.41	1.507	15	.91	2.484	25	10.0	22026.	.41	.664	-7	.91	.403	10.0	.000045
.42	1.522	15	.92	2.509	26	$\pi/2$	4.810	.42	.657	-6	.92	.399	$\pi/2$.208
.43	1.537	16	.93	2.535	25	$2\pi/2$	23.14	.43	.651	-6	.93	.395	$3\pi/2$.0432
.44	1.553	15	.94	2.560	26	$3\pi/2$	111.3	.44	.644	-6	.94	.391	$4\pi/2$.00898
0.45	1.568	16	0.95	2.586	26	$4\pi/2$	535.5	0.45	.638	-7	0.95	.387	$5\pi/2$.00187
.46	1.584	16	.96	2.612	26	$5\pi/2$	2576.	.46	.631	-6	.96	.383	$6\pi/2$.000388
.47	1.600	16	.97	2.638	26	$6\pi/2$	12392.	.47	.625	-6	.97	.379	$7\pi/2$.000081
.48	1.616	16	.98	2.664	27	$7\pi/2$	59610.	.48	.619	-6	.98	.375	$8\pi/2$.000017
.49	1.632	17	.99	2.691	27	$8\pi/2$	286751.	.49	.613	-6	.99	.372		.000003
0.50	1.649		1.00	2.718				0.50	0.607		1.00	.368		

*Note: Do not interpolate in this column.

NORMAL OR GAUSSIAN DISTRIBUTION

TABLE 10. NORMAL OR GAUSSIAN DISTRIBUTION

x	$p(x) = (2\pi)^{-1/2} \exp(-x^2/2)$	$P(x) = \int_{-x}^x p(t) dt$	x	$p(x) = (2\pi)^{-1/2} \exp(-x^2/2)$	$P(x) = \int_{-x}^x p(t) dt$
0.00	0.39894	0.50000	2.55	0.01545	0.9 ² 4614
0.05	0.39844	0.51994	2.60	0.01358	0.9 ² 5339
0.10	0.39695	0.53983	2.65	0.01191	0.9 ² 5975
0.15	0.39448	0.55962	2.70	0.01042	0.9 ² 6533
0.20	0.39104	0.57926	2.75	0.0 ² 9094	0.9 ² 7020
0.25	0.38667	0.59871	2.80	0.0 ² 7915	0.9 ² 7445
0.30	0.38139	0.61791	2.85	0.0 ² 6873	0.9 ² 7814
0.35	0.37524	0.63683	2.90	0.0 ² 5953	0.9 ² 8134
0.40	0.36827	0.65542	2.95	0.0 ² 5143	0.9 ² 8411
0.45	0.36053	0.67364	3.00	0.0 ² 4432	0.9 ² 8650
0.50	0.35207	0.69146	3.05	0.0 ² 3810	0.9 ² 8856
0.55	0.34294	0.70884	3.10	0.0 ² 3267	0.9 ³ 0324
0.60	0.33322	0.72575	3.15	0.0 ² 2794	0.9 ³ 1836
0.65	0.32297	0.74215	3.20	0.0 ² 2384	0.9 ³ 3129
0.70	0.31225	0.75804	3.25	0.0 ² 2029	0.9 ³ 4230
0.75	0.30114	0.77337	3.30	0.0 ² 1723	0.9 ³ 5166
0.80	0.28969	0.78814	3.35	0.0 ² 1459	0.9 ³ 5959
0.85	0.27798	0.80234	3.40	0.0 ² 1232	0.9 ³ 6631
0.90	0.26609	0.81594	3.45	0.0 ² 1038	0.9 ³ 7197
0.95	0.25406	0.82894	3.50	0.0 ³ 8727	0.9 ³ 7674
1.00	0.24197	0.84134	3.55	0.0 ³ 7317	0.9 ³ 8074
1.05	0.22988	0.85314	3.60	0.0 ³ 6119	0.9 ³ 8409
1.10	0.21785	0.86433	3.65	0.0 ³ 5105	0.9 ³ 8689
1.15	0.20594	0.87493	3.70	0.0 ³ 4248	0.9 ³ 8922
1.20	0.19419	0.88493	3.75	0.0 ³ 3526	0.9 ⁴ 1158
1.25	0.18265	0.89435	3.80	0.0 ³ 2919	0.9 ⁴ 2765
1.30	0.17137	0.90320	3.85	0.0 ³ 2411	0.9 ⁴ 4094
1.35	0.16038	0.91149	3.90	0.0 ³ 1987	0.9 ⁴ 5190
1.40	0.14973	0.91924	3.95	0.0 ³ 1633	0.9 ⁴ 6092
1.45	0.13943	0.92647	4.00	0.0 ³ 1338	0.9 ⁴ 6833
1.50	0.12952	0.93319	4.05	0.0 ³ 1094	0.9 ⁴ 7439
1.55	0.12001	0.93943	4.10	0.0 ⁴ 8926	0.9 ⁴ 7934
1.60	0.11092	0.94520	4.15	0.0 ⁴ 7263	0.9 ⁴ 8338
1.65	0.10226	0.95053	4.20	0.0 ⁴ 5894	0.9 ⁴ 8665
1.70	0.09405	0.95543	4.25	0.0 ⁴ 4772	0.9 ⁴ 8931
1.75	0.08628	0.95994	4.30	0.0 ⁴ 3854	0.9 ⁵ 1460
1.80	0.07895	0.96407	4.35	0.0 ⁴ 3104	0.9 ⁵ 3193
1.85	0.07206	0.96784	4.40	0.0 ⁴ 2494	0.9 ⁵ 4587
1.90	0.06562	0.97128	4.45	0.0 ⁴ 1999	0.9 ⁵ 5706
1.95	0.05959	0.97441	4.50	0.0 ⁴ 1598	0.9 ⁵ 6602
2.00	0.05399	0.97725	4.55	0.0 ⁴ 1275	0.9 ⁵ 7318
2.05	0.04879	0.97982	4.60	0.0 ⁴ 1014	0.9 ⁵ 7888
2.10	0.04398	0.98214	4.65	0.0 ⁵ 8047	0.9 ⁵ 8340
2.15	0.03955	0.98422	4.70	0.0 ⁵ 6370	0.9 ⁵ 8699
2.20	0.03547	0.98610	4.75	0.0 ⁵ 5030	0.9 ⁵ 8983
2.25	0.03174	0.98778	4.80	0.0 ⁵ 3961	0.9 ⁶ 2067
2.30	0.02833	0.98928	4.85	0.0 ⁵ 3112	0.9 ⁶ 3827
2.35	0.02522	0.9 ² 0613	4.90	0.0 ⁵ 2439	0.9 ⁶ 5208
2.40	0.02239	0.9 ² 1802	4.95	0.0 ⁵ 1907	0.9 ⁶ 6289
2.45	0.01984	0.9 ² 2857	5.00	0.0 ⁵ 1487	0.9 ⁶ 7133
2.50	0.01753	0.9 ² 3790			

Note: 0.0²9094 = 0.009094

$$P(-x) = 1 - P(x)$$

0.9³0324 = 0.9990324

$$\int_{-x}^x p(t) dt = 2P(x) - 1$$

BESSEL FUNCTIONS

TABLE 11. $J_0(z)$

z	0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9
0	1.0000	0.9975	0.9900	0.9776	0.9604	0.9385	0.9120	0.8812	0.8463	0.8075
1	0.7652	0.7196	0.6711	0.6201	0.5669	0.5118	0.4554	0.3980	0.3400	0.2818
2	0.2239	0.1666	0.1104	0.0555	0.0025	-0.0484	-0.0968	-0.1424	-0.1850	-0.2243
3	-0.2601	-0.2921	-0.3202	-0.3443	-0.3643	-0.3801	-0.3918	-0.3992	-0.4026	-0.4018
4	-0.3971	-0.3887	-0.3766	-0.3610	-0.3423	-0.3205	-0.2961	-0.2693	-0.2404	-0.2097
5	-0.1776	-0.1443	-0.1103	-0.0758	-0.0412	-0.0068	+0.0270	0.0599	0.0917	0.1220
6	0.1506	0.1773	0.2017	0.2238	0.2433	0.2601	0.2740	0.2851	0.2931	0.2981
7	0.3001	0.2991	0.2951	0.2882	0.2786	0.2663	0.2516	0.2346	0.2154	0.1944
8	0.1717	0.1475	0.1222	0.0960	0.0692	0.0419	0.0146	-0.0125	-0.0392	-0.0653
9	-0.0903	-0.1142	-0.1367	-0.1577	-0.1768	-0.1939	-0.2090	-0.2218	-0.2323	-0.2403
10	-0.2459	-0.2490	-0.2496	-0.2477	-0.2434	-0.2366	-0.2276	-0.2164	-0.2032	-0.1881
11	-0.1712	-0.1528	-0.1330	-0.1121	-0.0902	-0.0677	-0.0446	-0.0213	+0.0020	0.0250
12	0.0477	0.0697	0.0908	0.1108	0.1296	0.1469	0.1626	0.1766	0.1887	0.1988
13	0.2069	0.2129	0.2167	0.2183	0.2177	0.2150	0.2101	0.2032	0.1943	0.1836
14	0.1711	0.1570	0.1414	0.1245	0.1065	0.0875	0.0679	0.0476	0.0271	0.0064
15	-0.0142	-0.0346	-0.0544	-0.0736	-0.0919	-0.1092	-0.1253	-0.1401	-0.1533	-0.1650

TABLE 12. $J_1(z)$

z	0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9
0	0.0000	0.0499	0.0995	0.1483	0.1960	0.2423	0.2867	0.3290	0.3688	0.4059
1	0.4401	0.4709	0.4983	0.5220	0.5419	0.5579	0.5699	0.5778	0.5815	0.5812
2	0.5767	0.5683	0.5560	0.5399	0.5202	0.4971	0.4708	0.4416	0.4097	0.3754
3	0.3391	0.3009	0.2613	0.2207	0.1792	0.1374	0.0955	0.0538	0.0128	-0.0272
4	-0.0660	-0.1033	-0.1386	-0.1719	-0.2028	-0.2311	-0.2566	-0.2791	-0.2985	-0.3147
5	-0.3276	-0.3371	-0.3432	-0.3460	-0.3453	-0.3414	-0.3343	-0.3241	-0.3110	-0.2951
6	-0.2767	-0.2559	-0.2329	-0.2081	-0.1816	-0.1538	-0.1250	-0.0953	-0.0652	-0.0349
7	-0.0047	+0.0252	0.0543	0.0826	0.1096	0.1352	0.1592	0.1813	0.2014	0.2192
8	0.2346	0.2476	0.2580	0.2657	0.2708	0.2731	0.2728	0.2697	0.2641	0.2559
9	0.2453	0.2324	0.2174	0.2004	0.1816	0.1613	0.1395	0.1166	0.0928	0.0684
10	0.0435	0.0184	-0.0066	-0.0313	-0.0555	-0.0789	-0.1012	-0.1224	-0.1422	-0.1603
11	-0.1768	-0.1913	-0.2039	-0.2143	-0.2225	-0.2284	-0.2320	-0.2333	-0.2323	-0.2290
12	-0.2234	-0.2157	-0.2060	-0.1943	-0.1807	-0.1655	-0.1487	-0.1307	-0.1114	-0.0912
13	-0.0703	-0.0489	-0.0271	-0.0052	+0.0166	0.0380	0.0590	0.0791	0.0984	0.1165
14	0.1334	0.1488	0.1626	0.1747	0.1850	0.1934	0.1999	0.2043	0.2066	0.2069
15	0.2051	0.2013	0.1955	0.1879	0.1784	0.1672	0.1544	0.1402	0.1247	0.1080

TABLE 13. $J_2(z)$

z	0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9
0	0.0000	0.0012	0.0050	0.0112	0.0197	0.0306	0.0437	0.0588	0.0758	0.0946
1	0.1149	0.1366	0.1593	0.1830	0.2074	0.2321	0.2570	0.2817	0.3061	0.3299
2	0.3528	0.3746	0.3951	0.4139	0.4310	0.4461	0.4590	0.4696	0.4777	0.4832
3	0.4861	0.4862	0.4835	0.4780	0.4697	0.4586	0.4448	0.4283	0.4093	0.3879
4	0.3641	0.3383	0.3105	0.2811	0.2501	0.2178	0.1846	0.1506	0.1161	0.0813

TABLE 14. $J_3(z)$

z	0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9
0	0.0000	0.0000	0.0002	0.0006	0.0013	0.0026	0.0044	0.0069	0.0102	0.0144
1	0.0196	0.0257	0.0329	0.0411	0.0505	0.0610	0.0725	0.0851	0.0988	0.1134
2	0.1289	0.1453	0.1623	0.1800	0.1981	0.2166	0.2353	0.2540	0.2727	0.2911
3	0.3091	0.3264	0.3431	0.3588	0.3734	0.3868	0.3988	0.4092	0.4180	0.4250
4	0.4302	0.4333	0.4344	0.4333	0.4301	0.4247	0.4171	0.4072	0.3952	0.3811

TABLE 15. $J_4(z)$

z	0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9
0	0.0000	0.0000	0.0000	0.0000	0.0001	0.0002	0.0003	0.0006	0.0010	0.0016
1	0.0025	0.0036	0.0050	0.0068	0.0091	0.0118	0.0150	0.0188	0.0232	0.0283
2	0.0340	0.0405	0.0476	0.0556	0.0643	0.0738	0.0840	0.0950	0.1067	0.1190
3	0.1320	0.1456	0.1597	0.1743	0.1891	0.2044	0.2198	0.2353	0.2507	0.2661
4	0.2811	0.2958	0.3100	0.3236	0.3365	0.3484	0.3594	0.3693	0.3780	0.3853

49

Miscellaneous Data

Effect of Altitude on Breakdown Voltages 49-3

Pressure-Altitude Graph
Spark-Gap Breakdown Voltages

Weather Data 49-3

Temperature Extremes
World Temperatures
Precipitation Extremes
Wind-Velocity and Temperature Extremes in North America
Celsius Table of Relative Humidity or Percent of Saturation

World Time Chart 49-5

Materials and Finishes for Tropical and Marine Use 49-5

Corrosion
Fungus and Decay
Finish Application Table

Principal Low-Voltage Power Supplies in the World 49-8

Power Supply Wiring 49-11

Guide to Use of Tables

Wiring of Electronic Equipment and Chassis 49-12

Resistance Change with Temperature 49-15

Wire Identification 49-16

Color Coding

Diameter of Circle Enclosing a Given Number of Smaller
Circles 49-17

Arrangement of Center or Core Circles
Diameter of Enclosing Circle When Outer Layer of Circles Is
Complete

Diameter of Enclosing Circle When Outer Layer of Circles Is Not
Complete

Approximate Equation When Number of Enclosed Circles Is Large

- Torque and Horsepower 49-18
- Transmission-Line Sag Calculations 49-18
 - Supports at Same Level
 - Supports at Different Levels
- Structural Standards for Steel Radio Towers 49-20
 - Material
 - Loading
 - Unit Stresses
 - Foundations
 - Wind Velocities and Pressures
- Vibration and Shock Isolation 49-22
 - Symbols
 - Equations
 - Acceleration
 - Magnitudes
 - Natural Frequency
 - Resonance
 - Damping
 - Practical Application
 - Shock
- US Graphic Symbols 49-26
 - Diagram Types
 - Orientation
 - Detached Elements
 - Terminals
 - Associated or Future Equipment
- British Graphic Symbols 49-28
- Standardized Graphic Symbols for Logic Design 49-28
- Safety Levels for Exposure to EM Fields 49-28
 - References
 - Additional Reading
- Patent, Copyright, and Trade Secret Protection of Inventions 49-31
 - Intellectual Property Protection of Technology in the United States
 - Trade Secret Protection of Inventions
 - Patent Protection of Inventions
 - Copyright Protection of Inventions
- Summary of Military Nomenclature System 49-34
 - Nomenclature Policy
 - Modification Letters
 - Developmental Indicators
 - Examples of JETDS Type Numbers

EFFECT OF ALTITUDE ON BREAKDOWN VOLTAGES

Pressure-Altitude Graph

Design of electrical equipment for aircraft is somewhat complicated by the requirement of additional insulation for high voltages as a result of the decrease in atmospheric pressure. The extent of this effect may be determined from Figs. 1 and 2 and Table 1. (1 inch mercury = 25.4 millimeters mercury = 0.4912 pound/inch² = 3.38 × 10³ pascals.)

Spark-Gap Breakdown Voltages

Fig. 2 is for a voltage that is continuous or at a frequency low enough to permit complete deionization between cycles, between needle points or clean smooth spherical surfaces (electrodes ungrounded) in dust-free dry air. Temperature is 25 degrees Celsius and pressure is 760 millimeters (29.9 inches) of mercury. Peak kilovolts shown in the figure should be multiplied by the factors given in Table 1 for atmospheric conditions other than the above.

An approximate rule for uniform fields at all frequencies up to at least 300 megahertz is that the breakdown gradient of air is 30 peak kilovolts/centimeter or 75 peak kilovolts/inch at sea level (760 millimeters of mercury) and normal temperature (25 degrees Celsius). The breakdown voltage is approximately proportional to pressure and inversely proportional to absolute (Kelvin) temperature.

Certain synthetic gases have higher dielectric strengths than air. Two such gases that appear to be useful for electrical insulation are sulfur hexafluoride (SF₆) and Freon 12 (CCl₂F₂), which both have about

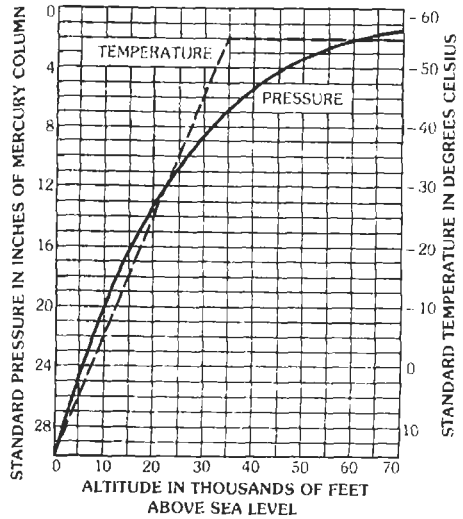


Fig. 2. Spark-gap breakdown voltages.

2.5 times the dielectric strength of air. Mixtures of sulfur hexafluoride with helium and of perfluoromethylcyclohexane (C₇F₁₄) with nitrogen have good dielectric strength as well as other desirable properties.

WEATHER DATA*

Temperature Extremes

United States (contiguous):

Lowest temperature:

Rodgers Pass, Montana (January 20, 1954)—-57 °C (-70 °F).

Highest temperature:

Greenland Ranch, Death Valley, California (July 10, 1913)—57 °C (134 °F).

Alaska:

Lowest temperature:

Prospect Creek Camp (January 23, 1971)—-62 °C (-79.8 °F).

Highest temperature:

Fort Yukon (June 27, 1915)—38 °C (100 °F).

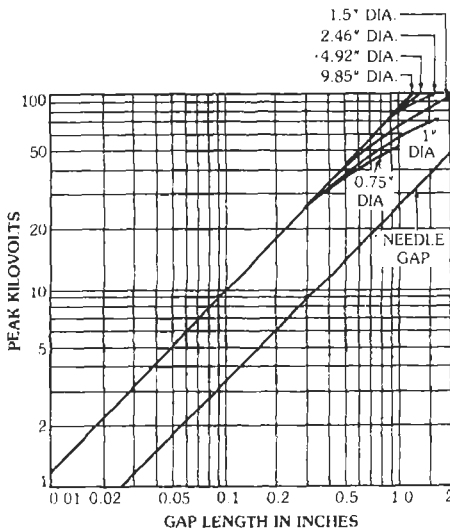


Fig. 1. Pressure as a function of altitude.

* Compiled in part from "Climate and Man," *Yearbook of Agriculture*, US Dept. of Agriculture. Obtainable from Superintendent of Documents, Government Printing Office, Washington, DC 20402. See also "Weather Extremes Around the World," a world map compiled (1984) by the Geographic Sciences Laboratory, US Army Engineer Topographic Laboratories, Belvoir, Virginia 22060.

For a comprehensive summary of available climatological information, refer to *Selective Guide to Climatic Data Sources* (No. 4.11), by Warren L. Hatch, July 1983, available from the National Climatic Data Center, Asheville, North Carolina 28801-2696; phone: (704) 259-0682.

TABLE 1. MULTIPLYING FACTORS

Pressure		Temperature in Degrees Celsius					
in Hg	mm Hg	-40	-20	0	20	40	60
5	127	0.26	0.24	0.23	0.21	0.20	0.19
10	254	0.47	0.44	0.42	0.39	0.37	0.34
15	381	0.68	0.64	0.60	0.56	0.53	0.50
20	508	0.87	0.82	0.77	0.72	0.68	0.64
25	635	1.07	0.99	0.93	0.87	0.82	0.77
30	762	1.25	1.17	1.10	1.03	0.97	0.91
35	889	1.43	1.34	1.26	1.19	1.12	1.05
40	1016	1.61	1.51	1.42	1.33	1.25	1.17
45	1143	1.79	1.68	1.58	1.49	1.40	1.31
50	1270	1.96	1.84	1.73	1.63	1.53	1.44
55	1397	2.13	2.01	1.89	1.78	1.67	1.57
60	1524	2.30	2.17	2.04	1.92	1.80	1.69

World:

Lowest temperature:

Vostok, Antarctica (July 21, 1983)—
-89 °C (-129 °F).

Highest temperature:

Azizia, Libya, North Africa (September 13,
1922)—58 °C (136 °F).

Lowest mean temperature (annual):

Framheim, Antarctica—-26 °C (-14 °F).

Highest mean temperature (annual):

Dalol Ethiopia—34 °C (94 °F).

World Temperatures

Territory	Maximum		Minimum	
	°C	°F	°C	°F
NORTH AMERICA				
Alaska	38	100	-62	-79.8
Canada	45	113	-63	-81
Canal Zone	36	97	-17	63
Greenland	30	86	-43	-46
Mexico	48	118	-12	11
USA	57	134	-57	-70
West Indies	39	102	7	45
SOUTH AMERICA				
Argentina	46	115	-33	-27
Bolivia	28	82	-4	25
Brazil	42	108	-6	21
Chile	37	99	-7	19
Venezuela	39	102	7	45
EUROPE				
British Isles	38	100	-16	4
France	42	107	-26	-14
Germany	38	100	-27	-16

Territory	Maximum		Minimum	
	°C	°F	°C	°F
EUROPE				
Iceland	22	71	-21	-6
Italy	46	114	-16	4
Norway	35	95	-32	-26
Russia	43	110	-52	-61
Spain	51	124	-12	10
Sweden	33	92	-45	-49
Turkey	38	100	-8	17
ASIA				
Saudi Arabia	51	123	2	35
China	44	111	-23	-10
East Indies	38	101	16	60
India	49	120	-28	-19
Iraq	52	125	-7	19
Japan	38	101	-22	-7
Malaysia	36	97	19	66
Philippine Islands	38	101	14	58
Russia	43	109	-68	-90
Thailand	41	106	11	52
Tibet	29	85	-29	-20
Turkey	44	111	-30	-22
Vietnam	45	113	1	33
AFRICA				
Algeria	56	133	-17	1
Angola	33	91	1	33
Egypt	51	124	-1	31
Ethiopia	44	111	0	32
Libya	58	136	2	35
Morocco	48	119	-15	5
Rhodesia	44	112	-8	18
Somalia	34	93	16	61
Sudan	52	126	-2	28
Tunisia	50	122	-2	28
Union of South Africa	44	111	-6	21
Zaire	36	97	1	34
AUSTRALASIA				
Australia	53	127	-7	19
Hawaii	33	91	11	51
New Zealand	34	94	-5	23
Samoan Islands	36	96	16	61
Solomon Islands	36	97	21	70

Precipitation Extremes*United States:*

Wettest state:

Louisiana—average annual rainfall 57.34 inches.

Driest state:

Nevada—average annual rainfall 8.60 inches.

Maximum recorded:

Camp Leroy, California (January 22-23,
1943)—26.12 inches in 24 hours.

Minimums recorded:

Bagdad, California (1909-1913)—3.93 inches
in 5 years.

Greenland Ranch, California—1.76 inches annual average.

World:

Maximums recorded:

Cherrapunji, India (July, 1861)—366 inches in 1 month. (Average annual rainfall of Cherrapunji is 450 inches).

Baguio, Luzon, Philippines, July 14–15, 1911—46 inches in 24 hours.

Minimums recorded:

Wadi Halfa (Sudan) and Aswan (Egypt) are in the “rainless” area; average annual rainfall is too small to be measured.

Wind-Velocity and Temperature Extremes in North America

Data regarding extremes of temperature and wind velocity for several locations in North America may be found in Table 2.

Celsius Table of Relative Humidity or Percent of Saturation

To find the relative humidity when the wet-bulb and dry-bulb temperatures are known, consult Table 3.

WORLD TIME CHART

Time differences between selected major world cities may be determined from Chart 1.

MATERIALS AND FINISHES FOR TROPICAL AND MARINE USE

Corrosion

Ordinary finishing of equipment fails to meet satisfactorily conditions encountered in tropical and marine use. Under these conditions, corrosive influences are greatly aggravated by prevailing higher relative humidity.

TABLE 2. WIND-VELOCITY AND TEMPERATURE EXTREMES IN NORTH AMERICA

Station	Wind* (miles/hour)	Temperature, Degrees Fahrenheit	
		Maximum	Minimum
UNITED STATES, 1871–1955			
Albany, New York	71	104	–26
Amarillo, Texas	84	108	–16
Buffalo, New York	91	99	–21
Charleston, South Carolina	76	104	7
Chicago, Illinois	87	105	–23
Bismarck, North Dakota	72	114	–45
Hatteras, North Carolina	110	97	8
Miami, Florida	132	95	27
Minneapolis, Minnesota	92	108	–34
Mobile, Alabama	87	104	–11
Mt. Washington, New Hampshire	188**	71	–46
Nantucket, Massachusetts	91	95	–6
New York, New York	99	102	–14
North Platte, Nebraska	72	112	–35
Pensacola, Florida	114	103	7
Washington, D.C.	62	106	–15
San Juan, Puerto Rico	149†	94	62
CANADA, 1955			
Banff, Alberta	52‡	97	–60
Kamloops, British Columbia	34‡	107	–37
Sable Island, Nova Scotia	64‡	86	–12
Toronto, Ontario	48‡	105	–46

*Maximum corrected wind velocity (fastest single mile).

**Gusts were recorded at 231 miles/hour (corrected).

†Estimated.

‡For a period of 5 minutes.

TABLE 3. CELSIUS TABLE OF RELATIVE HUMIDITY OR PERCENT OF SATURATION

Dry Bulb °C	Difference Between Readings of Wet and Dry Bulbs in Degrees Celsius																																							
	0.5	1.0	1.5	2.0	2.5	3.0	3.5	4.0	4.5	5	6	7	8	9	10	11	12	13	14	15	16	18	20	22	24	26	28	30	32	34	36	38	40							
4	93	85	77	70	63	56	48	41	34	28	15																													
8	94	87	81	74	68	62	56	50	45	39	28	17																												
12	94	89	84	78	73	68	63	58	53	48	38	30	21	12	4																									
16	95	90	85	81	76	71	67	62	58	54	45	37	29	21	14	7																								
20	96	91	87	82	78	74	70	66	62	58	51	44	36	30	23	17	11																							
22	96	92	87	83	79	75	72	68	64	60	53	46	40	34	27	21	16	11																						
24	96	92	88	85	81	77	74	70	66	63	56	49	43	37	31	26	21	14	10																					
26	96	92	89	85	81	77	74	71	67	64	57	51	45	39	34	28	23	18	13																					
28	96	92	89	85	82	78	75	72	68	65	59	53	47	42	37	31	26	21	17	13																				
30	96	93	89	86	82	79	76	73	70	67	61	55	50	44	39	35	30	24	20	16	12																			
32	96	93	90	86	83	80	77	74	71	68	62	56	51	46	41	36	32	27	23	19	15																			
34	97	93	90	87	84	81	77	74	71	69	63	58	53	48	43	38	34	30	26	22	18	10																		
36	97	93	90	87	84	81	78	75	72	70	64	59	54	50	45	41	36	32	28	24	21	13																		
38	97	94	90	87	84	81	79	76	73	70	65	60	56	51	46	42	38	34	30	26	23	16	10																	
40	97	94	91	88	85	82	79	76	74	71	66	61	57	52	48	44	40	36	32	29	25	19	13																	
44	97	94	91	88	86	83	80	77	75	73	68	63	59	54	50	47	43	39	36	32	29	23	17	12																
48	97	94	92	89	86	84	81	78	76	74	69	65	61	56	53	49	45	42	39	35	33	27	21	16	12															
52	97	94	92	89	87	84	82	79	77	75	70	66	62	58	55	51	48	44	41	38	35	30	25	20	16	11														
56	97	95	92	90	87	85	83	80	78	76	72	68	64	60	57	53	50	46	43	40	38	32	27	23	19	15	11													
60	98	95	93	90	88	86	83	81	79	77	73	69	65	62	58	55	52	48	45	43	40	35	30	26	21	18	14	11												
70	98	96	93	91	89	87	85	83	81	79	75	71	68	65	61	58	55	52	50	47	44	40	35	31	27	23	20	17	14	11										
80	98	96	94	92	90	88	86	84	83	81	77	74	71	67	64	61	58	56	53	50	48	43	39	35	31	28	24	22	19	16	14	11								
90	98	97	95	93	91	89	87	85	84	82	79	76	73	69	67	64	61	58	56	53	51	47	42	39	35	32	28	26	23	20	18	16	14							
100	99	97	95	93	92	90	88	86	85	83	80	77	74	71	68	66	63	60	58	56	54	49	45	42	38	35	32	29	26	24	22	19	17							

Example: Assume dry-bulb reading (thermometer exposed directly to atmosphere) is 20 °C and wet-bulb reading is 17 °C, or a difference of 3 °C. The relative humidity at 20 °C is then 74%.

CHART 1. WORLD TIME CHART

Tutuila, Samoa	Hawaiian Islands, Tahiti	Los Angeles, San Francisco, Seattle	Chicago, Central America (except Panama), Mexico, Winnipeg	New York, Montreal, Miami, Havana, Panama, Bogota, Lima, Quito	Bermuda, Puerto Rico, Caracas, La Paz, Asuncion	Buenos Aires, * Rio de Janeiro, Santos, Sao Paulo, Montevideo	Iceland	Lisbon, Dublin, Algiers, Dakar, Ascension Island	Greenwich Civil Time (GCT) or Universal Time (UT)	London, * Paris, * Madrid, * Brussels, Rome, Berlin, Vienna, Oslo, Stockholm, Copenhagen, Amsterdam, Tunis, Warsaw	Athens, Israel, Ankara, Cairo, Capetown	Moscow, * Ethiopia, Iraq, Malagasy Republic	Bombay, Sri Lanka, New Delhi	Bangkok, Chungking, Chengtu, Kunning	Hong Kong, Manila, Shanghai, Saigon, Taipeh, Celebes	Japan, Adelaide, Korea, Manchuria	Sydney, Melbourne, Brisbane, Guam, New Guinea, Khabarovsk	Solomon Islands, New Caledonia	Wellington, * Auckland*
1:00	2:00	4:00	6:00	7:00	8:00	9:00	11:00	Midnite	0000	1:00	2:00	3:00	5:30	7:00	8:00	9:00	10:00	11:00	11:30
2:00	3:00	5:00	7:00	8:00	9:00	10:00	Midnite	1:00	0100	2:00	3:00	4:00	6:30	8:00	9:00	10:00	11:00	Noon	12:30
3:00	4:00	6:00	8:00	9:00	10:00	11:00	Midnite	2:00	0200	3:00	4:00	5:00	7:30	9:00	10:00	11:00	Noon	1:00	1:30
4:00	5:00	7:00	9:00	10:00	11:00	Midnite	2:00	3:00	0300	4:00	5:00	6:00	8:30	10:00	11:00	Noon	1:00	2:00	2:30
5:00	6:00	8:00	10:00	11:00	Midnite	1:00	3:00	4:00	0400	5:00	6:00	7:00	9:30	11:00	Noon	1:00	2:00	3:00	3:30
6:00	7:00	9:00	11:00	Midnite	1:00	2:00	4:00	5:00	0500	6:00	7:00	8:00	10:30	Noon	1:00	2:00	3:00	4:00	4:30
7:00	8:00	10:00	Midnite	1:00	2:00	3:00	5:00	6:00	0600	7:00	8:00	9:00	11:30	1:00	2:00	3:00	4:00	5:00	5:30
8:00	9:00	11:00	1:00	2:00	3:00	4:00	6:00	7:00	0700	8:00	9:00	10:00	12:30	2:00	3:00	4:00	5:00	6:00	6:30
9:00	10:00	Midnite	2:00	3:00	4:00	5:00	7:00	8:00	0800	9:00	10:00	11:00	1:30	3:00	4:00	5:00	6:00	7:00	7:30
10:00	11:00	1:00	3:00	4:00	5:00	6:00	8:00	9:00	0900	10:00	11:00	Noon	2:30	4:00	5:00	6:00	7:00	8:00	8:30
11:00	Midnite	2:00	4:00	5:00	6:00	7:00	9:00	10:00	1000	11:00	Noon	1:00	3:30	5:00	6:00	7:00	8:00	9:00	9:30
Midnite	1:00	3:00	5:00	6:00	7:00	8:00	10:00	11:00	1100	Noon	1:00	2:00	4:30	6:00	7:00	8:00	9:00	10:00	10:30
1:00	2:00	4:00	6:00	7:00	8:00	9:00	11:00	Noon	1200	1:00	2:00	3:00	5:30	7:00	8:00	9:00	10:00	11:00	11:30
2:00	3:00	5:00	7:00	8:00	9:00	10:00	Noon	1:00	1300	2:00	3:00	4:00	6:30	8:00	9:00	10:00	11:00	Midnite	12:30
3:00	4:00	6:00	8:00	9:00	10:00	11:00	1:00	2:00	1400	3:00	4:00	5:00	7:30	9:00	10:00	11:00	Midnite	1:00	1:30
4:00	5:00	7:00	9:00	10:00	11:00	Noon	2:00	3:00	1500	4:00	5:00	6:00	8:30	10:00	11:00	Midnite	1:00	2:00	2:30
5:00	6:00	8:00	10:00	11:00	Noon	1:00	3:00	4:00	1600	5:00	6:00	7:00	9:30	11:00	Midnite	1:00	2:00	3:00	3:30
6:00	7:00	9:00	11:00	Noon	1:00	2:00	4:00	5:00	1700	6:00	7:00	8:00	10:30	Midnite	1:00	2:00	3:00	4:00	4:30
7:00	8:00	10:00	Noon	1:00	2:00	3:00	5:00	6:00	1800	7:00	8:00	9:00	11:30	1:00	2:00	3:00	4:00	5:00	5:30
8:00	9:00	11:00	1:00	2:00	3:00	4:00	6:00	7:00	1900	8:00	9:00	10:00	12:30	2:00	3:00	4:00	5:00	6:00	6:30
9:00	10:00	Noon	2:00	3:00	4:00	5:00	7:00	8:00	2000	9:00	10:00	11:00	1:30	3:00	4:00	5:00	6:00	7:00	7:30
10:00	11:00	1:00	3:00	4:00	5:00	6:00	8:00	9:00	2100	10:00	11:00	Midnite	2:30	4:00	5:00	6:00	7:00	8:00	8:30
11:00	Noon	2:00	4:00	5:00	6:00	7:00	9:00	10:00	2200	11:00	Midnite	1:00	3:30	5:00	6:00	7:00	8:00	9:00	9:30
Noon	1:00	3:00	5:00	6:00	7:00	8:00	10:00	11:00	2300	Midnite	1:00	2:00	4:30	6:00	7:00	8:00	9:00	10:00	10:30
1:00	2:00	4:00	6:00	7:00	8:00	9:00	11:00	Midnite	2400	1:00	2:00	3:00	5:30	7:00	8:00	9:00	10:00	11:00	11:30

- Notes:
- (1) Lightface figures designate am, bold figures pm.
 - (2) Time is that used at places indicated. In general, this is standard time, but for places marked with asterisks it is permanent daylight saving time. Temporary daylight saving time is commonplace but not indicated above.
 - (3) When passing the heavy line going down or to the right, add 1 day. When passing the heavy line going up or to the left, subtract 1 day.

ties, and temperature cycling causes alternate condensation on and evaporation of moisture from finished surfaces. Useful equipment life under adverse atmospheric influences depends largely on proper choice of base materials and finishes applied. Especially important in tropical and marine applications is avoidance of electrical contact between dissimilar metals.

Dissimilar metals widely separated in the galvanic series should not be bolted, riveted, etc., without separation by insulating material at the facing surfaces. The only exception occurs when both surfaces have been coated with the same protective metal, e.g., electroplating, hot dipping, galvanizing, etc.

Aluminum, steel, zinc, and cadmium should never be used bare. Electrical contact surfaces should be given copper-nickel-chromium or copper-nickel finish, and, in addition, they should be silver plated. Adjustable-capacitor plates should be silver plated.

An additional 0.000015-inch to 0.000020-inch electroplating of hard, bright gold over the silver greatly improves resistance to tarnish and oxidation and to attack by most chemicals, lowers electrical resistance, and provides long-term solderability.

Fungus and Decay

The value of fungicidal coatings or treatments is controversial. If equipment is to operate under tropical conditions, greater success can be achieved by the use of materials that do not provide a nutrient medium for fungus and insects. The following types or kinds of materials are examples of nonnutrient mediums that are generally considered acceptable.

- Metals
- Glass
- Ceramics (steatite, glass-bonded mica)
- Mica
- Polyamide
- Cellulose acetate
- Rubber (natural or synthetic)
- Plastic materials using glass, mica, or asbestos as a filler
- Polyvinylchloride
- Polytetrafluoroethylene
- Monochlorotrifluoroethylene

The following types or kinds of materials should not be used, except where such materials are fabricated into completed parts and their use is acceptable to the customer.

- Linen
- Cellulose nitrate
- Regenerated cellulose
- Wood
- Jute
- Leather
- Cork
- Paper and cardboard
- Organic fiberboard

Hair or wool felts

Plastic materials using cotton, linen, or wood flour as a filler

Wood should not be used as an electrical insulator, and its use for other purposes should be restricted to those parts for which a superior substitute is not known. When used, it should be pressure-treated and impregnated to resist moisture, insects, and decay with a waterborne preservative (as specified in Federal Specification TT-W-571), and it should also be treated with a suitable fire-retardant chemical.

Finish Application Table

For information regarding the application of finishes, see Table 4.

PRINCIPAL LOW-VOLTAGE POWER SUPPLIES IN THE WORLD

Territory (Frequency) Voltage

North America:

Alaska (60) 120/240
 Belize (60) 110/220
 Bermuda (60) 115/230; some 120/208
 Canada (60) 120/240; some 115/230
 Costa Rica (60) 110/220
 El Salvador (60) 110/220
 Guatemala (60) 110/240; some 220, 120/208
 Honduras (60) 110/220
 Mexico (50, 60) 127/220 and other voltages
 City (50) 125/216
 Nicaragua (60) 120
 Panama (60) 110/220; some 120/240, 115/230
 United States (60) 120/240 and 120/208

West Indies:

Antigua (60) 230/400
 Bahamas (60) 115/200; some 115/220
 Barbados (50) 120/208; some 110/200
 Cuba (60) 115/230; some 120/208
 Dominican Republic (60) 115/230
 Guadeloupe (50) 127/220
 Jamaica (50, some 60) 110/220
 Martinique (50) 127/220
 Puerto Rico (60) 120/240
 Trinidad (50) 115/230
 Virgin Islands (60) 120/240

South America:

Argentina (50) 220/380; also 220/440 dc
 Bolivia (50, also 60) 220 and other voltages
 Brazil (50, 60) 110, 220; also other voltages and dc Rio de Janeiro (50) 125/216
 Chile (50) 220/380; some 220 dc
 Colombia (60) 110/220; also 120/240 and others

TABLE 4. FINISH APPLICATION TABLE*

Material	Finish	Remarks
Aluminum alloy	Anodizing	An electrochemical-oxidation surface treatment, for improving corrosion resistance; not an electroplating process. For riveted or welded assemblies, specify chromic acid anodizing. Do not anodize parts with nonaluminum inserts. Colors vary: yellow-green, gray, or black.
	“Alrok”	Chemical-dip oxide treatment. Cheap. Inferior in abrasion and corrosion resistance to the anodizing process, but applicable to assemblies of aluminum and nonaluminum materials.
Copper and zinc alloys	Bright acid dip	Immersion of parts in acid solution. Clear lacquer applied to prevent tarnish.
Brass, bronze, zinc die-casting alloys	Brass, chrome, nickel, tin	As discussed under steel.
Magnesium alloy	Dichromate treatment	Corrosion-preventive dichromate dip. Yellow color.
Stainless steel	Passivating treatment	Nitric-acid immunizing dip.
Steel	Cadmium	Electroplate, dull white color, good corrosion resistance, easily scratched, good thread antiseize. Poor wear and galling resistance.
	Chromium	Electroplate, excellent corrosion resistance and lustrous appearance. Relatively expensive. Specify hard chrome plate for exceptionally hard abrasion-resistant surface. Has low coefficient of friction. Used to some extent on nonferrous metals particularly when die-cast. Chrome-plated objects usually receive a base electroplate of copper, then nickel, followed by chromium. Used for buildup of parts that are undersized. Do not use on parts with deep recesses.
	Blueing	Immersion of cleaned and polished steel into heated saltpeter or carbonaceous material. Part then rubbed with linseed oil. Cheap. Poor corrosion resistance.
	Silver plate	Electroplate, frosted appearance; buff to brighten. Tarnishes readily. Good bearing lining. For electrical contacts, reflectors.
	Zinc plate	Dip in molten zinc (galvanizing) or electroplate of low-carbon or low-alloy steels. Low cost. Generally inferior to cadmium plate. Poor appearance. Poor wear resistance: Electroplate has better adherence to base metal than hot-dip coating. For improving corrosion resistance, zinc-plated parts are given special inhibiting treatments.
	Nickel plate	Electroplate, dull white. Does not protect steel from galvanic corrosion. If plating is broken, corrosion of base metal will be hastened. Finishes in dull white, polished, or black. Do not use on parts with deep recesses.
	Black-oxide dip	Nonmetallic chemical black oxidizing treatment for steel, cast iron, and wrought iron. Inferior to electroplate. No buildup. Suitable for parts with close dimensional requirements as gears, worms, and guides. Poor abrasion resistance.
	Phosphate treatment	Nonmetallic chemical treatment for steel and iron products. Suitable for protection of internal surfaces of hollow parts. Small amount of surface buildup. Inferior to metallic electroplate. Poor abrasion resistance. Good paint base.

Continued on next page.

TABLE 4 (CONT.) FINISH APPLICATION TABLE*

Material	Finish	Remarks
Steel (cont.)	Tin plate	Hot dip or electroplate. Excellent corrosion resistance, but if broken will not protect steel from galvanic corrosion. Also used for copper, brass, and bronze parts that must be soldered after plating. Tin-plated parts can be severely worked and deformed without rupture of plating.
	Brass plate	Electroplate of copper and zinc. Applied to brass and steel parts where uniform appearance is desired. Applied to steel parts when bonding to rubber is desired.
	Copper plate	Electroplate applied before nickel or chrome plates. Also for parts to be brazed or protected against carburization. Tarnishes readily.

* By Z. Fox. Reprinted by permission from *Product Engineering*, Vol. 19, January 1948, p. 161.

Ecuador (60) 120/208; also 110/220 and others
 French Guiana (50) 127/220
 Guyana (50, 60) 110/220
 Paraguay (50) 220/440; some 220/440 dc
 Peru (60) 220; some 110
 Surinam (50, 60) 127/220; some 115/230
 Uruguay (50) 220
 Venezuela (60, some 50) 120/208, 120/240

Europe:

Austria (50) 220/380; Vienna also has 220/440 dc
 Azores (50) 220/380
 Belgium (50) 220/380 and many others; some dc
 Canary Islands (50) 127/220
 Denmark (50) 220/380; also 220/440 dc
 Finland (50) 220/380
 France (50) 120/240, 220/380, and many others
 Germany (Federal Republic) (50) 220/380; also others, some dc
 Gibraltar (50) 240/415
 Greece (50) 220/380; also others, some dc
 Iceland (50) 220; some 220/380
 Ireland (50) 220/380; some 220/440 dc
 Italy (50) 127/220, 220/380 and others
 Luxembourg (50) 110/190, 220/380
 Madeira (50) 220/380; also 220/440 dc
 Malta (50) 240/415
 Monaco (50) 127/220, 220/380
 Netherlands (50) 220/380; also 127/220
 Norway (50) 230
 Portugal (50) 220/380; some 110/190
 Spain (50) 127/220; also 220/380, some dc
 Sweden (50) 127/220, 220/380; some dc
 Switzerland (50) 220/380
 Turkey (50) 220/380; some 110/190
 United Kingdom (50) 240/415 and others, some dc
 Yugoslavia (50) 220/380

Asia:

Afghanistan (50) 220/380
 Burma (Myanmar) (50) 230
 Cambodia (50) 120/208; some 220/380

Cyprus (50) 240
 Hong Kong (50) 200/346
 India (50) 230/400 and others, some dc
 Indonesia (50) 127/220
 Iran (50) 220/380
 Iraq (50) 220/380
 Israel (50) 230/400
 Japan (50, 60) 100/200
 Jordan (50) 220/380
 Korea (60) 100/200
 Kuwait (50) 240/415
 Laos (50) 127/220; some 220/380
 Lebanon (50) 110/190; some 220/380
 Malaysia (50) 230/400; some 240/415
 Nepal (50) 110/220
 Okinawa (60) 120/240
 Pakistan (50) 230/400 and others, some dc
 Philippines (60) 110, 220, and others
 Saudi Arabia (50, 60) 120/208; also 220/380, 230/400
 Singapore (50) 230/400
 Sri Lanka (50) 230/400
 Syria (50) 115/200; some 220/380
 Taiwan (60) 100/200
 Thailand (50) 220/380; also 110/190
 Vietnam (50) 220/380 future standard
 Yemen Arab Republic (50) 220
 Yemen, Peoples Democratic Republic (50) 230/400

Africa:

Algeria (50) 127/220, 220/380
 Angola (50) 220/380
 Benin (50) 220/380
 Burkina Faso (50) 220/380
 Egypt (50) 110, 220 and others; some dc
 Ethiopia (50) 220/380; some 127/220
 Guinea (50) 220/380; some 127/220
 Kenya (50) 240/415
 Liberia (60) 120/240
 Libya (50) 125/220; some 230/400
 Madagascar (50) 220/380; some 127/220
 Mauritius (50) 230/400
 Morocco (50) 115/200; also 230/400 and others

Mozambique (50) 220/380
 Niger (50) 220/380
 Nigeria (50) 230/400
 Senegal (50) 127/220
 Sierra Leone (50) 230/400
 Somalia (50) 220/440; also 110, 230
 South Africa (50) 220/380; also others, some dc
 Sudan (50) 240/415
 Tanzania (50) 230/400
 Tunisia (50) 220/380; also others
 Uganda (50) 240/415
 Zaire (50) 220/380
 Zimbabwe (50) 220/380; also 230/400

Oceania:

Australia (50) 240/415; also others and dc
 Fiji Islands (50) 240/415
 Hawaii (60) 120/240
 New Caledonia (50) 220/440
 New Zealand (50) 230/400

Notes:

1. Abstracted from *Electric Power Abroad*, issued 1963 by the Bureau of International Commerce of the US Department of Commerce. This pamphlet is obtainable from the Superintendent of Documents, US Government Printing Office, Washington, D.C. 20402.
2. The listings show electric (residential) power supplied in each country; as indicated, in very many cases other types of supply also exist to a greater or lesser extent. Therefore, for specific characteristics of the power supply of particular cities, reference should be made to *Electric Power Abroad*. This pamphlet also gives additional details such as number of phases, number of wires to the residence, frequency stability, grounding regulations, and some data on types of commercial service.
3. In the United States in urban areas, the usual supply is 60-hertz 3-phase 120/208 volts; in less densely populated areas, it is usually 120/240 volts, single phase, to each customer. Any other supplies, including dc, are rare and are becoming more so. Additional information for the US is given in the current edition of *Directory of Electric Utilities*, published by McGraw-Hill Book Co., New York, N.Y.
4. All voltages in the table are ac except where specifically stated as dc. The latter are infrequent and in most cases are being replaced by ac. The lower voltages shown for ac, wye or delta ac, or for dc distribution lines, are used mostly for lighting and small appliances; the higher voltages are used for larger appliances.

boards, based in the USA primarily on the National Electrical Code (NEC)* and the National Electric Safety Codes.† Brief extracts from some NEC requirements are given here for convenient reference.

Many products such as wire and cable, fuses, outlet boxes, appliances, etc., are governed within the USA by Underwriters Laboratories (UL) Standards, which specify the terminology used for the various classes of an item as well as the safety requirements that must be met by UL approved items. Note that the overall performance of assemblies such as appliances, motors, radio equipment, or television equipment is not covered by NEC or UL standards, which are primarily for personnel safety.

The following tables are provided.

Tables 5 and 6: NEC standard types of insulated wires and cables.

Tables 7 and 8: Allowable currents for conductors.

Table 8: Derating factors to be applied for ambient temperatures above 30 °C (86 °F) and for more than three conductors in a cable or conduit.

Table 9: Motor starting currents, which determine the overcurrent protection requirements during the starting period.

Table 10: Motor full-load operating currents for usual conditions and speeds.

Guide to Use of Tables

Determine the total equipment load by adding the loads of the various individual items, estimating motor currents according to Table 10 if specific operating-current information is not otherwise available. Any load substantially bigger than this should be interrupted by an overload protective device; normally the next-larger standard fuse or circuit breaker is considered satisfactory.

Determine the total starting load by using the locked-rotor currents computed from Table 9 and the steady-state currents for resistive devices. Make an additional allowance for any large quantities of tungsten lamps, starting transients, and high-inertia loads that will increase the duration of the starting period. The circuit overload protection must be designed to carry this load for the entire starting period. Time-lag fuses or time-delay circuit breakers are usually desirable.

Using the starting currents, determine the voltage drops in the supply circuit; thus be sure that the terminal voltage of the motor or other device will be adequate at start. Increase the size of the supply conductor or reduce the source impedance if necessary. From the starting and running currents, determine the required size of supply conductors.

POWER SUPPLY WIRING

Electric power supply (mains) wiring is usually controlled for public safety by local or state government

* American National Standards Institute, Inc., ANSI Standard C1, prepared by the National Fire Protection Association.

† ANSI Standard C2.

TABLE 5. NEC CONDUCTOR APPLICATIONS AND INSULATIONS

Trade Name	Type Letter	Maximum Operating Temperature	Application	Insulation	Outer Covering
Heat-resistant rubber	RH	75 °C 167 °F	Dry locations	Heat-resistant rubber	Moisture-resistant flame-retardant non-metallic covering
Heat-resistant rubber	RHH	90 °C 194 °F	Dry locations	Heat-resistant rubber	Moisture-resistant flame-retardant non-metallic covering
Moisture- and heat-resistant rubber	RHW	75 °C 167 °F	Dry and wet locations	Moisture- and heat-resistant rubber	Moisture-resistant flame-retardant non-metallic covering
Heat-resistant latex rubber	RUH	75 °C 167 °F	Dry locations	90% unmilled grainless rubber	Moisture-resistant flame-retardant non-metallic covering
Moisture-resistant latex rubber	RUW	60 °C 140 °F	Dry and wet locations	90% unmilled grainless rubber	Moisture-resistant flame-retardant non-metallic covering
Thermoplastic	T	60 °C 140 °F	Dry locations	Flame-retardant thermoplastic compound	None
Moisture-resistant thermoplastic	TW	60 °C 140 °F	Dry and wet locations	Flame-retardant moisture-resistant thermoplastic	None
Moisture- and heat-resistant thermoplastic	THW	75 °C 167 °F	Dry and wet locations	Flame-retardant moisture- and heat-resistant thermoplastic	None
Heat-resistant thermoplastic	THHN	90 °C 194 °F	Dry locations	Flame-retardant heat-resistant thermoplastic	Nylon jacket
Moisture- and heat-resistant thermoplastic	THWN	75 °C 107 °F	Dry and wet locations	Flame-retardant, moisture- and heat-resistant thermoplastic	Nylon jacket
Moisture- and heat-resistant cross-linked synthetic polymer	XHHW	90 °C 194 °F	Dry locations	Flame-retardant cross-linked synthetic polymer	None
Extruded polytetrafluoroethylene	TFE	250 °C 482 °F	Dry locations	Extruded polytetrafluoroethylene	None
Silicone asbestos	SA	90 °C 194 °F	Dry locations	Silicone rubber	Asbestos
Fluorinated Ethylene Polypylene	FEP	90 °C 194 °F	Dry locations	Fluorinated Ethylene Propylene	None
Varnished Cambric	V	85 °C 185 °F	Dry locations	Varnished Cambric	Nonmetallic Covering or lead sheath

WIRING OF ELECTRONIC EQUIPMENT AND CHASSIS

There are few official standards for the internal wiring of electronic equipment and chassis. Nevertheless, the following points should be considered.

- (A) Probable maximum continuous ambient temperature where the wiring is located.
- (B) Allowable temperature rise of conductor surface under full-load conditions (determines minimum wire size).
- (C) Maximum voltage to ground or to surrounding metal parts (determines required insulation thickness).
- (D) Possibility of corona on high-voltage leads; some insulating materials deteriorate rapidly under corona conditions.
- (E) Need for shield braid on some conductors to reduce noise pickup. Shields must be insulated if positive single-point grounding is to be attained.
- (F) Skin effect on conductors carrying high radio-frequency currents.
- (G) Vibration, shock, or relative motion of conductors during normal use of the equipment. Stranded or flexible conductors and adequate clamping or other tie-down of conductors and cables may be essential.

TABLE 6. NEC FLEXIBLE-CORD DATA

Trade Name*	Type Letter†	Size Range (AWG)	No. of Conductors	Insulation	Outer Covering
All-rubber parallel cord	SP-3	18 – 12	2 or 3	Rubber	Rubber
All-plastic parallel cord	SPT-3	18 – 10	2 or 3	Thermoplastic	Thermoplastic
Lamp cord	C	18 – 10	2 or more	Rubber	None
Twisted portable cord	PD	18 – 10	2 or more	Rubber	Cotton or None
Vacuum-cleaner cord	SV	18	2 or 3	Rubber	Rubber
Vacuum-cleaner cord	SVT	18 – 17	2 or 3	Thermoplastic	Thermoplastic
Junior hard-service cord	SJ	18 – 14	2 – 4	Rubber	Rubber
Junior hard-service cord	SJO	18 – 14	2 – 4	Rubber	Oil-resistant compound
Junior hard-service cord	SJT	18 – 14	2 – 4	Rubber or thermoplastic	Thermoplastic
Junior hard-service cord	SJTO	18 – 14	2 – 4	Rubber or thermoplastic	Oil-resistant thermoplastic
Hard-service cord	S	18 – 2	2 or more	Rubber	Rubber
Hard-service cord	SO	18 – 2	2 or more	Rubber	Oil-resistant compound
Hard-service cord	ST	18 – 2	2 or more	Rubber or thermoplastic	Thermoplastic
Hard-service cord	STO	18 – 2	2 or more	Rubber or thermoplastic	Oil-resistant thermoplastic
Rubber-jacketed heat-resistant cord	AFSJ	18 – 16	2 or 3	Impregnated asbestos	Rubber
Rubber-jacketed heat-resistant cord	AFS	18 – 14	2 or 3	Impregnated asbestos	Rubber
Heater cord	HPD	18 – 12	2 – 4	Rubber or thermoplastic with asbestos or all neoprene	Cotton or rayon
Rubber-jacketed heater cord	HSJ	18 – 16	2 – 4	Rubber or thermoplastic and asbestos or all neoprene	Cotton and rubber
Jacketed heater cord	HSJO	18 – 16	2 – 4	Rubber with asbestos or all neoprene	Cotton and oil-resistant compound
Jacketed heater cord	HS	14 – 12	2 – 4	Rubber with asbestos or all neoprene	Cotton and rubber or neoprene
Jacketed heater cord	HSO	14 – 12	2 – 4	Rubber with asbestos or all neoprene	Cotton and oil-resistant compound
Parallel heater cord	HPN	18 – 12	2 or 3	Thermosetting	Thermosetting

*All types shown are recommended for use in damp locations.
 †“S” series cords may also be used in pendant applications.

TABLE 7. NEC CURRENT-CARRYING CAPACITY, IN AMPERES, OF FLEXIBLE CORDS

Size AWG	Rubber TP, TS Thermoplastic TPT, TSP	Rubber C, PD, E, EO, EN, S, SO, SRD, SS, SSO, SV, SVO, SP		AFS, AFSJ, HPD, HSJ, HSJO, HS HSO, HPN	Cotton* CFPD. Asbestos* AFC, AFPD
		Thermoplastic ET, ETT, ETLB, ETP, ST, STO, SRDT, SVT, SVTO, SPT			
27†	0.5	—	—	—	—
18		7‡	10§	10	6
16		10‡	13§	15	8
14		15‡	18§	20	17
12		20‡	25§	30	23
10		25‡	30§	35	28
8		35‡	40§	—	—
6		45‡	55§	—	—
4		60‡	70§	—	—
2		80‡	95§	—	—

* Generally used in fixtures exposed to high temperatures, derated accordingly.

† Tinsel.

‡ Three-conductor and other multiconductor cords connected so only three conductors are current-carrying.

§ Two-conductor and other multiconductor cords connected so only two conductors are current-carrying.

Notes:

1. For not more than three current-carrying conductors in a cord. If four to six conductors are used, allowable capacity of each conductor shall be reduced to 80% of values for not more than three current-carrying conductors.

2. A conductor used for equipment grounding and a neutral conductor which carries only the unbalanced current from other conductors shall not be considered as current-carrying conductors.

3. Based on room temperature of 30 °C (86 °F).

TABLE 8. NEC ALLOWABLE CURRENT-CARRYING CAPACITIES, IN AMPERES, OF CONDUCTORS

Size AWG	Copper-Conductor Insulation				Aluminum or Copper-Clad Aluminum Conductor Insulation			
	R UW (14 - 2) T, TW, UF	RH, RHW, RUH, (14 - 2) THWN, XHHW	TA, TBS, SA, FEP, FEPB, RHH, THHN, XHHW*	T FE†	R UW (12 - 2) T, TW, UF	RH, RHW, RUH (12 - 2), THWN, XHHW	TA, TBS, SA, RHH, THHN, XHHW*	
14	15	15	25‡	40	—	—	—	
12	20	20	30‡	55	15	15	25§	
10	30	30	40‡	75	25	25	30§	
8	40	45	50	95	30	40	40	
6	55	65	70	120	40	50	55	
4	70	85	90	145	55	65	70	
3	80	100	105	170	65	75	80	
2	95	115	120	195	75	90	95	
1	110	130	140	220	85	100	110	
0	125	150	155	250	100	120	125	
00	145	175	185	280	115	135	145	
000	165	200	210	315	130	155	165	
0000	195	230	235	370	155	180	185	

Correction Factors for Higher Room Temperatures									
°C	°F								
40	104	0.82	0.88	0.91	—	0.82	0.88	0.91	
45	113	0.71	0.82	0.87	—	0.71	0.82	0.87	
50	122	0.58	0.75	0.82	—	0.58	0.75	0.82	
55	131	0.41	0.67	0.76	—	0.41	0.67	0.76	
60	140	—	0.58	0.71	.95	—	0.58	0.71	

* Dry locations only.
 † Nickel or nickel-coated copper only.
 ‡ For types FEP, FEPB, RHH, THHN, and XHHW, sizes 14, 12, 10 shall be the same as designated for RH, RHW, etc.
 § For types RHH, THHN, and XHHW, sizes 12 and 10 shall be the same as designated for RH, RHW, etc.

Notes:

1. Not more than three conductors in raceway or cable.
2. Based on room temperature of 30 °C (86 °F). See correction factors for higher temperatures.
3. Derating factors—more than three conductors in raceway or cable:

Number of conductors	4-6	7-24	25-42	>42
% of current capacity	80	70	60	50

- (H) Wiring should be shielded from the direct heat radiation of high-temperature parts such as electron tubes and power resistors.
- (I) Wire identification may be required for convenience in manufacture, installation, or servicing.

As a matter of expediency, most electronic equipment employs the smallest conveniently handled wire size (usually 20 to 24 AWG) for most wiring, with the larger conductors being installed only for circuits carrying currents greater than that permitted for the "general-use" wire size. With the trend toward compact

solid-state integrated-circuit equipment, smaller wire sizes are being used. However, the reduction in wiring bundle size may be small unless the conductor insulation thickness (determined by voltage considerations) can be reduced.

Table 11 gives recommended current ratings for copper and some aluminum based on a 45 °C (40 °C for wires smaller than 22 AWG) conductor temperature rise due to load current. Table 11 may be used for the temperature conditions given in Table 12.

A 60 °C ambient temperature around the wiring (20 °C internal temperature rise from 40 °C [104 °F])

TABLE 9. NEC MOTOR STARTING-CURRENT DATA*

Code Letter	Kilovolt-Amperes per Horsepower with Locked Rotor
A	0 - 3.14
B	3.15 - 3.54
C	3.55 - 3.99
D	4.00 - 4.49
E	4.50 - 4.99
F	5.00 - 5.59
G	5.60 - 6.29
H	6.30 - 7.09
J	7.10 - 7.99
K	8.00 - 8.99
L	9.00 - 9.99
M	10.00 - 11.19
N	11.20 - 12.49
P	12.50 - 13.99
R	14.00 - 15.99
S	16.00 - 17.99
T	18.00 - 19.99
U	20.00 - 22.39
V	22.40 up

* Locked-rotor currents of motors are useful in determining branch-circuit overcurrent protection requirements and voltage drop at start. These values are indicated by code letters on motor nameplate. Note NEMA standard.

ambient around the equipment) is typical of some electronic equipment. If higher ambient temperatures, high power, or compact designs with electron tubes or

magnetic-core components (except for very low power) are a factor, the temperature in the wiring space should be specifically determined.

“Wiring confined” ratings are based on 15 or more wires in a bundle, with the sum of all the actual load currents of the bundled wires not exceeding 20% of the permitted “wiring confined” sum total carrying capacity of the bundled wires. These ratings approximate 60% of the free-air ratings (with some variations due to rounding). They should be used for wire in harnesses, cable, conduit, and general chassis conditions. Bundles of fewer than 15 wires may have the allowable sum of the load currents increased as the bundle approaches the single-wire condition.

RESISTANCE CHANGE WITH TEMPERATURE

The resistance of most conductor materials changes with temperature. Table 11 shows the copper-wire resistance at 100 °C. Correction factors must be applied to determine the resistance at other temperatures, or for other materials. Thus, from Table 11 determine the copper-wire resistance at 100 °C (multiply by conversion factor *m* of Table 13 for other materials). Use the equation

$$R_t = R_r m [1 + K(t - 100)]$$

where,

- R_t = resistance at desired temperature *t*,
- R_r = resistance at 100 °C for copper (Table 11),

TABLE 10. NEC MOTOR FULL-LOAD RUNNING CURRENTS IN AMPERES (USUAL CONDITIONS AND SPEEDS)

Horsepower	Single-Phase AC		3-Phase AC†			DC	
	115 V	230 V*	115 V	230 V*	460 V	120 V	240 V
1/6	4.4	2.2	—	—	—	—	—
1/4	5.8	2.9	—	—	—	3.1	1.6
1/3	7.2	3.6	—	—	—	4.1	2.0
1/2	9.8	4.9	4	2	1	5.4	2.7
3/4	13.8	6.9	5.6	2.8	1.4	7.6	3.8
1	16	8	7.2	3.6	1.8	9.5	4.7
1 1/2	20	10	10.4	5.2	2.6	13.2	6.6
2	24	12	13.6	6.8	3.4	17	8.5
3	34	17	—	9.6	4.8	25	12.2
5	56	28	—	15.2	7.6	40	20
7 1/2	80	40	—	22	11	58	29
10	100	50	—	28	14	76	38

* For 208 V, multiply by 1.1; for 200 V, multiply by 1.15.
 † Induction type, squirrel cage and wound rotor.

TABLE 11. RECOMMENDED CURRENT RATINGS (CONTINUOUS DUTY) FOR ELECTRONIC EQUIPMENT AND CHASSIS WIRING*

Wire Size		Copper Conductor (100 °C) Nominal Resistance (Ohms/1000 ft)	Maximum Current in Amperes			
AWG	Circular Mils		Copper Wire		Aluminum Wire	
			Wiring in Free Air	Wiring Confined	Wiring in Free Air	Wiring Confined
32	63.2	188.0	0.53	0.32		
30	100.5	116.0	0.86	0.52		
28	159.8	72.0	1.4	0.83		
26	254.1	45.2	2.2	1.3		
24	404.0	28.4	3.5	2.1		
22	642.4	22.0	7.0	5.0		
20	1022	13.7	11.0	7.5		
18	1624	6.50	16	10		
16	2583	5.15	22	13		
14	4107	3.20	32	17		
12	6530	2.02	41	23		
10	10 380	1.31	55	33		
8	16 510	0.734	73	46	60	36
6	26 250	0.459	101	60	83	50
4	41 740	0.290	135	80	108	66
2	66 370	0.185	181	100	152	82
1	83 690	0.151	211	125	174	105
0	105 500	0.117	245	150	202	123
00	133 100	0.092	283	175	235	145
000	167 800	0.074	328	200	266	162
0000	211 600	0.059	380	225	303	190

* See Table 12.

 m = material factor for 100 °C resistance value (Table 13), K = correction factor (Table 13), t = desired temperature (°C).

WIRE IDENTIFICATION

In a complex wiring assembly, or if both ends of a wire cannot be seen from one station, a means of

identifying each lead simplifies manufacture, installation, and servicing. Common identification methods are:

(A) Tag each end of a lead with an assigned designation (an alternative method is to print the designation at frequent intervals along the wire insulation).

(B) Color code the wires. The wire insulation may be a solid color, color stripes may be spiraled around the wire, or the name of the color (or its numerical code equivalent) may be stamped at frequent intervals along the wire.

TABLE 12. TEMPERATURE CONDITIONS FOR TABLE 11

Maximum Allowable Conductor Temperature °C	Maximum Ambient Temperature Around Wire °C	Typical Conductor and Insulation
105	60	Bare or tinned copper or aluminum; polyvinyl-chloride insulation
200	155	Silver-coated copper; FEP or PTFE with FEP jacket insulation*
260	215	Nickel-coated copper; PTFE insulation*

* FEP = Fluorinated ethylene propylene.
PTFE = Polytetrafluoroethylene.

TABLE 13. TEMPERATURE AND MATERIALS CORRECTION FACTORS

Conductor Material	Material Factor m	Correction Factor K
Soft copper	1.00	0.0039
Hard copper	1.03	0.0038
Copper-clad steel:		
30% conductivity	3.47	0.0044
40% conductivity	2.56	0.0041
Aluminum	1.64	0.0039
Nickel	5.28	0.0047
Nickel-clad copper:		
10%	1.07	0.0038
30%	1.35	0.0036
Silver	0.94	0.0038

Color Coding

The commonly used colors and their numerical codes are:

0 Black	5 Green
1 Brown	6 Blue
2 Red	7 Violet (purple)
3 Orange	8 Gray (slate)
4 Yellow	9 White

While spiral stripes can be applied on top of any basic insulation color, under less favorable viewing conditions it is difficult to distinguish some colors from the basic insulation color. Identification may be slow and subject to error. The preferred combination consists of one or two (sometimes three) colored stripes on a white basic insulation. To minimize identification errors, the first stripe is made wider than the second (or third), and some rules require that the second stripe be of higher numerical code than the first stripe. If the required variety of wire color codes is not great, the preceding guides should be followed.

Table 14 gives a standard color code used to distinguish by function the various leads in electronic circuits.

In manufacturing practice, it is preferred that, at any harness breakout point, all wires of the same color code be connected to the same terminal at that location. When this rule and the wire color coding of Table 14 are both applicable, additional tracers may be used to supplement the primary coding of Table 14.

DIAMETER OF CIRCLE ENCLOSING A GIVEN NUMBER OF SMALLER CIRCLES*

Four of many possible compact arrangements of circles within a circle are shown in Fig. 3. To determine the diameter of the smallest enclosing circle for a particular number of enclosed circles all of the same size, three factors that influence the size of the enclosing circle should be considered, as follows.

Arrangement of Center or Core Circles

The four most common arrangements of center or core circles are shown in cross section in Fig. 3. It may seem that Fig. 3A would require the smallest enclosing circle for a given number of enclosed circles, but this is not always the case since the most compact arrangement will depend in part on the number of circles to be enclosed.

TABLE 14. COLORS FOR WIRE IDENTIFICATION BY FUNCTION

Function	Color	Identification No.
Grounds, grounded elements	Black	0
Heaters or filaments	Brown	1
Power supply B+	Red	2
Screen grids	Orange	3
Cathodes and transistor emitters*†	Yellow	4
Control grids and transistor bases†	Green	5
Anodes (plates) and transistor collectors*†	Blue	6
Power supply, negative (-)	Violet (purple)	7
Ac power lines	Gray (slate)	8

* Applies to diodes, semiconductor elements, photoelectric cells, mercury-arc rectifiers, and other elements with operation similar to vacuum tubes and transistors.

† Applies to all types of gas tubes with operation similar to vacuum tubes.

Diameter of Enclosing Circle When Outer Layer of Circles Is Complete

Successive, complete "layers" of circles may be placed around each of the central cores of 1, 2, 3, or 4 circles. The number of circles contained in arrangements of complete "layers" around a central core of circles, as well as the diameter of the enclosing circle, may be obtained from Table 15. Thus, for example, Fig. 3A has a total of 18 circles arranged in two complete "layers" around a central core consisting of one circle; this agrees with the data shown in the left half of Table 15 for $n = 2$.

To determine the diameter of the enclosing circle, the data in the right half of Table 15 are used. Thus, for $n = 2$ and an "A" pattern, diameter D is 5 times the diameter, d , of the enclosed circles.

Diameter of Enclosing Circle When Outer Layer of Circles Is Not Complete

In most cases, it is possible to reduce the size of the enclosing circle from that required if the outer layer

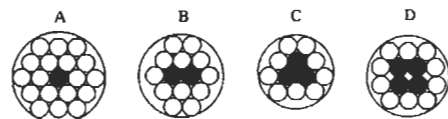


Fig. 3. Arrangements of circles within a circle. (Reprinted with permission from Machinery's Handbook, 17th Edition, Industrial Press, Inc., New York.)

* Dutka, J. "How Many Wires Can Be Packed Into a Circular Conduit." *Machinery's Handbook*. New York: Industrial Press, Inc., 1956.

TABLE 15. NUMBER OF CIRCLES CONTAINED IN COMPLETE LAYERS OF CIRCLES AND DIAMETER OF ENCLOSING CIRCLE

No. Complete Layers Over Core, <i>n</i>	Number of Circles in Center Pattern							
	1	2	3	4	1	2	3	4
	Arrangement of Circles in Center Pattern (see Fig. 3)							
	"A"	"B"	"C"	"D"	"A"	"B"	"C"	"D"
	Number of Circles, <i>N</i> , Enclosed				Diameter, <i>D</i> , of Enclosing Circle*			
0	1	2	3	4	<i>d</i>	2 <i>d</i>	2.155 <i>d</i>	2.414 <i>d</i>
1	7	10	12	14	3 <i>d</i>	4 <i>d</i>	4.055 <i>d</i>	4.386 <i>d</i>
2	19	24	27	30	5 <i>d</i>	6 <i>d</i>	6.033 <i>d</i>	6.379 <i>d</i>
3	37	44	48	52	7 <i>d</i>	8 <i>d</i>	8.024 <i>d</i>	8.375 <i>d</i>
4	61	70	75	80	9 <i>d</i>	10 <i>d</i>	10.018 <i>d</i>	10.373 <i>d</i>
5	91	102	108	114	11 <i>d</i>	12 <i>d</i>	12.015 <i>d</i>	12.372 <i>d</i>
<i>n</i>	†	†	†	†	†	†	†	†

* Diameter *D* is given in terms of *d*, the diameter of the enclosed circles.

† For *n* complete layers over core, the number of enclosed circles *N* for "A" center pattern is $3n^2 + 3n + 1$; for "B," $3n^2 + 5n + 2$; for "C," $3n^2 + 6n + 3$; for "D," $3n^2 + 7n + 4$; while the diameter *D* of the enclosing circle for "A" center pattern is $(2n + 1)d$; for "B," $(2n + 2)d$; for "C," $[1 + 2(n^2 + n + \frac{1}{2})^{1/2}]d$; and for "D," $[1 + (4n^2 + 5.644n + 2)^{1/2}]d$.

Reprinted with permission from *Machinery's Handbook*, 17th edition. New York: Industrial Press, Inc.

were complete. Thus, for example, Fig. 3B shows that the central core consisting of two circles is surrounded by one complete layer of eight circles and one partial outer layer of four circles so that the total number of circles enclosed is 14. If the outer layer were complete, then (from Table 15) the total number of enclosed circles would be 24 and the diameter of the enclosing circle would be 6*d*; however, since the outer layer is composed of only four circles out of a possible 14 for a complete second layer, a smaller diameter of enclosing circle may be used. Table 16 shows that for a total of 14 enclosed circles arranged in a "B" pattern with the outer layer of circles incomplete, the diameter for the enclosing circle is 4.606*d*.

Table 16 can be used to determine the smallest enclosing circle for a given number of circles to be enclosed by direct comparison of the "A," "B," and "C" columns. For data outside the range of Table 16, use the equations in Dr. Dutka's article.*

Approximate Equation When Number of Enclosed Circles is Large

When a large number of circles is to be enclosed, the arrangement of the center circles has little effect on the diameter of the enclosing circle. For numbers of circles

greater than 10 000, the diameter of the enclosing circle may be calculated within 2 percent from the equation

$$D = d[1 + (N/0.907)^{1/2}]$$

where,

- D* = diameter of enclosing circle,
- d* = diameter of enclosed circles,
- N* = number of enclosed circles.

TORQUE AND HORSEPOWER

Torque varies directly with power and inversely with rotating speed of the shaft, or

$$T = KP/N$$

where,

- T* = torque in pound-inches,
- P* = horsepower,
- N* = revolutions/minute,
- K* = 63 000 (constant).

TRANSMISSION-LINE SAG CALCULATIONS*

For transmission-line work, with towers on the same or slightly different levels, the cables are assumed to

* Dutka, J. "How Many Wires Can Be Packed Into a Circular Conduit." *Machinery's Handbook*. New York: Industrial Press, Inc., 1956.

* Reprinted by permission from *Transmission Towers*. Pittsburgh, Pa.: American Bridge Co., 1923; p. 70.

TABLE 16. FACTORS FOR DETERMINING DIAMETER, *D*, OF SMALLEST ENCLOSING CIRCLE FOR VARIOUS NUMBERS, *N*, OF ENCLOSED CIRCLES*

No. <i>N</i>	Center Circle Pattern			No. <i>N</i>	Center Circle Pattern			No. <i>N</i>	Center Circle Pattern		
	"A"	"B"	"C"		"A"	"B"	"C"		"A"	"B"	"C"
	Diameter Factor <i>K</i>				Diameter Factor <i>K</i>				Diameter Factor <i>K</i>		
2	3	2	—	34	7	7.083	7.110	66	9.718	9.544	9.326
3	3	2.732	2.155	35	7	7.245	7.110	67	9.718	9.544	9.326
4	3	2.732	3.309	36	7	7.245	7.110	68	9.718	9.544	9.326
5	3	3.646	3.309	37	7	7.245	7.429	69	9.718	9.660	9.326
6	3	3.646	3.309	38	7.928	7.245	7.429	70	9.718	9.660	10.018
7	3	3.646	4.055	39	7.928	7.557	7.429	71	9.718	9.888	10.018
8	4.464	3.646	4.055	40	7.928	7.557	7.429	72	9.718	9.888	10.018
9	4.464	4	4.055	41	7.928	7.557	7.429	73	9.718	9.888	10.018
10	4.464	4	4.055	42	7.928	7.557	7.429	74	10.165	9.888	10.018
11	4.464	4.606	4.055	43	7.928	8	8.024	75	10.165	10	10.018
12	4.464	4.606	4.055	44	8.211	8	8.024	76	10.165	10	10.238
13	4.464	4.606	5.163	45	8.211	8	8.024	77	10.165	10.539	10.238
14	5	4.606	5.163	46	8.211	8	8.024	78	10.165	10.539	10.238
15	5	5.359	5.163	47	8.211	8	8.024	79	10.165	10.539	10.452
16	5	5.359	5.163	48	8.211	8	8.024	80	10.165	10.539	10.452
17	5	5.359	5.163	49	8.211	8.550	8.572	81	10.165	10.539	10.452
18	5	5.359	5.163	50	8.211	8.550	8.572	82	10.165	10.539	10.452
19	5	5.583	5.619	51	8.211	8.550	8.572	83	10.165	10.539	10.452
20	6.292	5.583	5.619	52	8.211	8.550	8.572	84	10.165	10.539	10.452
21	6.292	5.583	5.619	53	8.211	8.810	8.572	85	10.165	10.644	10.866
22	6.292	5.583	6.033	54	8.211	8.810	8.572	86	11	10.644	10.866
23	6.292	6	6.033	55	8.211	8.810	9.083	87	11	10.644	10.866
24	6.292	6	6.033	56	9	8.810	9.083	88	11	10.644	10.866
25	6.292	6.196	6.033	57	9	8.937	9.083	89	11	10.849	10.866
26	6.292	6.196	6.033	58	9	8.937	9.083	90	11	10.849	10.866
27	6.292	6.568	6.033	59	9	8.937	9.083	91	11	10.849	11.214
28	6.292	6.568	6.773	60	9	8.937	9.083	92	11.392	10.849	11.214
29	6.292	6.568	6.773	61	9	9.185	9.083	93	11.392	11.149	11.214
30	6.292	6.568	6.773	62	9.718	9.185	9.083	94	11.392	11.149	11.214
31	6.292	7.083	7.110	63	9.718	9.185	9.083	95	11.392	11.149	11.214
32	7	7.083	7.110	64	9.718	9.185	9.326	96	11.392	11.149	11.214
33	7	7.083	7.110	65	9.718	9.544	9.326	97	11.392	11.440	11.214

* The diameter, *D*, of the enclosing circle is equal to the diameter factor, *K*, multiplied by *d*, the diameter of the enclosed circles, or $D = K \times d$. For example, if the number of circles to be enclosed, *N*, is 12, and the center circle arrangement is "C," then for $d = 1\frac{1}{2}$ inches, $D = 4.005 \times 1\frac{1}{2} = 6.008$ inches. Reprinted with permission from *Machinery's Handbook*, 17th Edition. New York: Industrial Press, Inc.

take the form of a parabola instead of their actual form of a catenary. The error is negligible, and the computations are much simplified. In calculating sags, the changes in cables due to variations in load and temperature must be considered.

Supports at Same Level

Refer to Fig. 4. The equations used in calculating sags are

$$H = WL^2/8S$$

$$S = WL^2/8H = [(L_c - L)3L/8]^{1/2}$$

$$L_c = L + 8S^2/3L$$

where,

L = length of span in feet,

L_c = length of cable in feet,

S = sag of cable at center of span in feet,

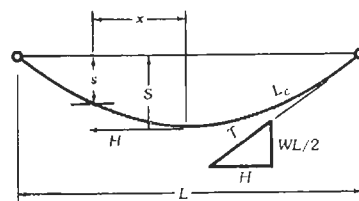


Fig. 4. Supports at same elevation.

H = tension in cable at center of span in pounds
 = horizontal component of the tension at any point,
 W = weight of cable in pounds per lineal foot.

If cables are subject to wind and ice loads, W = the algebraic sum of the loads. That is, for ice on cables, W = weight of cables plus weight of ice; for wind on bare or ice-covered cables, W = the square root of the sum of the squares of the vertical and horizontal loads.

For any intermediate point at a distance x from the center of the span, the sag is

$$S_x = S(1 - 4x^2/L^2)$$

Supports at Different Levels

Refer to Fig. 5.

$$S = S_0 = WL_0^2 \cos\alpha / 8T$$

$$= WL^2 / 8T \cos\alpha$$

$$S_1 = WL_1^2 / 8H$$

$$S_2 = WL_2^2 / 8H$$

$$L_1/2 = L/2 - (hH \cos\alpha / WL)$$

$$L_2/2 = L/2 + (hH \cos\alpha / WL)$$

$$L_c = L + \frac{4}{3}[(S_1^2/L_1) + (S_2^2/L_2)]$$

where,

W = weight of cable in pounds per lineal foot between supports or in direction of L_0 ,
 T = tension in cable direction parallel with line between supports.

The change, l , in length of cable, L_c , for varying temperature is found by multiplying the number of degrees, n , by the length of the cable in feet times the coefficient of linear expansion per foot per degree Fahrenheit, c . This is

$$l = L_c \times n \times c$$

A short approximate method for determining sags under varying temperatures and loadings that is close enough for all ordinary line work is as follows:

(A) Determine sag of cable with maximum stress under maximum load of lowest temperature occurring at the time of maximum load, and find length of cable with this sag.

(B) Find length of cable at the temperature for which the sag is required.

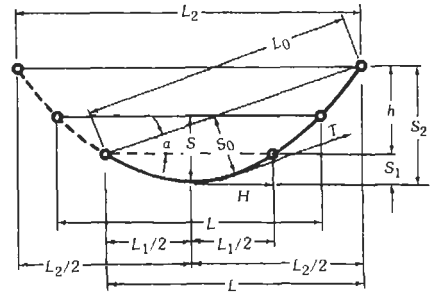


Fig. 5. Supports at different elevations.

(C) Assume a certain reduced tension in the cable at the temperature and under the loading combination for which the sag is required; then find the decrease in length of the cable due to the decrease of the stress from its maximum.

(D) Combine the algebraic sum of (B) and (C) with (A) to get the length of the cable under the desired conditions; from this length, the sag and tension can be determined.

(E) If this tension agrees with that assumed in (C), the sag in (D) is correct. If it does not agree, another assumption of tension in (C) must be made and the process repeated until (C) and (D) agree.

STRUCTURAL STANDARDS FOR STEEL RADIO TOWERS*

Material

(A) Structural steel shall conform to American Society for Testing Materials *Standard Specifications for Steel for Bridges and Buildings*, Serial Designation A-7, as amended to date.

(B) Steel pipe shall conform to American Society for Testing Materials standard specifications either for electric-resistance welded steel pipe, Grade A or Grade B, Serial Designation A-135, or for welded and seamless steel pipe, Grade A or Grade B, Serial Designation A-53, each as amended to date.

Loading

(A) 20-pound design: Structures up to 600 feet in height (unless they are to be located within city limits) shall be designed for a horizontal wind pressure of 20

* Abstracted from *American Standard Minimum Design Loads in Buildings and Other Structures*, A58.1-1955. New York: American National Standards Institute, Inc. Also from Electronic Industries Association Standard TR-116, October 1949. Sections on manufacture and workmanship, finish, and plans and marking of the standard are not reproduced here. The section "Wind Velocities and Pressures" is not part of the standard.

pounds/foot² on flat surfaces and 13.3 pounds/foot² on cylindrical surfaces.

(B) 30-pound design: Structures more than 600 feet in height and those of any height to be located within city limits shall be designed for a horizontal wind pressure of 30 pounds/foot² on flat surfaces and 20 pounds/foot² on cylindrical surfaces.

(C) Other designs: Certain structures may be designed to resist loads greater than those described in (A) and (B). Figure 1 of American Standard A58.1-1955 shows sections of the United States where greater wind pressures may occur. In all such cases, the pressure on cylindrical surfaces shall be computed as being $\frac{2}{3}$ of that specified for flat surfaces.

(D) For open-face (latticed) structures of square cross section, the wind pressure normal to one face shall be applied to 2.20 times the normal projected area of all members in one face, or 2.40 times the normal projected area of one face for wind applied to one corner. For open-face (latticed) structures of triangular cross section, the wind pressure normal to one face shall be applied to 2.00 times the normal projected area of all members in one face, or 1.50 times the normal projected area for wind parallel to one face. For closed-face (solid) structures, the wind pressure shall be applied to 1.00 times the normal projected area for square or rectangular shape, 0.80 for hexagonal or octagonal shape, and 0.60 for round or elliptical shape.

(E) Provisions shall be made for all supplementary loadings caused by the attachment of guys, antennas, transmission and power lines, ladders, etc. The pressure shall be as described for the respective designs and shall be applied to the projected area of the construction.

(F) The total load specified above shall be applied to the structure in the directions that will cause the maximum stress in the various members.

(G) The dead weight of the structure, and all materials attached thereto, shall be included.

R = corresponding radius of gyration, both in inches.

Maximum L/R for main leg members = 140

Maximum L/R for other compression members with calculated stress = 200

Maximum L/R for members with no calculated stress = 250

Bending on extreme fibres = 20 000 pounds/inch²

Single shear on bolts = 13 500 pounds/inch²

Double shear on bolts = 27 000 pounds/inch²

Bearing on bolts (single shear) = 30 000 pounds/inch²

Bearing on bolts (double shear) = 30 000 pounds/inch²

Tension on bolts and other threaded parts, on nominal area at root of thread = 16 000 pounds/inch²

Members subject to both axial and bending stresses shall be so designed that the calculated unit axial stress divided by the allowable unit axial stress, plus the calculated unit bending stress, divided by the allowable unit bending stress, shall not exceed unity.

(B) Minimum thickness of material for structural members:

Pinted structural angles and plates = $\frac{3}{16}$ inch

Hot-dip galvanized structural angles and plates = $\frac{1}{8}$ inch

Other structural members to mill minimum for standard shapes.

(C) Where materials of higher quality than specified under "Material" above are used, the above unit stresses may be modified. The modified unit stresses must provide the same factor of safety based on the yield point of the materials.

Unit Stresses

(A) All parts of the structure shall be so designed that the unit stresses resulting from the specified loads shall not exceed the following values:

Axial tension on net section = 20 000 pounds/inch²

Axial compression on gross section:

For members with values of L/R not greater than 120,

$$= 17\,000 - 0.485L^2/R^2 \text{ pounds/inch}^2$$

For members with values of L/R greater than 120,

$$= \frac{18\,000}{1 + L^2/18\,000R^2} \text{ pounds/inch}^2$$

where L = unbraced length of the member, and

Foundations

(A) Standard foundations shall be designed for a soil pressure not to exceed 4000 pounds/foot² under the specified loading. In uplift, the foundations shall be designed to resist 100 percent more than the specified loading, assuming that the base of the pier will engage the frustum of an inverted pyramid of earth whose sides form an angle of 30 degrees with the vertical. Earth shall be considered to weigh 100 pounds/foot³ and concrete 140 pounds/foot³.

(B) Foundation plans shall ordinarily show standard foundations as defined in (A). Where the actual soil conditions are not normal, requiring some modification in the standard design, and complete soil information is provided to the manufacturer by the purchaser, the foundation plan shall show the required design.

(C) Under conditions requiring special engineering such as pile construction, roof installations, etc., the manufacturer shall provide the necessary information so

that proper foundations can be designed by the purchaser's engineer or architect.

(D) In the design of guy anchors subject to submer- sion, the upward pressure of the water should be taken into account.

Wind Velocities and Pressures

Data regarding wind velocities and pressures are contained in Table 17.

VIBRATION AND SHOCK ISOLATION

Symbols

- b = damping factor
- d = static deflection in inches
- E = relative transmissibility
= (force transmitted by isolators)/(force transmitted by rigid mountings)
- F = force in pounds
- F_0 = peak force in pounds
- f = frequency in hertz
- f_0 = resonant frequency of system in hertz
- G = acceleration of gravity
≈ 386 inches per second²
- g = peak acceleration in dimensionless gravitational units

- = \ddot{X}_0/G
- j = $(-1)^{1/2}$, vector operator
- k = stiffness constant; force required to compress or extend isolators unit distance in pounds per inch
- r = coefficient of viscous damping in pounds per inch per second
- t = time in seconds
- W = weight in pounds
- x = displacement from equilibrium position in inches
- X_0 = peak displacement in inches
- \dot{x} = velocity in inches per second
= dx/dt
- \dot{X}_0 = peak velocity in inches per second
- \ddot{x} = acceleration in inches per second²
= d^2x/dt^2
- \ddot{X}_0 = peak acceleration in inches per second²
- ϕ = phase angle in radians
- ω = angular velocity in radians per second
= $2\pi f$

Equations

The following relations apply to simple harmonic motion in systems with one degree of freedom. Although actual vibration is usually more complex, the equations provide useful approximations for practical purposes.

TABLE 17. WIND VELOCITIES AND PRESSURES

Actual Velocity V_a^* (miles/hour)	Indicated Velocity V_i (miles/hour)		Pressure P (pounds/foot ²) Projected Areas [†]	
	3-cup Anemometer	4-cup Anemometer	Cylindrical Surfaces ($P = 0.0025V_a^2$)	Flat Surfaces ($P = 0.0042V_a^2$)
10	9	10	0.25	0.42
20	20	23	1.0	1.7
30	31	36	2.3	3.8
40	42	50	4.0	6.7
50	54	64	6.3	10.5
60	65	77	9.0	15.1
70	76	91	12.3	20.6
80	88	105	16.0	26.8
90	99	119	20.3	34.0
100	110	132	25.0	42.0
110	121	146	30.3	50.8
120	133	160	36.0	60.5
130	144	173	42.3	71.0
140	155	187	49.0	82.3
150	167	201	56.3	94.5

* Although wind velocities are measured with cup anemometers, all data published by the US Weather Bureau since January 1932 include instrumental corrections and are actual velocities. Prior to 1932, indicated velocities were published.

In calculating pressures on structures, the "fastest-single-mile velocities" published by the Weather Bureau should be multiplied by a gust factor of 1.3 to obtain the maximum instantaneous actual velocities. See "Wind-Velocity and Temperature Extremes in North America" in this chapter for fastest-single-mile records at various places in the United States and Canada.

† The American Bridge Company equations given here are based on a ratio of 25/42 for pressures on cylindrical and flat surfaces, respectively, while the Electronic Industries Association specifies a ratio of 2/3. The actual ratio varies in a complex manner with Reynolds number, shape, and size of the exposed object.

$$F = W(\ddot{x}/G) \tag{Eq. 1}$$

$$F_0 = Wg \tag{Eq. 2}$$

$$x = X_0 \sin(\omega t + \phi) \tag{Eq. 3}$$

$$X_0 = 9.77 g/f^2 \tag{Eq. 4}$$

$$\dot{X}_0 = \omega X_0 = 6.28fX_0 = 61.4g/f \tag{Eq. 5}$$

$$\ddot{X}_0 = \omega^2 X_0 = 39.5f^2 X_0 = 386g \tag{Eq. 6}$$

$$E = \left| \frac{r - j(k/\omega)}{r + j[(\omega W/G) - k/\omega]} \right| \tag{Eq. 7}$$

$$f_0 = 3.13 (k/W)^{1/2} \tag{Eq. 8}$$

$$b = 9.77r/(kW)^{1/2} \tag{Eq. 9}$$

$$E = \left| \frac{1}{(f/f_0)^2 - 1} \right| \tag{Eq. 10}$$

When damping is neglected

$$k = W/d \tag{Eq. 11}$$

$$f_0 = 3.13/d^{1/2} \tag{Eq. 12}$$

$$E = 9.77/(df^2 - 9.77) \tag{Eq. 13}$$

Acceleration

The intensity of vibratory forces is often defined in terms of *g* values. From Eq. 2, it is apparent, for example, that a peak acceleration of 10*g* on a body will result in a reactionary force by the body equal to 10 times its weight.

If an object is mounted on vibration isolators, the accelerations of the vehicle are transmitted to the object (or vice versa) in an amplitude and phase that depend on the elastic flexing of the isolators in the direction in which the accelerations (dynamic forces) are applied.

For critical damping, $b = 1$.
Neglecting dissipation ($b = 0$), or at $ff_0 = (2)^{1/2}$ for any degree of damping

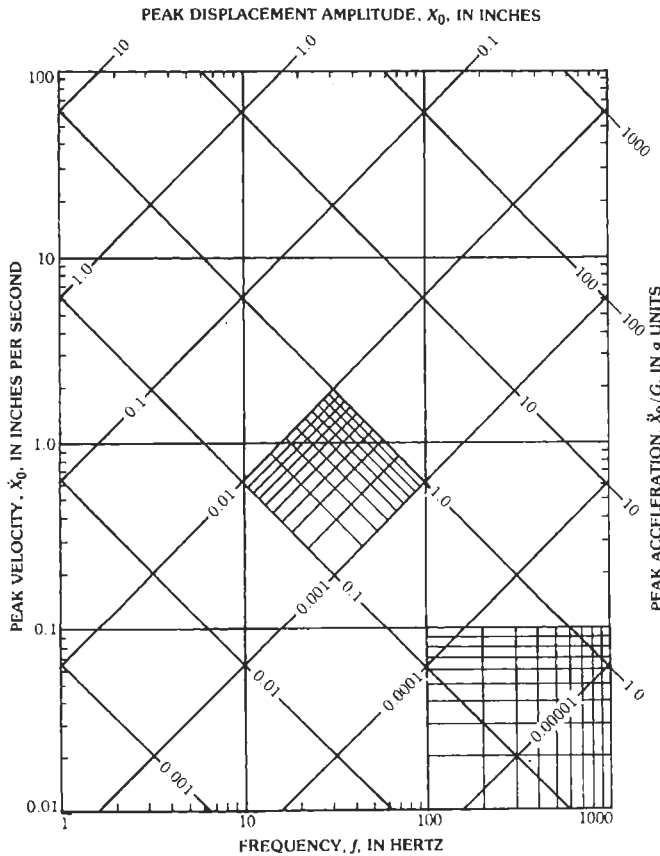


Fig. 6. Relation of frequency and peak values of velocity, displacement, and acceleration.

Magnitudes

The relations among X_0 , \dot{X}_0 , \ddot{X}_0 , and f are shown in Fig. 6. Any two of these parameters applied to the graph locate the other two. For example, suppose $f = 10$ hertz and peak displacement $X_0 = 1$ inch. From Fig. 6, peak velocity $\dot{X}_0 = 63$ inches per second and peak acceleration $\ddot{X}_0 = 10g$.

Natural Frequency

If damping is neglected, the natural frequency, f_0 , of vibration of an isolated system in the vertical direction can be calculated from Eq. 12 from the static deflection of the mounts. For example, suppose an object at rest causes a 0.25-inch deflection of its supporting springs. Then

$$f_0 = 3.13/(0.25)^{1/2} = 6.3 \text{ hertz}$$

Resonance

In Fig. 7, E is plotted against f/f_0 for various damping factors. Note that resonance occurs when $f_0 \approx f$ and that the vibratory forces are then increased by the isolators. To reduce vibration, f_0 must be less than $0.7f$, and it should be as small as $0.3f$ for good isolation.

It is not possible to secure good isolation at all vibrational frequencies in vehicles and similar environments where several different and varying exciting

frequencies are present and where the isolators may have to withstand shock as well as vibration. In such cases, f_0 is often selected as about 1.5 to 2 times the predominant f . Information regarding vibration in typical vehicles is shown in Table 18.

Although all supporting structures have compliance and may reduce the effects of vibration and shock, the apparent stiffness of many "rigid" mountings is merely a matter of degree, and in conjunction with the supported mass, they can also give rise to resonance effects, thus magnifying the amplitude of certain vibrations.

Damping

Damping is desirable to reduce vibration amplitude when the exciting frequency is in the vicinity of f_0 . This occurs occasionally in most installations. Any isolator that absorbs energy provides damping.

It is seldom practical to introduce damping as an independent variable in the design of vibration isolators for relatively small objects. The usual practice is to rely on the inherent damping characteristics of the rubber or other elastic material employed in the mounting. Damping achieved in this way seldom exceeds 5 percent of the amount needed to produce a critically damped system. In vibration isolators for large objects, such as variable-speed engines, the system often can be designed to produce nearly critical damping by employing fluid dashpots or similar devices.

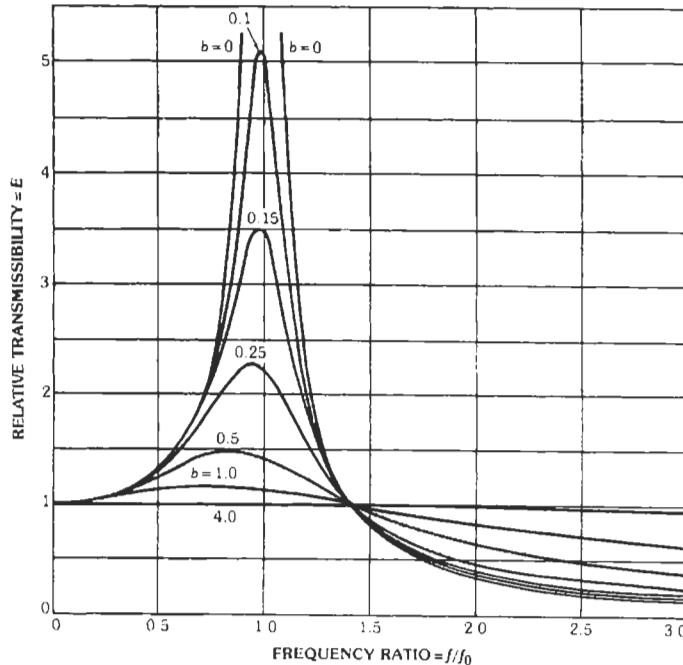


Fig. 7. Relative transmissibility E as a function of frequency ratio f/f_0 for various amounts of damping b . (By permission from *Vibration Analysis*, by N. O. Myklestad. ©1944, McGraw-Hill Book Co.)

TABLE 18. VIBRATION IN TYPICAL VEHICLES

Vehicle	Range of Frequencies (hertz)	Approximate Peak Amplitude (inches)	Nature of Excitation	Usual Choice of Isolator Resonant Frequency
Ships	0 to 15	0.02	Engine vibration in diesel or reciprocating steam drive	6 hertz for vibration isolation in commercial vessels, 27 to 30 hertz for shock isolation on naval vessels. These latter mounts amplify most vibrations to some extent.
	0 to 33	0.01	Propeller-blade frequency = (propeller rpm) × (number of blades)/60	
Piston-engine aircraft	0 to 60	0.01	Engine vibrations	Above 20 hertz. Amplitude of vibrations varies with location in aircraft. Landing shock can be neglected.
	0 to 100	0.01	Propeller vibrations. Aerodynamic vibrations due to buffeting	
Turboprop aircraft	0 to 60	0.01	Engine vibrations = (engine rpm)/60	9 hertz
	0 to 100	0.01	Propeller vibrations	
Jet aircraft	Up to 500	0.001	Audible noise frequencies due to jet wake and combustion turbulence; very little engine vibration	9 hertz
Passenger automobiles	1	6	Suspension resonance	25 hertz will usually avoid resonance with wheel hop and suspension resonant frequencies.
	8 to 12	0.02	Unsprung weight resonance (wheel hop)	
	20+	0.002	Irregular transient vibrations due to resonances of structural members with road roughnesses	
Automobile trucks	4	5	Suspension resonance	Above 20 hertz and should not correspond with any structural resonance. It is not advisable to attempt to isolate suspension and unsprung weight resonances.
	20	0.05	Unsprung weight resonance	
	80+	0.005	Structural resonances	
Military tanks	1 to 3	2	Suspension resonance	Similar to automobile truck
	Depends on speed	—	Track-laying frequency ≈ 17.6 (speed in mph)/(tread spacing in inches)	
	100+	0.001	Structural resonances	
Railroad trains	Broad and erratic		Similar to automobiles with additional excitations from rail joints and from side slop in rail trucks and draft gear	20 hertz has been successful in railroad applications. Shock with velocity changes up to 100 inches/second in direction of train occurs when coupling cars or starting freight trains.

Practical Application

Vibration can be accurately precalculated only for the simplest systems. In other cases, the actual vibration should be measured on experimental assemblies by using electrical vibration pickups. Complex vibration is often described by a plot of the g values against frequency. These plots usually show several frequencies at which the largest accelerations are present. The patterns will vary from place to place in a complicated structure and will also depend on the direction in which the acceleration is measured.

After vibration has been measured and plotted in this way, attention can be devoted to reduction of the predominant components. The equations and principles given above are used as guides in selecting the size, stiffness, damping characteristics, and location of isolators.

Shock

In many practical situations, vibration and shock occur simultaneously. The design of isolators for vibration should anticipate the effects of shock and vice versa.

When heavy shock is applied to a system using vibration isolators, there is usually a definite deflection at which the isolators snub or at which their stiffness suddenly becomes much greater. These actions may amplify the shock forces. To reduce this effect, it is generally desirable to use isolators that have smoothly increasing stiffness with increasing deflection.

Shock protection is improved by isolators that permit large deflections in all directions before the protected equipment is snubbed or strikes neighboring apparatus. The amplitude of vibration resulting from shock can be reduced by employing isolators that absorb energy and thus damp oscillatory movement.

Probabilities of damage to the apparatus itself from impact shock can be minimized by:

- (A) Making the weight of equipment components as small as possible and the strength of structural members as great as possible.
- (B) Distributing rather than concentrating the weights of equipment components and avoiding rigid connections between components.
- (C) Employing structural members that have high ratios of stiffness to weight, such as tubes, I beams, etc.
- (D) Avoiding, so far as is practical, stress concentrations at joints, supports, discontinuities, etc.
- (E) Using materials such as steel that yield rather than rupture under high stress.

US GRAPHIC SYMBOLS

The USA standard graphic symbols for electrical and electronics diagrams cover both the communication and power fields. Symbols of primary interest to communications workers are in Fig. 8. They have been abstracted

from IEEE Standard 315-1975.* Symbols that also agree with Recommendation No. 117 of the International Electrotechnical Commission are indicated by "IEC."

Diagram Types

Block diagrams consist of simple rectangles and circles with names or other designations within or adjacent to them to show the general arrangement of apparatus to perform desired functions. The direction of power or signal flow is often indicated by arrows near the connecting lines or arrowheads on the lines.

Schematic diagrams show all major components and their interconnections. A single-line diagram, as indicated by that name, uses single lines to interconnect components even though two or more conductors are actually required. It is a shorthand form of schematic diagram. It is always used for waveguide diagrams.

Wiring diagrams are complete in that all conductors are shown and all terminal identifications are included. The contact numbers on electron-tube sockets, colors of transformer leads, rotors of adjustable capacitors, and other terminal markings are shown so that a workman having no knowledge of the operation of the equipment can wire it properly.

Orientation

Graphic symbols are not considered to be coarse pictures of specific pieces of equipment but are true symbols representing the functions of parts in the circuit. Consequently, they may be rotated to any orientation with respect to each other without changing their meanings. Ground, chassis, and antenna symbols, for instance, may "point" in any direction that is convenient for drafting purposes.

Graphic symbols may be correlated with parts lists, descriptions, or instructions by means of reference designations (MIL-STD-16).

Detached Elements

Switches and relays often have many sets of contacts, and these sets may be separated and placed in the parts of the drawing to which they apply. Each separated element should be clearly labeled as part of the basic switch or relay.

Terminals

The terminal symbol need not be used unless it is needed. Thus, it may be omitted from relay and switch symbols. In particular, the terminal symbol often shown at the end of the movable element of a relay or

* Published by the Institute of Electrical and Electronics Engineers, 345 East 47th Street, New York, N.Y.

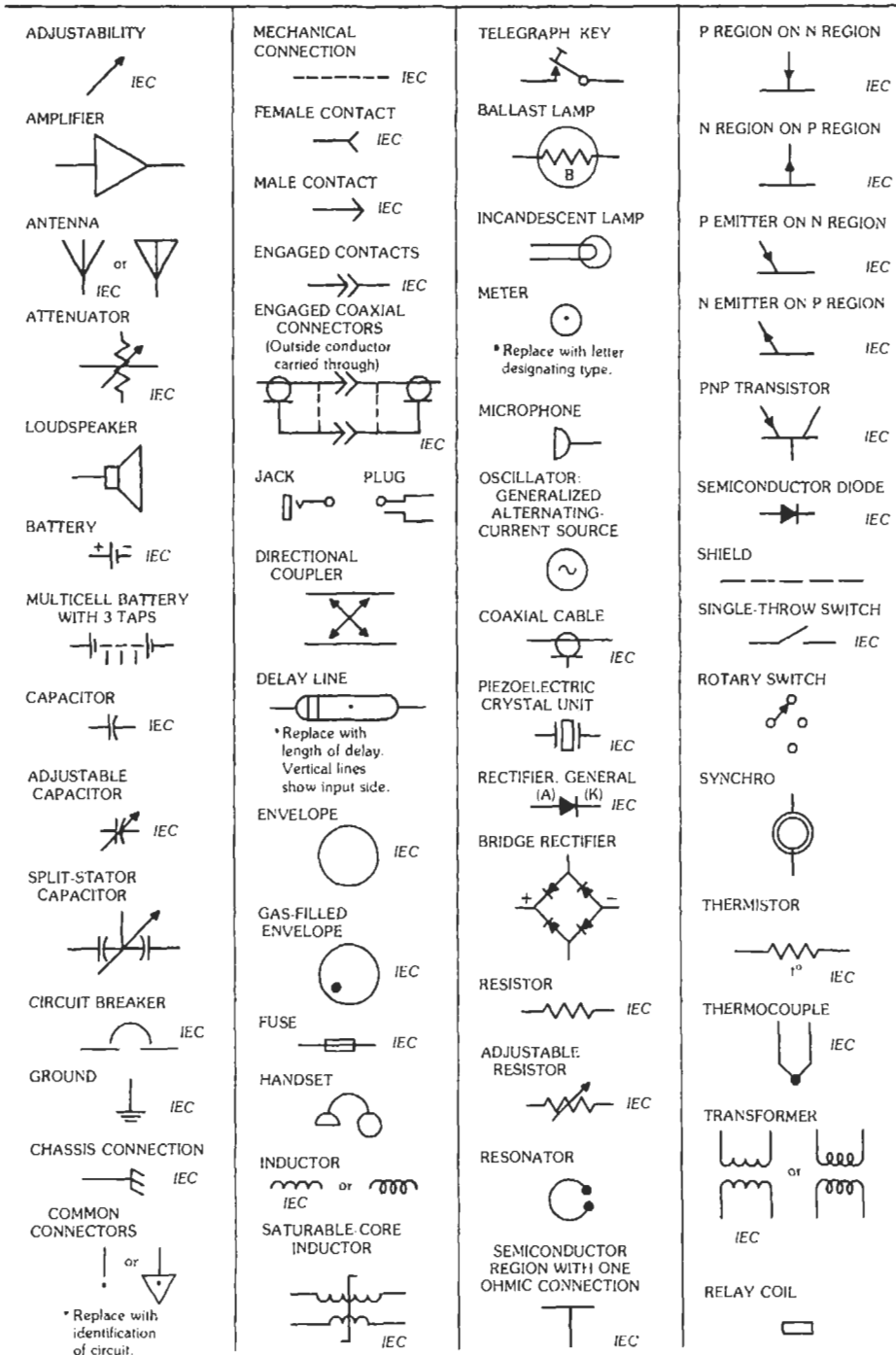


Fig. 8. Selected graphic symbols from IEEE Standard 315-1975.

switch should not be considered as the fulcrum or bearing but only as a terminal.

Associated or Future Equipment

Associated equipment, such as for measurement purposes, or additions that may be made later are identified as such by using broken lines for both symbols and connections.

BRITISH GRAPHIC SYMBOLS

Commonly used British block-diagram graphic symbols are in Fig. 9. They have been abstracted with permission from British Standard 530: 1948 and Supplement No. 5 (1962), and superseding BS 3939 (1966-1969). The issuing organization is the British Standards Institution, British Standards House, 2 Park Street, London, W. 1, England.

STANDARDIZED GRAPHIC SYMBOLS FOR LOGIC DESIGN

Symbols for logic design are so varied and in such a state of flux that many authors use their own symbols. Fig. 10 shows a short set of standardized graphic symbols for logic design selected from a revised draft of ANSI/IEEE Standard 91-1982.

SAFETY LEVELS FOR EXPOSURE TO EM FIELDS*

The American National Standards Institute (ANSI) has published recommendations concerning safety levels with respect to human exposure to electromagnetic fields in the frequency range from 300 kHz to 100 GHz.† The ANSI radio-frequency protection guides

* This section contributed by Charles A. Cain.
 † Reference 1.

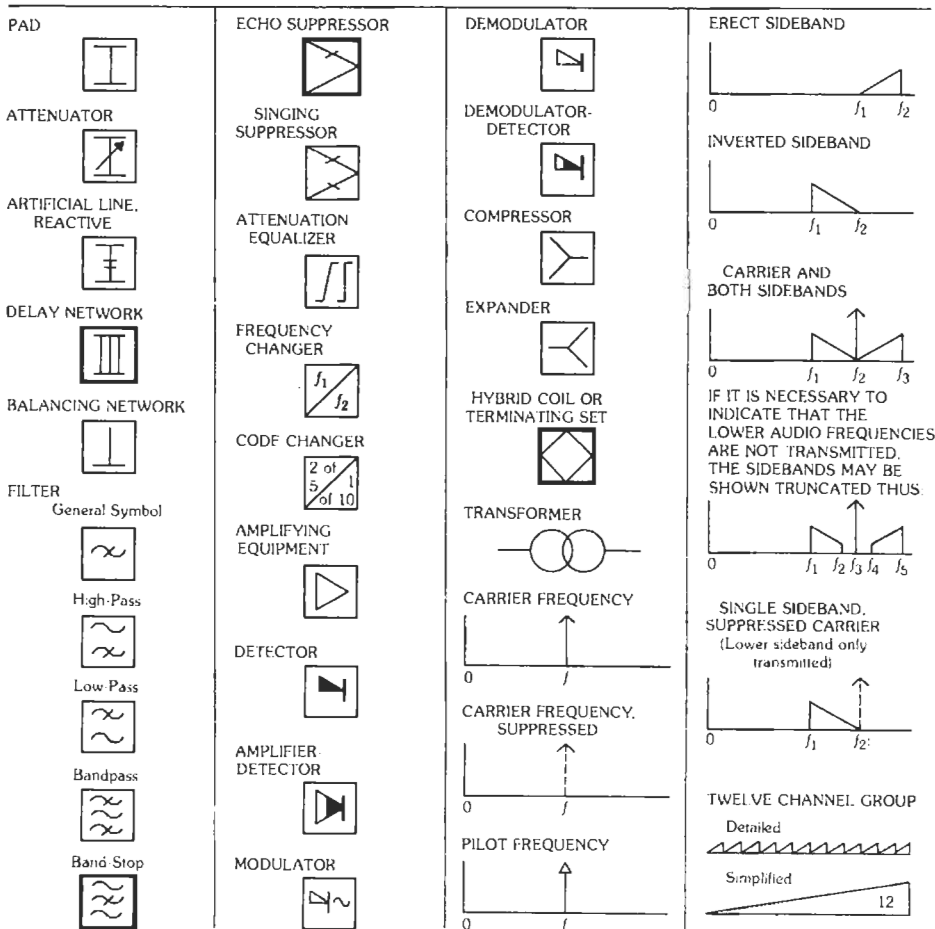


Fig. 9. Selected British block-diagram graphic symbols.




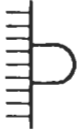










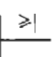







Description of Logic Functions	Rectangular-Shape Symbols	Distinctive-Shape Symbols
AND		
Accommodation of additional inputs		
Dot AND	  	 
OR		
Accommodation of additional inputs		
Dot OR	  	 
Exclusive OR (one and only one)		
Logical inversion		

Fig. 10. Selected standardized graphic symbols for logic diagrams. (Excerpted with permission from the revised draft of IEEE Standard Graphic Symbols for Logic Diagrams (Two State Devices). For a complete listing of all standardized symbols for logic diagrams, including the relationship between the new and superseded standards, see ANSI/IEEE Std. 91-1982, Appendix C.)

(RFPG) are intended to apply to both occupational and nonoccupational exposure, but not to exposure of patients for medical purposes. These guides specify

maximum allowable levels as a function of frequency in terms of the mean squared electric and magnetic fields, E and H , respectively, and the equivalent free-space plane-wave power density. To determine adherence to this RFPG, the power density or squares of the field strengths should be averaged over any 0.1-h period with all measurements to be made at a distance 5 cm or greater from any object. Table 19 is a summary of the ANSI C95.1-1982 RFPG.

A detailed rationale for the ANSI C95.1-1982 RFPG is given in reference 1. The frequency dependency of the maximum allowable equivalent plane-wave power density and of the mean squared field strengths is based on research results which show that the whole-body-average specific absorption rates* approach maximum values when the long axis of the body is parallel to the E -field vector and is approximately 0.4λ in length where λ is the free-space wavelength of the incident field.† This resonance occurs at a frequency of about 70 MHz for a man of average size. To account for the range of body dimensions from small infants to large adults, the RFPG specifies a maximum allowable power density of 1 mW/cm^2 over the frequency range from 30 MHz to 300 MHz. As indicated in Table 19, the allowable incident intensities increase below 30 MHz and above 300 MHz.

The ANSI C95.1-1982 standard specifies an exclusion from the recommended protection guides for low-power devices. The levels in Table 19 may be exceeded if the input power to a radiating device at radio frequencies between 300 kHz and 1 GHz is seven watts or less. The RFPG may also be exceeded between 300 kHz and 100 GHz if the exposure conditions result in specific absorption rates below 0.4 W/kg as averaged over the whole body and spatial peak specific absorption rates below 8 W/kg as averaged over any one gram of tissue for any 0.1-h time interval.

The type of situation that this exclusion is meant to address is illustrated by Table 20, which shows the field strengths and power densities in the near field of a 50-cm monopole antenna radiating at 150 MHz and 2 W of input power. It is clear from these data that the RFPG maximum recommended limits of field strength or of plane-wave equivalent power density are exceeded at locations near the antenna. However, this device radiates less than 7 watts and may be excluded.

References

1. American National Standard Safety Levels with Respect to Human Exposure to Radio Frequency Electromagnetic Fields, 300 kHz to 100 GHz, ANSI C95.1-1982. Institute of Electrical and Electronics Engineers, 1982.

* Specific absorption rate: the time rate per unit mass of radio-frequency electromagnetic energy deposition in biological tissues.

† Reference 2.

TABLE 19. RADIO-FREQUENCY PROTECTION GUIDES*

Frequency Range (MHz)	E^2 (V^2/m^2)	H^2 (A^2/m^2)	Power Density (mW/cm^2)
0.3- 3	400 000	2.5	100
3- 30	4000 ($900/f^2$)	0.025 ($900/f^2$)	$900/f^2$
30- 300	4000	0.025	1.0
300- 1500	4000 ($f/300$)	0.025 ($f/300$)	$f/300$
1500-100 000	20 000	0.125	5.0

* From ANSI C95.1-1982.
 Note: f = frequency in megahertz.

- Durney, C. H., Johnson, C. C., Barber, P. W., Massoudi, H., Iskander, M. F., Lords, J. L., Ryser, D. K., Allen, S. J., and Mitchell, J. C. *Radiofrequency Radiation Dosimetry Handbook*, Second Edition, May 1978, Report SAM-TR-78-22. USAF School of Aerospace Medicine, Brooks Air Force Base, Texas.
- Jordan, E. C. *Electromagnetic Waves and Radiating*

Systems. Radiating Systems. Englewood Cliffs, N.J.: Prentice-Hall, Inc., 1950, p. 323.

Additional Reading

The subject of safe levels of exposure to electromagnetic fields is one of great continuing interest, with much research being done worldwide. For an indication

TABLE 20. FIELD STRENGTHS (E AND H) AND POWER DENSITIES IN PROXIMITY TO A CURRENT-FED, QUARTER-WAVE, RADIATING MONOPOLE ANTENNA*

	Distance in Centimeters				
	0	10	20	30	40
$E = 64.3$ V/m		48.7	32.7	24.2	19.3
$H = 0.00$ A/m		0.0217	0.0317	0.0341	0.0348
$PD = 1.10$ mW/cm ²		0.628	0.287	0.156	0.0985
*		*	*	*	*
		95.1	45.7	29.8	22.4
		0.0524	0.0524	0.0489	0.0456
		2.40	0.553	0.236	0.133
		*	*	*	*
		121	53.4	33.3	24.5
		0.136	0.0856	0.0682	0.0533
		3.88	0.757	0.295	0.159
		*	*	*	*
		110	51.7	33.3	25.0
		0.229	0.124	0.089	0.0716
		3.72	0.708	0.300	0.193
		*	*	*	*
		83.6	43.0	30.2	24.1
		0.306	0.157	0.108	0.0531
		3.54	0.931	0.437	0.261
		*	*	*	*
		50.1	32.0	26.2	22.7
		0.358	0.180	0.120	0.0910
		4.82	1.22	0.548	0.312
		*	*	*	*
		27.7	26.3	24.3	22.1
		0.375	0.188	0.125	0.0938
		5.31	1.33	0.589	0.332
		*	*	*	*

* From ANSI C95.1-1982

A 50-cm antenna operating at 150 MHz and 2 W of input power is assumed, with an input impedance of 36 Ω . Each asterisk associated with a set of field measurements denotes the spatial point of measurement in relation to the antenna. The vertical and horizontal distance between adjacent points of measurement is 10 cm. (Calculations based on reference 3.)

of ongoing research on the biological effects of EM fields, and the efforts being made to promulgate regulations in the United States relating to EM-field standards, see summary articles by Eric J. Lerner: "Biological Effects of Electromagnetic Fields." *IEEE Spectrum*, May 1984, pp. 57-69; "The Drive to Regulate Electromagnetic Fields." *IEEE Spectrum*, March 1984, pp. 6370.

PATENT, COPYRIGHT, AND TRADE SECRET PROTECTION OF INVENTIONS

Intellectual Property Protection of Technology in the United States*

Depending on the nature of the invention or other development, the owner may seek protection in the United States by means of trade secret, patent, or copyright law.† Each nation is responsible for its own intellectual property laws, and on various points the law in many nations differs from that in the United States.

Trade Secret Protection of Inventions

What Is a Trade Secret?—One widely used definition is that a trade secret is "information, including a formula, pattern, computation, program, device, method, technique, or process, that: (i) derives independent economical value, actual or potential, from not being generally known to, and not being readily ascertainable by proper means, by other persons who can obtain economic value from its disclosure or use, and (ii) is the subject of efforts that are reasonable under the circumstances to maintain its secrecy." Basically, when a business has secret information giving it a commercial edge over competitors who don't have it, the business may seek to protect the information as a trade secret. Although trade secret protection often extends to such commercial items as customer lists and cost data, much of the information protected by trade secret law is technological.

For example, one may maintain as trade secrets chemical processes used in the trade secret owner's plant, or the structure and functioning of machines used by the owner's licensees who have agreed to keep these items confidential. And, along with copyright protection, trade secret law is quite commonly used to protect computer programs. However, a device marketed to consumers and embodying a "trade secret"

whose characteristics are obvious from a casual inspection, would not be a good candidate for trade secret protection after marketing commenced, because there would then be no secret. Any person who lawfully comes into possession of a product embodying a trade secret, and is not contractually precluded from doing so, is entitled (if clever enough) to "reverse engineer" the product and learn the trade secret.

The Law of Trade Secrets—The basic rule of trade secret law is as easy to state as it is difficult to apply: a person may not make unauthorized access to, or use or disclosure of, another's trade secret. Such conduct constitutes misappropriation of the trade secret. It is not necessary that any tangible object be taken for misappropriation to occur; rather, misappropriation can occur on the basis of use or disclosure of information existing in a person's memory.

Unlike patent and copyright law, trade secrets are governed largely by state (rather than federal) law. One result of this is that trade secret law often differs from state to state, although in recent years some 42 states have enacted the "Uniform Trade Secret Act," which has introduced more uniformity into trade secret law. Under certain circumstances, in many states trade secret misappropriation may constitute a crime, and it does constitute a federal crime.

However, unlike patent rights, trade secret rights do not extend to preclude another from use or disclosure of information independently developed or otherwise lawfully acquired free of restrictions. Accordingly, a trade secret owner is continually subject to loss of the trade secret through the lawful conduct of others.

General Requirements for a Trade Secret—In many lawsuits, the following six factors have been viewed as important in determining whether a trade secret exists:

- (1) The extent to which the information is known outside the owner's business. Obviously, if everyone in the industry knows it, there is no secret. But if one or two others know it, and maintain it as a secret, it may still be a trade secret.
- (2) The extent to which it is known by employees and others involved in the owner's business. A business that restricts internal proliferation of its secrets on a "need to know" basis is in a better posture to claim that it has trade secrets.
- (3) The extent of measures taken by the owner to guard the secrecy of the information. The owner will rarely be required to duplicate the security features of Fort Knox. But if the owner permits people to walk in off the street and roam around in areas where this information is available, it will be difficult to convince a court that it constitutes a trade secret.
- (4) The value of the information to the owner and his competitors. Value embodies the commercial importance and desirability of the information.

* This section contributed by David Bender, White & Case, LLP.

† In addition, a special type of protection, known as semiconductor mask work protection, is available for protecting the structure and design of semiconductor integrated circuits.

Without such value, it is difficult to see how the information gives its owner an advantage over competitors.

- (5) The amount of money or effort expended by the owner in developing the information. This may be looked to as one indicator of value. Moreover, as an equitable matter, a court or jury may be reluctant to permit a defendant to reap where he or she has not sown.
- (6) The ease or difficulty with which the information could be properly acquired or duplicated by others. If it is a trivial matter for another to acquire it, it should not be protected.

Relationships Giving Rise to Trade Secret Relationship—The owner of a trade secret has rights, not against the world, but only against one who unlawfully acquires, uses, or discloses the trade secret. Typically, the owner's right to restrain use or disclosure arises from a "confidential relationship" that exists between the owner and the other party, whereby the other agrees to maintain the information as confidential. This confidential relationship may arise from a contract expressly creating it, or it may be implied from the relationship of the parties. For example, a confidential relationship arises generally from the owner's employment of an employee to work on developing a process or product.

Also, employers sometimes require certain key employees to sign a "covenant not to compete." In such a covenant, the employee agrees that for a specified period following employment by the employer, he or she will not, perhaps in a particular geographic region, undertake activities of the type in which he or she engaged for the present employer.

Who Owns Rights in a Trade Secret?—It is not always easy to determine who owns rights to a trade secret. For example, in the employer-employee situation, some factors that a court would look to in adjudicating that issue might be: the terms of any contract between the two; whether the employee was hired to do research and/or development work; whether the secrets in question are the result only of the employee's efforts; whether the work resulting in the trade secret was assigned to the employee by the employer; and whether the work done by the employee was done on the employer's premises, during normal working hours, and with the use of the employer's facilities and materials.

Patent Protection of Inventions

Nature of the Patent Grant—A United States patent grants to its owner the right to exclude others from making, using, selling, offering for sale, or importing the claimed invention in the United States for the term of that patent.* A patent can be asserted even against one who invents independently and with-

out copying from the patent owner. Patents issuing on applications filed today will expire twenty years after the filing date.† Most importantly, a patent does not confer on its owner the right to practice that invention, for one or more other patents may also cover the invention. So, for example, if a patent for a carburetor coexisted with a patent for a dual-barrelled carburetor (an improvement on the basic carburetor), the owner of the latter could not lawfully manufacture his or her invention without a license from the owner of the former. However, when the patent expires, the owner loses all right to exclude. Once an application for a patent is filed, it is freely assignable as an item of personal property, as is an issued patent.

What Can be Protected by Patent?—A patent may be obtained on any new and useful process, machine, manufacture, or composition of matter, or any new and useful improvement thereof. The patent standard of novelty is a strict one.‡ One novelty requirement is that the inventor must have made the invention before it was used or known by others in the United States, and before it was patented or described by others in a printed publication anywhere on earth. Even beyond novelty itself, the claimed invention must be nonobvious, when made, to a person of ordinary skill in the area of technology to which it pertains.

Computer Programs and Methods of Doing Business—In recent years, many inventors of program-related inventions, and of methods of doing business, have obtained patents on them. Since at least 1981, it has been clear that the US Patent and Trademark Office (the PTO) will issue patents for program-related inventions. An important court case in 1998 has made it clear that a business method, at least when implemented by computer, is patentable subject matter. Since that case was decided, there has been a huge surge of patent applications from the financial services sector and others.

Recognizing Inventions—If an original development appears to be of commercial significance (or if there is some other reason why a patent is desired), the

* There are other less commonly exercised rights as well, including the right to preclude importation into, or sale, offering for sale or use in, the United States of a product made abroad by a process patented in the United States.

† Patents on applications filed prior to June 9, 1985 generally expire on the later of twenty years after the filing date or seventeen years after the issue date.

‡ The following hypothetical example shows just how strict the standard is. Suppose a person makes a simple, mechanical invention, yet one so ingenious that no person alive has contemplated it before. The inventor obtains a patent. A year later a cave is discovered in Massachusetts with wall paintings indicating that in 1795 the local tribe of native Americans publicly used that invention. Under this set of facts, the patent would be invalid.

inventor may wish to transmit the pertinent details to a patent attorney or patent agent in order to have a patent search done to assess whether the invention meets the patent standards of novelty and nonobviousness. The attorney or agent will generally compare the invention with issued U.S. patents and other items of "prior art" in reaching this determination.

Who Is the Inventor?—Patent law recognizes that invention comprises two stages: (1) conception of the idea, and (2) its reduction to practice (which amounts to either constructing the invention or making a detailed written description in a patent application). The person(s) actually responsible for the conception and reduction to practice is (are) the inventor(s). If a person's contribution to an invention does not rise above the level of ordinary mechanical skill, that person is not an inventor. Supervisors should not, as a matter of course, be named on patent applications as inventors. They (and anyone else) should be named only where they actually participated in the invention.

The Inventive Process—The question of whether an invention is patentable does not depend on the manner in which it is made. In particular, the fact that an invention is made by accident, or serendipity, will not negate patentability. Often the inventive process comprises the following steps: someone identifies either a problem to be solved or a result to be achieved; the inventor conceives a way of solving the problem or effecting the result; and the inventor reduces the conception to reality, either by actually building or practicing the invention, or by making a detailed written description of it.

Obtaining a Patent—One seeks a patent by filing and prosecuting a patent application with the PTO. There are many requirements that must be met in the application. The application must include: a written description of the invention, and the manner and process of making and using it, in such full, clear, concise, and exact terms as to enable any person skilled in the appropriate technology to make and use it; a statement of the best mode contemplated by the inventor at the time of filing; and one or more "claims" distinctly pointing out the subject matter regarded as the invention. The preparation and prosecution of a patent application is a very technical and detailed task. By law, an application can be filed and prosecuted only by the inventor or by a registered patent attorney or patent agent. Each year, over 200,000 patent applications are filed in the United States, and some 140,000 patents issued. The time between filing of the application and issuance of the patent varies considerably from application to application and currently averages between 1 1/2 and 2 years.

Interferences—On occasion, different inventors independently file patent applications claiming the same invention. (The PTO maintains all patent applications in secrecy for at least eighteen months, so that

when an application is filed, the applicant may be unaware of other applications.) In the United States, the first person to invent is entitled to any resulting patent. When the PTO identifies applications by different inventors claiming substantially the same invention, the PTO declares an "interference" proceeding to determine the first to invent. Similarly, if two interfering patents inadvertently issue, the first inventor may be determined by suit in federal court. In most foreign countries, the first inventor to file a patent application is entitled to any resulting patent. Thus, the US patent owner may be a different person than the owner of foreign patents.

Publication—For a valid United States patent to issue, an application must be filed within one year of publication or public use of the invention. However, it is generally unwise publicly to disclose an invention before a patent application is filed. The law in some foreign countries is such that publication before filing for a patent will extinguish the right to obtain a patent there.

Importance of Records—In certain instances, the patent owner may be called upon to prove various facts associated with his or her patent rights. In an interference, he or she will be expected to prove when the invention was conceived, when it was reduced to practice, and perhaps that due diligence was exercised in the interim. In licensing negotiations, the owner may be challenged to show similar facts. And in litigation, the owner will be expected to show these and other facts relating to the invention. Interferences, licensing negotiations, or litigation may occur years—even decades—after the invention was made. Memories dim and witnesses die. With clear and detailed contemporaneous documents and models, it is less difficult to reconstruct precisely what happened and when.

Such documents may consist of the lab or engineering notebooks of the inventors and others, memoranda, and sketches. The notebooks in particular are often of prime importance, and certain characteristics may render them more credible as evidence:

- (1) They should be bound, rather than loose-leaf.
- (2) Entries should be in chronological order, written clearly and completely in ink, without blank pages or spaces.
- (3) Changes should be made by drawing one line through material to be deleted (but not obliterating it).
- (4) Erasures and removal of pages should be avoided.
- (5) Where other materials must be included in the notebook (e.g., small lab samples, graphs), these should be fastened securely to the appropriate page, referenced on that page, dated, signed, and witnessed.
- (6) Each entry should be made at or shortly after the events recorded and should be dated and signed when made.

- (7) When an event of perceived significance is recorded, the entry should be shown to someone of discretion who will understand it. This person should then write the date and his or her signature below the entry, subscribing above the signature, "read and understood."

Copyright Protection of Inventions

Copyright law protects the expression of an "original work of authorship," as opposed to the tangible embodiment of a concept or idea. For this reason, copyright is not generally an appropriate vehicle for protecting technology. The single major exception to this generality occurs in connection with the computer program. Copyright springs into being automatically as soon as the work achieves fixed tangible form, without the necessity of filing any application. The copyright may be registered with the US Copyright Office, but registration is not a requirement for protection (although it offers certain advantages and is generally necessary in order to sue for infringement). Copyright provides protection against the unauthorized copying and distribution of programs. A copyright precludes only one who copies the owner's work, and not one who independently developed a similar (or even identical) work.

SUMMARY OF MILITARY NOMENCLATURE SYSTEM*

In the Joint Electronics Type Designation System (JETDS), formerly called the "AN" system, nomenclature for electronic equipment consists of a name, followed by a type number.

A type designation assignment for equipment such as a definitive system, subsystem, center, central, set, etc., shall consist of at least an AN, a slant bar, a three-letter equipment designation (see Table 21), a dash, and a number. Example: AN/VRC-12 would be a radio communication set installed in a vehicle designed for functions other than carrying electronic equipment.

All groups, including commercial off-the-shelf equipment, are identified by a two-letter indicator from Table 22. Applicable equipment indicator letters (see Table 21) follow the slant bar to indicate the potential of the group for multiple or peculiar application. Example: OE-162/ARC indicates an antenna for aircraft radiocommunication equipment. Equipment indicators with a specific model number (e.g., OK-450/TRC-26) are used following the slant bar when the

group is peculiar to specific equipment (e.g., AN/TRC-26) with no known potential for other use.

The type designation for units having one end use consists of an indicator (Table 23), a dash, a number, a slant bar, and the equipment the unit is a part of or used with. Example: the receiver portion of the AN/VRC-12 is identified as R-40/VRC-12. If the unit has multiple usage, only those indicators that are common or appropriate are included after the slant bar.

Examples: A power supply, part of or used with the AN/VRC-12 and AN/VRC-19 would be identified as PP-50/VRC. A power supply, "part of" the AN/VRC-12 and "used with" the AN/VRR-40 would be identified as PP-60/VR.

The system indicator (AN) does not mean that the Army, Navy, and Air Force use the equipment, but simply that the type number was assigned in the JETDS system.

Nomenclature Policy

JETDS nomenclature will be assigned to:

- (A) Complete systems, subsystems, centers, centrals, sets, groups, kits, and units of military design, either definitive or variable in configuration.
- (B) Groups of articles of either commercial or military design that are grouped for a military purpose.
- (C) Electronic articles of military design that are part of or used with an item not identified in the JETDS.
- (D) Commercial articles requiring military identification for use by US Government.
- (E) Electronic materials of military design which are not part of or used with a set.

JETDS nomenclature will not be assigned to:

- (A) Articles cataloged commercially except in accordance with paragraph (D) above.
- (B) Minor components of military design for which other adequate means of identification are available.
- (C) Small parts such as capacitors and resistors.
- (D) Articles having other adequate identification in joint military specifications.

Nomenclature assignments will remain unchanged regardless of later changes in installation and/or application.

Modification Letters

Component modification suffix letters will be assigned for each modification of a component when detail parts and subassemblies used therein are no longer interchangeable, but the component itself is interchangeable physically, electrically, and mechanically.

* Adapted from MIL-STD-196C, Joint Electronics Type Designation System, 22 April 1971; Notice 1, 20 April 1972; Notice 2, 8 June 1972; Notice 3, 14 November 1972; and Notice 4, 17 July 1977. Available from the Superintendent of Documents, Washington, DC 20402.

TABLE 21. SET OR EQUIPMENT INDICATOR LETTERS

1st Letter (Type of Installation)	2nd Letter (Type of Equipment)	3rd Letter (Purpose)
A Piloted aircraft	A Invisible light, heat radiation	A Auxiliary assemblies (not complete operating sets used with or part of two or more sets or sets series) (inactivated. do not use)
B Underwater mobile submarine	B Pigeon (do not use)	B Bombing
C Air transportable (inactivated, do not use)	C Carrier	C Communications (receiving and transmitting)
D Pilotless carrier	D Radiac	D Direction finder, reconnaissance, and/or surveillance
F Fixed ground	E Nupac (inactivated, do not use)	E Ejection and/or release
G General ground use (includes two or more ground-type installations)	F Photographic*	G Fire control or searchlight directing
	G Telegraph or Teletype	H Recording and/or reproducing (graphic meteorological and sound)
	I Interphone and public address	
	J Electromechanical or inertial wire covered	
K Amphibious	K Telemetering	K Computing
	L Countermeasures	L Searchlight control (inactivated, use G)
M Ground, mobile (installed as operating unit in a vehicle which has no function other than transporting the equipment)	M Meteorological	M Maintenance and/or test assemblies (including tools)
	N Sound in air	N Navigational aids (including altimeters, beacons, compasses, beacons, depth sounding, approach, and landing)
P Pack or portable (animal or man)	P Radar	P Reproducing (inactivated, use H)
	Q Sonar and underwater sound	Q Special, or combination of purposes
	R Radio	R Receiving, passive detecting
S Water surface craft	S Special types, magnetic, etc., or combinations of types	S Detecting and/or range and bearing, search
T Ground, transportable	T Telephone (wire)	T Transmitting
U General utility (includes two or more general installation classes, airborne, shipboard, and ground)		
V Ground, vehicular (installed in vehicle designed for functions other than carrying electronic equipment, etc., such as tanks)	V Visual and visible light	
W Water surface and underwater combination	W Armament (peculiar to armament, not otherwise covered)	W Automatic flight or remote control
	X Facsimile or television	X Identification and recognition
		Y Data processing surveillance (search, detect, and multiple target tracking) and control (both fire control and air control)
Z Piloted and pilotless airborne vehicle combination		

* Not for US use except for assigning suffix letters to previously nomenclatured items.

TABLE 22. GROUP INDICATORS

Indicator	Family Name	Indicator	Family Name
OA	Miscellaneous groups	ON	Interconnecting groups
OB	Multiplexer and/or demultiplexer groups	OP	Power supply groups
OD	Indicator groups	OQ	Test set groups
OE	Antenna groups	OR	Receiver groups
OF	Adapter groups	OT	Transmitter groups
OG	Amplifier groups	OU	Converter groups
OH	Simulator groups	OV	Generator groups
OI	Consoles and console groups	OW	Terminal groups
OK	Control groups	OX	Coder, decoder, interrogator, transponder groups
OL	Data analysis and data processing groups	OY	Radar set groups
OM	Modulator and/or demodulator groups	OZ	Radio set groups

Set modification letters will be assigned for each modification not affecting interchangeability of the sets or equipment as a whole, except that in some special cases these will be assigned to indicate functional interchangeability and not necessarily complete electrical and mechanical interchangeability. Modification letters will only be assigned if the frequency coverage of the unmodified equipment is maintained.

The suffix letters X, Y, and Z will be used only to designate a set or equipment modified by changing the power input voltage, phase, or frequency. They will be used as follows: X will indicate the first change, Y the second, Z the third, XX the fourth, etc. These letters will be in addition to other modification letters applicable.

Developmental Indicators

Experimental Sets—To identify a set or equipment of an experimental nature with the development organization concerned, the following indicators are used within the parentheses:

XA	Aeronautical Systems Division, Wright-Patterson Air Force Base, OH.
XB	Naval Research Laboratory, Washington, DC.
XC	US Army Signal Engineering Laboratories, The Hexagon, Fort Monmouth, NJ (inactivated, use XE).
XD	Electronic Systems Division, Laurence G. Hanscom Field, Bedford, MA.
XE	US Army Electronics Laboratories, Fort Monmouth, NJ.
XF	Frankford Arsenal, Philadelphia, PA.
XG	USN Electronics Laboratory, San Diego, CA.

*XH	Aerial Reconnaissance Laboratory, Wright-Patterson Air Force Base, OH.
XI	Air Force Armament Laboratory, Eglin Air Force Base, FL.
XJ	Naval Air Development Center, Johnsville, PA.
*XK	Flight Control Laboratory, Wright-Patterson Air Force Base, OH.
XL	US Army Signal Electronics Research Unit, Mountain View, CA.
XM	US Army Signal Engineering Laboratories, The Hexagon, Fort Monmouth, NJ (inactivated, use XE).
XN	Department of the Navy, Washington, DC.
XO	US Army Missile Command, Redstone Arsenal, AL.
XP	Canadian Department of National Defence, Ottawa, Ontario, Canada.
*XQ	Aeronautical Accessories Laboratory, Wright-Patterson Air Force Base, OH.
XR	National Security Agency, Fort George G. Meade, MD.
*XS	Electronic Components Laboratory Wright-Patterson Air Force Base, OH.
XT	US Army Security Agency, Arlington Hall Station, Arlington, VA.
XU	USN Underwater Sound Laboratory, Fort Trumbull, New London, CT.
XV	Air Force Weapons Laboratory, Kirtland Air Force Base, NM.
XW	Rome Air Development Center, Rome, NY.
*XY	Weapons Guidance Laboratory, Wright-Patterson Air Force Base, OH.

* Not for Air Force use except for assigning additional developmental designations to previously type-designated items. Use XA for all new equipment.

- XZ** USN Bureau of Naval Weapons Activities.
- XAA** Air Force Ballistic Systems Division, Norton Air Force Base, CA.
- XAE** US Army Electronics Research and Development Activity, Fort Huachuca, AZ.
- XAN** Naval Avionics Facility, Indianapolis, IN.
- XBB** US Army Electronics Command, Proc and Prod Div., Fort Monmouth, NJ.
- XCA** US Naval Ammunition Depot, Crane, IN.
- XCC** Air Force Missile Test Center, Patrick Air Force Base, FL.
- XCL** Naval Weapons Center, China Lake, CA.
- XCR** Naval Weapons Center, Corona Laboratory, Corona, CA.
- XDD** US Army Signal Air Defense Engineering Agency, Fort George G. Meade, MD.
- XDV** US Naval Weapons Laboratory, Dahlgren, VA.
- XGS** Grand Support Equipment Division, Naval Air Engineering Center, Philadelphia, PA.
- XIH** US Naval Ordnance Station, Indianhead, MD.
- XLW** US Army Limited War Laboratory, Aberdeen Proving Ground, MD.
- XMG** Naval Missile Center, Point Mugu, CA.
- XPM** US Army, Project Michigan, Ypsilanti, MI.
- XSC** US Army Satellite Communications Agency, Fort Monmouth, NJ.
- XWH** US Naval Weapons Handling Laboratory, Naval Weapons Station Earle, Colts Neck, NJ.
- XWO** Naval Ordnance Laboratory, White Oak, Silver Spring, MD.

Example: Radio Set AN/ARC-3 might be assigned for a new airborne radio communication set under development. The cognizant development organization might then assign AN/ARC-3(XA-1), AN/ARC3(XA-2), etc., type numbers to the various sets developed for test. When the set was considered satisfactory for use, the experimental indicator would be dropped, and procurement nomenclature **AN/ARC-3** would be officially assigned thereto.

- (A) A set to train for a specific basic set will be assigned the basic-set type number followed by a dash, the letter T, and a number. Example: Radio Training Set AN/ARC-6A-T1 would be the first training set for Radio Set AN/ARC-6A.
- (B) A set to train for general types of sets will be assigned the usual set indicator letters followed by a dash, the letter T, and a number. Example: Radio Training Set AN/ARC-T1 would be the first training set for general airborne radio communication sets.

Parenthesis Indicators: A series of a basic item, i.e., all production and/or nonproduction versions, may be identified by a type designation with an empty parenthesis. Examples: AN/APS-25 () or R-275()/APS-25. Such an assignment is all inclusive and does not refer to any specific version within the series.

Systems, subsystems, centers, centrals, sets, groups, or units with variable parts lists are assigned type designations in the same manner, except a parenthetical V expressed as (V) is added to the type designation. Units designed to accept plug-ins that change the function, frequency, or technical characteristics of the type-designated unit are designated with a (P) preceding the slant bar. The plug-in is not considered a part of the unit.

Examples of JETDS Type Numbers

- AN/SRC-3()** General reference set nomenclature for water surface craft radio communication set number 3.
- AN/SRC-3** Original procurement set nomenclature applied against AN/ SRC-3().
- AN/SRC-3A** Modification set nomenclature applied against AN/SRC-3.
- AN/APQ-13-TI()** General reference training set nomenclature for the AN/APQ-13 set.
- AN/APQ-13-TI** Original procurement training set nomenclature applied against AN/APQ-13-TI().
- AN/APQ-13-TIA** Modification training set nomenclature applied against AN/APQ-13-TI.
- AN/UPT-T3()** General reference training set nomenclature for general utility radar transmitting training set number 3.
- AN/UPT-T3** Original procurement training set nomenclature applied against AN/UPT-T3().
- AN/UPT-T3A** Modification training set nomenclature applied against AN/ UPT-T3.
- T-51()/ARQ-8** General reference component nomenclature for transmitter number 51, part of or used with airborne radio special set 8.
- T-51/ARQ-8** Original procurement component nomenclature applied against T-51()/ARQ-8.
- T-51A/ARQ-8** Modification component nomenclature applied against T51/ARQ-8.
- RD-31()/U** General reference component nomenclature for recorder-reproducer number 31 for general utility use, not part of a specific set.
- RD-31/U** Original procurement component nomenclature applied against RD-31()/U.
- RD-31A/U** Modification component nomenclature applied against RD31/U.

TABLE 23. UNIT INDICATORS

Indicator	Family Name	Indicator	Family Name
AB	Supports, antenna	MF*	Magnets or magnetic-field generators
AM	Amplifiers	MK	Miscellaneous kits
AS	Antennas, complex and simple	ML	Meteorological devices
AT*	Antennas, simple	MT	Mountings
BA	Batteries, primary type	MU	Memory units
BB	Batteries, secondary type	MX	Miscellaneous
BZ	Alarm units	O	Oscillators
C	Controls	OC*	Oceanographic devices
CA*	Commutator assemblies, sonar	OS	Oscilloscopes, test
CB*	Capacitor banks	PD*	Prime drivers
CG	Cable assemblies, RF	PF*	Fittings, pole
CK*	Crystal kits	PG*	Pigeon articles
CM	Comparators	PH*	Photographic articles
CN	Compensators	PL	Plug-in units
CP	Computers	PP	Power supplies
CR*	Crystals	PT	Mapping and plotting units
CU	Couplers	PU	Power equipments
CV	Converters (electronic)	R	Receivers
CW	Radomes	RC*	Reels
CX	Cable assemblies, non-RF	RD	Recorder-reproducers
CY	Cases and cabinets	RE	Relay assemblies
D	Dispensers	RF*	Radio-frequency components
DA	Loads, dummy	RG*	Cables, RF bulk
DT	Detecting heads	RL	Reeling machines
DY*	Dynamotors	RO	Recorders
E*	Hoists	RP	Reproducers
F	Filter units	RR	Reflectors
FN*	Furniture	RT	Receiver and transmitter
FR	Frequency-measuring devices	S	Shelters
G	Generators, power	SA	Switching units
GO*	Goniometers	SB	Switchboards
GP*	Ground rods	SG	Generators, signal
H	Head, hand, and chest sets	SM	Simulators
HC*	Crystal holders	SN	Synchronizers
HD	Environmental apparatus (heating, cooling, etc.)	ST*	Straps
ID	Indicators, noncathode-ray tube	SU	Optical devices
IL*	Insulators	T	Transmitters
IM	Intensity-measuring devices	TA	Telephone apparatus
IP	Indicators, cathode-ray tube	TB	Towed bodies
I	Interface units	TC*	Towed cables
KY	Keying devices	TD	Timing devices
LC*	Tools, line-construction	TF	Transformers
LS	Loudspeakers	TG	Positioning devices
M	Microphones	TH	Telegraph apparatus
MA*	Magazines	TK*	Tool kits
MD	Modulators, demodulators, and discriminators	TL*	Tools
ME	Meters	TN	Tuning units

TABLE 23 (CONT). UNIT INDICATORS

Indicator	Family Name	Indicator	Family Name
TR	Transducers	VS*	Signaling equipment, visual
TS	Test units	WD*	Cables, two-conductor
TT	Teletypewriter and facsimile apparatus	WF*	Cables, four-conductor
TV*	Testers, tube	WM*	Cables, multiple-conductor
TW	Tape units	WS*	Cables, single-conductor
U*	Connectors, audio and power	WT*	Cables, three-conductor
UG*	Connectors, RF	ZM	Impedance-measuring devices
V	Vehicles		

* Not for US use except for assigning suffix letters to previously nomenclatured items.

Index

- A
- A- and B-Bit signaling, 38-17
 - Abampere, 3-13
 - Abcoulomb, 3-13
 - Abfarad, 3-13
 - Abhenry, 3-13
 - Abmho, 3-13
 - Abohm, 3-13
 - Absorptance, optical waves, 21-9
 - Absorption, 41-2
 - atmospheric, 33-20-21
 - on optical communications, 22-14-16
 - coefficients, 40-14-15
 - building materials, 40-16
 - furnishings, 40-16
 - cross section, 41-4
- Abvolt, 3-14
- Ac
- filter reactors, 13-19
 - motor, drive for adjustable-speed, 14-48
 - output power supplies, 14-46-51
 - voltage regulators, 14-46
- Ac-dc converter circuits, 14-12-39
- Acceleration
- displacement, and vclocity magnitudes, 49-23
 - error constant, unity feedback, 15-11
 - vibratory forces, 49-23
- Acceptor impurities, 18-4
- Access time, ICs, 20-3
- Accumulators, computer, 42-11
- Åckerberg-Mossberg (AM) biquad filters, 10-21-22
- Acoustic impedance, specific, 40-3
- Acoustic power levels, of typical noise sources, 40-9
- Acoustic properties, materials, 4-31-33
- Acoustical standards, 40-25-28
- Acoustical systems
- electrical analogies for, 40-5-7
 - mechanical analogies for, 40-5-7
- Acre, 3-14
- Actinium, 4-4-5, 4-8-9
- Active filter(s)
- Åckerberg-Mossberg (AM) biquad, 10-21-22
 - building blocks, 10-10-15
 - circuit elements, 10-4-7
 - design
 - fundamentals of, 10-7-10
- Active filter(s)—cont
- techniques of, 10-7-10
- example, low-pass notch, 10-44-45
- general realization methods, 10-8
- introduction to, 10-3-4
- operational amplifiers in, 10-5-7
- transconductance-C (g_m -C), 10-22-23
- zero placement formulas, 10-43
- Active loads, transistor amplifiers, 19-17-19
- Active switched-capacitor filters, 10-34-45
- Adaptive arrays, antenna, 32-45-47
- Adaptive digital signal processing, 28-22-29
- Adaptive filters, 28-23-29
- Adder, definition of, 20-3
- Addition
 - nondecimal numbers, 42-46
 - using twos complement method, 42-8
- Address register
 - computer processor, 42-15
 - indirect addressing, 42-19
- Addressing
 - absolute short, computer, 42-17, 42-19
 - based, computer, 42-19
 - immediate, computer, 42-17
 - indirect, 42-25-26
 - modes, computer, 42-17-20
 - relative, computer, 42-19
- Adler-Hassner code, 25-7
- Admittance
 - formulas, 6-6, 6-9-13
 - input, of resonant line, 29-16-17
 - measurement of, 12-8
 - measuring circuit, twin-T, 12-10
 - transmission lines, 29-6-7
- ADS, 37-8
- ADSL, 38-7-8
- Advanced Data Communication Control Protocols (ADCCP), 26-16
- Advanced mobile phone service (AMPS), 26-20
- Advanced television (ATV) systems, 35-43-46
- Aeronautical Mobile Satellite Service (AMSS), 27-3
- Aeronautical Radio, Inc., 37-2
- AF amplifier
 - class-A, design of, 17-7
 - class-B, design of, 17-7-8
- Air traffic control radar beacon system, 37-5
- Aircraft
 - jet, vibration in, 49-25
 - piston engine, vibration in, 49-25
 - turboprop, vibration in, 49-25
- Airlines Electronic Engineering Committec (AEEC), 37-2
- Algebra
 - Boolean, 43-2
 - matrix, 47-31-35
 - switching, 43-2
- Algebraic equations, 47-6-9
- Algorithm(s)
 - Berlekamp-Massey, 25-15
 - distributed control, local area networks, 26-22-23
 - distributed-tree retransmission, 26-21-22
 - Euclidean, 25-15
 - Fano, 25-16
 - FFT, 28-9-11
 - Gauss, 47-32
 - Levinson, 25-25
 - overlap-add, 28-9
 - overlap-save, 28-8-9
 - stack, 25-16
 - Viterbi, 25-16
- Aliasing, 35-36
 - distortion, 28-6
- Allocation(s)
 - frequency, radio, 1-3-141, 1-3-5
 - to services, 1-3-141
- Alloys
 - resistivity of, 4-15-16
 - superconductivity of, 4-17
- ALOHA, 26-19, 27-29
 - reservation-, 26-21
- Alternating current bridge, fundamental, 12-3
- Altitude, effect on breakdown voltages, 49-3
- Aluminum, 4-4-5, 4-8-9
 - electrolytic capacitors, 5-20-21
- AM; See also Amplitude modulation
 - broadcast stations
 - classification of, 35-3-4
 - clear-channel, 35-3
 - coverage data, 35-5
 - field-strength requirements, 35-4
 - local, 35-3
 - regional, 35-3
 - transmission system requirements, 35-5
 - detection, 23-6

- AM—cont
 stereo
 independent sideband method, 35-6
 Kahn system, 35-6
 mixed mode method, 35-5
 Motorola system, 35-6
 quadrature method, 35-6
 transmission, 35-5-6
- Ambient conditions, for measurements, 5-5, 5-7
- Ambient temperature, test conditions, 5-5
- American National Standards Institute, 5-4
- American Standard Code for Information Interchange (ASCII), 42-10, 42-11
- Americium, 4-4-5, 4-8-9
- Ampere, 3-2
 per meter, 3-2
 per square meter, 3-2
- Amplification factor (μ), of vacuum tube, 16-20
- Amplification, in quantum system, 41-3-4
- Amplifier(s)
 AF
 class-A, design of, 17-7
 class-B, design of, 17-7-8
 bandpass, 19-22-24
 broad-band distributed, 17-13
 circuits
 RF, 17-8-13
 tetrode, 17-9, 17-11
 classifications, 19-20
 common-base, 19-12, 19-14
 common-collector, 19-11, 19-14
 common-drain, 19-11, 19-14
 common-emitter, 19-10, 19-13
 common-gate, 19-12, 19-14
 common-source, 19-10, 19-13
 crossed-field, 16-37-38
 current-differencing, 20-43
 differential, 19-17-18, 20-42-43
 Doherty, 17-11
 double-tuned interstage, 19-23-24
 emitter-follower, 19-11
 feedback, analysis, 19-22
 topologies, 19-20-21
 transistor, 19-18-22
 Klystron/TM, 17-10, 17-18
 linear, 17-10-11
 low-noise, 16-34
 noise, 21-20
 Norton, 20-43
 operating data, typical, 17-2
 operational, in active filters, 10-5-7
 operational transconductance, 20-43
 in active filters, 10-7
 operational voltage, 20-43, 20-44
 RF
 bandwidth, 17-17
 class-B, design of, 17-6-7
 class-C, design of, 17-4-6
 power gain, 17-17
- Amplifier(s)—cont
 sense, definition of, 20-10
 single-tuned interstage, 19-22-23
 source-follower, 19-11
 transistor
 single stage, 19-6-17
 frequency response, 19-6, 19-15
- Amplitude measurement, range, 12-18-19
- Amplitude-modulated pulses
 spectra of
 double polarity, 23-12
 single polarity, 23-13
- Amplitude modulation, 23-4-5
 comparison of systems, 23-6
 conventional, 23-4
 demodulation, 23-5-6
 detection, 23-5-6
 double sideband, 23-4
 phase deviation, 23-7
 single sideband, 23-4-5
 vestigial sideband, 23-4
- Amplitude response, measurement of, 12-26-27
- Amplitude-shift keying (ASK), 24-5
- Amplitude, traveling wave, 31-2
- Analog, ASICs, 20-53-54
- Analog switching systems, cellular systems, 46-21
- Analog technologies, discrete-time, 28-29-37
- Analog transmission systems, satellite systems, 27-21-22
- Analysis
 logic circuits, 43-3-4
 microwave-link, 12-26-33
 network, 12-13-17
 signal, 12-17-20
 Weilbull, 45-22-23
- Analytic geometry, solid, 47-19-21
- Analytic signal, representation of modulated waveforms, 23-3
- Analyzer
 network, 12-13
 spectrum, 12-17-19
 vector impedance, 12-11, 12-13
- Ancillary data format, digital television, 35-43
- AND gate, 43-3
 definition of, 20-6
- Angle(s)
 azimuth, geostationary satellite, 27-11-12
 divergence, collimated beam, 22-7
 elevation, geostationary satellite, 27-11-12
 modulation, multitone, 23-8
 small, approximation of, 47-11
 sum of, 47-11
 between two lines, equation for, 47-17
- Angstrom, 3-14, 18-4
- Angular modulation, 23-6
 frequency spectrum of single-tone, 23-7-8
- Anodes, electron tube, 16-22
- ANSI, 5-4
- Antenna(s)
 adaptive arrays, 32-45-47
 aperiodic loop, 32-6
 aperture distributions, 32-47-54
 design of, 32-47
 array, 32-32-47
 electronic scanning, 32-44-45
 multiple-beam, 32-41-44
 mutual impedance, 32-38-40
 thinned, 32-40-41
 tolerances, 32-41
 backfire, 32-25, 32-27-28
 bandwidth, 32-4, 32-6
 Bayliss \bar{n} one-parameter difference line source, 32-50-51
 beam angle, 33-26
 beam width, 33-26-27
 Cassegrain, 32-54-57
 definitions, 32-3-4
 dipole
 directivity, 32-14
 open sleeve, 32-15
 self-impedance, 32-13
 directional, TV broadcast stations, 35-12
 directivity, 32-3
 earth stations, 27-40-43
 effective areas of typical, 33-22-23, 33-26
 errors, radar, 36-13
 far-field region, 32-4
 field regions, 32-4
 field strength in proximity with, 49-30
 fractal, 32-32
 frequency-independent, 32-28-29
 gain, 32-3
 circularly polarized, 32-11
 equivalent, optical, 22-10
 radar, 36-12
 relative to hypothetical isotropic, 33-26
 ground systems, radial-wire, 32-15-16
 half-wave dipole, 32-13-15
 Hansen one-parameter circular source distribution, 32-53
 helix, 32-23-24
 high frequency, angles of departure and arrival, 33-8
 horn, 32-22-23, 32-24
 log-periodic, 32-30-32
 loop, 32-18-19
 spheroidal core, 32-19-20
 loss, cosecant-squared, radar, 36-12-13
 low-gain, 32-13-22
 low side-lobe patterns, measurement of, 32-52-53
 medium-gain, 32-22-32
 microstrip patch, 32-17-18
 monopole, 32-14
 near-field
 field region, 32-4
 power density, 32-7-8

- Antenna(s)—cont
 noise factor, 34-2
 noise temperature, 32-8–11, 34-5
 parasitic patch, 32-18
 path loss, 32-3
 power density, 32-7–8
 in proximity with, 49-30
 printed-circuit, 32-17–18
 pyramidal horn, 32-23
 radar, 36-11–13
 radiation resistance, 32-5
 rectangular patch, 32-18
 reflectors, 32-54–58
 regions, 32-4
 separation requirements, cellular system, 46-21
 slot, 32-19–22
 small, 32-4–6
 spacecraft, 27-31–32, 27-33
 spiral, 32-28–29
 stripline slot, 32-17
 superconducting, 32-6
 Taylor
 \bar{n} circular source distribution, 32-53–54
 \bar{n} line source distribution, 32-49–50
 one-parameter line source distribution, 32-47–48
 terminals, low-power, very small, 27-4
 tolerances, 32-41
 Yagi-Uda, 32-24–28
- Antialiasing filter, 35-37
- Antimony, 4-4–5, 4-8–9
- Antipodal signal set, 24-3
- APCVD, definition of, 20-3
- Aperiodic loop antenna, 32-6
- Aperture distributions
 antenna, 32-47–54
 design of, 32-47
- Applebaum-Howells antenna, 32-45–46
- Application specific integrated circuits (ASICs), 20-49–56
- Application specific memory, 20-54
- ARC, definition of, 20-3
- Arc lamps, 21-10
- Archimedean spiral antenna, 32-28–29
- Architecture, spacecraft, 27-14–15
- Are, 3-14
- Areas of plane figures, equations for, 47-4–5
- Argon, 4-4–5, 4-8–9
- Arithmetic and logic unit, computer, 42-15
- Arithmetic operations, error in, 47-44–45
- Arithmetic progression, 47-8
- Army-navy preferred list, radio-frequency cables, 29-30–39
- ARPANET, 26-9, 26-10
- Array(s)
 adaptive, antenna, 32-45–47
 antenna(s), 32-32–47
- Array(s)—cont
 Applebaum-Howells, 32-45–46
 directivity, 32-34–35
 electronic scanning, 32-44–45
 electronically scanned, 32-39
 grating lobes, 32-35–36
 linear feeds, 32-36–38
 LMS, 32-45–46
 mutual impedance, 32-38–40
 positioning beam of, 32-32–33
 power inversion, 32-45
 quantization lobes, 32-35–36
 traveling-wave, 32-36–37
 two-element, 32-32–33
 field emitter, 16-10–12
 gate, 20-49–54
 field programmable, 20-51, 20-53, 20-55
 programmable logic, 20-51, 20-53
 multiple-beam, 32-41–44
 thinned, antenna, 32-40–41
- Arsenic, 4-4–5, 4-8–9
- Artwork, printed circuit, 5-30–31
- ASCII codes, 42-10, 42-11
- Ash, definition of, 20-3
- ASICs
 analog, 20-53–54
 future trends, 20-54–56
 merged analog and digital, 20-55–56
- ASK *See* Amplitude-Shift Keying
- Assembler
 computer, 42-23–24
 cross, computer, 42-23
- Assembly language, computer, 42-3, 42-23–24
- Assigned-slot listen-before-transmission protocol, 26-22
- Astable multivibrator, 19-29–30
- Astatine, 4-4–5, 4-8–9
- Astronomical unit, 3-14
- Asymmetric DSL, 38-7–8
- Asynchronous
 DLC protocols, 26-13–14
 transfer mode, network, 26-36, 39-10
- ATCRBS, 37-5
- ATM networks, 39-10
 bandwidth and traffic allocations, 26-37–38
- ATM-VP allocation, network, 26-38
- Atmosphere
 refractive index, 33-3
 standard, 3-14
 technical, 3-14
- Atmospheric
 absorption, 33-20–21
 on optical communications, 22-14–16
 effects, optical communications, 22-13–19
 multipath fading, 33-19
 noise, 32-9, 34-3–4
 transmission, over sea water, 41-2
- Atomic
 mass, relative, 4-4, 4-6
 number, 4-3, 4-4, 4-6, 4-8, 4-10
- Atomic—cont
 radius, 4-3, 4-4, 4-6
 time, 1-143
- Attenuation
 cables, 29-29–30
 comparison of measurement systems, 11-8
 curve, for quadratic factor, 15-7, 15-8
 distortion, telephony, 2-10
 free-space, 33-22
 gaseous, 33-20
 lines, 29-29–30
 measurement, 11-5–10
 broadband, 11-7
 fixed frequency, 11-5–6
 swept or stepped frequency, 11-9–10
 microstrip lines, 29-26–27
 rain, 33-20–21
 satellite communications, 27-36–37
 transmission lines, at UHF, 29-15
 in waveguide beyond cutoff, 30-7
- Attenuator(s)
 connectors, 11-5, 11-8
 definition of, 11-2
 pads, resistance values, 11-4
 resistance networks for, 11-3–5
 resistive, typical designs, 11-2–5
- Audio
 frequency
 inductors, 13-21–23
 laminated-core, 13-21
 toroidal core, 13-21–23
 response
 AM, 35-5
 FM broadcast, 35-10
 transformers, 13-12–14
 cores for, 13-12–13
 design of, 13-13–14
 recording, digital, 35-27
 signals, WWV and WWVH, by telephone, 1-157
 transformers
 impedance match, 13-14
 insertion loss, 13-14
- Audio Engineering Society, 35-28
- Audio frequencies, standard, WWV and WWVH, 1-151–152
- Auditory system, 40-6–7
- Aural intercity stations, 35-23
- Aural transmitter, TV, 35-13–14
- Automatic dependence surveillance system, 37-8
- Automatic detection and track (ADT), 36-22
- Automatic repeat request (ARQ), 25-14
- Automobiles, vibration in, 49-25
- Auxiliary broadcast services, 35-23–24
 frequency sharing, 35-24
- Auxiliary stations
 low power, 35-24
 television, 35-23
- Avalanche breakdown, 18-11
 definition of, 20-3

- Average power, periodic waveform, 7-8
 Average terrain, antenna site, 35-12
 Avogadro's number, 18-4
 Azimuth angle, geostationary satellite, 27-11-12
- B**
- Backfire, antenna, 32-25, 32-27-28
 Background, definition of, 20-3
 Background radiation, optical communications, 22-12-13
 Backlash-type nonlinearity, 15-25-27
 Backward transition matrix, 25-12
 Backward-wave oscillators, 16-33-34
 Backward-wave tubes, gyro, 16-38-39
 Balance
 capacitance, mutual-inductance, 12-7
 Felici mutual-inductance, 12-7
 Ball bond, definition of, 20-3
 Bandpass
 amplifier, 19-22-24
 filter, 8-1-9
 gain at resonance, 8-2-3
 resonance selectivity far from, 8-3
 selectivity near, 8-3, 8-6
 Bands, radio frequency, nomenclature, 1-2, 1-3
 Bandwidth
 allocation, ATM networks, 26-37-38
 antenna, 32-4, 32-6
 increasing, 32-6
 considerations in multitone FM, 23-9
 detector, spectrum analyzer, 12-20
 fiber, 22-24
 IF filter, spectrum analyzer, 12-19-20
 limitations, high frequency transmissions, 33-9-11
 resonant arrays, 32-36
 RF amplifiers, 17-17
 transponders, 27-17-18
 Yagi-Uda antennas, 32-25
 Banyan network, 39-9
 Bar, 3-14
 Barium, 4-4-5, 4-8-9
 Barn, 3-14
 Barrel, 3-14
 Barrier layer, definition of, 20-3
 Bartlett window function, 28-11
 Barye, 3-14
 Base 2, logarithms of, table of, 48-5
 Base, number system, 42-4-5
 Base-to-collector current amplification factor, 18-19
 Base transport factor, transistor, 18-19
 Baseband measurements, 12-30-33
 Baseband signal
 in-phase component, 24-5
 quadrature component, 24-5
 sets, 24-2-4
 Baseline network, 39-9
 Bayes criterion, 24-14
 Bayesian statistics, 45-23-25
 Bayliss \bar{n} one-parameter difference line source, antenna, 32-50-51
 BCD
 -decimal equivalents, 1-157
 definition of, 20-3
 representation, 42-10
 time code, WWV and WWVH, 1-152-154
 Beacon system, air traffic control radar, 37-5
 Beacons, radio, 37-3-4
 Beads, dielectric, characteristic impedance, 29-22
 Beam
 angle, antenna, 33-26
 deviation factor, 32-57
 -forming networks, antenna array, 32-41-44
 modes, Hermite-Gaussian, 41-8
 width, antenna, 33-26-27
 Beamwidth, conversion, universal horn, 32-55-56
 Becquerel, 3-2
 BER *See* Bit Error Rate
 Berkelium, 4-4-5, 4-8-9
 Berlekamp-Massey algorithm, 25-15
 Beryllium, 4-4-5, 4-8-9
 BESOI, 20-103
 Bessel functions, 47-41-43
 table of, 48-9-10
 β , of transistor, 18-19
 Beta ratio, definition of, 20-3
 Bias
 forward, definition of, 20-3
 reverse, definition of, 20-3
 techniques
 bipolar transistor, 19-2-5
 field-effect transistor, 19-5
 Bilinear z-transformation, 28-14
 Binary
 baseband data transmission, receivers for, 24-13-16
 -coded decimal (BCD) number system, 42-10
 coded decimal, definition of, 20-3
 direct-sequence, spread-spectrum communications, 24-24-26
 division, 42-9
 entropy function, 25-9
 frequency-shift keying (FSK), 24-8-9
 logic, definition of, 20-3
 multiplication, 42-8-9
 orthogonal signaling, 22-5
 phase shift keying (BPSK), 24-5-6, 27-23
 radix, 42-5
 RF signals, coherent receivers for, 24-16
 Binomial distribution, 44-4, 45-11
 Bipolar junction transistors, 18-16-21
 Bipolar logic families, 20-30-35
 Bipolar power transistors, 14-29-30, 14-32
 Bipolar technology, definition of, 20-3
 Bipolar transistors
 bias techniques, 19-2-5
 insulated-gate, 14-7-8
 Biquads, 10-15-23
 Bird beak, definition of, 20-3
 Bismuth, 4-4-5, 4-8-9
 Bistable element, definition of, 20-3-4
 Bistable multivibrator, 19-29
 Bit
 computer, definition of, 42-4
 definition of, 20-4
 energy, 25-19-20
 error
 probability of, 22-4, 24-23
 rate (BER), 25-19-20
 satellite systems, 27-21
 ratio (BER), definition of, 38-11
 -oriented DLC protocols, 26-16-18
 -parallel interfaces, 35-40
 -serial interfaces, 35-40
 BJT *See* Bipolar Junction Transistor
 Blackbody radiation, 21-3, 21-5-8
 basic laws and constants, 21-6
 Blackbody responsivity, 21-17
 Blackbody temperature, 22-12
 Blackman window function, 28-11
 Blackouts, radio, 1-155-156
 Block codes, 25-4, 25-15
 fixed-length, 25-10
 variable-length, 25-10
 Block coding, modulation technique, 22-4
 Block diagrams, 49-26
 British, 49-28
 Blocking networks, 39-6
 traffic-carrying capability, 39-6
 Board foot, 3-14
 Boat, definition of, 20-4
 Bode diagrams, 15-6-7
 Boiling point, of elements, 4-8, 4-10
 Bolometer, 21-21
 Boltzmann's constant, 18-4, 27-14
 Boolean algebra, 43-2
 Boost converter, 14-28-29
 Boron, 4-4-5, 4-8-9
 Bouguer's law, 22-14
 Bounded-input bounded-output (BIBO) system, 15-4
 BPSK *See* Binary Phase-Shift Keying
 Breakdown diode, 18-12
 Breakdown, reverse, semiconductor, 18-10-11
 Breakdown voltages
 effect of altitude on, 49-3
 spark-gap, 49-3
 Bridge
 alternating current, fundamental, 12-3
 circuit, capacitor balance, 12-4
 four-terminal-pair, 12-11
 Hay, 12-5
 impedance, 12-3-13
 Maxwell, 12-5
 Owcn, 12-5
 ratio-arm, 12-10

- Bridge—cont
 rectifier, 14-13
 full-wave, single phase, 14-13
 three-phase full-wave, 14-13
 resonance, 12-5
 Schering, 12-5
 series-resistance-capacitance, 12-4
 substitution method, high impedances, 12-5-6
 Wayne-Kerr, 12-10-11
 Wheatstone, 12-3
 Wien, 12-4
- Brightness, typical approximate values, 16-48
- British thermal unit, 3-14
- Broadband distributed amplifiers, 17-13
 integrated services digital network (B-ISDN), 26-36
- Broadcast
 bus, 26-6
 channels, 25-26-27
 receivers, noise measurements, 34-8-9
 routing, networks, 26-25
 services, auxiliary, 35-23-24
 stations
 classification of, 35-3-4
 remote pickup, 35-23
- Broadcast Television Sound Committee (BTSC), 35-14
- Broadcasting
 FM, 35-6-11
 geographical zones for, ITU, 35-24-25
 international standards for, 35-46-49
 international, United States, 35-24-26
 stations, international, frequencies, 35-24
 television, 35-11-14
- Broadcasting Satellite Service (BSS), 27-13
- Broadcasts, frequency, 1-148-162
- Broadcasts, time, 1-148-162
- Bromine, 4-4-5, 4-8-9
- BTS, definition of, 20-4
- Buck converter, 14-19, 14-28
- Buffer memory, definition of, 20-4
- Building materials, absorption coefficients, 40-16
- Bursty, communications, 26-3
- Bus
 broadcast, 26-6
 definition of, 20-4
- Buses
 computer, 42-30
 linear, network, 26-33-34
- Bushel, 3-14
- Busy-tone
 multiple access (BTMA), 26-20
 proposed standard, 2-18
- Butler beam-forming network, 32-41-44
- Butterworth low-pass filter, 9-14
 design, 9-18, 9-20-21
- Byte, definition of, 20-4, 42-4
- C
- Cable, 3-14
 ready, TV receivers, 35-15
 systems
 bidirectional, 35-15
 fiber optic usage, 35-16-17
 television, 35-14-17
 channel frequencies, 35-15, 35-21-22
- Cables
 attenuation, 29-29-30
 power rating, 29-30-31
 radio-frequency, army-navy preferred list, 29-30-39
 telephone, characteristics of, 38-12
- CAD/CAE/CAM, definition of, 20-4
- Cadmium, 4-4-5, 4-8-9
- Calcium, 4-4-5, 4-8-9
- Calculus, differential, 47-21
- Caliber, 3-14
- Calibration factor, K_B , definition of, 12-25
- Californium, 4-4-5, 4-8-9
- Call blocking, 46-8
- Call completion, 46-8
- Call drops, 46-8
- Call intensity, 39-11
- Call termination, 46-8
- Call-waiting tone, proposed standard, 2-18
- Caloric, 3-14
- Camera tube
 image orthicon, 16-50-52
 vidicon, 16-53-54
- Canceller, N-pulse, 36-16
- Candela, 3-2
 per square meter, 3-2
- Canonical
 expressions, 43-5-6
 product, 43-5-6
 sum, 43-5-6
- Capacitance
 balance, mutual-inductance, 12-7
 coaxial cylindrical capacitor, 6-15
 direct, measurement of, 12-6-7
 parallel-plate capacitor, 6-14-15
 semiconductor junction, 18-11
- Capacitive loss, inductors, 13-27
- Capacitively coupled RF plasma discharges, 20-92-93
- Capacitor(s), 5-16-26
 in active filters, 10-5
 balance, bridge circuit, 12-4
 ceramic, 5-22-23
 color code, 5-23, 5-24
 temperature coefficient, 5-23
 charge and discharge, 6-15-16
 classes of, 5-17
 coaxial cylindrical, capacitance of, 6-15
- Capacitor(s)—cont
 definitions, 5-16-17
 electrolytic, 5-18-21
 energy stored, 6-15
 film, 5-18
 input, filter design, 14-41, 14-42
 mica, 5-25-26
 paper-foil, 5-23
 parallel-plate, capacitance of, 6-14-15
 plastic film, 5-17-18
 polypropylene, 5-17
 polystyrene, 5-17
 in series with unknown, measurement of, 12-6
 tantalum electrolytic, 5-21
 voltage variable, 18-12, 18-13
- Carat, 3-14
- Carbon, 4-4-5, 4-8-9
 composition resistors, 5-15
- Carcinotron, 16-36-37
- Carrier
 amplitude/carrier shift, AM, 35-5
 concentrations, 18-5
 -derived channel signaling, 38-17
 digital loop, 38-5-6
 excess minority, definition of, 20-5
 -frequency transformers, 13-15
 design example, 13-15-16
 noise, measurement of, 12-33
 offset, cochannel TV stations, 35-12
 sense multiple access (CSMA), 26-19-20
 storage, rectifiers, 14-4
 systems, telephone, 38-13-16
 -to-noise power ratio, 27-19-21
- Carriers
 drift of, 18-5-6
 excess, 18-6-8
 minority and majority, 18-5
- Cascade realization, filter sections, 10-23-25
- Cascade, rinse, definition of, 20-9
- Cascaded networks, noise factor, 34-11
- Cassegrain reflector system, 32-54-57
- Cathode
 electron tube, 16-21
 follower circuit, 17-8-9
 impregnated tungsten dispenser, 17-18
 materials, tube, 16-3-4
 -ray tubes, 16-55-57
 beam focusing, 16-55
 deflection of, 16-55
 operation of, 16-55
 storage, 16-55-57
 bistable, 16-57-58
- Cauer-parameter
 low-pass filters, 9-4-5
 design, 9-5-13
- Cavity(ies)
 coupling techniques, 30-22-24
 coupling through, equations for, 30-22
 equations, 30-21-22
 optical, 41-6-8

- Cavity(ies)—cont
 stability of, 41-7
 resonant, 30-20-24
 cylindrical, 30-21
 rectangular, 30-20
 spherical, 30-21
 tuning, effect of temperature and humidity on, 30-21
 waveguide, simple, 30-24
- CCIR; See also International Radio Consultative Committee recommendations, 2-3
- CCITT; See also International Telephone and Telegraph Consultative Committee recommendations, 2-3
 telephone, 38-8
 and telegraphy, 2-13
 telephone signaling systems, 2-14
 and telephony, 2-13
 X-20bis recommendation, 26-11-13
 X-21 protocol, 26-11-13, 26-28-29
 X-21bis recommendation, 26-11-13
 X-24 recommendation, 26-11-13
- CDI, definition of, 20-4
- CDMA, 27-20-21
- Cell
 definition of, 20-4
 ECL storage, definition of, 20-5
 multiple-emitter, definition of, 20-7
 site, cellular phone, 46-4-5
 solar, 18-13
 thermopneumatic, 21-21
- Cells
 photoconductive, 21-28-29
 photovoltaic, 21-28-29
 standard, ASICs, 20-50-51, 20-52
- Cellular
 CDMA system, 46-24
 signals, combining techniques, 46-21
 system(s)
 antenna separation requirements, 46-21
 call blocking, 46-8
 call completion, 46-8
 call drops, 46-8
 call termination, 46-8
 cochannel interference reduction factor, 46-18
 converters, 46-7
 data-link network, 46-5
 design aspect, 46-18-19
 digital, 46-22-24
 Digital AMPS, 46-22
 elements, 46-4-7
 enhancr, 46-6
 4G, 46-24
 frequencies, 46-3
 frequency reuse, 46-17-19
 factor, 46-20
 frequency spectrum, 46-3-4
 handoff procedure, 46-8
 high-speed data channel, 46-25
 history of, 46-4
 hooks and extensions, 46-7
 major world, 46-3-4
- Cellular—cont
 mobile-originated call, 46-7
 mobile-unit initialization, 46-7
 network-originated call, 46-7-8
 operation of, 46-7-8
 radio capacity, 46-19-21
 switching equipment, 46-21-22, 46-23
 3G, 46-5, 46-6, 46-24
 traffic models, 46-21-22
- Celsius
 degree, 3-2
 temperature, 3-2
- Centimeter of mercury, 3-14
- Centimeter of water, 3-14
- Centosymmetric materials, 41-10
- Central processing unit (CPU), 42-10
 microprocessor, 20-37-38
- Centralized routing, network, 26-24-25
- Centrally controlled demand assignment, multiaccess link control, 26-20-21
- Ceramic capacitors, 5-22-23
- Ceramic-envelope tubes, advantages of, 16-19-20
- Ceramic flat package, ICs, 20-29
- Ceramic packages, ICs, 20-28-29
- Cerium, 4-4-5, 4-8-9
- Cermet-element resistors, 5-15
- Cesium, 4-4-5, 4-8-9
- Cesium oscillators, 1-146-147
- Chain, 3-14
- Channel(s)
 average mutual information, 25-13-14
 broadcast, 25-26-27
 capacity, 25-14
 codeword, 25-3
 coding, for discrete noisy, 25-11-12
 continuous, 25-17-22
 capacity of, 25-18
 decoder, 25-4
 designations, television, 35-11-12
 discrete, 25-12
 Gaussian noise, additive, 25-18-19
 memoryless, 25-12
 noisy, 25-12
 transition probability, 25-12
 two-way, 25-26
 waveform, 25-19
- Character codes, 42-10
- Character generator, definition of, 20-4
- Character oriented, DLC protocols, 26-15-16
- Characteristic impedance, of transmission lines, 29-21-25
 insertion loss, 11-2
- Characteristics, components, 5-4-5
- Charge carriers, semiconductors, 18-3-5
- Charge coupled device (CCD), 20-46, 20-47, 20-48
 definition of, 20-4
- Charge injection device (CID), 20-46
- Charge, mobile, definition of, 20-7
- Charge potential, image, 16-8
- Chart, projective, 31-6
- Charts, conformal, 31-6
- Chebyshev low-pass filters, 9-14
 design, 9-15-17
- Chemical vapor deposition, thin films, 20-13
- Chi-Square distribution, 44-4-5, 45-11
- Chi-Square test, 44-4
 fitting a distribution using, 45-21
- Chip enable, definition of, 20-4
- Chlorine, 4-4-5, 4-8-9
- Chokes, swinging, 14-41
- Chromatic dispersion, 22-24
- Chromium, 4-4-5, 4-8-9
- Circle
 enclosing smaller circles, 49-17-18
 arrangement of center, 49-17
 diameter of, 49-17, 49-19
 large number of, equation, 49-18
 equations for, 47-4, 47-18
 sector of, equations for, 47-4
 segment of, equations for, 47-4
- Circuit(s)
 characteristics, telephone, 2-9-10
 coupled, coefficient of coupling, 8-4-5
 elements, waveguide, 30-18-19
 noise, telephony, 2-11-13
 phase shift, single-tuned, 8-6-7
 power grid tube, classification, 17-8
 printed, 5-26-31
 pulse, 19-26-30
 rectifier, 14-12-13
 selectivity, single-tuned, 8-6-7
 sequential, 43-16-18
 switch-mode converter control, 14-34-35
 switched networks, 26-5-6
 switching, 26-2
 networks, 39-9
 two-stage gate, 43-12
- Circular mil, 3-14
- Circular patch, antenna, 32-18
- Circular polarization, 32-11-12
- Circular waveguides, 30-5-6
- Clamps, 14-32, 14-34
- Class-A AF amplifier, design of, 17-7
- Class-A operation, power grid tubes, 17-2
- Class-AB AF amplifier, design of, 17-7-8
- Class-AB operation, power grid tubes, 17-2
- Class-B AF amplifier, design of, 17-7-8
- Class-B operation, power grid tubes, 17-2
- Class-B RF amplifier, design of, 17-6-7
- Class-C operation, power grid tubes, 17-2
- Class-C oscillator, design of, 17-4-6
- Class-C RF amplifier, design of, 17-4-6
- Class N, definition of, 20-4
- Classification
 of the elements, periodic, 4-3, 4-12
 power grid tube, circuit, 17-8

- Cleaning
 metallic, 20-91
 organic, 20-89-91
 procedure, substrate, 20-12
 reagents, IC processing, 20-12
 surface
 chemical oxide, by vapor-phase HF, 20-86-87
 IC manufacturing, 20-82-85
 low-temperature, germane-based, 20-85-86
 native oxide, by vapor-phase HF, 20-86-87
- Clear-channel AM broadcast stations, 35-3
- Clear sky noise temperature, 27-35-36
- Climate conditions, rain rate, US and Canada, 33-20, 33-22
- Climate regions, global rain rate, 33-20, 33-22
- Clo, 3-14
- Clock, definition of, 20-4
- Closed-loop system, evaluation of, 15-4-5
- Clutter
 fluctuations, radar, 36-17
 radar, 36-18, 36-20-21
 surface, radar, 36-8
 volume, radar, 36-8-9
- CMOS well formation, 20-100
- CMTT *See* Joint Committee on Television Transmission
- C_p , definition of, 20-4
- Coaxial
 connectors, list of, 11-8
 cylindrical capacitor, capacitance of, 6-15
 line, single, characteristic impedance of, 29-21
- Cobalt, 4-4-5, 4-8-9
- Cochannel interference reduction factor (CIRT), cellular system, 46-18
- Code division multiple access (CDMA) multiple system, 46-24
- CODECs, switched capacitor filters for, 28-31-32
- Coded trellis signaling, 27-24
- Code(s)
 Adlcr-Hassner, 25-7
 block, 25-4, 25-15
 fixed-length, 25-10
 variable-length, 25-10
 character, 42-10
 convolutional, 25-4, 25-15-16
 country, international telephony, 2-13
 Elias, 25-9
 error-control, 25-14
 Franaszek, 25-6
 Hamming, 25-15
 Huffman, 25-10
 Lempel-Ziv, 25-9
 Morse, 25-4-6
 multifrequency numerical
 four-by-four, 2-14, 2-16
- Code(s)—cont
 two-out-of-six, 2-14, 2-16
 quality assurance, 45-25
 Reed-Solomon, 25-15
 sliding block, 25-4
 source compaction, 25-7-11
 source compression, 25-25-26
 trec, 25-4
 variable length, 25-10-11
 universal, 25-9
- Codeword, channel, 25-3
- Coding
 color, wire, 49-17
 component value, 5-8
 for discrete noisy channels, 25-11-12
 for noiseless channels, 25-4-7
 semiconductor diode type number, 5-7
 signal, digital television, 35-38, 35-40
- Coefficient(s)
 absorption, 40-14-16
 coupled circuits, 8-4-5
 of coupling, 8-2
 reflection
 definition, 31-2
 measurement of, 31-2
- Coherent
 receivers, 24-16
 sources, optical communications, 22-7-8
- Coils
 audio frequency, core material characteristics, 13-22
 loosely coupled, measurement of coupling coefficient, 12-9-10
 loss mechanisms, 13-26-28
 unshielded, Q of, 6-5
- Cold-cathode gas tubes, 16-42
- Collector junction, 18-18
- Collimated beam
 divergence angle, 22-7
 signal intensity, 22-7
- Collision ionization, electron tube, 16-19
- Color code
 ceramic capacitors, 5-23, 5-24
 components, 5-4, 5-5
 resistors, 5-10
 wire, 49-17
- Color television system, NTSC-M, 35-40
- Colpitts oscillator, 19-25
- Column decoder, definition of, 20-4
- Combinational switching circuits, 43-3
- Combinations, table of, logic circuit, 43-4
- Common
 -base amplifier, 19-12, 19-14
 -base configuration, 18-18
 -carrier hierarchies, 38-13-16
 carrier transmissions, definitions, 38-13
 -channel signaling, telephone, 38-16
 -collector amplifier, 19-11, 19-14
- Common—cont
 -drain amplifier, 19-11, 19-14
 -emitter amplifier, 19-10, 19-13
 -emitter circuit, 18-20
 -gate amplifier, 19-12, 19-14
 lasers, 41-8-9
 -source amplifier, 19-10, 19-13
- Communications
 binary direct-sequence spread-spectrum, 24-24-26
 bursty, 26-3
 digital audio radio service (DARS), 27-3, 27-5
 direct to home (DTH), 27-3, 27-5
 earth-space, 33-33
 frequency-hop spread-spectrum, 24-26
 games, 25-27-31
 networks, computer, types of, 26-6-7
 satellite(s)
 choice of frequency, 27-13
 systems, classification of, 27-3-6
 transmission system
 considerations, 27-19-21
 transponders for, 27-15-19
 spread-spectrum, 24-24-26
 systems, satellite vs fiber-optic, 27-5
- Communications Act of 1934, 1-2-3
- Comparators, IC, 20-45
- Compatibility, electromagnetic, 12-35-36
- Compiler, computer, 42-3
- Compilers, silicon, 20-56
- Complement
 logic, definition of, 43-3
 number systems, 42-7
- Complementary MOS, definition of, 20-4
- Complexity, definition of, 20-4
- Components
 characteristics, 5-4-5
 coding of signals, digital television, 35-41-42
 color code, 5-4, 5-5
 preferred values, 5-4, 5-6
 reliability, 45-9-11
 tolerance, 5-4
 value coding, 5-8
 voltage rating, 5-4
- Composite encoded signals (NTSC & PAL), digital television, 35-40-41
- Composite television signal, 35-18-19
- Composite triple beat interference, 35-16
- Compounds, superconductivity of, 4-17
- Compression
 codes, source, 25-25-26
 pulse, radar, 36-18, 36-19
- Computer
 accumulators, 42-11
 address-register indirect addressing, 42-19
 addressing modes, 42-17-20
 -aided analysis, feedback control systems, 15-33

- Computer—cont
- aided design, feedback control systems, 15-33
 - arithmetic and logic unit, 42-15
 - assembly language, 42-3, 42-23-24
 - based addressing, 42-19
 - bit, definition of, 42-4
 - buses, 42-30
 - byte, definition of, 42-4
 - communications networks, types of, 26-6-7
 - compiler, 42-3
 - control of instruments, 12-33-35
 - data types, 42-4
 - debugger, 42-3
 - direct memory access, 42-14
 - general-register processor, organization of, 42-14-16
 - hardware, 42-3
 - high-level language, 42-3
 - immediate addressing, 42-17
 - input/output, 42-14, 42-30-33
 - instruction formats, 42-16-17
 - instruction set, 42-17
 - instructions, 42-20-22
 - interfaces, 42-30
 - interpreter, 42-3
 - machine language, 42-3
 - program, 42-22-24
 - mass storage devices, 42-14
 - memory, 42-12-14
 - networks
 - physical layer, 26-9-13
 - structure of, 26-2
 - nibble, definition of, 42-4
 - organization, basic, 42-10-14
 - peripheral devices, 42-14, 42-30
 - ports, 42-31-32
 - processors, 42-10-12, 42-14-22
 - effective address register, 42-15
 - general register, 42-15
 - instruction register, 42-15
 - program counter, 42-15
 - registers, 42-15
 - program
 - operation of simple, 42-25-30
 - subroutines, 42-25-26
 - programming, 42-14-22
 - programs, patent protection, 49-32
 - register-direct addressing, 42-17
 - relative addressing, 42-19
 - simulator, 42-3
 - software, 42-3-4
 - stack machine, 42-12
 - text editor, 42-3
 - use in reliability, 45-25
- Conditional
- branch instruction, computer, 42-21
 - code register, computer processor, 42-15
 - probability function, 44-2-3
- Conducting materials, 4-3
- Conduction cooling, electron tube, 16-17-18
- Conductivity, thermal, of elements, 4-9, 4-11
- Conductors
- applications and insulations, NEC, 49-12
 - current capacity, printed circuit, 5-27, 5-29
 - NEC allowable current-carrying capacities, 49-14
- Cone, equations for, 47-6
- Confidence level, 45-15
- Confidence limits, 45-15-21
- Conformal charts, 31-6
- Conformal model, hyperbolic trigonometry, 47-14
- Congestion control
- ATM, 26-38-39
 - networks, 26-25-27
- Conic frustrum, equations for, 47-6
- Conic surface, development of, equations for, 47-6
- Conical spiral antenna, 32-29
- Connectors
- attenuators, 11-5, 11-8
 - coaxial, list of, 11-8
 - optical fiber, 22-24-25
- Constant humidity tests, 5-6
- Constant-voltage transformers, 13-10-11
- Constant(s)
- Boltzmann's, 18-4
 - fundamental physical, 3-7-13
 - physical, 18-4
 - Planck's, 18-4, 21-2
- Continuous channels and sources, 25-17-22
- Continuous count, 12-23
- Continuous modulation, 23-3-10
- Continuous-phase frequency-shift keying (CPFSK), 24-9-10
- Continuous time (CT) signal, 28-2
- Control system(s)
- component models, 15-2-4
 - definition of, 15-2
 - digital, 15-15-21
 - state variable analysis, 15-17-18
- feedback
- computer-aided analysis, 15-33
 - computer-aided design, 15-33
 - disturbance inputs, 15-11
 - nonlinear, 15-22-32
 - robustness, 15-11-12
 - sensitivity, 15-11-12
 - tracking accuracy, 15-10-11
 - system performance, 15-4-12
- Controller
- current-mode control, 14-36, 14-37
 - design
 - LQ-optimal control, 15-15
 - methods, classical design, 15-12
 - pole placement, 15-14
 - lag, 15-12
 - lead, 15-12
 - proportional-integral-derivative, 15-12
 - switch-mode converter, 14-34-35
- Controller—cont
- voltage-mode, 14-35-36
- Conversion
- factors, 3-13-17
 - physical constants, 18-4
 - NPR to SNR, 12-32-33
 - number systems, 42-6
 - wavelength-frequency, 1-2
- Converter(s)
- A/D, 35-35-36
 - basic, 14-19
 - boost, 14-28-29
 - buck, 14-19, 14-28
 - cellular system, 46-7
 - circuits, ac-dc, 14-12-39
 - current-mode control, 14-36, 14-37
 - D/A, 35-35-36
 - definition of, 13-11-12
 - flyback, 14-28-29
 - forward-mode, 14-19, 14-28
 - general impedance, 10-12-13
 - /inverter transformers, 13-11-12
 - step-down, 14-19, 14-28
 - step-up, 14-28-29
 - switch-mode, 14-34-35
 - thermal, 12-24
 - voltage-mode control, 14-35-36
- Convolution, finite, equation, 28-3
- Convolutional codes, 25-4, 25-15-16
- Cooling
- electron tube
 - conduction, 16-17-18
 - evaporative, 16-17
 - forced-air, 16-16-17
 - radiation, 16-16
 - water, 16-16
- Coordinate switch, network, 39-2
- Coordinate switches, single-stage, 39-2-3
- Coordinate switching networks, 39-2-10
- definitions for, 39-2
 - properties of, 39-2
- Coordinated Universal Time (UTC), 1-142-144, 1-151
- Coordinates
- rectangular, 47-35-36
 - transformation of, 47-17-18
- Coplanar transmission lines, 29-28-29
- Copper, 4-4-5, 4-8-9
- regulation, power transformer, 13-6
- Copyright, protection of inventions, 49-34
- Cord, 3-14
- Cordless Telephone-2 (CT-2), 46-24
- Cords
- flexible
 - current-carrying capacity, NEC, 49-13
 - NEC data, 49-13
- Core(s)
- for audio-frequency transformers, 13-12-13
 - inductor, magnetic, 13-3
 - loss, magnetic, 13-26
 - magnetic, types of, 13-3-5

- Core(s)—cont
material
 audio-frequency coils, 13-21
 wideband transformers, 13-15
 reactor, magnetic, 13-3
- Corona, transformers, 13-10
- Corporate feeds, linear arrays, 32-38
- Correlation receiver, 24-15-16
- Corrosion, finishes, 49-5, 49-8
- Cosecant-squared antenna loss, radar, 36-12-13
- Cosine integrals, 47-39-40
- Cosines
 for angles, law of, 47-12
 direction, 47-19-20
 hyperbolic, table of, 48-3
 law of, 47-11, 47-15
 of sides, law of, 47-12
- Cosmic noise, 27-12-13, 34-4-5
- Coulomb, 3-2
 per cubic meter, 3-2
 per kilogram, 3-2
 per square meter, 3-2
- Counter, frequency, 12-21-22
- Country codes, international telephony, 2-13
- Coupling
 to cavities and loaded Q, 30-22
 coefficient, 8-2
 loosely coupled coils,
 measurement of, 12-9-10
 techniques, cavity, 30-22-24
 through cavity, equations for, 30-22
- Coverage
 data, AM broadcast stations, 35-5
 estimated, FM broadcasting, 35-9
 radar, 36-13-14
 step, definition of, 20-10
- CPFSK *See* Continuous-Phase Frequency-Shift Keying
- Cramer-Rao bound, 25-23-24
- Cramer's rule, 47-32
- Criterion
 Bayes, 24-14
 minimax, 24-14
- Cross assembler, computer, 42-23
- Cross talk
 definition of, 38-10-11
 linear, telephony, 2-10
- Crossed-field
 amplifiers, 16-37-38
 tubes, 16-34-38
 voltage-tuned oscillator, 16-36-37
- Crosspoints, network, 39-2-3
- CRT, digital storage, 12-14
- Cryogenic magnetron etch, 20-97-99
- Cryptography, 25-30-31
- Crystal, single, definition of, 20-10
- Crystal, uniaxial, 41-16-17
- Cu-low k, definition of, 20-4
- Cubic
 meter, 3-2
 per kilogram, 3-2
- Cubit, 3-14
- Cumulative distribution function (CDF), 44-2, 45-2
- Cumulative distribution function (CDF)—cont
 multipath fading, 46-12, 46-13, 46-14
- Cumulative law, vector analysis, 47-35
- Cup, 3-14
- Curie, 3-14
- Curium, 4-4-5, 4-8-9
- Current
 capacity, printed circuit conductors, 5-27, 5-29
 -carrying capacities
 conductors, NEC allowable, 49-14
 flexible cords, NEC, 49-13
 density, macroscopic, 16-10
 -differencing amplifier, 20-43
 drift, definition of, 20-5
 equations, transmission lines, 29-5-6
 forward, definition of, 20-4
 gain, dc, power transistor, 14-5
 mirrors, 19-17-19
 -mode control, 14-36, 14-37
 motor full-load running, NEC, 49-15
 protective resistance, 16-13
 ratings, rectifiers, 14-3
 sources, transistor amplifiers, 19-17-19
 standards, TTL, 20-34
 -voltage relationship, 16-12
- Curve, V, 12-32
- Curves, space, vector equations for, 47-37
- CVD, definition of, 20-4
- CW lasers, 41-8
 common, 41-13-14
- CW oscillation, 41-4-5
- Cycle time, definition of, 20-4
- Cyclic redundancy check, 26-15
- Cycloconverters, 14-50-50
- Cyclotron resonance microwave tubes, 16-38-39
- Cylinder, equations for, 47-5
- Cylindrical coordinates, vector equations for, 47-36
- Czocharlski, definition of, 20-4

D

- D/A converter, 35-35-36
- D-I-T FFT, 28-9-10, 28-12
- D layer, wave propagation in, 33-7
- D**, optical detectors, 21-15
- D*, optical detectors, 21-14-15
- DAMA, 27-21
- Damascene processing
 definition of, 20-4
 ICs, 20-22-26
- Damping, vibration reduction, 49-24
- Data
 circuit terminating equipment (DCE), 26-9
 format, ancillary, digital television, 35-43
 link control
 high-level, 26-16
- Data—cont
 layer, network, 26-13-18
 synchronous, 26-16
 link network, cellular system, 46-5-6
 pulse, energy per, 24-4
 rate, maximum, optical fiber, 22-25
 register, computer processor, 42-15
 sources, generic, 45-13
 terminal equipment (DTE), 26-9
 transfer modes, 26-16
 transmission, binary baseband,
 receivers for, 24-13-16
 types, computer, 42-4
- Day, 3-14
- Day-night
 sound exposure (DNSE), 40-11
 sound level (DNL), 40-11-12
- dBrn, definition of, 38-9
- DBS-TV, 27-3
- DBSA, 27-3
- Dc-ac inverters, 14-47-50
- Dc motors, 15-2-3
- Dc-output power supplies, 14-43-46
- DC signaling, 38-16
- De-emphasis, signal-to-noise improvement through, 23-10
- Deadlock, store-and-forward networks, 26-26
- Debugger, computer, 42-3
- Decay, materials, 49-8
- Decibels, 3-6
- Decimal-BCD equivalents, 1-157
- Decimation-in-frequency FFT, 28-9
- Decision statistic, 24-2
- Decision theory, 25-22-25
- Decoder
 channel, 25-4
 definition of, 20-4
- Defense Electronics Supply Center, 5-4
- Definite integrals, 47-27-29
- Deglaze, definition of, 20-4
- Degraded diversity system, 25-27
- Degrec, 3-14
- Del operator, 47-35
- Delay distortion, group, telephony, 2-9
- Delay equations, traffic networks, 39-13-14
- Delta
 modulation, 23-11, 23-18-19
 signal-to-noise ratio, 23-18-19
 routing, network, 26-25
- Demand-assigned multiple access (DAMA) systems, 27-24
- Demand assignment with distributed control, network, 26-21
- Demodulation, amplitude modulation, 23-5-6
- DeMorgans theorem, 43-10
- Denier, 3-14
- Densities
 of elements, 4-8
 various metals, 4-33
- Density, 4-8, 4-10
- Department of Transportation, 37-2
- Dependent data, remote compaction, 25-27

- Depletion region, definition of, 20-4
 Depletion region, semiconductor, 18-9
 Depolarization, 27-40
 Deposition
 spin-on, thin films, 20-14
 sputter, thin films, 20-14
 Derivatives, list of, 47-21
 DESC, 5-4
 Describing function
 calculation of, 15-24-25
 technique, nonlinear elements, 15-24-29
 Descum, definition of, 20-4
 Design examples, filter, low-pass, 9-19, 9-22, 9-24-26
 Design, IC, 20-26-28
 Design reviews, 45-6
 Destination indicators, international, telegraph, 2-13
 Detection
 amplitude modulation, 23-5-6
 DSB, 23-5-6
 envelope, 23-6
 SSB, 23-6
 synchronous, 23-6
 Detectivity, optical detectors, 21-14
 Detector(s)
 bandwidth, spectrum analyzer, 12-20
 classes, laser, 41-2
 envelope, 24-18
 optical, 21-14-30, 22-2-3
 characterization of, 21-16-18
 d^* , 21-14-15
 d^{**} , 21-15
 detectivity, 21-14
 noise equivalent power, 21-14
 noise measurements, 21-18
 quantum efficiency, 21-15-16
 responsivity, 21-14-15, 21-16-18
 spectral response, 21-16
 ultimate sensitivity of, 21-18-21
 photovoltaic, 21-18
 pyroelectric, 21-21
 quantum, 21-21-31
 sensitivity, ultimate, 21-20-21
 spectral response of, 21-30
 thermal, 21-21
 Determinant, matrix algebra, 47-31-32
 Developmental indicators, military equipment, 49-36-37
 Deviation
 from the mean
 RMS, 44-3
 standard, 44-3
 mean absolute, 44-4
 ratio, FSK, 24-9
 Diagnostic measurement, microwave systems, 12-33
 Diagrams
 block, 49-26
 Bode, 15-6-7
 schematic, 49-26-27
 state, 25-5-7
 wiring, 49-26
 Dial pulse signal, telephonic, 2-14
 Dial pulsing, 38-16
 Dial tone, proposed standard, 2-18
 Die-by-die alignment, definition of, 20-5
 Dielectric
 constant, various terrains, 33-4, 33-6
 guides, rectangular, 30-27-28
 insulation, transformers, 13-10
 resonators, 30-28-29
 -rod waveguides, 30-26-27
 Difference
 amplifiers, 19-17-18, 20-42-43
 calculus, 47-21
 entropy, 25-17-18
 linear, equation, 28-3
 phase-shift keying, 24-19-21
 Differentiation, 47-44
 Diffraction propagation, knife-edge, 33-29
 Diffusion current
 definition of, 20-5
 semiconductor, 18-9-10
 Diffusion, excess carriers, 18-6-7
 Diffusion length, definition of, 20-5
 Digital
 AMPS, cellular system, 46-22
 audio recording, 35-27
 cellular systems, 46-22-24
 cellular systems, high-speed data channel, 46-25
 communications
 messages, 24-2
 optimum receivers for, 24-13-18
 system, error probability, 24-18-24
 control system, 15-15-21
 state variable analysis, 15-17-18
 filter(s), 28-2
 adaptive, 28-23-29
 design, 28-13-16
 finite wordlength effects, 28-20-22
 implementation, 28-16-22
 impulse invariance, 28-13-14
 LMS algorithm, 28-23-24
 network structures, 28-16-20
 overflow oscillations, 28-22
 quantization error
 coefficient error, 28-20
 input error, 28-20
 overflow oscillations, 28-22
 roundoff error, 28-20-22
 integrated circuits, 20-30-39
 loop carrier, 38-5-6
 modulation, 27-23
 performance, 38-11
 signal processing, 28-2
 signal processing, adaptive, 28-22-29
 storage CRT, 12-14
 subscriber line, 38-7-8
 switching systems, cellular systems, 46-21
 systems, linear time-invariant, stability of, 15-18-20
 Digital—cont
 television
 ancillary data format, 35-43
 blanking interval considerations, 35-41-43
 component coding of signals, 35-41-42
 composite and component signal coding, 35-38, 35-40
 composite encoded signals (NTSC & PAL), 35-40-41
 encoding parameter for component television, 35-42
 quantization levels for component television, 35-42
 signal
 coding, 35-38, 35-40
 dynamic range, 35-38
 quantizing, 35-38
 synchronization considerations, 35-41-43
 systems, 35-34-46
 basics, 35-34-35
 sampling, 35-36-38
 sampling pulse frequency, 35-37
 spectra, 35-36-38
 timing references, 35-41-42
 transmission systems, satellite communications, 27-23-24
 Digital audio radio service (DARS), 27-3, 27-5
 Digital Equipment Corporation network (DECNET), 26-9, 26-10
 Digital European Cordless Telephone (DECT), 46-24
 Digits, random, table of, 48-6
 Dimensionless quantities, 3-5-6
 Diode(s)
 breakdown, 18-12
 Gunn, 18-24-25
 IMPATT, 18-24-25
 isolation, definition of, 20-5
 light-emitting, 18-13, 18-23-24
 reference, 18-13, 20-45
 requirements, switching power supplies, 14-20-27
 Schottky barrier, definition of, 20-9
 semiconductor, 18-12
 type number coding, 5-7
 semiconductor, definition of, 20-5
 sensors, 12-24-25
 superluminescent, 21-12
 zener, 14-4-5, 18-12-13
 DIP, definition of, 20-5
 Dipole antennas
 array, directivity, 32-34-35
 directivity, 32-14
 half-wave, 32-13-15
 open sleeve, 32-15
 self-impedance, 32-13
 Direct Broadcast Satellite (DBS), 35-23
 Direct-detection receivers, optical, 22-9-10
 Direct field intensity, sound source, 40-12-13

- Direct mapping, definition of, 20-5
- Direct memory access (DMA), computer, 42-14
- Direct to home (DTH), 27-3, 27-5
- Direction cosines, 47-19-20
- Direction numbers, 47-19-20
- Directional antennas, TV broadcast stations, 35-12
- Directivity
antenna, 32-3
arrays, antenna, 32-34-35
dipole antennas, 32-14
- Directory routing, network, 26-24
- Discrete channel, 25-12
- Discrete Fourier transform, 28-7-13
- Discrete-time analog technologies, 28-29-37
- Discrete-time (DT) signal, 28-2
- Discrete-time Fourier transform, 28-5-6
- Discrete-time systems
definitions, 28-2-3
fundamentals, 28-2-7
Z-transform, 28-3-5
- Discrimination, minimum, 25-24
- Disk recordings, video, 35-33-34
- Dispersion
chromatic, 22-24
-flattened fibers, 22-24
material, optical, 22-22
optical fiber, 22-22-24
-shifted fibers, 22-24
waveguide, 22-22-24
- Displacement, velocity, and acceleration magnitudes, 49-23
- Display devices, semiconductor, 18-23-24
- Displays, liquid crystal, 18-23
- Dissipation
electrode, 16-16-18
factor, insulating materials, 4-13
power, within tee network, 11-4, 11-7
- Distance
effect on field strength, 33-18-19
frequency reuse, cellular system, 46-17-19
hyperbolic, measurement of, 31-6-9
measuring equipment, 37-5
from point to a line, equation for, 47-17
between two points, 47-20
- Distance separation, minimum, FM broadcast stations, 35-9
- Distortion
aliasing, 28-6
attenuation, telephony, 2-10
envelope-delay, 12-27-28
harmonic, audio transformers, 13-13
- Distortion-rate function, 25-25-26
- Distributed
algorithms, network, 26-25
amplifiers, broadband, 17-13
control algorithms, local area networks, 26-22-23
- Distributed—cont
power supplies, 14-45-46
queue dual bus (DQDB), 26-33
tree retransmission algorithms, 26-21-22
- Distribution(s), 44-2-5
binomial, 44-4, 45-11
Chi-Square, 44-4-5, 45-11
exponential, 44-5
fitting, using Chi-Squared test, 45-23
-free tests of goodness of fit, 45-23
frequency, common, 45-11
gamma, 45-11
Gaussian, table of, 48-8
joint, 44-2
lognormal, 45-11
mean, 44-3
median, 44-3
mode, 44-3
normal, 44-5, 45-11
table of, 48-8
Poisson, 44-5, 45-11
Rayleigh, 24-23
Rician, 24-23
student-*t*, 45-11
variance, 44-3
Weibull, 45-11
- Distributive law
for scalar multiplication, 47-35
for vector multiplication, 47-35
- Disturbance inputs, feedback control system, 15-11
- Divergence angle, collimated beam, 22-7
- Diversity, 33-19-20
frequency, 37-3
schemes, 46-19-21
site, 27-39-40
space, 37-3
systems, 33-11
degraded, 25-27
- Division, binary, 42-9
- DLC protocols
asynchronous, 26-13-14
bit-oriented, 26-16-18
character-oriented, 26-15-16
synchronous, 26-13-14
- DM *See* Delta Modulation
- DME, 37-5
- Documents, Federal Communications Commission, obtaining, 35-26
- Doherty amplifier, 17-11
- Domestic Satellite Service, 35-20
- Donor impurities, 18-4
- Doors, sound isolation, 40-19
- Dope, definition of, 20-5
- Doping procedures, substrate, 20-21-22
- Doppler
frequency, measurement error, 36-21
frequency shift, 36-14
radar, pulse, 36-14-16
signal processing, radar, 36-14-16
- Dose, definition of, 20-5
- Dot matrix, definition of, 20-5
- Double-rail logic, 43-10
- Double sideband, amplitude modulation, 23-4
- Dram, 3-14
- DRAM, characteristics, 20-39
- DRAM/sRAM, definition of, 20-5
- Drift current
definition of, 20-5
semiconductor, 18-10
- Drift, of carriers, 18-5-6
- Drive-in, definition of, 20-5
- Dry etching, ICs, 20-19-20
- DSB, 23-5-6
- DSL, 38-7-8
- DSW, definition of, 20-5
- Dual-tone multifrequency signaling (DTMF), 38-16, 38-17
- Duality, principle of, 43-3
- Duty cycle, effect on transformer design, 13-9
- Dynamic
logic, 20-36
RAM, 20-42
range, digital TV signal, 35-38
storage elements, definition of, 20-5
vs static routing, network, 26-24
- Dyne, 3-14
- Dysprosium, 4-4-5, 4-8-9
- ## E
- E and M lead signaling, 38-16-17
- E-beam, definition of, 20-5
- $e^{\lambda u} - \mu U$, exponentials of, table of, 48-7
- E_{u2UPROM} *See* Electrically Erasable Programmable Read-Only Memory
- $e^{\lambda u} - \mu U$, exponentials of, table of, 48-7
- EAROM, definition of, 20-5
- Earth
effective radius, 33-29-30
-space communication, 33-33
stations, 27-40-43
antennas, 27-40-43
categories, 27-40-43
power ranges, 27-40-43
surface, illumination of, approximate, 16-48
-to-spacecraft communications, 27-3
- E_c , definition of, 20-5
- Echo
control, 38-8
listener, definition of, 38-10
return loss, definition of, 38-10
talker, definition of, 38-10
- ECL circuits, definition of, 20-5
- ECL storage cell, definition of, 20-5
- Eclipses, solar, effect on satellite communications, 27-9-11
- EEPROM, 20-40-41; *See also* Electrically Erasable Programmable Read-Only Memory
definition of, 20-5
 E_c , definition of, 20-5

- Effective address register, computer processor, 42-15
- Effective efficiency, η , definition of, 12-25
- Effective radiated power, 35-12
- Efficiency
 effective, η , definition of, 12-25
 transmission lines, 29-10-12
- E_g , definition of, 20-5
- E_i , definition of, 20-5
- EIA, 5-4
 RS-232-C, 26-11-13
 RS-422-A, 26-11-13
 RS-423-A, 26-11-13
 RS-449, 26-11-13
 RS-449-C, 26-11-13
 television standards, 35-34
- Eigenvalues, 47-32
- Eigenvectors, 47-32
- Einstein relation, 18-8
- Einsteinium, 4-4-5, 4-8-9
- Einstein's law, 21-22
- Elasticity modulus, of elements, 4-9, 4-11
- Electrical insulating materials, classification of, 13-11
- Electrically erasable programmable read-only memory, 42-14
- Electrochemical equivalents, 4-3, 4-5, 4-7
- Electrode deposition
 gate oxide, 20-100-101
 polysilicon, 20-100-101
- Electrode dissipation, 16-16-18
- Electrolytic capacitors, 5-18-21
 aluminum, 5-20-21
 tantalum, 5-21
- Electromagnetic
 compatibility, 12-35-36
 interference, 12-36
 measurements, 12-36, 12-38
 regulations, 12-36-37
 susceptibility, 12-36
 measurements, 12-38
 waves, propagation in waveguides, 30-3
- Electromigration, definition of, 20-5
- Electromotive force, thermal, of alloys, 4-15
- Electron
 -beam lithography, definition of, 20-5
 cyclotron resonance (ECR) discharges, 20-96-97
 emission, 16-3-16
 laser, free, 16-24, 16-39-41
 power tubes, comparison with solid-state devices, 17-13-17
 tube
 collision ionization, 16-19
 flicker effect, 16-18-19
 grid temperature, 16-18
 induced noise, 16-19
 partition noise, 16-18
 work function, 4-3, 4-5, 4-7
- Electronegativity, 4-3, 4-5, 4-7
- Electronic
 charge, 18-4
 equipment, wiring, 49-12, 49-14-15, 49-16
 Industries Association, 5-4, 35-28
 logic circuit, latch, 43-16
 rest mass, 18-4
 scanning, antenna array, 32-44-45
 -tube materials, thermal emissivity of, 16-16
- Electronvolt, 3-14, 18-3
- Elements
 atomic mass, 4-4, 4-6
 atomic number, 4-3, 4-4, 4-6, 4-8, 4-10
 atomic radius, 4-3, 4-4, 4-6
 boiling point, 4-8, 4-10
 density, 4-8, 4-10
 elasticity modulus, 4-9, 4-11
 electrochemical equivalents, 4-3, 4-5, 4-7
 electron work function, 4-3, 4-5, 4-7
 electronegativity, 4-3, 4-5, 4-7
 evaporation of, 20-15-17
 first ionization potential, 4-3, 4-5, 4-7
 general properties, 4-3, 4-4-7
 gram atomic volume, 4-3, 4-4, 4-6
 heat, specific of, 4-9
 latent heat of fusion, 4-9, 4-11
 linear expansion, 4-9, 4-11
 mass number, 4-3, 4-4, 4-6
 melting point, 4-8, 4-10
 periodic classification of, 4-3, 4-12
 physical properties, 4-3, 4-8-11
 relative hardness, 4-8, 4-10
 specific heat, 4-9, 4-11
 static storage, definition of, 20-10
 tensile strength, 4-9, 4-11
 thermal conductivity, 4-9, 4-11
- Elevation angle, geostationary satellite, 27-11-12
- Elias code, 25-9
- Ellipse, equations for, 47-5, 47-18-19
- Ellipsoid, equations for, 47-6, 47-20
- Elliptic function network, 9-4
- Elliptic integrals, 47-41
- Elliptical polarization, 32-11-12
- ELO, 20-103
- EM fields, exposure to, safety levels for, 49-28-30
- Emission(s)
 electron, 16-3-16
 field, 16-7-10
 current-voltage relationship, 16-12
 displays, 16-13-14
 macroscopic current density, 16-10
 metals, 16-7-9
 modulation, 16-12
 protective resistance, 16-13
 semiconductors, 16-9-10
 transconductance, 16-12
 secondary, 16-5-7
 spontaneous, 41-42
 stimulated, 41-2
- Emission(s)—cont
 stimulated Raman, 41-19
 thermionic, 16-3-5
- Emitter(s)
 arrays, field, 16-10-12
 ccll, multiple, definition of, 20-7
 -coupled logic (ECL), 20-32-33, 20-35
 -follower amplifier, 19-11
 infrared, 21-11
 injection efficiency, 18-19
 junction, 18-18
 visible light, 21-11
- Emu, 3-14
- Encoders, incremental, in error measuring systems, 15-2
- Encoding, source, 25-10
- End-of-active-video reference, 35-41-42
- Energy
 bands, semiconductors, 18-3-5
 bit, 25-19-20
 per data pulse, 24-4
 photon, 21-2
 stored
 capacitor, 6-15
 inductor, 6-15
- Engset equation, 39-12
- Enhanced definition television (EDTV), 35-43
- Enhancer, cellular system, 46-6
- Entropy
 differential, 25-17-18
 function, 25-8-10
 maximum, 25-24
 power, 25-18
- Entry-to-exit flow control, networks, 26-26-27
- Envelope
 -delay distortion (EDD), 12-27-28
 detection, 23-6
 detector, 24-18
- Environmental
 noise, 40-10-12
 test methods, 5-5
- Ephemeris time, 1-143
- Epitaxial growth, definition of, 20-5
- EPROM, 20-40-41; See also Erasable Programmable Read-Only Memory
 definition of, 20-5
- Equal loudness contours, 40-7
- Equated busy-hour call (EBHC), definition of, 39-10
- Equation(s)
 cavity, 30-21-22
 delay, traffic networks, 39-13-14
 Engset, 39-12
 Erlang-B, 39-12
 Erlang-C, 39-12
 finite convolution, 28-3
 half-angle, 47-11-12, 47-13, 47-16
 half-side, 47-13, 47-16
 linear difference, 28-3
 mensuration, 47-4-6
 Poisson, 39-12

- Equation(s)—cont
 Sabine, 40-13-14
 selectivity, 8-4-5
 traffic, 39-12-13
 transmission lines, 29-4-6
 wave, 40-3
- Equiangular spiral antenna, 32-28-29
- Equipment indicator letters, 49-35
- Equiripple approximation, FIR filter
 computer design, 28-16
- Equivalent noise sideband input (cnsi),
 34-8
- Equivocation, 25-13
- Erasable programmable read-only
 memory (EPROM), 42-14
- Erbium, 4-4-5, 4-8-9
- Erg, 3-14
- Erlang-B equation, 39-12
- Erlang-C equation, 39-12
- Erlang, definition of, 39-10-11
- Error
 in arithmetic operations, 47-44-45
 code, satellite digital systems, 27-23
 -control codes, 25-14
 function, 47-40-41
 mean absolute, 44-4
 -measuring systems, 15-2
 probability
 coherent systems with white
 Gaussian noise channels,
 24-18-19
 digital communications system,
 24-18-24
 noncoherent systems with
 nonselective fading channels,
 24-22-24
 noncoherent systems with white
 Gaussian noise, 24-19-22
 propagation, reduction of, 37-2
 range measurement, rms, radar,
 36-21
 rms Doppler frequency
 measurement, 36-21
- Error-free second, definition of, 38-11
- Errored second, definition of, 38-11
- Estimation, spectral, 25-24-25
- Estimation theory, 25-22-25
- Estimator, matched-filter, 25-23-24
- Esu, 3-14
- Etch, cryogenic magnetron, 20-97-99
- Etch filament, definition of, 20-5
- Etch, plasma, definition of, 20-8
- Etch selectivity, definition of, 20-5
- Etchants, used with films, 20-20
- Etching
 Damascene processing, 20-20-21
 ICs, 20-19-21
 dry, 20-19-20
 wet, 20-19
 plasma, gases used with, 20-20
- ETHERNET, 26-19
- Euclidean algorithm, 25-15
- Euler transformation, 28-15
- Europe, cellular systems in, 46-4
- European Broadcasting Union (EBU),
 35-48-49
- Europium, 4-4-5, 4-8-9
- E_v , definition of, 20-5
- Evaporation, thin films, 20-13-14,
 20-15-17
- Evaporative cooling, electron tube,
 16-17
- Event probability, calculation of, 45-7,
 45-9
- Event tree, 45-7
- Excess carriers, 18-6-8
- Excess minority carriers, definition of,
 20-5
- Exchange-access trunks, 38-6
- Exchange plant, 38-3-6
- Excitation mechanisms, lasers, 41-9
- Exclusive-OR gate, 43-11
- Execution cycle, computer, 42-15
- Expected value, 44-3
- Experiment, random, 44-2
- Exponential
 distribution, 44-5
 integrals, 47-39
 integrands, 47-25
 modulation, 23-6-10
- Exponentials of $e^{i\omega t}$ and $e^{-i\omega t}$,
 table of, 48-7
- Exposure to EM fields, safety levels
 for, 49-28-30
- Express-Net, 26-33
- Expressions, canonical, 43-5-6
- F
- F_1 layer, wave propagation in, 33-7
- F_2 layer, wave propagation in, 33-7
- Fabrication
 IC, 20-11
 process, 20-22-26
- Fabry-Perot cavity, 41-6-8
- Factors, conversion, 3-13-17
- Factory-programmed ROM, definition
 of, 20-5
- Fades
 average duration of, mobile signal,
 46-12
 distribution of durations, 46-12
- Fading, 33-19-20
 atmospheric multipath, 33-19
 multipath, 46-14
 reflection multipath, 33-19
 Rican, slow nonselective, 24-22
 selective, mobile signal, 46-12,
 46-17
- Failure data collection and assessment,
 45-9
- Failure mode analysis, 45-6-8
- Failure modes and effects analysis
 (FMEA), 45-6
- Failure rate, definition of, 20-5
- Fall time, definition of, 20-5
- False-pulse noise, in pulse-code
 modulation, 23-17-18
- Fan-out, definition of, 20-5
- Fano algorithm, 25-16
- Far-field region, antenna, 32-4
- Far-zone noise sources, 34-6-7
- Farad, 3-2
 per meter, 3-2
- Faraday, 3-14
- Faraday rotation, effect on propagation,
 27-32-33
- FasNet, 26-33
- Fast Fourier transform *See* FFT
- Fast packet switching, networks, 39-10
- Fathoni, 3-14
- Fault tree, 45-7, 45-8
- FCC, 1-2-3
 Record, 35-26
- FDM/FM/FDMA, multiple carrier
 case, 27-22
- FDM/FM, single carrier case,
 27-21-22
- FDMA, 27-20-21
- FDNR circuit, in ladder circuit,
 10-29-30
- FDNR element, 10-14
- Federal Aviation Administration
 (FAA), 37-2
- Federal Communications Commission
 (FCC), 37-2
 obtaining documents, 35-26
- Feedback
 amplifier
 analysis, 19-22
 topologies, 19-20-21
 transistor, 19-18-22
 control system(s)
 computer-aided analysis, 15-33
 computer-aided design, 15-33
 disturbance inputs, 15-11
 nonlinear, 15-22-32
 characteristics of, 15-22-23
 robustness, 15-11-12
 sensitivity, 15-11-12
 tracking accuracy, 15-10-11
 follow-the-leader, 10-23
 multiple-topologies, 10-25-27
 negative, 19-20
 positive, 19-20
 sensors, process control, 20-60-61
 unity
 acceleration error constant, 15-11
 positive error constant, 15-11
 velocity error constant, 15-11
- Felici mutual-inductance balance, 12-7
- Fermi, 3-14
- Fermium, 4-4-5, 4-8-9
- Ferrites, 13-4
 properties of, 4-19, 4-25-27
- Ferromagnetic materials, 13-3
- Ferroresonant transformers, 13-10-11
- Ferroresonant voltage regulators,
 13-10-11
- FET devices, operating modes, 19-6
- Fetch cycle, computer, 42-15
- FFT
 algorithms, 28-9-11
 leakage and picket-fence effect,
 28-11, 28-13
 mixed radix, 28-10
 in spectral analysis, 28-11-13

- FFT—cont
 windowing, 28-11
 zero padding, 28-11
- Fiber(s)
 bandwidth, 22-24
 dispersion-flattened, 22-24
 dispersion-shifted, 22-24
 distributed data interface (FDDI), 26-33
 graded-index, 22-19
 in the loop (FITL), 38-6
 -optic
 facilities, telephone, 38-13
 scatter sensor, 20-73-74
 usage, cable systems, 35-16-17
 vs satellite communications systems, 27-5
- optical
 connectors, 22-24-25
 dispersion, 22-22-24
 group delay, 22-20-21
 group index, 22-20-21
 losses, 22-21-22
 numerical aperture, 22-20
 splicing, 22-24
 terminology, 22-19
 types of, 22-19-20
 V-number, 22-21
 single-mode step-index, 22-19
 waveguide transmission, 22-19-26
- Fiducial, definition of, 20-5
- Field
 effect transistors, 14-8
 bias techniques, 19-5
 junction, 18-14-15
 MOS, 18-15-16
- emission, 16-7-10
 current-voltage relationship, 16-12
 displays, 16-13-14
 macroscopic current density, 16-10
 metals, 16-7-9
 modulation, 16-12
 protective resistance, 16-13
 semiconductors, 16-9-10
 transconductance, 16-12
- emitter arrays, 16-10-12
- programmable gate arrays (FPGA), 20-51, 20-53
- regions, antenna, 32-4
- strength
 noise, 34-4
 in proximity to antenna, 49-30
 requirements, AM broadcast stations, 35-4
 variation with distance, 33-18-19
 of view (FOV), optical receiver, 22-9
- Figure
 noise, 34-9
 calculation of, 34-10-11
- Filaments, tungsten, 16-3-5, 16-21
- Film capacitors, 5-18
- Film resistors, 5-12-13
- Film systems, television, international standards, 35-34
- Filming, definition of, 20-5
- Films
 etchants used with, 20-20
 thin, 20-13-14
- Filter(s)
 Åckerberg-Mossberg (AM) biquad, 10-21-22
- active
 building blocks, 10-10-15
 general realization methods, 10-8
 low-pass notch, 10-44-45
 operational transconductance amplifiers in, 10-7
 switched-capacitor, 10-34-45
 types of, 10-3
 zero placement formulas, 10-43
- adaptive, 28-23-29
- antialiasing, 35-37
- autoregressive, 25-24
- bandpass, 8-1-9
 gain at resonance, 8-2-3
 selectivity far from resonance, 8-3
 selectivity near resonance, 8-3, 8-6
- Butterworth low-pass, 9-14
- canonical dual-mode, 27-18
- Chebyshev low-pass, 9-14
- circuits, 14-39-43
 switched capacitor, 28-29-32
 basic principles, 28-30-31
 voice quality, 28-31-32
- coupled dual-mode circular-cavities, 27-18
- design
 Butterworth low-pass, 9-18, 9-20-21
 capacitor input, 14-41
 Cauer-parameter low-pass, 9-5-13
 Chebyshev low-pass, 9-15-17
 examples, low-pass, 9-19, 9-22, 9-24-26
 inductor input, 14-40-41
 resonant, 14-43
- digital, 28-2
 adaptive, 28-23-29
 design, 28-13-16
 finite wordlength effects, 28-20-22
 implementation, 28-16-22
 impulse invariance, 28-13-14
 LMS algorithm, 28-23-24
 network structures, 28-16-20
 overflow oscillations, 28-22
 quantization error
 coefficient error, 28-20
 input error, 28-20
 overflow oscillations, 28-22
 roundoff error, 28-20-22
- FIR, 28-15-16
 computer design, equiripple approximation, 28-16
- IF, bandwidth, spectrum analyzer, 12-19-20
- IIR, 28-13-15
 integrated, 10-3
- interference, optical receiver, 22-9
- low-pass
 to band-stop transformation, 9-19
- Filter(s)—cont
 to bandpass transformation, 9-14
 Cauer-parameter, 9-4-5
 to high-pass transformation, 9-14
 matched, Gaussian noise, 24-14-15
 network, requirements for, 9-2
 nonrecursive, 36-16
 reference low-pass, 9-3
 SC, 10-34-45
 operation of ideal, 10-36-38
 single-amplifier, 10-16-18
 single-mode cascaded-cavities, 27-18
 surface-acoustic-wave, 28-33-37
 response, 28-34
 switched capacitor (SCF), 28-29
 systems, sampled-data, 10-34-36
 three-amplifier, 10-21-22
 transconductance- C (g_m - C), 10-22-23
 transfer function, 10-7-8
 transversal, 36-16
 two-amplifier, 10-18-21
- Fin-lines, 30-31
- Finish application table, 49-9-10
- Finishes, for tropical and marine use, 49-5, 49-8
- Finite convolution equation, 28-3
- Finite impulse response (FIR), 28-3
- FIPOS, 20-103
- FIR filter, 28-15-16
 computer design, equiripple approximation, 28-16
 LMS algorithm, 28-23-24
 window designs, 28-15-16
- First-in first-out reservation scheme, 26-21
- First ionization potential, 4-3, 4-5, 4-7
- First-order transfer functions, 10-14-15
- Fixed assignment techniques, multiaccess link control, 26-18
- Fixed composition resistors, 5-10-11
- Fixed-point representations, 42-9-10
- Fixed-program storage (read-only memory or read-only store), definition of, 20-5
- Fixed Satellite Service (FSS), 27-3, 27-13
- Flange
 waveguide
 insertion loss, 30-17
 leakage, 30-16
 resistance, 30-16-17
- Flanking, sound isolation, 40-21, 40-23
- Flexible cords, current-carrying capacity, NEC, 49-13
- Flexible waveguides, 30-8, 30-14-15
- Flicker effect, electron tube, 16-18-19
- Flip-chip package, ICs, 20-29-30
- Flip-flop
 D, definition of, 20-6
 definition of, 20-5-6
 JK, definition of, 20-6
 RS, definition of, 20-6
 RST, definition of, 20-6

- Flip-flop—cont
 T, definition of, 20-6
- Float zone, definition of, 20-6
- Floating-point representations, 42-9-10
- Flow control
 entry-to exit, networks, 26-26-27
 hop-level, networks, 26-26
 network access, 26-27
- Fluctuations, clutter, radar, 36-17
- Fluid ounce, 3-14
- Fluorescent lamps, 21-10
- Fluorine, 4-4-5, 4-8-9
- Flux density, core materials, 13-5
- Flux, units, 16-45-47
- Flyback converter, 14-28-29
- FM
 broadcast stations
 minimum distance separation, 35-9
 performance required, 35-9-10
 broadcast, stereophonic
 transmission, 35-10-11
 broadcasting, 35-6-11
 channel availability, 35-8-9
 coverage, estimated, 35-9
 frequency allocations, 35-7-8
 interference ratios, 35-8
 station classification, 35-8-9
 multitone, bandwidth considerations
 in, 23-9
 systems, signal-to-noise
 improvement, 23-9-10
 translators, 35-11
- FM/FM signal, spectral distribution of, 23-8
- FMX stereo, 35-11
- Focusing, cathode-ray tube, 16-55
- Folded network, 39-2
- Foot, 3-14
 US survey, 3-14
 of water, 3-15
- Footcandle, 3-15
- Footlambert, 3-15
- Force/pressure-current analogy, 40-6-7
- Force/pressure-voltage analogy, 40-6
- Forced-air cooling, electron tube, 16-16-17
- Format, ancillary data, digital
 television, 35-43
- Formats, instruction, computer, 42-16-17
- Forward bias
 definition of, 20-3
 semiconductor, 18-9-10
- Forward, current, definition of, 20-4
- Forward-mode converter, 14-19, 14-28
- Foundation, steel radio towers, 49-21-22
- 4G, 46-24
- Four-by-four, multifrequency
 numerical code, 2-14, 2-16
- Four-terminal-pair bridge, 12-11
- Four-wire repeaters, 38-11
- Fourier
 discrete transform, 28-7-13
 expansion, harmonics, 7-8-9
- Fourier—cont
 series, 7-2, 7-8-9
 complex form, 7-2, 7-8-9
 periodic waveforms and, 7-11-13
 real form, 7-2
 transfer of function, 7-2, 7-3-8
 transform, discrete-time, 28-5-6
- Fowler Nordheim equation, 16-8
- Fractal antennas, 32-32
- Frame format, bit-oriented protocols, 26-16-17
- Franaszek code, 25-6
- Francium, 4-4-5, 4-8-9
- Free electron laser (FEL), 16-24, 16-39-41
- Free fall, 3-15
- Free space
 permittivity of, 18-4
 propagation, equations for, 33-27
 transmission equations, 33-21-22
- Frequency(ies)
 accuracy, international telephony, 2-10
 allocations
 FM broadcasting, 35-7-8
 radio, 1-3-141, 1-3-5
 allotments, 1-141-142
 bands, letter designation, 1-3
 bands, radio, nomenclature, 1-2, 1-3
 broadcasts, 1-148-162
 cable television channels, 35-15, 35-21-22
 cellular system, 46-3
 counter, 12-21-22
 -dependent negative resistance. *see*
 FDNR
 distributions, common, 45-11
 diversity, 37-3
 functions, pulse shapes, 7-6-8
 -hop spread-spectrum
 communications, 24-26
 -independent antennas, 32-28-29
 international broadcasting stations, 35-24
 lowest useful high, 33-8
 maximum usable, 33-7
 measurement, 12-21-23
 modulation, 23-7; *See also* FM
 normalized, 9-3
 optimum working, 33-8
 repair, definition of, 20-9
 response, single-stage transistor
 amplifier, 19-6, 19-15
 reuse, cellular system, 46-17-19
 reuse factor, cellular system, 46-18
 satellite communications, choice of, 27-12-13
 sharing, auxiliary broadcast services, 35-24
 shift, Doppler, 36-14
 shift keying
 binary, 24-8-9
 continuous-phase, 24-9-10
 multiple, 24-10
 spectrum, cellular system, 46-4
 standards, 1-144-148
- Frequency(ies)—cont
 and time-domain relationships, 12-15-16
 and time-interval analysis, 12-23-24
 tracking loop, 12-27
 vibration, natural, 49-24
 -wavelength, conversions, 1-2
- Fresnel
 integrals, 47-41
 -Kirchhoff theory, 33-16
 zones, 33-16-17
- Frustrum, equations for, 47-6
- FSK *See* Binary Frequency-Shift
 Keying
- Full-wave rectifier
 single-phase
 bridge, 14-13
 center-tap, 14-12-13
- Function(s)
 Bessel, 47-41-43
 table of, 48-9-10
 conditional probability, 44-2-3
 cumulative distribution, 44-2
 entropy, 25-8-10
 error, 47-40-41
 Fourier transfer of, 7-2, 7-3-8
 gamma, 47-40
 hyperbolic, 47-13-14
 marginal probability, 44-2
 network, 26-7-9
 polygamma, 47-40
 probability density, 44-2
 ψ , 47-40
 transcendental, 47-21
- Fundamental
 products, 43-5
 quantities, transmission lines, 29-4-5
 sums, 43-5
- Fundamental physical constants, 3-7-13
- Fungus, on materials, 49-8
- Furlong, 3-15
- Furnishings, absorption coefficients, 40-16
- Fusible metallization, definition of, 20-6
- Fusion, latent heat of, of elements, 4-9, 4-11
- G
- GaAs LEDs, 21-11
- Gadolinium, 4-4-5, 4-8-9
- Gain
 antenna, 32-3
 radar, 36-12
 relative to hypothetical isotropic, 33-26
 circularly polarized antenna, 32-11
 current, power transistor, 14-5
 insertion, 12-26
 loop, feedback, 15-11
 luminous, 16-49
 -phase plot, nonlinear system, 15-29

- Gain—cont
 power
 apparent, 33-26
 RF amplifiers, 17-17
 at resonance, bandpass filter, 8-2-3
- Gal, 3-15
- Galactic noise, 32-9, 34-5
- Galileo System, 37-8
- Gallium, 4-4-5, 4-8-9
- Gallium arsenide LEDs, 21-1
- Gallon, 3-15
- Galois field, 25-15
- Galvanic series, in sea water, 4-3, 4-13
- Games, communications, 25-27-31
- Gamma, 3-15
- Gamma distribution, 45-11
- Gamma function, 47-40
- Gap noise, temporal structure, 34-8
- Gas photodiodes, 16-50
- Gas tubes, 16-41-45
 amplification of, 16-42-43
 cold-cathode, 16-42
 hot-cathode, 16-42
 microwave applications for, 16-44-45
 as noise source, 16-44-45
 voltage regulator, 16-43
- Gaseous attenuation, 33-20
- Gases, sound in, 40-3-4
- Gate
 AND, 43-3
 definition of, 20-6
 arrays, 20-49-54
 field programmable, 20-51, 20-53, 20-55
 circuit, two-stage, 43-12
 Exclusive-OR, 43-11
 Inclusive-OR, 43-11
 NAND, definition of, 20-6
 networks, general, 43-10-11
 NOR, definition of, 20-6
 OR, 43-3
 definition of, 20-6
 oxide electrode deposition, 20-100-101
 reduced, 43-11
 set, complete, 43-10
 structure formation, 20-101
 sum modulo two, 43-11
 turn-off, thyristors, 14-11
 XOR, definition of, 20-6
- Gauss, 3-15
 algorithm, 47-32
- Gaussian distribution, table of, 48-8
- Gaussian noise, 32-9
 channel, additive, 25-18-19
 general, optimum threshold for, 24-14
- General impedance converter (GIC), 10-12-13, 10-19
- General purpose interface bus (GPIB), 12-33-34
- General Radio Type 21-3, 12-10
- General-register
 computer processor, 42-15
 organization of, 42-14-16
- General-register—cont
 processors, 42-11
- Generation-recombination noise, 21-19
- Generators
 character, definition of, 20-4
 rate, 15-4
- Geographical zones for broadcasting, ITU, 35-24-25
- Geometric progression, 47-8
- Geometry
 analytic, solid, 47-19-21
 plane analytic, 47-16-19
- Geophysical alerts, WWV and WWVH, 1-155-156
- Geostationary orbit, 27-11
- Geostationary satellite, elevation and azimuth angles, 27-11-12
- Germanium, 4-4-5, 4-8-9
- Germany, cellular systems in, 46-3
- Getters, vacuum tube, 16-22
- GIC *See* General Impedance Converter
- Gilbert, 3-15
- Gill, 3-15
- Glass, metallic, 13-4
- Glass spray, definition of, 20-6
- Global alignment, definition of, 20-6
- Global positioning system, 37-6-8
 measurement techniques, 1-160-162
 status announcement, WWV and WWVH, 1-156-157
- GLONASS, 37-8
- g_m , definition of, 20-6
- GOI, definition of, 20-6
- Gold, 4-4-5, 4-8-9
- Goodness of fit, 44-4-5
- GPIB, 12-33-34
- GPS, 37-6-8
- Grad, 3-15
- Grade of service, switching system, 39-11-12
- Graded-index fibers, 22-19
 bandwidth, 22-24
- Grain, 3-15
- Gram, 3-15
- Gram atomic volume, 4-3, 4-4, 4-6
- Graphic design, power grid tube circuits, 17-4-8
- Graphic symbols
 British, 49-28
 logic design, 49-29
 US, 49-26-27
- Grating lobes, antenna array, 32-35-36
- Gray, 3-2
 per second, 3-2
- Great circle
 calculations, 33-11-12
 distance, 33-8
- Greek alphabet, 3-6
- Grid drive power, power grid tube, 17-3
- Grid temperature, electron tube, 16-18
- Grids
 electron tube, 16-21-22
 network, 26-34
- Ground
 conductivity
 Canada, 35-5, 35-8
- Ground—cont
 US, 35-5, 35-7
 various terrains, 33-4, 33-6
 connection, Wagner, 12-3
 packet radio, 26-6
 systems, radial-wire, 32-15-16
- Grounded-cathode circuit, 17-8-9
- Grounded-grid circuit, 17-8-9
- Group delay
 distortion
 measurement of, 12-28-29
 telephony, 2-9
 optical fiber, 22-20-21
- Group index, optical fiber, 22-20-21
- Group indicators, military equipment, 39-35
- GSM digital cellular system, 46-22, 46-24
- Gunn diode, 18-24-25
- Gyration resistance, 10-13-14
- Gyator, 10-13-14
- Gyro
 backward-wave tubes, 16-38-39
 klystron, 16-38-39
 traveling-wave tubes, 16-38-39
- Gyrotron, 16-38-39
- H
- Hadamard matrix, 24-3
- Hafnium, 4-4-5, 4-8-9
- Half-angle equation, 47-11-12, 47-13, 47-16
- Half-side equations, 47-13, 47-16
- Half-wave dipole antennas, 32-13-15
- Half-wave rectifier, single-phase, 14-12
- Hamming code, 25-15
- Hamming window function, 28-11
- Hand, 3-15
- Handoff procedure, 46-8
- Hanning window function, 28-11
- Hansen one-parameter circular source distribution, antennas, 32-53
- Hardness, relative, of elements, 4-8, 4-10
- Hardware, computer, 42-3
- Harmonic currents, power supply filters, 14-41
- Harmonic distortion
 audio transformers, 13-13
 transistor amplifiers, 19-16-17
- Harmonic motion, equations, 49-22-23
- Harmonics, Fourier expansion, 7-8-9
- Hartley oscillator, 19-25
- Hay bridge, 12-5
- Haze, Hilton, definition of, 20-6
- HDSL, 38-8
- Hearing, 40-6-7
- Heat
 latent, of elements, 4-9
 specific, of elements, 4-9
- Heat of fusion, latent, of elements, 4-9, 4-11
- Heat, specific, of elements, 4-9, 4-11
- Heaviside expansion theorem, 6-21

- Hectare, 3-15
Helium, 4-4-5, 4-8-9
Helix antennas, 32-23-24
Helix, traveling-wave tube, 16-30, 16-33
HEMTs, 27-17
Henry, 3-2
 per meter, 3-2
Hermite-Gaussian beam modes, 41-8
Hermitian forms, matrix algebra, 47-34
Hertz, 3-2
Heterodyne-detection, optical
 receivers, 22-10-12
Heterostructures, laser, 41-9
Hewlett-Packard, Q meter, 12-8-10
Hexadecimal, number system, 42-5
Hierarchical routing, network, 26-24
High
 altitude tests, 5-6, 5-7
 definition television (HDTV), 35-18, 35-23, 35-43-46
 frequencies, propagation of, 33-6-12
 frequency
 antenna, angles of departure and arrival, 33-8
 propagation, forecasts, 33-8-10
 transmissions, bandwidth limitations, 33-9-11
 intensity visible LEDs, 21-11-12
 level data link control, 26-16
 level language, computer, 42-3
 order transfer functions, 10-23-34
 speed
 data channel systems, 46-25
 digital services, 38-6
 DSL, 38-8
 local area network, 26-30-36
 wide area networks, 26-36-39
 tone, telephony, proposed standard, 2-18
Higher-level protocols, network, 26-29-30
Hillocks, definition of, 20-6
Hilton haze, definition of, 20-6
Hogshead, 3-15
Hold time, definition of, 20-6
Hollow waveguides, propagation in, 30-3
Holmium, 4-4-5, 4-8-9
Homojunction lasers, 41-9
Hooks and extensions, 46-7
Hop-level flow control, networks, 26-26
Horizon, radio, distance, 33-13-15
Horizontal resolution, television, 35-35-36
Horn antennas, 32-22-23, 32-24
Horsepower, 3-15, 49-18
Hot-cathode, gas tubes, 16-42
Hour, 3-15
Huffman code, 25-10
Humidity
 constant, tests, 5-6
 effect on cavity tuning, 30-21
 relative, Celsius table of, 49-6
Hundredweight, 3-15
Hybrid
 -coil method, 12-7
 junctions, waveguide, 30-19-20
 repeaters, 38-11, 38-13
 routing, network, 26-25
Hydrogen, 4-4-5, 4-8-9
Hydrogen masers, 1-147-148
Hydrogen thyratron, 16-42-43
Hydrometeors, 27-36
Hydrophilic, definition of, 20-6
Hydrophobic, definition of, 20-6
Hyperbola, equations for, 47-19
Hyperbolic
 cosines, table of, 48-3
 distance, measurement of, 31-6-9
 functions, 47-13-14
 protractor, 31-7
 sines, table of, 48-2
 tangents, table of, 48-4
 trigonometry, 47-14-16
Hyperboloid, of revolutions, equations for, 47-21
Hypothesis
 alternate, 25-22
 null, 25-22
 testing, 25-22-23

I
I/O, memory mapped, 42-32
I/O section, microprocessor, 20-36-37
IBM System Network Architecture (SNA), 26-9, 26-10
IC(s)
 access time, 20-3
 Damascene processing, 20-22-26
 design, 20-26-28
 modeling, 20-26
 simulation, 20-26
 topological realization of, 20-26-28
 etching, 20-19-21
 fabrication, 20-11
 process, 20-22-26
 yield, 20-23, 20-26
 image-sensing, 20-46
 micro-electro-mechanical systems (MEMS), 20-49
 for microcomputers, 20-36-39
 for microprocessors, 20-36-39
 processing, 20-11-26
 speech-synthesizer, 20-46-47
 system in a package (SIP), 20-47-48
 system on a chip (SOC), 20-47-48
 three-dimensional, 20-108
 voltage reference, 20-43, 20-45
 voltage regulators, 20-45
IDSL, 38-7
IEC, 5-4, 35-29
 standard
 625, 12-33-34
 625-2, 12-34-35
IEEE standard,
 488.1, 12-33-34
 488.2, 12-34-35
IF filter, bandwidth, 12-19-20
IGBT *See* Insulated-Gate Bipolar Transistors
IGFET; *See also* Insulated-Gate Field-Effect Transistor
 definition of, 20-6
Ignitrons, 16-42
IIR filters, 28-13-15
 LMS algorithm, 28-24
Illumination, earth surface,
 approximate, 16-48
Illuminators, multizone RTP, 20-67-68, 20-69
ILS, 37-4
Image
 charge potential, 16-8
 guides, 30-27-28
 intensifier tube, 16-48-50
 orthicons, 16-50-52
 -sensing ICs, 20-46
 tube, 16-48-50
Imaging, optical, 16-47-48
Imaging telescope, 22-6
Immittance, vs normalized frequency,
 double-tuned circuits,
 8-9-10
IMPATT diode, 18-24-25
Impedance(s)
 acoustic, specific, 40-3
 and admittance
 matrix of a junction, 31-9-10
 waveguide, normalized, 31-9
 analyzer, vector, 12-13
 bridges, 12-3-13
 characteristic, of transmission lines,
 29-21-25
 formulas, 6-6, 6-9-13
 input, of resonant line, 29-16-17
 lower than directly measurable,
 measurement of, 12-8-9
 match, audio transformers, 13-14
 matching
 with open stub, 29-19
 with shorted stub, 29-19
 meters, automatic, 12-11-13
 mutual, antenna array, 32-38-40
 scan, 32-44
 telephony trunks, 2-9
 transformation on high SWR lines,
 29-12-13
 transmission lines, 29-6-7
Implicants, prime, 43-13-14
Impregnated tungsten dispenser
 cathodes, 17-18
Impulse invariance, digital filter,
 28-13-14
Impulse noise, definition of, 38-9
Impurities
 acceptor, 18-4
 donor, 18-4
Inch, 3-15
 of mercury, 3-15
 of water, 3-15
Inclusive-OR gate, 43-11
Incoherent sources, optical
 communications, 22-6-7

- Incremental encoders, in error measuring systems, 15-2
- Indefinite integrals, 47-21-24
- Independent sideband method, AM stereo, 35-6
- Index, refractive, typical mean values, 33-29, 33-31
- Indicator letters, equipment, 49-35
- Indicators
developmental, military equipment, 49-36-37
unit, military equipment, 49-38-39
- Indirect addressing, 42-25-26
- Indium, 4-4-5, 4-8-9
- Indoor design, noise levels, 40-15, 40-17
- Inductance
self-, circular ring, 6-14
of single-layer solenoids, 6-3-5
- Induction-heating supply, 14-48-49
- Inductively coupled RF plasma discharges, 20-93-95
- Inductor(s)
audio-frequency, 13-21-23
laminated core, 13-21
toroidal core, 13-21-23
charge and discharge, 6-16-17
energy stored, 6-15
frequency ranges, 13-3
input, filter design, 14-40-41
loss equations, 13-26
magnetic core, 13-3
magnetic-core, 13-19, 13-21
pot-core, 13-23
precision adjustable, 13-23, 13-26-28
Q, 13-26
simulation, 10-29, 10-30-32
temperature stability, 13-28
types of, 13-3
- Industrial noise, assessment of, 40-12
- Inequalities, 47-11
- Inference, probability and statistical, 45-15-25
- Infinite impulse response (IIR), 28-3
- Information networks, multiterminal, 25-26-27
- Information theory, fundamental theorem of, 25-14
- Infrared, emitters, 21-11
- Infrared spectrum, 21-2
- Injection laser, semiconductor, 21-12-14
- INMARISAT communications system, 27-5
- Input loading factor (ILF), definition of, 20-6
- Input/output
computer, 42-14, 42-30-33
programming, 42-32
- Insertion
characteristic loss, 11-2
gain, 12-26
loss, 12-26
audio transformers, 13-14
definition of, 38-9-10
- Insertion—cont
resonant lines, 29-17-18
- Instantaneous sampling, pulse amplitude modulation, 23-12
- Instantaneous sound pressure, definition of, 40-9
- Instruction register, computer processor, 42-15
- Instruction set, computer, 42-17
- Instructional Television Fixed Service (ITFS), 35-18, 35-20
- Instrument landing system, 37-4
- Instruments, computer control of, 12-33-35
- Insulated-gate bipolar transistors, 14-7-8, 14-30, 14-32
- Insulated-gate field-effect transistor (IGFET), 18-15-16
- Insulating materials, 4-12, 4-13
characteristics of, 4-20-23
transformers
classes of, 13-10-11
operating temperatures, 13-10-11
- Insulation, dielectric, transformers, 13-10
- Insulator, definition of, 20-6
- Integer(s)
reciprocals of, sums of powers of, 47-8-9
signed, 42-4
sums of powers of, 47-8
unsigned, 42-4
- Integrals
cosine, 47-39-40
definite, 47-27-29
elliptic, 47-41
exponential, 47-39
Fresnel, 47-41
indefinite, 47-21-24
sine, 47-39-40
table of, 47-21-29
- Integrands
exponential, 47-25
logarithmic, 47-24-25
trigonometric, 47-25-27
- Integrated circuit(s); See also IC
application specific, 20-49-56
ceramic flat package, 20-29
ceramic packages, 20-28-29
Damascene processing, 20-22-26
definitions, 20-3-11
digital, 20-30-39
flip-chip package, 20-29-30
leadless ceramic chip carrier, 20-29
linear, 20-42-45
memory, 20-39-42
micro-electro-mechanical systems (MEMS), 20-49
miscellaneous, linear, 20-45
packaging, 20-28-30
plastic packages, 20-28
processing, single wafer, 20-56-102
special purpose linear, 20-45
system in a package (SIP), 20-47-48
system on a chip (SOC), 20-47-48
- Integrated circuit(s)—cont
terminology, 20-3-11
trends in, 20-45-49
- Integrated CMOS processing, RTP for, 20-74-76
- Integrated device manufacturing, 20-58-66
- Integrated Digital Services Network (ISDN), 27-5
- Integrated filters, 10-3
- Integrated injection (I^2L), 20-33-34
- Integrated services digital network (ISDN), 38-3
- Integration, 47-44
- Integrators, 10-11-12
Miller, 10-11-12
SC, 10-38-39
- INTELSAT, 27-28-29, 27-32
- INTELSAT IV, 27-24
- INTELSAT VI, 27-18
- Intensity
call, 39-11
modulation, 22-3
signal, LED, 22-7
sound, 40-5
traffic, 39-11
- Intercept-intercept form, equation, 47-17
- Interdepartment Radio Advisory Committee (IRAC), 1-3
- Interface protocols, 26-8
- Interfaces
bit-parallel, 35-40
bit-serial, 35-40
computer, 42-30
- Interference
composite triple beat, 35-16
between direct and reflected waves, 33-17-18
electromagnetic, 12-36
regulations, 12-36-37
filter, optical receiver, 22-9
levels, speech, 40-15
ratios, FM broadcasting, 35-8
signals from other sources, 34-11-12
tone, definition of, 38-9
- Interferences, patents applications, 49-33
- Interior design, noise levels, 40-17
- Interlevel isolation, 20-101-102
- Intermediate Data Rate (IDR), 27-28
- Intermodulation noise, transponder, 27-15-17
- International
broadcasting stations, 35-3
frequencies, 35-24
schedules, 35-24
broadcasting, United States, 35-24-26
sound recording standards, 35-29
standards
for broadcasting, 35-46-49
telephony, 2-3
television film systems, 35-34
telephony country codes, 2-13
television recording standards, 35-34

- International Air Transportation Association (IATA), 37-2
- International Business Service (IBS), 27-28
- International Civil Aviation Organization (ICASO), 37-2
- International Electrotechnical Commission, 5-4
publications, 35-47-48
- International Organization for Standardization
recommendations, 35-48
- International Radio Consultative Committee (CCIR), 35-46
recommendations, 35-47
- International Standards Organization, 5-4
- International Telecommunications Union (ITU), 35-46, 37-2
- International Telephone and Telegraph Consultative Committee (CCITT), 35-46
- Internet time signals, 1-158, 1-160
- Interoffice signaling, common carrier, 38-16-18
- Interpolation, at equidistant points, 47-44
- Interpolation polynomial, 47-44
- Interpreter, computer, 42-3
- Intersatellite Satellite Service (ISS), 27-13
- Intertoll trunks, telephone, 38-6-7
- Intrinsic concentration, definition of, 20-6
- Intrinsic material, 18-5
- Inventions
copyright protection of, 49-34
patent protection of, 49-32-34
recognizing, 49-32-33
- Inventive process, 49-33
- Inventor, who is, 49-33
- Inverse baseline network, 39-9
- Inverse transforms, 47-39
- Inverse z-transfer function, 15-17
- Inversion formula method, inverse z-transfer evaluation, 15-17
- Inverter, 43-2
dc-ac, 14-47-50
definition of, 13-11-12
ratioless, definition of, 20-9
- Iodine, 4-4-5, 4-8-9
- Ion-beam lithography, definition of, 20-6
- Ion implantation
definition of, 20-6
substrate doping procedures, 20-21-22
- Ion milling, definition of, 20-6
- Ionization potential, first, 4-3, 4-5, 4-7
- Ionospheric scatter propagation, 33-12
- Ionospheric scintillations, 27-33-35
- Iridium, 4-4-5, 4-8-9
- IRIDIUM, 27-11
- Irises, resonant, 30-24
- Iron, 4-4-5, 4-8-9
- Iron powder, 13-4
- Irradiance, solar spectral, 22-13
- ISDN, 27-5, 38-3
- ISDN DSL, 38-7
- ISO, 5-4
reference network, 26-9, 26-10
- Isocline method, phase portrait construction, 15-31
- Isolated routing, network, 26-25
- Isolation
sound, 40-17, 40-19-23
doors, 40-19
flanking, 40-21, 40-23
resilient channels, 40-19
walls, 40-19, 40-20-21
weighting of components, 40-19, 40-21
- Isolators, shock, 49-26
- Isoplanar, definition of, 20-7
- ITU
geographical zones for broadcasting, 35-24-25
Radio Regulations, 1-2
- J
- Jammer saddle point, 25-27-29
- Jansky, 3-15
- Japan, cellular systems in, 46-3-4
- Jaynes maximum entropy principle, 25-24
- JEDEC, 5-4
- JETDS
nomenclature policy, 39-35
type numbers, examples, 49-37
- JFET
large symbol model, 19-2-3
transistor, symbol, 19-2-3
- Johnson noise, 21-18
- Joint Committee on Television Transmission (CMTT), 35-46
- Joint distribution, 44-2
- Joint Electron Device Engineering Council, 5-4
- Joint Precision Approach and Landing System, 37-8
- Joint Tactical Information Distribution System, 37-8
- Joule, 3-2
per cubic meter, 3-2
per kelvin, 3-2
per kilogram, 3-2
per kilogram kelvin, 3-2
per mole, 3-2
per mole kelvin, 3-2
- JPALS, 37-8
- JTIDS, 37-8
- Jump resonance, 15-23
- Junction(s)
avalanche, definition of, 20-7
barrier, definition of, 20-7
capacitor, definition of, 20-7
collector, 18-18
emitter, 18-18
- Junction(s)—cont
field-effect transistors. *see* JFET
photovoltaic, 18-13
pn, 18-8-12
at equilibrium, 18-8-9
transistor, definition of, 20-7
two-port, 31-3
- Jury table, 15-18-19
- K
- Kahn system, AM stereo, 35-6
- Karnaugh map, 43-13
- Kayser, 3-15
- Kelvin, 3-3
temperature, 3-3
- Kepler's laws, 27-6-7
- Keying
amplitude-shift, 24-5
binary frequency-shift, 24-8-9
binary phase-shift, 24-5-6
continuous-phase frequency-shift, 24-9-10
differential phase-shift, 24-19-21
minimum-shift (MSK), 24-8
multiple frequency shift, 24-10-11
offset quadriphase shift, 24-8
phase-shift, 24-11-13
quadriphase shift, 24-6-8
- Kilocalorie, 3-15
- Kilogram, 3-3
per cubic meter, 3-3
- Kilogram-force, 3-15
- Kilopond, 3-15
- Kip, 3-15
- Kirchhoff's law, 21-3
- Klystron™
amplifier, 17-10, 17-18
tube, 16-22
- Klystrons, 16-25-30
gyro, 16-38-39
multicavity, 16-27-29
multiple-beam, 16-29-30
reflex, 16-27
classes of, 16-28
two-cavity, 16-26
- Knife-edge diffraction propagation, 33-29
- Knot, 3-15
- KOOI effect, definition of, 20-7
- Krypton, 4-4-5, 4-8-9
- K_u band, satellites, 35-23
- Kullback principle of minimum discrimination, 25-24
- L
- L, definition of, 20-7
- Laboratory microphones, 40-8
- Ladder
low-pass, L-C, simulation of, 10-27-30
simulation, 10-27-34
leapfrog topology, 10-29
using FDNR circuits, 10-29-30

- Ladder—cont
 using inductor simulation, 10-29
 using OTAs, 10-30-34
- Lag controller, 15-12
- Lambert, 3-15
- Lambert's law, 21-9
- Laminations, transformer, 13-3-4
- Lamps
 arc, 21-10
 fluorescent, 21-10
 tungsten, 21-10
- LAN *See* Local Area Network
- Land Mobile Satellite Service (LMSS), 27-3
- Landing system, instrument, 37-4
- Langley, 3-15
- Lanthanum, 4-4-5, 4-8-9
- Laplace transforms, 6-19-21, 15-17, 47-38-39
 table of, 47-38-39
- Large-scale integration, definition of, 20-7
- Large-signal measurements, 12-16-17
- Laser anncal, definition of, 20-7
- Laser(s)
 common, 41-8-9
 CW, 41-8
 common, 41-13-14
 detector classes, 41-2
 excitation mechanisms, 41-9
 free electron, 16-24, 16-39-41
 heterostructures, 41-9
 homojunction, 41-9
 losses, 41-4
 pulsed, 41-8
 high power, 41-12
 second-harmonic generation, 41-12, 41-14, 41-16-19
 semiconductor, 21-12-14, 41-9
 spectral region of operation, 41-2
 threshold, 41-4-6
 titanium-doped sapphire, 41-9
 transient operation, 41-5-6
 transverse modulators, 41-18
 tuning curves for various dyes, 41-10-11
- Latch
 circuit, set-reset, 43-16
 electronic logic circuit, 43-16
 SR circuit, 43-16
- Latent heat of fusion, of elements, 4-9, 4-11
- Law
 Bouguer's, 22-14
 Einstein's, 21-22
 Kirchhoff's, 21-3
 Lambert's, 21-9
 Planck's, 21-3, 21-6
 Snell's, 22-19-20
 Stefan-Boltzmann, 21-3, 21-5
 of trade secrets, 49-31-32
 Wien's displacement, 21-3
- Lawrencium, 4-4-5, 4-8-9
- Layer, network, 26-23-29
- LCR Meter, automatic, 12-11-12
- Lead, 4-4-5, 4-8-9
- Lead controller, 15-12
- Leadless ceramic chip carrier, ICs, 20-29
- Leaguc, 3-15
- Lcakage, flange, waveguide, 30-16
- Leaks, sound, 40-22-23
- Leap seconds, 1-143-144
- Lcap seconds, WWV and WWVH, 1-154
- Leapfrog topology, ladder simulation, 10-29
- Least mean square algorithm, 28-23-24
- LEDs, 21-10-11; *See also* Light Emitting Diodes
 high-intensity visible, 21-11-12
 radiation pattern, 22-6
- Lempel-Ziv code, 25-9
- Length, of transmission line, 29-21
- Lengths, associated with plane figures, equations for, 47-4-5
- LEO satellites, 27-11
- Letter designation
 frequency bands, 1-3
 radar bands, 36-11
- Level
 confidence, 45-15
 crossing rate (LCR), mobile signal, 46-12
 reference level point, 38-9
 relative, 38-9
 sound, 40-8
 power, 40-8-9
 pressure, 40-5-6
- Levinson algorithm, 25-25
- L'Huilier's theorem, 47-13
- Life expectancy, power grid tubes, 17-17-18
- Life testing, 45-13-14
- Lifetime, definition of, 20-7
- Lifetime, photon, 41-7
- Light
 -emitting diodes, 18-13, 18-23-24
 -emitting tubes, 16-45-48, 16-55-58
 -sensing tubes, 16-45-54
 speed of, 18-4
 visible, emitters, 21-11
- Light year, 3-15
- Limit cycle, 15-23
- Limits, confidence, 45-15-21
- Line
 parameters, transmission lines, 29-4-5
 shape function, laser, 41-2
 of sight transmissions, 33-13-29
 signal, telephone, 2-14-17
 straight, equations of a, 47-20
- Linear
 amplifier, 17-10-11
 array feeds, 32-36-38
 -beam tube, 16-25-34
 buses, network, 26-33-34
 difference equation, 28-3
 electro-optic effect, 41-17-19
 expansion, of elements, 4-9, 4-11
 integrated circuits, 20-42-45
 miscellaneous, 20-45
- Linear—cont
 special purpose, 20-45
 regulated power supplies, 14-43-44
 time-invariant digital systems, stability of, 15-18-20
 time-invariant systems
 control system, 15-12-13
 stability of, 15-5-10
 transformations, matrix algebra, 47-32
- Lines
 attenuation, 29-29-30
 fin-, 30-31
 microstrip, 29-25-27
 power rating, 29-30-31
 resonant, 29-15-19
 slot, 30-29-31
 transmission
 coplanar, 29-28-29
 strip, 29-27-28
- Link, 3-15
- Link budgets, satellite communications, 27-13-14
- Link control, multiaccess, 26-18-23
- Liquid crystal displays, 18-23
- Liquids, sound in, 40-4-5
- Listener echo, definition of, 38-10
- Liter, 3-15
- Literal, logic, definition of, 43-3
- Lithium, 4-4-5, 4-8-9
- Lithography, 20-14, 20-17-19
 mask generation, 20-14
 resist casting, 20-17-18
 resist development, 20-18-19
 resist exposure, 20-18
- LMS algorithm, 28-23-24
- Load line, definition of, 20-7
- Load reflection coefficient, Γ_L , definition of, 12-25
- Load resistance, power grid tube, 17-3
- Load, traffic carrying, for grade of service, 39-13-14
- Loading
 factor, input, definition of, 20-6
 factors, definition of, 20-7
 noise, 12-31
 steel radio towers, 49-20-21
 telephony circuit, 2-5-6
- Loads, active, transistor amplifiers, 19-17-19
- Local AM broadcast stations, 35-3
- Local area networks
 classes of, 26-31-32
 distributed control algorithms, 26-22-23
 high speed, 26-30-36
 low speed, 26-32
 medium speed, 26-32
- Local interoffice trunks, telephone, 38-6
- LOCOS, definition of, 20-7
- LOCSTAR, 27-11
- Log (0.4343), multiples of, table of, 48-4
- Log-periodic antennas, 32-30-32
- Logarithmic integrands, 47-24-25

- Logarithmic plots, 15-6-7
 Logarithms, base 2, table of, 48-5
 Log₁₀ (2.3026), multiples of, table of, 48-5
- Logic
 array, programmable, 20-51, 20-53
 binary, definition of, 20-3
 circuit, electronic latch, 43-16
 circuits, analysis of, 43-3-4
 complement, definition of, 43-3
 design
 graphic symbols, 49-29
 map method, 43-12-15
 synthesis, 43-4-8
 devices, programmable, 20-51, 20-53
 double-rail, 43-10
 dynamic, 20-36
 emitter-coupled (ECL), 20-32-33, 20-35
 families
 bipolar, 20-30-35
 MOS, 20-35-36
 integrated injection (I^2L), 20-33-34
 literal, definition of, 43-3
 negative, 43-2
 definition of, 20-8
 networks, 43-6-8
 positive, 43-2
 ratioed, 20-35-36
 ratioless, 20-36
 Schottky transistor (STL), 20-34-35
 single-rail, 43-10
 standard terminology, 20-32
 static, 20-35-36
 swing, definition of, 20-7
- Logical circuit routing, networks, 26-25
- Lognormal distribution, 45-11
- Long-cable method, return loss measurement, 12-30
- Longitudinal modulators, laser, 41-18
- Longitudinal redundancy check, 26-15
- Longitudinal slot
 resonant conductance, 32-20-21
 resonant length, 32-20-21
- Loop(s)
 antennas, 32-18-19
 spheroidal core, 32-19-20
 frequency tracking, 12-27
 gain, feedback, 15-11
 reverse-battery signaling, 38-16
 subscriber
 telephone, 38-5-6
 resistance design, 38-5
- Loran-C, 37-3
- Loss(es), 38-8
 characteristic insertion, 11-2
 cosecant-squared antenna, radar, 36-12-13
 echo return, definition of, 38-10
 flange insertion, waveguide, 30-17
 inductors, 13-26-28
 insertion, 12-26
 definition of, 38-9-10
 resonant lines, 29-17-18
- Loss(es)—cont
 magnetic core, 13-26
 nct, 38-9
 optical fiber, 22-21-22
 path, radio waves, 33-4
 power semiconductors, 14-31-32
 propagation, mobile radio, 46-9
 return, 12-29-30
 definition of, 38-10
 singing margin, definition of, 38-10
 singing return, definition of, 38-10
 transducer
 definition of, 38-10
 transmission lines, 29-13-14
 transmission, radio waves, 33-4
 turn-off, 14-32, 14-33
 turn-on, 14-32, 14-33
 winding-resistance, inductors, 13-27
- Loudspeaker, placement in distributed system, 40-24
- Low frequencies, propagation of, 33-3-6
- Low-gain antennas, 32-13-22
- Low-pass
 to band-stop transformation, filters, 9-19
 to bandpass transformation, filters, 9-14
- filter
 Butterworth, 9-14
 Cauer-parameter, 9-4-5
 Chebyshev, 9-14
 design
 Butterworth, 9-18, 9-20-21
 Cauer-parameter, 9-5-13
 Chebyshev, 9-15-17
 examples, 9-19, 9-22, 9-24-26
 reference, 9-3
 to high-pass transformation, filters, 9-14
 ladder, L-C, simulation of, 10-27-30
 notch, active filter, 10-44-45
- Low power auxiliary stations, 35-24
- Low-power television (LPTV) stations, 35-17-18
- Low voltage power supplies, world, 49-8, 49-10-11
- Lowest useful high frequency (luf), 33-8
- LPCVD, definition of, 20-7
- LQ-optimal control, controller design, 15-15
- Lumen, 3-3
- Luminous gain, 16-49
- Lutetium, 4-4-5, 4-8-9
- Lux, 3-3
- Magnet-wire data, 6-3-4
- Magnetic core
 inductors, 13-3, 13-19, 13-21
 loss, 13-26
 reactors, 13-3, 13-18-19
 types of, 13-3-5
- materials, 4-17, 4-19
 ferrites, 4-19, 4-25, 4-25-27
 permanent-magnet, 4-19, 4-25
 permeability, 13-4
 soft, 13-4
 soft magnetic metals, 4-19, 4-24
- metals, soft, properties of, 4-19, 4-24
- Magnetostriction, 4-27-29
- Magnetrans, 16-34-36
 -enhanced discharges, 20-95-96
 etch, cryogenic, 20-97-99
 oscillator, 16-34-36
 pulsed, 16-35
 tuning, 16-35
 voltage-tuned, 16-36-37
- Magnitude, formulas, series and parallel combinations, 6-9-13
- Mainframe memory, definition of, 20-7
- Majority carriers, 18-5
- Man-made radio noise, 34-6
- Manganesec, 4-4-5, 4-8-9
- MANs, 26-32
- Manufacturing
 device, intelligent cluster-tool-based, 20-63, 20-66
 integrated device, 20-58-66
 semiconductor, 20-56-58
- Map, Karnough, 43-13
- Map method, logic design, 43-12-15
- Marginal probability function, 44-2
- Marine storm warnings, WWV and WWVH, 1-156
- MARISAT communications system, 27-5
- Maritime Mobile Satellite Service (MMSS), 27-3
- Mark, FSK signal, 24-9
- Markov source, 25-8
- Maser, 41-2
- Masers, hydrogen, 1-147-148
- Mask generation, lithography, 20-14
- Mass number, 4-3, 4-4, 4-6
- Mass, relative atomic, 4-4, 4-6
- Mass storage devices, computer, 42-14
- Matched filter, 24-14-15
 estimator, 25-23-24
- Matching sections, quarter-wave, 29-19
- Material(s)
 acoustic properties, 4-31-33
 conducting, 4-3
 dispersion, optical, 22-22
 ferromagnetic, 13-3
 insulating, 4-12, 4-13
 characteristics of, 4-20-23
 dissipation factor, 4-13
 permittivity of, 4-12, 4-13
 transformers
 classes of, 13-10-11

M

- M-ary signaling, 24-10-11
- Machine language, computer, 42-3
- Machine-language program, computer, 42-22-24
- Macroscopic current density, 16-10
- Magnesium, 4-4-5, 4-8-9

- Material(s)—cont**
 operating temperatures, 13-10-11
 magnetic, 4-17, 4-19
 permeability, 13-4
 microwave, 4-27
 permanent-magnet, 4-19, 4-25
 semiconducting, 4-3, 4-18
 soft magnetic, 4-19, 4-24, 13-4
 steel radio towers, 49-20
 for tropical and marine use, 49-5, 49-8
- Mathematical expectation, 44-3-4**
- Matrix**
 algebra, 47-31-35
 Hadamard, 24-3
 of a junction, impedance and admittance, 31-9-10
 scattering
 of a junction, 31-3
 measurement of, 31-4-6
 state transition, 15-13-14
 transformation, 31-3-4, 31-10-11
- Matter, interaction with optical waves, 21-9-10**
- Maximum**
 available power, P_{avs} , definition of, 12-25
 data rate, optical fiber, 22-25
 entropy, 25-24
 usable frequency (muf), 33-7
- Maxterms, 43-5**
- Maxwell, 3-15**
 bridge, 12-5
- Mean**
 absolute deviation, 44-4
 absolute error, 44-4
 noise power, 12-31
 root mean square deviation from, 44-3
 standard deviation from, 44-3
 time before failure (MTBF), 45-3
 time between failure (MTBF), 45-12
 time to repair (MTTR), 45-3, 45-13
- Measurement(s)**
 admittance, 12-8
 ambient conditions for, 5-5, 5-7
 amplitude
 range, 12-18-19
 response, 12-26-27
 attenuation, 11-5-10
 broadband, 11-7
 fixed frequency, 11-5-6
 swept or stepped frequency, 11-9-10
 baseband, 12-30-33
 C_i and true L_c , with Q-meter, 12-8-10
 capacitor in series with unknown, 12-6
 carrier/noise, 12-33
 coupling coefficient, loosely coupled coils, 12-9-10
 diagnostic, microwave systems, 12-33
 direct capacitance, 12-6-7
- Measurement(s)—cont**
 electromagnetic interference, 12-36
 electromagnetic susceptibility, 12-38
 frequency, 12-21-23
 group-delay distortion, 12-28-29
 hybrid-coil method, 12-7
 impedances lower than directly measurable, 12-8-9
 large-signal, 12-16-17
 microphones, 40-8
 microwave frequencies, 12-22
 noise, 34-8-11
 broadcast receivers, 34-8-9
 figure, 34-9-10
 optical detectors, 21-18
 Q of resonant circuit, by bandwidth, 12-7
 reflection coefficient, 31-2
 return loss, 12-30
 RF, 12-24-26
 sound, 40-8-9
 time, 12-21
 and frequency, 12-20-24
- Measures, traffic, 39-11**
- Mechanically programmable semiconductor ROM, definition of, 20-7**
- Median, 44-3**
- Medium frequencies, propagation of, 33-3-6**
- Medium-scale integration, definition of, 20-7**
- Melting point, of elements, 4-8, 4-10**
- Memory**
 application specific, 20-54
 buffer, definition of, 20-4
 computer, 42-12-14
 integrated circuits, 20-39-42
 main, microprocessor, 20-38
 main, two-level, definition of, 20-11
 mainframe, definition of, 20-7
 mapped I/O, 42-32
 random-access (RAM), 20-41-42
 read-only, nonvolatile, 20-40-41
 read-only (ROM), 20-40
- Memoryless channel, 25-12**
- MEMS, 20-49**
- MEMS, definition of, 20-7**
- Mendelevium, 4-4-5, 4-8-9**
- Mensuration equations, 47-4-6**
- MEO satellites, 27-11**
- Mercury, 4-6-7, 4-8-9**
- Mercury-vapor rectifier, 16-42**
- Merged analog and digital ASICs, 20-55-56**
- MESFET, definition of, 20-7**
- Message processor, 26-2**
- Message set, digital communications, 24-2**
- Message switching, 26-2**
- Messages, digital communications, 24-2**
- Metal-oxide semiconductor transistor (MOST), 18-15-16**
- Metal sheet resistance sensor, 20-74**
- Metallic cleans, 20-91**
- Metallic glass, 13-4**
- Metallization, definition of, 20-7**
- Metallization isolation, 20-101-102**
- Metallization, two-level, definition of, 20-11**
- Metals**
 characteristic impedances of, 4-33
 densities of, 4-33
 field emission, 16-7-9
 resistivity of, 4-15-16
 soft magnetic, properties of, 4-19, 4-24
 superconductivity of, 4-17
 velocities of, 4-33
- Meteor-burst propagation, 33-13**
- Meter, 3-3, 3-15, 3-16**
 cubic, 3-3
 impedance, automatic, 12-11-13
 LCR, automatic, 12-11-12
 per second, 3-3
 per second squared, 3-3
 sound level, 40-8
- Metropolitan Area Network, 26-32**
- MFSK. See Multiple Frequency Shift Keying**
- Mho, 3-16**
- Mica capacitors, 5-25-26**
- Micro-electro-mechanical systems (MEMS), 20-49**
- Micro-electro-mechanical systems (MEMS), definition of, 20-7**
- Microcomputers, ICs for, 20-36-39**
- Micron, 3-16, 18-4**
- Microphones**
 laboratory, 40-8
 measurement and calibration, 40-8
- Microprocessor(s)**
 central processor unit (CPU), 20-37-38
 characteristics, 20-38, 20-39
 i/o section, 20-36-37
 ICs for, 20-36-39
 memory, main, 20-38
- Microstrip lines, 29-25-27**
- Microstrip patch antenna, 32-17-18**
- Microwave**
 applications for gas tubes, 16-44-45
 devices, semiconductor, 18-24-26
 discharges, 20-96-97
 frequencies, measurement of, 12-22
 landing system, 37-5
 -link analysis, 12-26-33
 materials, 4-27
 power modic (MPM), 16-24, 16-32-33
 propagation, profile paper, 33-16
 systems
 path plotting, 33-15
 performance, 12-33
 tubes, 16-23-38
 cyclotron resonance, 16-38-39
 noise in, 16-19
 terminology, 16-23-24
 trends, 16-24-25
- Mil, 3-16**

- Mile
 international, 3-16
 nautical, 3-16
 statute, 3-16
- Military equipment
 developmental indicators, 49-36-37
 group indicators, 39-35
 training sets, 49-37
 unit indicators, 49-38-39
- Military nomenclature system,
 summary of, 49-34-39
- Miller integrators, 10-11-12
- Millibar, 3-16
- Millimeter, of mercury, 3-16
- Miniaturization, trends in, 20-45-46
- Minimal sums, formation of, 43-14-15
- Minimax criterion, 24-14
- MINIMOS, definition of, 20-7
- Minimum-loss pads, resistance values,
 11-6
- Minimum-shift keying (MSK), 24-8
- Minimum sum, of function, 43-12
- Minislotted alternating priorities
 (MSAP), 26-22
- Minority carriers, 18-5
- Minterms, 43-5
- Minute, 3-16
- Mirrors, current, 19-17-19
- MIS, definition of, 20-7
- Mismatch error, multiple, M_{ii} ,
 definition of, 12-25-26
- Mismatch, transmission lines,
 29-13-14
- Mixed mode method, AM stereo, 35-5
- MLS, 37-5
- Mobile charge, definition of, 20-7
- Mobile-originated call, 46-7
- Mobile radio
 environment, 46-8-17
 path loss, 46-8-12
 prediction of received signal,
 46-9-12
 propagation loss, 46-9
- Mobile-Satellite Service, 27-31
- Mobile services, satellite, 27-28-31
- Mobile signal
 average duration of fades, 46-12
 level-crossing rate, 46-12
 multipath fading, 46-12-17
 selective fading, 46-12, 46-17
 time delay spread, 46-12, 46-17
- Mobile subscriber units, 46-7
- Mobile Telephone Switching Office
 (MTSO), 46-5
- Mobile-unit initialization, 46-7
- Mobility, definition of, 20-7
- Modal dispersion, optical fiber, 22-22
- Mode, 44-3
- Modeling, IC design, 20-26
- Modelocking, 41-6
- Models, source, 25-7-8
- Modern network theory, 9-2
- Modulated waveforms, analytic signal
 representation of, 23-3
- Modulation
 AM requirements, 35-5
- Modulation—cont
 amplitude, 23-4-5
 conventional, 23-4
 detection, 23-5-6
 double sideband, 23-4
 single sideband, 23-4-5
 vestigial sideband, 23-4
 angular, 23-6
 capability, FM broadcast, 35-10
 continuous, 23-3-10
 delta, 23-11, 23-18-19
 digital, 27-23
 exponential, 23-6-10
 field emission, 16-12
 frequency, 23-7
 intensity, 22-3
 phase, 23-7
 pulse, 23-10-11
 pulse-amplitude, definition of, 23-10
 pulse-code, definition of, 23-11
 pulse duration, 17-13
 pulse-duration, definition of, 23-11
 pulse-frequency
 definition of, 23-11
 natural sampling, 23-14
 uniform sampling, 23-14-15
 pulse-length, definition of, 23-11
 pulse-phase
 definition of, 23-11
 natural sampling, 23-15
 uniform sampling, 23-15
 pulse-position, 22-3-4
 definition of, 23-11
 natural sampling, 23-15
 uniform sampling, 23-15
 pulse-time, 23-13-14
 definition of, 23-10
 spectra of pulses, 23-13-14
 pulse-width, definition of, 23-11
 quaternary direct sequence, 24-24
 ratio, FSK, 24-9
 techniques, optical communications,
 22-3-5
 two-tone angle, 23-8
- Modulators
 longitudinal, laser, 41-18
 polarization, 22-3
 power grid tube, 17-12-13
 transverse, laser, 41-18
- Modules, multiple-chip, 20-54-55
- Mole, 3-3
 per cubic meter, 3-3
- MOLNIYA satellites, 27-11
- Molybdenum, 4-6-7, 4-10-11
- Moment of order r , 44-3
- Monochromatic radiation,
 photocathode response to,
 21-22
- Monolithic integrated circuit, definition
 of, 20-7
- Monopole antennas, 32-14
- Monostable multivibrator, 19-29-30
- Month, 3-16
- Morse code, 25-4-6
- MOS
 capacitor, definition of, 20-7
- MOS—cont
 complementary, definition of, 20-4
 definition of, 20-7
 logic families, 20-35-36
- MOSFET, 18-15-16
 definition of, 20-7
 large symbol model, 19-2-3
 transistor symbol, 19-2-3
- MOSFETs, 14-29, 14-30, 14-32
- MOST *See* Metal-Oxide
 Semiconductor Transistor
- Motion, harmonic, equations, 49-22-23
- Motorola system, AM stereo, 35-6
- Motors
 ac, drive for adjustable-speed, 14-48
 dc, 15-2-3
 full-load running currents, NEC,
 49-15
 starting current data, NEC, 49-15
 step, 15-3-4
 synchronous-, drive, 14-49-50
- Moving target radar, 36-14-18
- MSK *See* Minimum-Shift Keying
- MTBF, definition of, 20-7
- MTI radar, 36-14-18
 staggered PRF, 36-17
- μ , of electron tube, 16-20
- Multiaccess
 -broadcast, packet communication
 network, 26-6
 link control, 26-18-23
 centrally controlled demand
 assignment, 26-20-21
 random access techniques,
 26-18-20
 protocols, 26-16
- Multichannel Multipoint Distribution
 Service (MMDS), 35-20
- Multichannel television sound system
 (MTS), 35-14
- Multifrequency numerical code
 four-by-four, 2-14, 2-16
 two-out-of-six, 2-14, 2-16
- Multifrequency signaling, 38-17-18
- Multilevel oxide, definition of, 20-7
- Multipath fading, mobile signal,
 46-12-17
- Multipath reduction factor, 33-10-11
- Multiple-access schemes, satellite
 transponders, 27-19-21
- Multiple-beam arrays, 32-41-44
- Multiple-beam reflector, 32-57-58
- Multiple-chip modules (MCM),
 20-54-55
- Multiple-emitter cell, definition of,
 20-7
- Multiple-feedback topologies,
 10-25-27
- Multiple frequency shift keying
 (MFSK), 24-10
- Multiple mismatch error, M_{ii} , definition
 of, 12-25-26
- Multiplexing, telephone systems, 38-13
- Multiplication, binary, 42-8-9
- Multipliers, voltage, 14-13

- Multipoint Distribution Service (MDS), 35-20
- Multitone angle modulation, 23-8
- Multitone FM, bandwidth considerations in, 23-9
- Multivariate, 44-2
- Multivibrators, 19-29-30
- Multizone RTP illuminators, 20-67-68, 20-69
- Mutual impedance, antenna array, 32-38-40
- Mutual-inductance, capacitance balance, 12-7
- N**
- N , definition of, 20-8
- N-type semiconductor, definition of, 20-8
- n -variable theorems, switching algebra, 43-9-10
- N_A , definition of, 20-8
- NAND gate, definition of, 20-6
- Napier's analogies, 47-16
- National Association of Broadcasters, 35-27
- National Electrical Manufacturers Association, 5-4
- National Telecommunications and Information Administration (NTIA), 1-2-3
- Native oxide removal techniques, comparison of, 20-87-89
- Natural frequency, vibration, 49-24
- Natural noise, 34-2-3
- Natural sampling, pulse amplitude modulation, 23-11
- Nautical mile, 3-16
- Navigation agencies, major, 37-2
aids, radio, 37-3-8
system, satellite, 37-5-8
- NAVSTAR, 37-6-8
- N_D , definition of, 20-8
- Near-field exact, dipoles, 32-13
power density, antenna, 32-7-8
region, antenna, 32-4
- Near-zone noise sources, 34-6-7
- NEC allowable current-carrying capacities of conductors, 49-14
conductor applications and insulations, 49-12
current-carrying capacity, flexible cords, 49-13
flexible-cord data, 49-13
motor full-load running currents, 49-15
- Negative feedback, 19-20
- Negative logic, 43-2
definition of, 20-8
- Negative numbers, representation of, 42-7
- NEMA, 5-4
- Neodymium, 4-6-7, 4-10-11
- Neon, 4-6-7, 4-10-11
- NEP *See* Noise Equivalent Power
- Nepers, 3-4, 3-6
- Neptunium, 4-6-7, 4-10-11
- Net loss, 38-9
- Net power flow, waveguide, 31-9
- Network(s)—cont
logic, 43-6-8
nonblocking, 39-2
nonblocking 3-stage, 39-4-6
comparison with single stage, 39-5, 39-7
folded, 39-6
-originated call, 46-7-8
packet communications, 26-5-6
packet switching, 39-9
presentation layer, 26-30
RC, response of, 19-26-27
ring, 26-7-8
RL, response of, 19-26
RLC, 6-17-19
equations for, 6-14-15
response of, 19-28-29
Rotman-Gent beam-forming, 32-43-44
sensitivity, 10-8-10
services, 26-23-24
routing, 26-24-25
session layer, 26-29-30
six-port, 12-15
sorting stages, 39-9-10
space-division switching, 39-2
STM allocation, 26-37
structures, digital filters, 28-16-20
supercomputer, 26-34
switching, 39-2
coordinate, 39-2-10
T- π transformation, 6-15-16
theory, modern, 9-2
time constant, 6-15
time division switching, 39-2
transformation, T- π , 6-15-16
transients, elementary cases, 6-15-19
transport layer, 26-29-30
trees, 26-34
ultragigabit, 26-35-36
Y- Δ transformation, 6-15-16
- Newton, 3-3
meter, 3-3
per meter, 3-3
- Neyman-Pearson theorem, 25-22-23
- Nibble, definition of, 20-8, 42-4
- Nickel, 4-6-7, 4-10-11
- Niobium, 4-6-7, 4-10-11
- Nitrogen, 4-6-7, 4-10-11
- NMOS, definition of, 20-8
- NMOS process, 20-22-26
- Nobelium, 4-6-7, 4-10-11
- Noise
1/f, 21-19-20
amplifier, 21-20
atmospheric, 32-9, 34-3-4
circuit, telephony, 2-11-13
cosmic, 27-12-13, 34-4-5
definition of, 38-9
environmental, 40-10-12
equivalent power, optical detectors, 21-14
exposure, permissible, 40-12
factor
antenna, 34-2

- Noise**—cont
- cascaded networks, 34-11
 - receiver, 34-9
 - falsc-pulse, in pulse-code modulation, 23-17-18
 - field strength, 34-4
 - figure, 34-9
 - calculation of, 34-10-11
 - measurement of, 34-9-10
 - galactic, 32-9, 34-5
 - Gaussian, 32-9
 - general, optimum threshold for, 24-14
 - generation-recombination, 21-19
 - immunity, definition of, 20-8
 - impulse, definition of, 38-9
 - induced, electron tube, 16-19
 - jamming radar equation surveillance, 36-9
 - tracking, 36-9
 - Johnson, 21-18
 - levels
 - indoor design, 40-15, 40-17
 - interior
 - design, 40-17
 - planning, 40-18
 - loading, 12-31
 - man-made, effect on mobile communications, 46-8
 - measurements, 34-8-11
 - broadcast receivers, 34-8-9
 - optical detectors, 21-18
 - microwave tubes, 16-19
 - natural, 34-2-3
 - Nyquist, 21-18
 - optical detectors, 22-2-3
 - photon, 21-20
 - power
 - line, 34-7-8
 - mean, 12-31
 - psophometric, 2-4-5
 - ratio, 12-31-32
 - psophometric, 2-4-5
 - quantization, in PCM system, 23-16-17
 - quantizing, 38-9
 - radio, man-made, 34-6
 - receiver, theoretical, 33-27
 - reference, definition of, 38-9
 - resistors, 5-9, 5-11
 - shot, 21-19
 - solar, 27-12
 - sources
 - far-zone, 34-6-7
 - gas-discharge tube as, 16-44-45
 - near-zone, 34-6-7
 - sound power levels of typical, 40-9
 - temperature, 21-20, 27-12, 27-15, 32-8-11
 - antenna, 32-8-11
 - clear sky, 27-36
 - effective, receiver, 34-10
 - losses and, 30-17-18
 - sky, with rain, 27-37-39
 - thermal, 21-18, 34-3
- Noise**—cont
- calculations, 34-8
 - definition of, 38-9
 - in tubes, 16-18-19
 - workplace-industrial, assessment of, 40-12
- Noiseless channels**
- capacity of discrete, 25-4-7
 - coding for, 25-4-7
- Noisy channel, 25-12**
- Non-SI units, tables of, 3-4**
- Nonblocking network, 39-2**
- 3-stage, 39-4-6
 - comparison with single stage, 39-5, 39-7
 - folded, 39-6
- Noncoherent correlation detector, 24-17**
- Noncoherent receivers, 24-18-26**
- Nondestructive readout, definition of, 20-8**
- Nongcostationary, orbits, 27-11**
- Nonlinear feedback control systems, 15-22-32**
- characteristics of, 15-22-23
- Nonlinearity**
- backlash-type, 15-25-27
 - saturation-type, 15-25
- Nonrecursive filter, 36-16**
- Nonvolatile, read-only memories, 20-40-41**
- NOR gate, definition of, 20-6**
- Normal distribution, 44-5, 45-11**
- table of, 48-8
- Normal form, equation, 47-17**
- Normalization, frequencies, 9-2-3**
- Normalized**
- frequencies, 9-3
 - resistances, 9-3
- North America, cellular systems in, 46-3-4**
- North American Digital Cellular (NADC) system, 46-22**
- Norton amplifier, 20-43**
- NOT operation, definition of, 20-8**
- Notation, matrix algebra, 47-31**
- Npn transistor**
- large symbol model, 19-2
 - symbol, 19-2
- NPR to SNR, conversion, 12-32-33**
- NTSC-M color television system, 35-40**
- Nuclear explosions, effect on radio propagation, 33-12**
- Null hypothesis, 25-22**
- Number(s)**
- atomic, 4-3, 4-4, 4-6, 4-8, 4-10
 - Avogadro's, 18-4
 - direction, 47-19-20
 - mass, 4-3, 4-4, 4-6
 - negative, representation of, 42-7
 - nondecimal, addition and subtraction of, 42-6
 - per meter, 3-3
 - system(s), 42-4-10
 - base, 42-4-5
- Number(s)**—cont
- binary, 42-5
 - complement, 42-7
 - conversion, 42-6
 - decimal, 42-4
 - hexadecimal, 42-5
 - octal, 42-5
 - positional, 42-4-5
 - radix, 42-4-5
- Numerical**
- analysis, 47-43-45
 - aperture, optical fiber, 22-20
- Numerical values, 3-6**
- NVRAM, definition of, 20-8**
- Nyquist**
- criterion, 35-36-37
 - diagram, nonlinear system, 25-27
 - loci, typical, 15-6
 - noise, 21-18
 - stability criterion, 15-5-6
- O**
- Object module, computer, 42-24**
- Oblate spheroid, equations for, 47-20**
- Oblique hyperbolic triangle, equations for, 47-15-16**
- Oblique triangle, equations for, 47-11-13**
- OC curve, definition of, 20-8**
- Octal, number system, 42-5**
- ODE, definition of, 20-8**
- Oersted, 3-16**
- Offset carrier, cochannel TV stations, 35-12**
- Offset quadruphase shift keying (OQPSK), 24-8**
- Ohm, 3-3**
- Ohmic contact, definition of, 20-8**
- Omega, 37-3**
- Omnidirectional range, VHF, 37-4**
- Onboard processing asynchronous transfer mode (OBP ATM), 27-19**
- One-minute rainfall rate, various cities, 33-24-25**
- 1/f noise, 21-19-20**
- Ones-complement, system, 42-7**
- Opcode word, 42-16**
- Operand, 42-16**
- Operands, computer instructions, 42-20-22**
- Operating data, typical amplifier, 17-2**
- Operation, parallel, definition of, 20-8**
- Operational amplifiers, in active filters, 10-5-7**
- Operational transconductance amplifiers (OTA), 20-43**
- in active filters, 10-7
- Operational voltage amplifier, 20-43, 20-44**
- Operations, matrix algebra, 47-31**
- Operator, Del, 47-35**
- Optical**
- cavities, 41-6-8

- Optical—cont
 stability of, 41-7
 coefficients, nonlinear, 41-17
 communications
 atmospheric effects, 22-13-19
 background radiation, 22-12-13
 coherent sources, 22-7-8
 modulation techniques, 22-3-5
 detectors, 21-14-30, 22-2-3
 characterization of, 21-16-18
 d^* , 21-14-15
 d^{**} , 21-15
 detectivity, 21-14
 noise, 22-2-3
 noise equivalent power, 21-14
 noise measurements, 21-18
 quantum efficiency, 21-15-16
 responsivity, 21-14-18
 spectral response, 21-16
 ultimate sensitivity of, 21-18-21
 excitation, semiconductors, 18-6
 fiber
 connectors, 22-24-25
 dispersion, 22-22-24
 group delay, 22-20-21
 group index, 22-20-21
 losses, 22-21-22
 numerical aperture, 22-20
 splicing, 22-24
 terminology, 22-19
 types of, 22-19-20
 V-number, 22-21
 imaging, 16-47-48
 receivers, 22-8-12
 heterodyne-detection, 22-10-12
 SNR, 22-11
 signal
 scattering, 22-16-17
 turbulence, 22-17-19
 sources, 21-10-14
 spectrum, 21-2
 thickness, propagation path, 22-14
 transmitters, 22-5-8
 waves
 absorptance, 21-9
 interaction with matter, 21-9-10
 reflectance, 21-9
 refraction, 21-10
 transmittance, 21-9-10
 Optics, nonlinear, 41-9-12
 Optimum threshold, Gaussian noise,
 general, 24-14
 Optimum working frequency (f_o),
 33-8
 OQSK *See* Offset Quadrature Shift
 Keying
 OR, definition of, 20-8
 OR gate, 43-3
 definition of, 20-6
 Orbit
 geostationary, 27-11
 nongeostationary, 27-11
 satellite, 27-6-7
 choice of, 27-8-9
 Organic cleans, 20-89-91
 Orthogonal curvilinear coordinates,
 vector equations for, 47-37
 Orthogonal polynomials, 47-43
 Orthogonal signal, 24-9
 set, 24-3
 Orthogonal signaling, binary, 22-5
 Oscillation, CW, 41-4-5
 Oscillator(s)
 backward-wave, 16-33-34
 cesium, 1-146-147
 class-C, design of, 17-4-6
 Colpitts, 19-25
 equations, 19-25
 Hartley, 19-25
 magnetron, 16-34-36
 phase-shift, 19-25
 quartz, 1-144
 rubidium, 1-144, 1-146
 sinusoidal, 19-24-26
 sweep, 12-14
 voltage-tuned crossed-field,
 16-36-37
 Wien bridge, 19-25
 Osmium, 4-6-7, 4-10-11
 OTA *See* Operational
 Transconductance Amplifiers
 Ounce, 3-16
 avoirdupois, 3-16
 force, 3-16
 mass (troy or apothecary), 3-16
 Out of band emissions, AM, 35-5
 Output stages, transistor, 19-15-16
 Over-the-horizon radar equation, HF,
 36-10
 Over-the-horizon radar, HF, 36-16
 Overflow oscillations, digital filter,
 28-22
 Overlap-add algorithm, 28-9
 Overlap-save algorithm, 28-8-9
 Owen bridge, 12-5
 Oxidation, definition of, 20-8
 Oxidation number, elements, 4-3
 Oxide, multilevel, definition of, 20-7
 Oxide, native, comparison of removal
 techniques, 20-87-89
 Oxygen, 4-6-7, 4-10-11
- P**
- P , definition of, 20-8
 P -type semiconductor, definition of,
 20-9
 Pace, 3-16
 Packaging, integrated circuits,
 20-28-30
 Packet
 communications network, 26-5-6
 multiaccess-broadcast, 26-6
 switching, 26-3
 networks, 39-9
 Pads
 attenuator, resistance values, 11-4
 minimum-loss, resistance values,
 11-6
 tee, power dissipation within, 11-7
 Palladium, 4-6-7, 4-10-11
 PAM *See* Pulse-Amplitude Modulation
 Paper-foil capacitors, 5-23
 Paper, probability, 45-22
 Parabola, equations for, 47-4, 47-18
 Parabolic reflectors, 32-54-57
 Paraboloid
 equations for, 47-6
 of revolution, equations for, 47-20
 Parallel, definition of, 20-8
 Parallel operation, definition of, 20-8
 Parallel-plate capacitor, capacitance of,
 6-14-15
 Parallelogram, equations for, 47-4
 Parasitic-insensitive SC biquads, 10-38,
 10-39, 10-40-44
 Parasitic patch, antenna, 32-18
 Parasitics, definition of, 20-8
 Parsec, 3-16
 Partial-fraction expansion method,
 inverse z -transfer equation,
 15-17
 Partition noise, electron tube, 16-18
 Pascal, 3-3
 second, 3-3
 Passivation, definition of, 20-8
 Passive elements, definition of, 20-8
 Patent
 importance of records, 49-33
 obtaining, 49-33
 protection, of inventions, 49-32-34
 Path clearance, radio transmission,
 33-17
 Path loss
 antenna, 32-3
 radio waves, 33-4
 Path plotting, microwave system, 33-15
 PCM *See* Pulse-Code Modulation
 PDM *See* Pulse-Duration Modulation
 Peak sound pressure level, definition of,
 40-10
 Peck, 3-16
 Pennyweight, 3-16
 Pentode RF amplifier circuit, 17-8-9
 Perch, 3-16
 Periodic classification of the elements,
 4-3, 4-12
 Periodic waveform
 average power, 7-8
 and Fourier series, 7-11-13
 Peripheral devices, computer, 42-14,
 42-30
 Permalloy powder, 13-4
 Permeability, magnetic materials, 13-4
 Permittivity
 of free space, 18-4
 of insulating materials, 4-12, 4-13
 Personal Communication Network
 (PCN), 46-24
 PFM *See* Pulse-Frequency Modulation
 Phase
 angles, series and parallel
 combinations, formulas,
 6-9-13
 -control power supplies, 14-44
 -locked loop servo systems,
 15-21-22

- Phase—cont
 modulation, 23-7
 narrow band, phase deviation, 23-7
 wideband, phase deviation, 23-7-8
 -plane method, nonlinear system analysis, 15-30-32
 portrait, construction of, 15-31
 Phase shift
 keying, 24-19-21
 binary, 24-5-6
 differential, 24-19-21
 oscillator, 19-25
 pairs of coupled tuned circuits, 8-7-9
 single-tuned circuits, 8-6-7
 Phased array radar tracking, 36-22
 Phaser quantization, 32-35-36
 Phasor, velocity, sound wave, 40-3
 ϕ_r , definition of, 20-9
 ϕ_s , definition of, 20-9
 Phosphorus, 4-6-7, 4-10-11
 Phot, 3-16
 Photocathode response, 16-14-15
 to monochromatic radiation, 21-22
 responsivities, typical peak, 21-23
 Photocathodes, 16-14-16
 Photoconductive cells, 21-28-29
 Photoconductivity, 21-26-28
 Photodiodes, 21-28-30
 gas, 16-50
 vacuum, 21-22-25
 Photoemissive devices, 21-22
 Photomasking, definition of, 20-8
 Photometric equivalents, 16-46
 Photometry, 16-45-46
 Photomultipliers, 21-25
 as scintillation and single-electron counters, 21-25-26
 Photon
 energy, 21-2
 lifetime, 41-7
 noise, 21-20
 Photosurfaces, characteristics of
 typical, 21-24
 Photovoltaic
 detector, 21-18
 junction, 18-13
 Photovoltaic cells, 21-28-29
 Physical properties of elements, 4-3, 4-8-11
 Physical constants, 18-4
 Physical constants, fundamental, 3-7-13
 Physical layer, computer network, 26-9-13
 Physical properties of elements, 4-3, 4-8-11
 Pica, 3-16
 Piezoelectric, strain coefficients, 4-29
 Piezoelectricity, 4-28-31
 Pinholes, definition of, 20-8
 Pint, 3-16
 Pipe, steel, towers, 49-20
 Pirahna, definition of, 20-8
 PLA, definition of, 20-8
 Planar transistor, definition of, 20-8
 Planar triode, 16-22
 Planck's constant, 18-4, 21-2
 Planck's law, 21-3, 21-6
 Plane analytic geometry, 47-16-19
 Plane, equations of, 47-20
 Plane trigonometry, 47-11-12
 Plane waves, 40-3
 Plant, exchange, 38-3-6
 Plasma
 definition of, 20-8
 deposition, thin films, 20-13
 etching, gases used with, 20-20
 processing, single-wafer, 20-91-99
 Plasma etch, definition of, 20-8
 Plastic film capacitors, 5-17-18
 Plastic packages, ICs, 20-28
 Plate dissipation, power grid tube, 17-3
 Plate resistance
 ac (r_p) of electron tube, 16-20
 total (R_p) of electron tube, 16-20
 Platinum, 4-6-7, 4-10-11
 Plug bar, definition of, 20-8
 Plutonium, 4-6-7, 4-10-11
 PMOS, definition of, 20-8
 Pn junction, definition of, 20-8
 Pn junctions, 18-8-12
 at equilibrium, 18-8-9
 Pockels effect, 41-17-19
 Point, 3-16
 Point-point form, equation, 47-17
 Point-slope form, equation, 47-17
 Poise, 3-16
 Poisson
 distribution, 44-5, 45-11
 equation, 39-12
 Polar-plot, nonlinear system, 25-27
 Polarization
 circular, 32-11-12
 diversity, 33-11
 elliptical, 32-11-12
 FM broadcast, 35-10
 modulators, 22-3
 Pole, 3-16
 Pole placement, controller design, 15-14
 Polling systems, network, 26-20-21
 Polonium, 4-6-7, 4-10-11
 Poly, definition of, 20-8
 Polygamma function, 47-40
 Polygon, regular, equations for, 47-4
 Polynomial
 interpolation, 47-44
 orthogonal, 47-43
 Polypropylene capacitors, 5-17
 Polysilicon electrode deposition, 20-100-101
 Polystyrene capacitors, 5-17
 Population, 44-2
 Ports, computer, 42-31-32
 Positioning servo, 15-4
 Positive
 error constant, unity feedback, 15-11
 feedback, 19-20
 logic, 43-2
 Postulates, 43-2-3
 Pot-core inductor, 13-23
 Potassium, 4-6-7, 4-10-11
 Potential barrier, definition of, 20-8
 Potential, ionization, first, 4-3, 4-5, 4-7
 Potentiometers, 5-13
 in error measuring systems, 15-2, 15-3
 terminal identification, 5-16
 Pound
 -force, 3-15, 3-16
 -mass, 3-15, 3-16
 Poundal, 3-16
 Powder
 iron, 13-4
 permalloy, 13-4
 Power
 average, periodic waveform, 7-8
 densities, in proximity to antenna, 49-30
 density
 antenna, 32-7-8
 signal, optical, 22-7
 dissipated in load, P_L , definition of, 12-25
 dissipation
 resonant lines, 29-19
 within tee network, 11-4, 11-7
 within tee pad, 11-7
 effective radiated, 35-12
 entropy, 25-18
 factor
 correction, active, 14-38-39
 correction circuits, 14-36
 definition, 14-38
 power supply filters, 14-41
 flow, waveguide, 31-2
 -flux density at satellite, up-link, 27-14
 gain
 apparent, 33-26
 RF amplifiers, 17-17
 grid tube
 amplifiers, classes of operation, 17-2
 circuit, classification, 17-8
 development trends, 17-18
 general design considerations, 17-2-4
 life expectancy, 17-17-18
 modulators, 17-12-13
 performance comparison with solid-state devices, 17-13-17
 special interest circuits, 17-10-13
 UHF operation, 17-10
 incident upon load, P_i , definition of, 12-25
 levels, telephony, 2-6-7
 line, noise, 34-7-8
 mean noise, 12-31
 -measurement definitions, 12-25-26
 module, microwave, 16-24, 16-32-33
 noise equivalent, optical detectors, 21-14
 noise, ratio, 12-31-32
 output, under pulse conditions, 17-17

- Power—cont
 phopsometric, 2-5
 ranges, earth stations, 27-40
 rating
 cables, 29-30-31
 lines, 29-30-31
 reflected by load, P_r , definition of, 12-25
 semiconductor
 characteristics, 14-3-12
 losses, 14-31-32
 -series method, inverse z -transfer evaluation, 15-17
 spectral density, digital communications, 24-11-13
 supplies
 ac-output, 14-46-51
 dc-output, 14-43-46
 distributed, 14-45-46
 linear regulated, 14-43-44
 low voltage, world, 49-8, 49-10-11
 phase-control, 14-44
 switching, 14-45
 transistor and diode requirements, 14-20-27
 unregulated, 14-43
 supply
 filters
 harmonic currents, 14-41
 phase-controlled thyristor, 14-41, 14-43
 power factor, 14-41
 uninterruptible, 14-47-48
 wiring, 49-11-12
 television broadcasting, maximum, 35-11
 transformer
 copper regulation, 13-6
 design for rectifiers, 13-5-9
 determining number of turns, 13-6
 determining wire size, 13-6
 effect of duty cycle on design, 13-9
 VA rating vs lamination sizes, 13-7
 wire table for, 13-8
 transformers, 13-3-10
 transistor, 14-5-7
 breakdown sustaining voltage, 14-5
 dc current gain, 14-5
 reverse voltage operation, 14-6
 saturation voltage, 14-5
 second breakdown, 14-5-6
 switching times, 14-6
 transmission lines, 29-10-12
 transmitter, for required output signal/noise ratio, 33-27-29
- Powers
 of 2, table of, 48-5
 of integers, sums of, 47-8
 of reciprocals of integers, sums of, 47-8-9
- PPM *See* Pulse-Position Modulation
 Praseodymium, 4-6-7, 4-10-11
- Precharge time (reset time), definition of, 20-9
- Precipitation
 attenuation due to, 33-20-21
 extremes, 49-4-5
 static, 34-8
- Precision adjustable inductors, 13-23, 13-26-28
- PREDEP, definition of, 20-9
- Predeposition and drive-in, substrate doping procedures, 20-21
- Preemption tone, proposed standard, 2-18
- Preferred values, components, 5-4, 5-6
- Prefixes, SI units, 3-3
- Presentation layer, network, 26-30
- Pressure-altitude graph, 49-3
- Pressure, on steel tower by wind, 49-22
- Prime implicants, 43-13-14
- Principle of duality, 43-3
- Printed circuit antennas, 32-17-18
- Printed circuits, 5-26-31
 materials, 5-27-28
- Priority-oriented demand assignment (PODA), 26-23
- PRO channel, 35-14
- Probability
 of bit error (PBE), 22-4
 definitions for, 44-2-3
 density function, 44-2
 distributions, continuous, 44-2, 44-5
 distributions, discrete, 44-2, 44-4
 distributions, theoretical, 44-4-5
 function, conditional, 44-2-3
 marginal, 44-2
 paper, 45-22
 and statistical inference, 45-15-25
 of word error (PWE), 22-4
- Process and device integration, SWP-based, 20-99-102
- Process control
 device fabrication, 20-63, 20-64-65
 feedback sensors, 20-60-61
- Process development, RTP based, overview, 20-76-78
- Processes, flow, definition of, 20-9
- Processor, message, 26-2
- Processors
 computer, 42-10-12, 42-14-22
 general-register, 42-11
 organization of, 42-14-16
- Production standards
 program, 35-26-34
 sound, 35-27-29
- Products
 canonical, 43-5-6
 fundamental, 43-5
- Profile paper, microwave propagation, 33-16
- Program
 computer, operation of simple, 42-25-30
 counter, computer processor, 42-15
 production standards, 35-26-34
- Program—cont
 signals, network distribution, 35-20, 35-22-23
 sound, 35-27-29
- Programmable logic array (PLA), 20-51, 20-53
- Programmable logic devices (PLDs), 20-51, 20-53
- Programmable read-only memory (PROM), 42-13-14
- Programming
 computer, 42-14-22
 input/output, 42-32
- Projective chart, 31-6
- Projective model, hyperbolic trigonometry, 47-14
- Prolate spheroid, equations for, 47-20
- PROM, 20-40-41; *See also* Programmable Read-Only Memory
- Promethium, 4-6-7, 4-10-11
- Propagation
 delay, definition of, 20-9
 errors, reduction of, 37-2
 freq-space, equations for, 33-27
 high-frequency forecasts, 33-8-10
 ionospheric scatter, 33-12
 knife-edge diffraction, 33-29
 loss, mobile radio, 46-9
 meteor-burst, 33-13
 path, optical thickness, 22-14
 satellite communications, 27-32-40
 time, one-way, telephony, 2-9
 tropospheric scatter, 33-29, 33-32
 wave, above 30 MHz, 33-13-29
 in waveguides, 30-3
- Properties, general, the elements, 4-3, 4-4-7
- Properties, physical, the elements, 4-3, 4-8-11
- Proportional-integral-derivative (PID) controller, 15-12
- Protactinium, 4-6-7, 4-10-11
- Protective resistance, 16-13
- Protocol(s)
 assigned-slot
 listen-before-transmission, 26-22
 asynchronous DLC, 26-13-14
 CCITT, X-21, 26-11-13
 DLC
 bit-oriented, 26-16-18
 character-oriented, 26-15-16
 EIA, RS-232-C, 26-11-13
 higher-level, network, 26-29-30
 multiaccess, 26-16
 network, 26-8
 synchronous DLC, 26-13-14
- Protractor, hyperbolic, 31-7
- Proximity-effect loss-tangent formula, 13-28
- PSG, definition of, 20-9
 ψ function, 47-40
- PSK *See* Phase-Shift Keying
- Psophometric
 noise, 2-4-5
 power, 2-5

- Psophometric—cont
 weighting factor, 2-5
 weighting, telephony circuits, 2-5
 PTM *See* Pulse-Time Modulation
 Public Service Telephone Network (PSTN), 46-5
 Publication, of patents, 49-33
 Pulse
 -amplitude modulation, 23-11-13
 definition of, 23-10
 signal-to-noise ratio in, 23-13
 circuits, 19-26-30
 -code modulation, 23-16-18
 definition of, 23-11
 false-pulse noise in, 23-17-18
 quantization noise in, 23-16-17
 compression, radar, 36-18, 36-19
 conditions, power output under, 17-17
 Doppler radar, 36-14-16
 duration modulation, 17-13
 definition of, 23-11
 frequency modulation
 definition of, 23-11
 natural sampling, 23-14
 uniform sampling, 23-14-15
 sampling, digital television system, 35-37
 -length modulation, definition of, 23-11
 modulation, 23-10-11
 sampling
 of bandpass function, 23-11
 in frequency domain, 23-11
 in time domain, 23-11
 -phase modulation
 definition of, 23-11
 uniform sampling, 23-15
 -position modulation (PPM), 22-3-4
 definition of, 23-11
 natural sampling, 23-15
 uniform sampling, 23-15
 rate, optical fiber, 22-25
 shapes, time and frequency functions, 7-6-8
 shaping, 19-26-29
 -time modulation, 23-13-14
 definition of, 23-10
 spectra of pulses, 23-13-14
 -train
 analysis, 7-9-10
 spectrum for, 7-9
 transformers, 13-17-18
 transmission, 37-3
 -width modulation, definition of, 23-11
 Pulsed lasers, 41-8
 high power, 41-12
 Pulses
 amplitude-modulated spectra of
 double polarity, 23-12
 single polarity, 23-13
 Purge, definition of, 20-9
 PWM *See* Pulse-Width Modulation
 pWp, definition of, 38-9
 Pyramid, equations for, 47-5
 Pyramidal frustrum, equations for, 47-6
 Pyramidal horn antennas, 32-23
 Pyroelectric detector, 21-21
 Pyrolytic graphite, tube grids, 17-18

 Q
 determination, by 3-decibel points, 8-7
 of inductor, 13-26
 meter, Hewlett-Packard, 12-8-10
 of resonant circuit, measurement of
 by bandwidth, 12-7
 of resonant line, 29-15-16
 switching, 41-5
 of unshielded coil, 6-5
 q, definition of, 20-9
 QPSK *See* Quadrature Shift Keying
 Quad, 3-16
 Quadrature method, AM stereo, 35-6
 Quadrature phase shift keying (QPSK), 24-6-8
 Quantitative reliability assessment, 45-8-9
 Quantities of dimension 1, 3-5-6
 Quantity, traffic, 39-11
 Quantization error
 digital filters
 coefficient error, 28-20
 input error, 28-20
 overflow oscillations, 28-22
 roundoff error, 28-20-22
 Quantization lobes, antenna array, 32-35-36
 Quantization noise, in PCM system, 23-16-17
 Quantizing, digital TV signal, 35-38
 Quantizing noise, definition of, 38-9
 Quantum
 detectors, 21-21-31
 efficiency, optical detectors, 21-15-16
 electronic device, elementary relationships, 41-2
 system, amplification in, 41-3-4
 Quart, 3-16
 Quarter-wave matching sections, 29-19
 Quartz oscillators, 1-144
 Quasi-TEM characteristics, microstrip lines, 29-25-26
 Quaternary direct sequence modulation, 24-24
 Quaternary phase shift keying (QPSK), 27-23-24
 Quaternary RF signals, coherent receivers for, 24-16

 R
 Rad, 3-16
 Radar
 antennas, 36-11-13
 errors, 36-13
 Radar—cont
 gain, 36-12
 Taylor illumination, 36-11
 automatic detection and track (ADT), 36-22
 beacon system, air traffic control, 37-5
 clutter, 36-18, 36-20-21
 fluctuations, 36-17
 coscant-squared antenna loss, 36-12-13
 coverage, 36-13-14
 Doppler, signal processing, 36-14-16
 fluctuating target models, 36-5-7
 fundamentals of, 36-3-4
 letter bands, 36-11
 measurement accuracy, 36-21
 MTI, 36-14-18
 noise jamming equation, surveillance, 36-9
 over-the-horizon equation, HF, 36-10
 over-the-horizon, HF, 36-16
 phased array tracking, 36-22
 pulse compression, 36-18, 36-19
 pulse Doppler, 36-14-16
 range measurement error, rms, 36-21
 range, prediction of, 36-4-5
 rms angle measurement accuracy, 36-21
 rms Doppler frequency measurement error, 36-21
 self-screening range equation, 36-9
 single-target tracker (STT), 36-21
 staggered PRF MTI, 36-17
 surface clutter range equation, 36-8
 surveillance equation, 36-8
 symbol definitions, 36-10-11
 synthetic aperture, 36-4, 36-16
 synthetic aperture equation, 36-10
 target cross section, 36-7, 36-10
 examples, 36-7
 track while scan (TWS), 36-22
 tracking, 36-21-22
 tracking equation, 36-8
 volume clutter equation, 36-8-9
 weather, 36-16
 weather equation, 36-9
 Radial-wire ground systems, 32-15-16
 Radian, 3-3
 per second, 3-3
 per second squared, 3-3
 Radiation
 blackbody, 21-3, 21-5-8
 cooling, electron tube, 16-16
 monochromatic, photocathode response to, 21-22
 pattern, LED, 22-6
 resistance, 32-5
 units, compatible system of, 16-46
 Radiative recombination, semiconductors, 18-6
 Radio
 beacons, 37-3-4
 blackouts, 1-155-156
 capacity, cellular system, 46-19-21

- Radio**—cont
 frequencies, international distribution, 1-2
 frequency
 allocations, 1-3-141, 1-3-5
 bands, nomenclature, 1-2, 1-3
 cables, army-navy preferred list, 29-30-39
 exposure protection guides, 49-30
 spectrum engineering software, 1-3
 spectrum management, 1-2-142
 US, 1-2-3
 horizon distance, 33-13-15
 navigation aids, 37-3-8
 noise, man-made, 34-6
 propagation, effect of nuclear explosions of, 33-12
 refraction, 33-13-15
 towers
 loading, 49-20-21
 steel
 foundation, 49-21-22
 material, 49-20
 pipe, 49-20
 standards, 49-20-22
 unit stresses, 49-21
 structural steel, 49-20-22
 transmission, path clearance, 33-17
 waves
 interference between direct and reflected, 33-17-18
 penetration into ground, 33-6
 transmission loss, 33-4
- Radio Technical Commission for Aeronautics (RTCA)**, 37-2
- Radio Technical Commission for Maritime Services (RTCM)**, 37-2
- Radiometric**, terminology, 21-3-4
- Radiometry**, 16-45-46, 21-3-8
- Radium**, 4-6-7, 4-10-11
- Radius**, atomic, 4-3, 4-4, 4-6
- Radix**, number system, 42-4-5
- Radon**, 4-6-7, 4-10-11
- RADSL**, 38-8
- Rain**
 attenuation, 33-20-21
 satellite communications, 27-36-37
 rate, 33-20
 climate conditions, US and Canada, 33-20, 33-22
 distribution values vs percent of year expected, 33-24
 global climate regions, 33-20, 33-22
 one-minute, various cities, 33-24-25
 sky noise temperature, with, 27-37-39
- RAM**; See also **Random Access Memory**
 dynamic, 20-42
 static, 20-41-42
- Raman emission**, stimulated, 41-19
- Random**
 access memory (RAM), 20-41-42, 42-13-14
 definition of, 20-9
 access techniques, multiaccess link control, 26-18-20
 digits, table of, 48-6
 experiment, 44-2
 variable, 44-2
- Range**
 measurement error, rms, radar, 36-21
 self-screening, radar, 36-9
 surface clutter, radar, 36-8
- Range**, definition of, 20-9
- Rapid thermal processing**
 definition of, 20-9
 integrated device, 20-66-81
 thin films, 20-14
- Rate adaptive DSL**, 38-8
- Rate generators**, 15-4
- Ratio-type shift register**, definition of, 20-9
- Ratioed logic**, 20-35-36
- Ratioless inverter**, definition of, 20-9
- Ratioless logic**, 20-36
- Ratioless-type shift register**, definition of, 20-9
- Ratio(s)**
 -arm bridge, 12-10
 deviation, FSK, 24-9
 inverter, definition of, 20-9
 modulation, FSK, 24-9
 noise power, 12-31-32
 relative dumping, 15-7
- Rayleigh**, 3-16
- Rayleigh distribution**, 24-23
- RC networks**, response of, 19-26-27
- Reactance charts**, 6-5-8
- Reactive ion etching**, definition of, 20-9
- Reactors**
 ac-filter, 13-19
 magnetic core, 13-3, 13-18-19
 rectifier-filter, 13-19-20
- Read-only memories (ROM)**, 20-40, 42-13
- Read-only memory**
 definition of, 20-9
 electrically erasable programmable, 42-14
 crasable programmable, 42-14
 factory-programmed, definition of, 20-5
 programmable, 42-13-14
- Read/write memory (RWM)**, 42-13
- Readout**, nondestructive, definition of, 20-8
- Received signal**, mobile radio, prediction of, 46-9-12
- Receiver(s)**
 coherent, 24-16
 correlation, 24-15-16
 direct-detection, optical, 22-9-10
 noise factor, 34-9
 noise, theoretical, 33-27
 noncoherent, 24-16
 optical, 22-8-12
- Receiver(s)**—cont
 heterodyne-detection, 22-10-12
 optimum, for digital communications, 24-13-18
- Receiving objective loudness rating (ROLR)**, 38-4
- Reception**, space-diversity, 33-18
- Recording**
 digital audio, 35-27
 standards
 international
 sound, 35-29
 television, 35-34
 systems
 sound, production standards, 35-27-29
 television, production standards, 35-29-34
 video disk, 35-33-34
- Records**, patent, importance of, 49-33-34
- Rectangular**
 coordinates, 47-35-36
 transformation of, 47-17-18
 dielectric guides, 30-27-28
 patch, antenna, 32-18
 waveguides, 30-3-4
 window function, 28-11
- Rectifier(s)**, 14-3-4, 14-31, 14-32, 18-12-13
 bridge, 14-13
 full-wave, single phase, 14-13
 carrier storage, 14-4
 circuits, 14-12-13
 chart, 14-14-15
 current ratings, 14-3
 filter reactors, 13-19-20
 forward characteristics, 14-3-4
 full-wave, single-phase bridge, 14-13
 half-wave, single-phase, 14-12
 mercury-vapor, 16-42
 power transformer design for, 13-5-9
 Schottky, 14-4, 14-32
 silicon controlled, 14-8-11
 star, 14-13
 three-phase
 diametric half-wave, 14-13
 double-Wye, half-wave, 14-13
 full-wave bridge, 14-13
 Wye, half-wave, 14-13
 ultra-fast, 14-32
 voltage ratings, 14-3
 Wye, 14-13
- Reduced gate**, 43-11
- Reduced instruction-set computing (RISC)**, 42-10-11
- Reed-Solomon code**, 25-15
- Reference**
 circuits, telephony, hypothetical, 2-7-9
 diodes, 18-13, 20-45
 level point, 38-9
 noise, definition of, 38-9
- Reflectance**, optical waves, 21-9
- Reflection**
 charts, geometry of, 31-6

- Reflection—cont
 - coefficient
 - definition, 31-2
 - load, Γ_L , definition of, 12-25
 - measurement of, 31-2
 - source, Γ_p , definition of, 12-25
 - waveguide, 31-2
 - multipath fading, 33-19
- Reflector(s)
 - antenna, 32-54–58
 - multiple-beam, 32-57–58
 - parabolic, 32-54–57
 - scanning, 32-57–58
- Reflux klystrons, 16-27
 - classes of, 16-28
- Reflow, definition of, 20-9
- Refraction
 - optical waves, 21-10
 - radio, 33-13–15
- Refractive index
 - atmosphere, 33-3
 - typical mean values, 33-29, 33-31
- Refresh, definition of, 20-9
- Regional AM broadcast stations, 35-3
- Regions, antenna, 32-4
- Register, definition of, 20-9
- Register-direct addressing, computer, 42-17
- Regulations, EMI, 12-36
- Regulators
 - ac voltage, 14-46
 - voltage
 - gas tube, 16-43
 - IC, 20-45
- Relative
 - atomic mass, 4-4, 4-6
 - damping ratio, 15-7
 - hardness of elements, 4-8, 4-10
 - humidity, Celsius table of, 49-6
 - level, 38-9
 - level points, 2-3–4
- Reliability
 - assessment
 - elements of, 45-6
 - quantitative, 45-8–9
 - component, 45-9–11
 - data, sources of, 45-13–15
 - definitions and terminology, 45-2–5
 - organizations for, 45-5
 - specifications, 45-20
 - US government, 45-20
 - standards, 45-25
 - system, 45-11–13
 - use of computers in, 45-25
- Rem, 3-16
- Remote pickup broadcast stations, 35-23
- Reorder tone, proposed standard, 2-18
- Reox, definition of, 20-9
- Repair frequency, definition of, 20-9
- Repeaters, telephone, 38-11, 38-13
- Reservation-ALOHA, 26-21
- Reservation scheme
 - first-in first-out, 26-21
 - round-robin, 26-21
- Reset, definition of, 20-9
- Resilient channels, sound isolation, 40-19
- Resist
 - casting, lithography, 20-17–18
 - development, lithography, 20-18–19
 - exposure, lithography, 20-18
- Resist, definition of, 20-9
- Resistance(s)
 - change with temperature, 49-15–16
 - design, telephone, 38-5
 - flange, waveguide, 30-16–17
 - gyration, 10-13–14
 - networks, attenuators, 11-3–5
 - normalized, 9-3
 - radiation, 32-5
 - sensor, metal sheet, 20-74
 - sheet, definition of, 20-10
 - temperature characteristics, resistors, 5-9
 - transmission lines, at UHF, 29-15
 - value, critical, 5-10
 - values, minimum-loss pads, 11-6
- Resistive attenuators, typical designs, 11-2–5
- Resistivity
 - definition of, 20-9
 - metals and alloys, 4-15–16
 - rocks, 4-17
 - soils, 4-17
- Resistors
 - in active filters, 10-4–5
 - adjustable, 5-13–16
 - carbon composition, 5-15
 - cermet-element, 5-15
 - color code, 5-10
 - definitions, 5-9–10
 - fixed composition, 5-10–11
 - fixed film, 5-12–13
 - low-power insulated, 5-11
 - noise, 5-11
 - package styles, 5-10–11
 - power, 5-12
 - precision wirewound, 5-11–12
 - RF effects, 5-11
 - temperature rating, 5-11
 - tolerance, 5-10
 - wirewound, 5-15
 - fixed, 5-11
- Resolution, television, 35-35–36
- Resonance
 - bridge, 12-5
 - gain at, bandpass filter, 8-2–3
 - jump, 15-23
 - selectivity far from, bandpass filter, 8-3
 - selectivity near, bandpass filter, 8-3, 8-6
 - vibratory forces, 49-24
- Resonant
 - cavities, 30-20–24
 - filter design, 14-43
 - irises, 30-24
 - lines, 29-15–19
 - input admittance of, 29-15–16
 - input impedance of, 29-15–16
 - insertion loss, 29-17–18
- Resonant—cont
 - power dissipation, 29-19
 - Q of, 29-15–16
 - transistor circuits, 14-36
- Resonators, dielectric, 30-28–29
- Response
 - amplitude, measurement of, 12-26–27
 - finite impulse (FIR), 28-3
 - infinite impulse (IIR), 28-3
 - spectral, semiconductor detectors, 21-30
- Responsivity
 - blackbody, 21-17
 - optical detectors, 21-14–18
 - photocathode, typical peak, 21-23
- Retransmission algorithms, distributed-tree, 26-21–22
- Return difference, feedback, 15-11
- Return loss, 12-29–30
 - definition of, 38-10
 - measurement of, 12-30
- Reverberation time, 40-13
 - calculation of, 40-13–14
 - measurement of, 40-13
 - Schroeder-Kuttruff method, 40-13
 - optimum, 40-15
- Reverse bias
 - definition of, 20-3
 - semiconductor, 18-9–10
- Reverse breakdown, semiconductor, 18-10–11
- Reverse voltage operation, power transistor, 14-6
- RF
 - amplifier(s)
 - bandwidth, 17-17
 - circuits, 17-8–13
 - triode, 17-8
 - class-B, design of, 17-6–7
 - class-C, design of, 17-4–6
 - high efficiency, 17-10
 - high linearity, 17-10–12
 - power gain, 17-17
 - measurements, 12-24–26
 - plasma discharges
 - capacitively coupled, 20-92–93
 - inductively coupled, 20-93–95
- RFI, additive, 24-25
- Rhe, 3-16
- Rhenium, 4-6–7, 4-10–11
- Rheostats, 5-14
- Rhodium, 4-6–7, 4-10–11
- Rician fading, slow nonselective, 24-22
- Rician distribution, 24-23
- Ridged waveguides, 30-8–13
- Right triangles
 - equations for, 47-11
 - hyperbolic, 47-15
 - spherical, 47-12
- Ring network, 26-7–8
- Ring of circular cross section, equations for, 47-5
- Ringdown signal, telephone, 2-17
- Ringtone, audible, proposed standard, 2-18

- Rinse cascade, definition of, 20-9
 RL networks, response of, 19-26
 RLC circuit, charge and discharge, 6-17-19
 RLC networks
 equations for, 6-14-15
 response of, 19-28-29
 Robustness, feedback control system, 15-11-12
 Rocks, resistivity of, 4-17
 Rod, 3-16
 Roentgen, 3-16
 ROM *See* Read-Only Memory
 Rooms, enclosed, sound in, 40-12-15
 Root loci examples, 15-8, 15-9, 15-10
 Root-locus method, 15-7-10
 Root mean square, 44-3
 Root mean square deviation from the mean, 44-3
 Rotman-Gent beam-forming network, 32-43-44
 Round-robin reservation scheme, 26-21
 Routh-Hurwitz criterion, 15-5
 Routing
 control center, 26-24-25
 network services, 26-24-25
 Row decoder, definition of, 20-9
 RTCVD process, 20-74-76
 RTP
 based process development, overview, 20-76-78
 for integrated CMOS processing, 20-74-76
 results, 20-78-81
 Rubidium, 4-6-7, 4-10-11
 Rubidium oscillators, 1-144, 1-146
 Rule, Cramer's, 47-32
 Ruthenium, 4-6-7, 4-10-11
 Rutherford, 3-16
- S**
- Safety levels, exposure to EM fields, 49-28-30
 Sag calculations, transmission-line, 49-18-20
 SAG, definition of, 20-9
 Samarium, 4-6-7, 4-10-11
 SAMPLE, definition of, 20-9
 Sampled-data filter systems, 10-36-38
 Sampling
 digital television system, 35-36-38
 pulse frequency, digital television system, 35-37
 pulsed modulation
 of bandpass function, 23-11
 in frequency domain, 23-11
 in time domain, 23-11
 theorem, 25-17
 Satellite(s)
 communications
 choice of frequency, 27-13
 DAMA systems, 27-24
 digital transmission systems, 27-21-24
 measurement of, 31-4-6
 Scattering, optical signal, 22-16-17
 Schematic symbols, 49-26-27
 Schematics
 diversity, 46-19-21
 URN, 26-22
 Schering bridge, 12-5
 Schmitt trigger, definition of, 20-9
 Schottky barrier diode, definition of, 20-9
 Schottky rectifiers, 14-4, 14-32
 Schottky transistor logic (STL), 20-34-35
 Schroeder-Kuttruff method, reverberation time measurement, 40-13
 Scintillation counters, photomultipliers as, 21-25-26
 SCPI, 12-35
 SCR, 18-26
 Scruple, 3-16
 SDSL, 38-8
 Sea water, galvanic series in, 4-3, 4-13
 Second breakdown, power transistors, 14-5-6
 Second harmonic generation, laser, 41-12, 41-14, 41-16-19
 Second order transfer functions, 10-15-23
 Second (plane angle), 3-16
 Second (time), 1-143-144, 3-3, 3-16
 Secondary emission, 16-5-7
 Section, 3-16
 Secular variations, 27-7
 Segregation, definition of, 20-9
 Selective fading, mobile signal, 46-12, 46-17
 Selectivity
 equations, 8-4-5
 far from resonance, bandpass filter, 8-3
 near resonance, bandpass filter, 8-3, 8-6
 pairs of coupled-tuned circuits, 8-7-9
 single-tuned circuits, 8-6-7
 Selenium, 4-6-7, 4-10-11
 Self-impedance, dipoles, 32-13
 Self-refresh, definition of, 20-9
 Self-screening range equation, radar, 36-9
 SEM, definition of, 20-9
 Semiconducting materials, 4-3, 4-18
 Semiconductor(s)
 charge carriers, 18-3-5
 controlled rectifier, 18-26
 definition of, 20-9-10
 depletion region, 18-9
 diode, definition of, 20-5
 diode, type number coding, 5-7
 diodes and rectifiers, 18-12-13
 display devices, 18-23-24
 energy bands, 18-3-5
 field emission, 16-9-10
 forward bias, 18-9-10
 impacts, 20-58-59
 Satellite(s)—cont
 downlink, 27-3
 effect of solar eclipse on, 27-9-11
 evolution, 27-4-6
 ionospheric scintillations, 27-33-35
 propagation, 27-32-40
 rain attenuation, 27-36-37
 transmission system considerations, 27-19-21
 transponders for, 27-15-19
 tropospheric effects, 27-35
 uplink, 27-3
 defining position, 27-7
 direct tv broadcast from, 27-4-5
 geostationary, elevation and azimuth angles, 27-11-12
 K_u band, 35-23
 mobile communication, 27-5
 MOLNIYA, 27-11
 navigation system, 37-5-8
 network program distribution, 35-20, 35-22-23
 orbits, 27-6-7
 choice of, 27-8-9
 polar, 27-11
 SPOT-4, 27-19
 -switched time-division multiple access. *see* SS-TDMA
 systems
 analog transmission systems, 27-21-24
 communications, classification of, 27-3
 transponders, multiple-access schemes, 27-19-21
 vs fiber-optic communications systems, 27-5
 SATO, definition of, 20-9
 Saturated gain coefficient, 41-4
 Saturation
 percent of, 49-6
 -type nonlinearity, 15-25
 voltage, power transistor, 14-5
 SC
 biquads, parasitic-insensitive, 10-38, 10-39, 10-40-44
 filter, 10-34-45
 operation of ideal, 10-36-38
 integrators, 10-38-39
 z-domain biquadratic transfer functions, 10-39-40
 Scalar multiplication, distributive law for, 47-35
 Scalar triple product, 47-35
 Scales, time, 1-142-144
 Scaling, ICs, trends in, 20-45-46
 Scan impedance, 32-44
 Scandium, 4-6-7, 4-10-11
 Scanning reflector, 32-57-58
 Scatter propagation
 ionospheric, 33-12
 tropospheric, 33-29, 33-32
 Scatter sensor, fiber-optic, 20-73-74
 Scattering matrix
 of a junction, 31-3

- Semiconductor(s)—cont
 injection laser, 21-12-14
 lasers, 21-12-14, 41-9
 manufacturing, single wafer, 20-56-58
 materials, 18-2-3
 properties of, 18-7
 motivations, 20-58-59
 n-type, definition of, 20-8
 optical excitation, 18-6
 p-type, definition of, 20-9
 power, characteristics, 14-3-12
 power losses, 14-31-32
 processing, sensor hierarchy, 20-59-61
 properties of, 4-18
 radiative recombination, 18-6
 reverse bias, 18-9-10
 reverse breakdown, 18-10-11
 ROM, mechanically programmable, definition of, 20-7
 switching devices, 18-26
 transition region, 18-9
- Sense amplifier, definition of, 20-10
- Sensitivity
 detector, ultimate, 21-20-21
 feedback control system, 15-11-12
 network, 10-8-10
- Sensors
 diode, 12-24-25
 thermistor, 12-24
 thermocouple, 12-24
- Separate audio program (SAP), 35-14
- Sequential circuits, 43-16-18
- Sequential switching circuits, 43-3
- Sequentially accessed memory, definition of, 20-10
- Serial accumulator, definition of, 20-10
- Serial operation, definition of, 20-10
- Series
 Fourier, 7-2, 7-8-9
 complex form, 7-2, 7-8-9
 real form, 7-2
 miscellaneous, 47-30
 -resistance-capacitance bridge, 12-4
 Taylor
 single variable, 47-29
 two variables, 47-30
- Servics, network, 26-23-24
- Servo
 positioning, 15-4
 systems, phase-locked loop, 15-21-22
- Servomotors, two-phase, 15-3
- Session
 layer, network, 26-29-30
 routing, networks, 26-25
- Set-reset latch electronic circuit, 43-16
- Severely errored second, definition of, 38-11
- Shake, 3-16
- Shannon's formula, 25-18-19
- Sheet resistance, definition of, 20-10
- Ships, vibration in, 49-25
- Shock and vibration, isolation of, 49-22-26
- Shock, isolators, 49-26
- Shot effect, electron tube, 16-18
- Shot noise, 21-19
- SI units, 3-2-4
 prefixes, 3-3
 symbols for, 3-2-3
 table of, 3-3
 tables of, 3-2-3
 uncertainties, 1-143
- Side-lobe
 distributions, low, antenna, 32-51-52
 patterns, low, antenna, measurement of, 32-52-53
- Siemens, 3-3
- Sievert, 3-3
- Signal-to-noise
 improvement, FM systems, 23-9-10
 through de-emphasis, 23-10
 improvement ratio in PTM, 23-15-16
 ratio
 delta modulation, 23-18-19
 in pulse-amplitude modulation, 23-13
 transmitter power for required, 33-27-29
- Signaling
 A- and B-Bit, 38-17
 with bandwidth constraints, 25-21-22
 binary orthogonal, 22-5
 carrier-derived channels, 38-17
 coded trellis, 27-24
 DC, 38-16
 dual-tone multifrequency, 38-16, 38-17
 E and M lead, 38-16-17
 interoffice, common carrier, 38-16-18
 loop reverse-battery, 38-16
 multifrequency, 38-17-18
 single-frequency, 38-17
 subscriber-loop, 38-16
 telephone, 2-13-19
 network control, 38-16-18
 without bandwidth constraints, 25-21
- Signal(s)
 analysis, 12-17-20
 characteristics, 12-19
 coding, digital television, 35-38, 35-40
 composite television, 35-18-19
 -flow graph (SFG) simulation, 10-32-34
 FM/FM, spectral distribution of, 23-8
 ideal sampled, 28-6
 intensity
 collimated beam, 22-7
 LED, 22-7
 Internet time, 1-158, 1-160
 mobile, multipath fading, 46-12-17
 power
 density, optical, 22-7
 received, optical, 22-10
- Signal(s)—cont
 program, network distribution, 35-20, 35-22-23
 set
 antipodal, 24-3
 orthogonal, 24-3
 for RF channels, 24-4-13
 sets, baseband, 24-2-4
 sky wave, 35-5, 35-8
 stereophonic, FM broadcast, 35-10-11
 telephone time, 1-158, 1-160
 telephony, conventional, 2-5
 test, on TV signal, 35-13
 VIR, TV, 35-13
- Signed integer, 42-4
- Signed-magnitude system, numbers, 42-7
- SILEX optical system, 27-19
- Silicon, 4-6-7, 4-10-11
- Silicon compilers, 20-56
- Silicon controlled rectifiers, 14-8-11
- Silicon-on-insulator: See also SOI
 technology and devices, 20-102-108
- Silver, 4-6-7, 4-10-11
- SIM, definition of, 20-10
- SIMOX, 20-103
- Simulation, IC design, 20-26
- Simulator, computer, 42-3
- Sine integrals, 47-39-40
- Sines
 hyperbolic, table of, 48-2
 law of, 47-11, 47-12, 47-15
- Singing margin, definition of, 38-10
- Singing return loss, definition of, 38-10
- Single
 channel per carrier (SCPC)
 transmissions, 27-24
 -degree-of-freedom structure, unity feedback, 15-11
 -frequency signaling, 38-17
 -loop system, unity feedback, 15-5
 -mode step-index fiber, 22-19
 -mode step-index fibers
 bandwidth, 22-24
 -rail logic, 43-10
 sideband, amplitude modulation, 23-4-5
 -stage coordinate switches, 39-2-3
 -target tracker (STT), 36-21
 -variable theorems, switching algebra, 43-8
 wafer epitaxy, 20-91
 wafer plasma processing, 20-91-99
 wafer surface cleaning, 20-81-91
- Single crystal, definition of, 20-10
- Singular points, types of, 15-32
- Sink tree, 26-24
- Sinter, definition of, 20-10
- Sinusoidal oscillators, 19-24-26
- SIP, 20-47-48
- Site diversity, 27-39-40
- Six-port network, 12-15
- Skein, 3-16
- Skin effect, 6-6-8, 6-14
 factor, 13-27

- Sky**
 noise temperature, with rain, 27-37-39
 wave, 33-5-6
 signals, 35-5, 35-8
SLD *See* Superluminescent Diodes
Slice, definition of, 20-10
Sliding block code, 25-4
Slope-intercept form, equation, 47-16
Slot(s)
 antennas, 32-19-22
 lines, 30-29-31
 longitudinal, 32-20-21
 waveguide, 32-20-21
Slug, 3-16
Small antennas, 32-4-6
Smirnov-Kolmogorov test, 45-23
SMPTTE, 35-29
Snell's law, 22-19-20
SNR, optical receiver, 22-11
Snubbers, 14-32, 14-34
SOC, 20-47-48
Society of Motion Picture and Television Engineers, 35-30
Sodium, 4-6-7, 4-10-11
Soft magnetic materials, 13-4
Software, computer, 42-3-4
SOI
 circuits, 20-107-108
 definition of, 20-10
 device structures using, 20-104-106
 devices on, 20-106
 fully depleted devices, 20-107
 material options, 20-102-104
 nonfully depleted devices, 20-106-107
Soils, resistivity of, 4-17
Solar
 cell, 18-13
 eclipses, effect on satellite communications, 27-9-11
 noise, 27-12
 spectral irradiance, 22-13
Solenoids
 effect of shielding on, 6-3-4
 inductance of, single-layer, 6-3-5
Solid
 analytic geometry, 47-19-21
 figures
 surface areas of, equations for, 47-5-6
 volumes of, equations for, 47-5-6
 sound in, 40-5
 state devices, comparison with electron power tubes, 17-13-17
 state diffusion, definition of, 20-10
 -state power amplifiers (SSPA), 27-17
SOS, definition of, 20-10
Sound
 in enclosed rooms, 40-12-15
 exposure
 day-night, 40-11
 level, definition of, 40-10
 in gases, 40-3-4
 breakdown voltages, 49-3
 triggered, 16-43
Special Mobile Group (GSM), cellular system, 46-22, 46-24
Specific heat, of elements, 4-9, 4-11
Specifications, reliability, 45-20
Spectra
 matching factors, 16-49
 of PSK and MSK signals, 24-11-13
Spectral analysis
 FFT in, 28-11-13
 Fourier transform, 7-10
 bit rate, 25-20
 distributions, 16-47
 FM/FM signal, 23-8
 estimation, 25-24-25
 irradiance, solar, 22-13
 response
 optical detectors, 21-16
 semiconductor detectors, 21-30
 vidicon, 16-54
Spectrum
 analyzer, 12-17-19
 detector bandwidth, 12-20
 IF filter bandwidth, 12-19-20
 laser operation, 41-2
 management
 engineering software, 1-3
 radio frequency, 1-2-142
 US, 1-2-3
 optical, 21-2
 spreading, 25-29-30
 Speech interference levels, 40-15
 Speech synthesizer ICs, 20-46-47
Speed, of light, 18-4
Sphere
 equations for, 47-5
 sector of, equations for, 47-5
 segment of, equations for, 47-5
Spherical
 coordinates, vector equations for, 47-36-37
 trigonometry, 47-12-13
 waves, 40-3
Spheroid
 oblate, equations for, 47-20
 prolate, equations for, 47-20
Spheroidal core, loop antenna, 32-19-20
SPICE 2, definition of, 20-10
Spin-on deposition, thin films, 20-14
Spiral antenna, 32-28-29
SPL, 40-5-6
Splicing, optical fiber, 22-24
Split-channel reservation multiple access (SMRA), 26-21
Spontaneous emission, 41-2
Sporadic E transmission, 33-7
SPOT-4, 27-19
Spread-spectrum
 communications, 24-24-26
 binary direct-sequence, 24-24-26
 frequency-hop, 24-26
 multiple access (SSMA), 26-20
 system, 25-29-30
Sound—cont
 intensity, 40-5
 isolation, 40-17, 40-19-23
 doors, 40-19
 flanking, 40-21, 40-23
 resilient channels, 40-19
 walls, 40-19, 40-20-21
 weighting of components, 40-19, 40-21
 leaks, 40-22-23
 level, 40-8
 day-night, 40-11-12
 definition of, 40-9-10
 meters, 40-8
 in liquids, 40-4-5
 mmeasurement of, 40-8-9
 power level, 40-8-9
 definition of, 40-9
 of typical noise sources, 40-9
 pressure
 instantaneous, definition of, 40-9
 level, 40-5-6
 definition of, 40-9, 40-10
 recording standards, international, 35-29
 recording systems, production standards, 35-27-29
 reinforcement systems, 40-23-24
 in solids, 40-5
 source
 direct field of, 40-12
 intensity, 40-12-13
 waves
 standing, 40-13
 theory of, 40-3-5
Source(s)
 compaction codes, 25-7-11
 compression codes, 25-25-26
 encoding, 25-10
 -follower amplifier, 19-11
 Markov, 25-8
 models, 25-7-8
 optical, 21-10-14
 reflection coefficient, Γ_g , definition of, 12-25
Space
 curves, vector equations for, 47-37
 diversity, 33-11, 37-3
 reception, 33-18
 division switching network, 39-2
 FSK signal, 24-9
 systems, classification of, 27-3
 transmissions in, equations for, 27-7-8
 weather, 1-155
Spacecraft
 antennas, 27-33
 zone coverage, 27-31-32
 architecture, 27-14-15
 -to earth communications, 27-3
 -to spacecraft communications, 27-3
Spaced-tapered arrays, 32-40
SPADE communications system, 27-24-27
Span, 3-16
Spark gaps

- Sputter deposition, thin films, 20-14
- Square
 - meter, 3-3
 - topped sampling, pulse amplitude modulation, 23-12
 - waveguides, 30-7
- SR latch circuit, 43-16
- sRAM cell, 20-41-42
- SS-TDMA
 - basic concept, 27-27-29
 - block diagram, 27-30
 - on-board signal processing, 27-28
 - operation of, 27-25
- SSB; See also Single Sideband detection, 23-6
- Stability, optical cavities, 41-7
- Stack
 - algorithm, 25-16
 - machine, computer, 42-12
 - registers, computer processor, 42-15
- Staggered PRF MTI radar, 36-17
- Standard
 - broadcast stations, classification of, 35-3-4
 - commands for programmable instruments (SCPI), 12-35
 - deviation from the mean, 44-3
 - mismatch method, return loss measurement, 12-30
 - time intervals, WWV and WWVH, 1-151
 - waveguides, 30-7, 30-10-11
- Standards
 - acoustical, 40-25-28
 - frequency, 1-144-148
 - international
 - broadcasting, 35-46-49
 - sound recording, 35-29
 - telephony, 2-3
 - television recording, 35-34
 - organizations for, 5-4
 - program production, 35-26-34
 - sound, 35-27-29
 - reliability, 45-25
 - steel radio towers, 49-20-22
 - television, 35-34
 - time, 1-144-148
 - transmission, TV broadcast stations, 35-12-13
 - video tape, 35-30-32
- Standing sound waves, 40-13
- Standing-wave ratio, transmission lines, 29-8-10
- Star rectifier, 14-13
- Start-of-active-video reference, 35-41-42
- Statampere, 3-16
- Statcoulomb, 3-17
- State
 - diagrams, 25-5-7
 - transition matrix, 15-13-14
 - variable analysis, digital control system, 15-17-18
 - variable design methods, digital case, 15-20-21
- State, definition of, 20-10
- Statfarad, 3-17
- Stathenry, 3-17
- Static
 - logic, 20-35-36
 - precipitation, 34-8
 - RAM, 20-41-42
 - vs dynamic routing, network, 26-24
- Static storage elements, definition of, 20-10
- Stations
 - aural intercity, 35-23
 - broadcast, classification of, 35-3-4
 - earth, 27-40-43
 - FM broadcasting, classification of, 35-8-9
- Statistical inference, probability and, 45-15-25
- Statistics, Bayesian, 45-23-25
- Statohm, 3-17
- Status register, computer processor, 42-15
- Statute mile, 3-17
- Statvolt, 3-17
- Steel
 - towers
 - foundation, 49-21-22
 - loading, 49-20-21
 - material, 49-20
 - structural standards, 49-20-22
 - unit stresses, 49-21
 - wind velocities and pressure on, 49-22
- Stefan-Boltzmann law, 21-3, 21-6
- Step coverage, definition of, 20-10
- Step-down converter, 14-19, 14-28
- Step motors, 15-3-4
- Step-up converter, 14-28-29
- Steradian, 3-3
- Stere, 3-17
- Stereo, FMX, 35-11
- Stereo transmission
 - AM broadcast, 35-5-6
 - FM broadcast, 35-10-11
- Stilb, 3-17
- Stimulated emission, 41-2
 - cross section, 41-4
 - Raman, 41-19
- STL stations, 35-23
- STM allocation, network, 26-37
- Stoke, 3-17
- Stokes wave, 41-19
- Storage
 - bistable, cathode-ray tubes, 16-57-58
 - cathode-ray tubes, 16-55-57
 - devices, mass, computer, 42-14
 - elements, static, definition of, 20-10
 - fixed-program, definition of, 20-5
- Store-and-forward
 - deadlock, networks, 26-26
 - transmission, 26-3, 26-4, 26-6
- Stored
 - charge programmable semiconductor ROM, 20-10
 - energy
 - capacitor, 6-15
- Stored—cont
 - inductor, 6-15
- Straggle, definition of, 20-10
- Straight line, equations for, 47-16, 47-20
- Strain coefficients, piezoelectric, 4-29
- Stresses, unit, steel radio towers, 49-21
- Strip transmission lines, 29-27-28
- Stripline slot antennas, 32-17
- Strontium, 4-6-7, 4-10-11
- Structural standards, steel radio towers, 49-20-22
- Structures, man-made, effect on mobile communications, 46-8
- Stub
 - open, impedance matching with, 29-19
 - shorted, impedance matching with, 29-19
- Student-*t* distribution, 45-11
- Subroutines, computer program, 42-26-30
- Subscriber
 - loop signaling, 38-16
 - loops
 - telephone, 38-5-6
 - resistance design, 38-5
 - sets, telephone, 38-3-4
 - units, mobile, 46-7
- Subsidiary Communications Authorization (SCA), 35-11
- Substitution method, high impedances, bridge, 12-5-6
- Substrate
 - cleaning procedure, 20-12
 - doping procedures, 20-21-22
 - diffusion from doped oxides, 20-21
 - by ion implantation, 20-21-22
 - predeposition and drive-in, 20-21
- Substrate, definition of, 20-10
- Subtraction
 - nondecimal numbers, 42-6
 - using twos complement method, 42-8
- Sulfur, 4-6-7, 4-10-11
- Sum, canonical, 43-5-6
- Sum modulo two gate, 43-11
- Summers, 10-10-11
- Sums
 - fundamental, 43-5
 - minimal, formation of, 43-14-15
 - of powers of integers, 47-8
 - of powers of reciprocals of integers, 47-8-9
- Supercomputer networks, 26-34
- Superconducting antennas, 32-6
- Superconductivity, metals, alloys, and compounds, 4-17
- Superluminescent diodes, 21-12
- Supervisor stack register, computer processor, 42-15
- Supervisory audio tone (SAT), 46-4
- Supervisory signal, telephone, 2-14-17
- Suppressors, transient-voltage, 14-12
- SUPREM, definition of, 20-10

- Surface
- acoustic wave (SAW), 28-29
 - filters, 28-33-37
 - response, 28-34
 - cleaning
 - chemical oxide, by vapor-phase HF, 20-86-87
 - IC manufacturing, 20-82-85
 - low-temperature, germane-based, 20-85-86
 - native oxide, by vapor-phase HF, 20-86-87
 - single wafer, 20-81-91
 - clutter range equation, radar, 36-8
 - wave transmission line, 30-24-26
- Surveillance radar
- equation, 36-8
 - noise jamming, 36-9
- Susceptibility, electromagnetic, 12-36
- Sweep oscillator, 12-14
- Swinging chokes, 14-41
- Switch, coordinate, network, 39-2
- Switch-mode converters, 14-34-35
- Switched
- capacitor filters (SCF), 28-29
 - circuits
 - active, 10-34-45
 - basic principles, 28-30-31
 - voice quality, 28-31-32
 - telecommunication system, 38-3
- Switches, coordinate, single-stage, 39-2-3
- Switching
- \s1 algebra, 43-2
 - theorems, 43-8-10
 - n -variable, 43-9-10
 - single-variable, 43-8
 - circuits, 26-2
 - combinational, 43-3
 - sequential, 43-3
 - devices, semiconductor, 18-26
 - equipment, cellular systems, 46-21-22, 46-23
 - fast packet, networks, 39-10
 - message, 26-2
 - networks, 39-2
 - coordinate, 39-2-10
 - packet, 26-3
 - networks, 39-9
 - power supplies, 14-45
 - bipolar power transistors, 14-29-30, 14-32
 - insulated-gate, 14-30, 14-32
 - MOSFETs, 14-29, 14-30, 14-32
 - power semiconductors, 14-29, 14-31-32, 14-33
 - transistor and diode requirements, 14-20-27
 - system, grade of service, 39-11-12
 - techniques, 26-2-5
 - time-division, 39-6-10
 - times, power transistor, 14-6
 - transistor circuits, 14-19
- SWP-based process and device integration, 20-99-102
- Symbol error, probability of, 24-23
- Symbols
- graphic
 - British, 49-28
 - US, 49-26-27
 - npn transistor, 19-2
 - schematic, 49-26-27
 - SI units, 3-2-3
 - transmission lines, 29-3-4
 - vibration and shock, 49-22
- Symmetric DSL, 38-8
- Synchronous
- data link control, 26-16
 - detection, 23-6
 - DLC protocols, 26-13-14
 - motor drive, 14-49-50
 - transfer mode, network, 26-37
- Synchros, in error measuring systems, 15-2, 15-3
- Synthesis, logic design, 43-4-8
- Synthetic aperture radar, 36-16
- Synthetic aperture radar (SAR), 36-4
- equation, 36-10
- System
- error-measuring, 15-2
 - performance
 - control system, 15-4-12
 - microwave systems, 12-33
 - reliability, 45-11-13
- System in a package (SIP), 20-47-48
- System in a package (SIP), definition of, 20-10
- System Network Architecture (SNA), 26-9-10
- System on a chip (SOC), 20-47-48
- System on a chip (SOC), definition of, 20-10
- T
- T- π transformation, network, 6-15-16
- Table of combinations, logic circuit, 43-4
- Tablespoon, 3-17
- TACAN, 37-5
- Tactical air navigation, 37-5
- Talker echo, definition of, 38-10
- Tangents
- hyperbolic, table of, 48-4
 - law of, 47-11
- Tanks, military, vibration in, 49-25
- Tantalum, 4-6-7, 4-10-11
- electrolytic capacitors, 5-21
- Tape formats, video, 35-29-30
- Target cross section
- radar, 36-7, 36-10
 - examples, 36-7
- Taylor illumination, radar antennas, 36-11
- Taylor \bar{n} circular source distribution, antennas, 32-53-54
- Taylor \bar{n} line source distribution, antenna, 32-49-50
- Taylor one-parameter line source distribution, antenna, 32-47-48
- Taylor series
- for a single variable, 47-29
 - for two variables, 47-30
- TCAS, 37-5-6
- Tchebyscheff *See* Chebyshev
- TDMA, 27-20-21, 27-24-27
- acquisition, 27-26
 - frame structure, 27-28
 - synchronization, 27-26
 - system operation, 27-24-25
 - systems characteristics, 27-25
- TDRSS, 27-11
- TE waves, 30-3
- Teaspoon, 3-17
- Technetium, 4-6-7, 4-10-11
- Tee network, power dissipation within, 11-4, 11-7
- Tee pad, power dissipation within, 11-7
- Telecommunication system, switched, 38-3
- Telephone, 38-17-18
- busy tone, proposed standard, 2-18
 - cables, characteristics of, 38-12
 - call-waiting tone, proposed standard, 2-18
 - carrier systems, 38-13-16
 - CCITT recommendations, 38-8
 - circuit characteristics, 2-9-10
 - dial pulse signal, 2-14
 - dial tone, proposed standard, 2-18
 - facilities, 38-11-16
 - fiber-optic facilities, 38-13
 - high tone, proposed standard, 2-18
 - intertoll trunks, 38-6-7
 - line signal, 2-14-17
 - local interoffice trunks, 38-6
 - network control signaling, 38-16-18
 - preemption tone, proposed standard, 2-18
 - reorder tone, proposed standard, 2-18
 - repeaters, 38-11
 - ringdown signal, 2-17
 - ringing tone, audible, proposed standard, 2-18
 - signaling, 2-13-19
 - systems, CCITT, 2-14
 - tones, 2-17-18
 - subscriber loops, 38-5-6
 - subscriber sets, 38-3-4
 - supervisory signal, 2-14-17
 - time signals, 1-158, 1-160
 - toll transmission plant, 38-6-7
 - wire, 38-11
 - wireless access, 38-6
- Telephony
- circuit loading, 2-5-6
 - busy hour, 2-5-6
 - conventional, 2-5-6
 - circuit noise, 2-11-13
 - circuits, psophometric weighting, 2-5
 - group delay distortion, 2-9
 - hypothetical reference circuits, 2-7-9
 - international standards, 2-3

- Telephony—cont
 linear cross talk, 2-10
 power levels, 2-6-7
 propagation time, one-way, 2-9
 signal, conventional, 2-5
- Telescope
 imaging, 22-6
 receiving, optical communications, 22-8-9
 transmitting, optical communications, 22-5-6
- Television
 auxiliary stations, 35-23
 broadcast stations
 directional antennas, 35-12
 transmission standards, 35-12-13
 broadcasting, 35-11-14
 aural transmitter, 35-13-14
 coverage data, 35-11
 grades of service, 35-11
 maximum power, 35-11
 minimum power, 35-11
 MTS, 35-14
 SAP, 35-14
 translators, 35-17-18
 utility audio-data channel, 35-14
 zones, 35-11
- cable, 35-14-17
- channels
 designations, 35-11-12
 instructional TV, 35-20, 35-23
- 865digital
 ancillary data format, 35-43
 blanking interval considerations, 35-41-43
 component coding of signals, 35-41-42
 composite encoded signals (NTSC & PAL), 35-40-41
 encoding parameter for component television, 35-42
 quantization levels for component television, 35-42
 signal coding, 35-38, 35-40
 synchronization considerations, 35-41-43
 timing references, 35-41-42
- enhanced definition (EDTV), 35-43
- film systems, international standards, 35-34
- high definition, 35-18, 35-23, 35-43-46
 wide-screen, encoding parameters, 35-44
- recording systems, 35-29-34
 production standards, 35-29-34
- signal, composite, 35-18-19
- standards, 35-34
- stations, low power, 35-17-18
- systems
 advanced, 35-43-46
 digital, 35-34-46
 NTSC-M, 35-40
 resolution, 35-35-36
- wide screen, encoding parameters, 35-45
- Telex, destination codes, international, 2-13
- Tellurium, 4-6-7, 4-10-11
- TEM, definition of, 20-10
- TEM waves, 30-3
- Temperature(s)
 antenna-noise, 34-5
 blackbody, 22-12
 coefficient
 ceramic capacitors, 5-23
 resistors, 5-9
 effect on cavity tuning, 30-21
 -EMF characteristics, thermocouples, 4-3, 4-12
 extremes, 49-3-4
 US, 49-3, 49-5
 world, 49-4
 noise, 21-20, 27-15, 32-8-11
 clear sky, 27-36
 operating, transformer insulating materials, 13-10-11
 resistance change with, 49-15-16
 sensors, multipoint, 20-68, 20-70-73
 stability, inductor, 13-28
 world, 49-4
- Temporary register, computer processor, 42-15
- Tensile strength, of elements, 4-9, 4-11
- Terbium, 4-6-7, 4-10-11
- Terminal plane, changing of, 31-3
- Terrain, effect on mobile communications, 46-8
- Terrestrial network program distribution, 35-20
- Tesla, 3-3
- Testing
 definitions and terminology, 45-2-5
 life, 45-13-14
- Test(s)
 Chi-square, 44-4
 constant humidity, 5-6
 high-altitude, 5-6, 5-7
 methods, environmental, 5-5
 sets, 12-14
 signals, on TV signal, 35-13
 Smirnov-Kolmogorov, 45-23
 vibration, 5-6, 5-7
- Tetrode amplifier circuits, 17-9, 17-11
 RF, 17-8-9
- Tex, 3-17
- Text editor, computer, 42-3
- Thallium, 4-6-7, 4-10-11
- Theorem(s)
 DeMorgans, 43-10
 Heaviside expansion, 6-21
 L'Huilier's, 47-13
 Neyman-Pearson, 25-22-23
 sampling, 25-17
 switching algebra, 43-8-10
n-variable, 43-9-10
 single-variable, 43-8
 two- & three-variable, 43-8-9
- Theory
 decision, 25-22-25
 estimation, 25-22-25
- Therm, 3-17
- Thermal
 compression bond, definition of, 20-10
 converters, 12-24
 detectors, 21-21
 electromotive force, of alloys, 4-15
 emissivity, electron-tube materials, 16-16
 generation, definition of, 20-10
 growth, thin films, 20-13
 noise, 21-18, 34-3
 calculations, 34-8
 definition of, 38-9
 processing, rapid, definition of, 20-9
 processing, rapid, thin films, 20-14
 resistance, IC package, 20-30
 unit, British, 3-14
- Thermal conductivity, of elements, 4-9, 4-11
- Thermionic emission, 16-3-5
- Thermistor sensors, 12-24
- Thermocouple(s), 21-21
 properties of, 4-3, 4-14
 sensors, 12-24
 temperature-EMF characteristics, 4-3, 4-12
- Thermopile, 21-21
- Thermopneumatic cell, 21-21
- Thin films, 20-13-14
 chemical vapor deposition, 20-13
 evaporation, 20-13-14, 20-15-17
 plasma deposition, 20-13
 rapid thermal processing (RTP), 20-14
 spin-on deposition, 20-14
 sputter deposition, 20-14
 thermal growth, 20-13
- Thinned arrays, antenna, 32-40-41
- Thorium, 4-6-7, 4-10-11
- 3G, 46-5, 46-6, 46-24
- Three-state output, definition of, 20-10
- Threshold, laser, 41-4-6
- Throughput, definition of, 20-10
- Thulium, 4-6-7, 4-10-11
- Thyatron, 16-42
 hydrogen, 16-42-43
- Thyristor(s), 14-8-11
 characteristics, 14-10-11
 circuits, phase controlled, 14-16-18
 circuits, phase-controlled, 14-13
 gate turn-off, 14-11
 MOS-controlled, 14-11
 phase-controlled, power supply filter, 14-41, 14-43
 ratings, 14-9-10
- Time
 atomic, 1-143
 -average sound level, definition of, 40-10
 broadcasts, 1-148-162
 chart, world, 49-5, 49-7
 code
 BCD, WWV and WWVH, 1-152-154
 WWVB, 1-157-158
 constant, network, 6-15

- Time—cont
 Coordinated Universal, 1-142–144, 1-151
 delay spread, mobile signal, 46-12, 46-17
 division multiple access. *see* TDMA
 division switching, 39-6–10
 network, 39-2
 time-slot interchange, 39-9
 domain, and frequency relationships, 12-15–16
 ephemeris, 1-143
 and frequency measurement, 12-20–24
 functions, pulse shapes, 7-6–8
 invariant digital systems, linear, stability of, 15-18–20
 invariant systems, stability of linear, 15-5–10
 measurement, 12-21
 modulated pulses, spectra of, 23-13–14
 reverberation, 40-13
 calculation of, 40-13–14
 measurement of, 40-13
 optimum, 40-15
 scales, 1-142–144
 universal, 1-143
 signal
 continuous, 28-2
 discrete, 28-2
 Internet, 1-158, 1-160
 telephone, 1-158, 1-160
 slot interchange, time-division switching, 39-9
 standards, 1-144–148
 zones, world, 1-144, 1-145
 Timing references, digital television, 35-41–42
 Tin, 4-6–7, 4-10–11
 TIRAM, definition of, 20-10
 Titanium, 4-6–7, 4-10–11
 -doped sapphire laser, 41-9
 TM waves, 30-3
 Toggle, definition of, 20-10
 Tolerance
 antennas, 32-41
 components, 5-4
 resistors, 5-10
 Toll transmission plant, telephone, 38-6–7
 Ton, 3-17
 Tone interference, definition of, 38-9
 Tones, telephone signaling, 2-17–18
 Tonne, 3-17
 Top sampling, pulse amplitude modulation, 23-11
 Topological realization of ICs, 20-26–28
 Toroidal cores
 audio-frequency inductors, 13-21–23
 design example, 13-23
 winding methods, 13-27
 Torque, 49-18
 bolt, waveguide, 30-16–17
 Torr, 3-17
 Torus, equations for, 47-5
 Towers
 radio
 loading, 49-20–21
 material, 49-20
 steel
 foundation, 49-21–22
 pipe, 49-20
 standards, 49-20–22
 structural, 49-20
 unit stresses, 49-21
 Township, 3-17
 t_{on} , definition of, 20-10
 TR tubes, 16-45
 Track while scan (TWS), 36-22
 Tracking
 accuracy, feedback control system, 15-10–11
 radar, 36-21–22
 equation, 36-8
 noise jamming equation, 36-9
 Trade secret, protection of, 49-31–32
 Traffic
 allocation, ATM networks, 26-37–38
 carrying, load for grade of service, 39-13–14
 concepts, 39-10–14
 definitions for, 39-10–11
 equations, 39-12–13
 intensity, 39-11
 measures, 39-11
 models, cellular system, 46-21–22
 quantity, 39-11
 Traffic alert and collision avoidance system, 37-5–6
 Training sets, military equipment, 49-37
 Trains, vibration in, 49-25
 Transcendental functions, 47-21
 Transconductance- C (g_m-C) filters, 10-22–23
 Transconductance (S_m) of electron tube, 16-20
 Transducer loss
 definition of, 38-10
 transmission lines, 29-13–14
 Transfer
 function
 filter, 10-7–8
 first-order, 10-14–15
 generic biquadratic, 10-39–40, 10-44
 high-order, 10-23–34
 inverse z -, 15-17
 second-order, 10-15–23
 z -, 15-16–17
 modes, data, 26-16
 storage tubes, 16-57
 Transform
 Fourier
 discrete, 28-7–13
 discrete-time, 28-5–6
 fast, 28-9–13
 Laplace, 15-17
 z -, 15-16, 15-17
 Transformation(s)
 filters
 low-pass to band-stop, 9-19
 low-pass to bandpass, 9-14
 low-pass to high-pass, 9-14
 linear, matrix algebra, 47-32–33
 matrix, 31-3–4, 31-10–11
 T- π , network, 6-15–16
 Y- Δ network, 6-15–16
 Transformer(s)
 audio, 13-12–14
 cores for, 13-12–13
 design of, 13-13–14
 impedance match, 13-14
 insertion loss, 13-14
 carrier-frequency, 13-15
 design example, 13-15–16
 constant-voltage, 13-10–11
 converter/inverter, 13-11–12
 definition of, 13-3
 dielectric insulation and corona, 13-10
 equivalent circuit, 13-3
 ferrocresonant, 13-10–11
 frequency ranges, 13-3
 insulating materials
 classes of, 13-10–11
 operating temperatures, 13-10–11
 laminations, 13-3–4
 methods of winding, 13-9–10
 power, 13-3–10
 copper regulation, 13-6
 design for rectifiers, 13-5–9
 determining number of turns, 13-6
 determining wire size, 13-6
 effect of duty cycle on design, 13-9
 VA rating vs lamination sizes, 13-7
 wire table for, 13-8
 pulse, 13-17–18
 types of, 13-3
 volt-ampere rating, 13-5–7
 wideband, 13-15–17
 core material, 13-15
 design example, 13-15–16
 high frequency, 13-16–17
 Transforms
 inverse, 47-39
 Laplace, 6-19–21, 47-38–39
 Transient operation, laser, 41-5–6
 Transient voltage suppressors, 14-12
 Transients
 elementary cases, networks, 6-15–19
 operational calculus, and Laplace transforms, 6-19–22
 Transistor(s), 14-32, 18-13–21
 amplifier(s)
 active loads, 19-17–19
 current sources, 19-17–19
 harmonic distortion, 19-16–17
 single stage, 19-6–17
 frequency response, 19-6, 19-15
 β of, 18-19
 bipolar
 bias techniques, 19-2–5

- Transistor(s)—cont
 insulated-gate, 14-7-8, 14-30, 14-32
 junction, 18-16-21
 power, 14-29-30, 14-32
 circuits
 resonant, 14-36
 switching, 14-19
 field-effect, 14-8
 bias techniques, 19-5
 insulated-gate field-effect, 18-15-16
 JFET, symbol, 19-2-3
 junction
 definition of, 20-7
 field-effect, 18-14-15
 large-signal characteristics, 19-15
 metal-oxide semiconductor, 18-15-16
 MOS field-effect, 18-15-16
 MOSFET, symbol, 19-2-3
 MOSFETs, 14-32
 npn
 large symbol model, 19-2
 symbol, 19-2
 output stages, 19-15-16
 planar, definition of, 20-8
 power, 14-5-7
 breakdown sustaining voltage, 14-5
 dc current gain, 14-5
 reverse voltage operation, 14-6
 saturation voltage, 14-5
 second breakdown, 14-5-6
 switching times, 14-6
 requirements, switching power supplies, 14-20-27
 small signal
 equations, 19-5-8
 design, 19-9
 models, 19-5-8
 transistor logic. *see* TTL
 voltage regulator, linear, 14-18-19
 Transit, 37-5
 Transition matrix, 25-12
 Transition region, semiconductor, 18-9
 Translators
 FM, 35-11
 television broadcasting, 35-17-18
 Transmission
 coefficient, in waveguide, 31-3
 free-space, equations, 33-21-22
 line of sight, 33-13-29
 lines
 admittance, 29-6-7
 attenuation at UHF, 29-15
 characteristic impedance of, 29-21-25
 coplanar, 29-28-29
 current equations, 29-5-6
 fundamental quantities, 29-4-5
 impedance, 29-6-7
 transformation, high SWR, 29-12-13
 length of, 29-21
 mismatch, 29-13-14
 open at far end, 29-7-8
 Transmission—cont
 parameters, 29-4-5
 power and efficiency, 29-10-12
 resistance at UHF, 29-15
 rule of subscripts, 29-3
 sag calculations, 49-18-20
 short-circuited at far end, 29-7-8
 sign conventions, 29-3
 standing-wave ratio, 29-8-10
 strip, 29-27-28
 surface-wave, 30-24-26
 symbols, 29-3-4
 transducer loss, 29-13-14
 voltage equations, 29-5-6
 voltage reflection coefficient, 29-8-10
 loss, radio waves, 33-4
 pulse, 37-3
 in space, equations for, 27-7-8
 standards, TV broadcast stations, 35-12-13
 store-and-forward, 26-3, 26-4, 26-6
 system considerations, satellite communications, 27-19-21
 system requirements, AM broadcast stations, 35-5
 Transmittance, optical waves, 21-9-10
 Transmitter(s)
 optical, 22-5-8
 power, for required output signal/noise ratio, 33-27-29
 Transmitting objective loudness rating (TOLR), 38-4
 Transponders
 bandwidth, 27-17-18
 communications satellites, 27-15-19
 onboard processing asynchronous transfer mode (OBP ATM), 27-19
 satellite, multiple-access schemes, 27-19-21
 transparent, 27-4, 27-15
 Transport layer, network, 26-29-30
 Transversal filter, 36-16
 Transverse modulators, laser, 41-18
 Trapezium, equations for, 47-5
 Trapezoid, equations for, 47-4
 Traveling wave
 amplitude, 31-2
 antenna array, 32-36-37
 tubes, 16-24-25, 16-30-32, 27-15
 gyro, 16-38-39
 Tree codes, 25-4
 variable length, 25-10-11
 Trees, network, 26-34
 Trellises, 25-5-7
 Triacs, 14-11
 Triangle(s)
 equations for, 47-4
 oblique, 47-11-13
 hyperbolic, 47-15-16
 right, 47-11
 hyperbolic, 47-15
 spherical, 47-12
 Trigger, definition of, 20-11
 Triggered spark gaps, 16-43
 Trigonometric
 equations, 47-9-12
 integrands, 47-25-27
 Trigonometry
 hyperbolic, 47-14-16
 plane, 47-11-12
 spherical, 47-12-13
 Trimmer resistors, 5-13-14
 Triode
 grounded-grid, 17-8
 neutralized, 17-8
 RF amplifier circuit, 17-8
 Tropospheric
 effects, satellite communications, 27-35
 scatter propagation, 33-29, 33-32
 Trucks, vibration in, 49-25
 Trunks
 exchange-access, 38-6
 local interoffice, telephone, 38-6
 Truth table, 43-4
 Truth table, definition of, 20-11
 TTL, 20-30-32
 definition of, 20-11
 performance characteristics, 20-34
 voltage and current standards, 20-34
 Tube(s)
 cathode materials, 16-3-5
 cathode-ray, 16-55-57
 ceramic-envelope, advantages of, 16-19-20
 crossed-field, 16-34-38
 electron
 amplification factor (μ) of, 16-20
 characteristics equations, 16-21
 collision ionization, 16-19
 cooling of, 16-16-18
 flicker effect, 16-18-19
 geometry of, 16-22
 grid temperature, 16-18
 induced noise, 16-19
 partition noise, 16-18
 total plate resistance (R_p) of, 16-20
 transconductance (S_m) of, 16-20
 gas, 16-41-45
 image, 16-48-50
 intensifier, 16-48-50
 light-emitting, 16-45-48, 16-55-58
 light-sensing, 16-45-54
 linear-beam, 16-25-34
 microwave, 16-23-38
 cyclotron resonance, 16-38-39
 trends, 16-24-25
 noise in, 16-18-19
 TR, 16-45
 traveling-wave, 16-24-25, 16-30-32, 27-15
 Tuned circuits
 phase shift
 coupled, 8-7-9
 single-, 8-6-7
 Tungsten, 4-6-7, 4-10-11
 filaments, 16-3-5, 16-21
 lamps, 21-10
 Turbulence, atmospheric, on optical signal, 22-17-19

Turn-off losses, 14-32, 14-33
 Turn-on losses, 14-32, 14-33
 TV, direct broadcast from satellite, 27-4-5
 TVRO, 27-5
 Twin-T admittance-measuring circuit, 12-10
 Two
 -out-of-six, multifrequency numerical code, 2-14, 2-16
 -port junctions, 31-3
 powers of, table of, 48-5
 -stage gate circuit, 43-12
 and three-variable theorems, switching algebra, 43-8-9
 -tone angle modulation, 23-8
 -way channels, 25-26
 Two-level main memory, definition of, 20-11
 Two-level metallization, definition of, 20-11
 Twos complement
 addition using, 42-8
 subtraction using, 42-8
 system, 42-7

U

U-Net, 26-33
 UHF operation, power grid tube, 17-10
 ULSI, definition of, 20-11
 Ultra-fast rectifiers, 14-32
 Ultragigabit networks, 26-35-36
 Ultraviolet spectrum, 21-2
 Uniaxial crystal, 41-16-17
 Uninterruptible power supply, 14-47-48
 Unit impulse, circuit response related to, 6-20
 Unit indicators, military equipment, 49-38-39
 Unit pole, 3-17
 Unit step, circuit response related to, 6-20-21
 Units
 dimensionless quantities, 3-5-6
 non-SI, tables of, 3-4
 numerical values, 3-6
 rules and style conventions, 3-4-5
 SI units, 1-143, 3-2-4
 Unity feedback
 acceleration error constant, 15-11
 positive error constant, 15-11
 single-deg-of-freedom structure, 15-11
 single-loop system, 15-5
 velocity error constant, 15-11
 Unity-power-factor circuits, 14-38
 Universal codes, 25-9
 Universal time scales (UT), 1-143
 Unregulated, power supplies, 14-43
 Unsigned integer, 42-4
 Up-link power-flux density at satellite, 27-14
 UPS *See* Uninterruptible Power Supply

Uranium, 4-6-7, 4-10-11
 URN scheme, 26-22
 US Coast Guard, 37-2
 US Government Printing Office, 35-26
 US government, reliability specifications, 45-20
 User stack register, computer processor, 42-15
 UT1 time corrections, WWV and WWVH, 1-154

V

V curve, 12-32
 V-number, optical fiber, 22-21
 Vacancy, definition of, 20-11
 Vacuum photodiodes, 21-22-25
 Vanadium, 4-6-7, 4-10-11
 Varactor, 18-13
 Variable, random, 44-2
 Variance, 44-3
 Variate, 44-2
 V_{BG} , definition of, 20-11
 V_{cc} , definition of, 20-11
 V_{dth} , definition of, 20-11
 VDSL, 38-8
 Vector
 analysis equations, 47-35-37
 generation, definition of, 20-11
 impedance analyzer, 12-11, 12-13
 multiplication, distributive law for, 47-35
 triple product, 47-35
 Vehicle motion, effect on mobile communications, 46-8
 Vehicles, typical, vibration, 49-25
 Velocity(ies)
 displacement, and acceleration magnitudes, 49-23
 error constant, unity feedback, 15-11
 phasor, sound wave, 40-3
 of sound in gases, 40-3-4
 sound, in liquids, 40-4-5
 various metals, 4-33
 wind, pressure on steel tower, 49-22
 Vertical redundancy check, 26-15
 Vertical resolution, television, 35-35-36
 Very-High-Speed DSL, 38-8
 Very low frequencies, propagation of, 33-3
 Very-small-aperture terminal *See* VSAT
 Vestigial sideband, amplitude modulation, 23-4
 VFB, definition of, 20-11
 V_{GG} , definition of, 20-11
 V_{GS} , definition of, 20-11
 VHF omnidirectional range, 37-4
 VHPIC, definition of, 20-11
 VHSIC, definition of, 20-11
 Vibration
 damping, reduction, 49-24
 natural frequency, 49-24
 and shock isolation, 49-22-26

Vibration—cont
 tests, 5-6, 5-7
 typical vehicles, 49-25
 Vibratory forces, resonance, 49-24
 Video
 disk recordings, 35-33-34
 recordings, recommended practices, 35-32-33
 tape
 formats, 35-29-30
 standards, 35-30-32
 Vidicons, 16-53-54
 VIR signal, TV, 35-13
 Visible light spectrum, 21-2
 Visual power, maximum, TV broadcasting, 35-11
 Viterbi algorithm, 25-16
 VLSI, definition of, 20-11
 Voice of America, 35-24
 Volt, 3-3
 per meter, 3-3
 Voltage
 breakdown, spark-gap, 49-3
 clamps, 14-32, 14-34
 current, relationship, 16-12
 doublers, 14-13
 cascade, 14-13, 14-16
 conventional, 14-13, 14-16
 equations, transmission lines, 29-5-6
 -mode control, 14-35-36
 multiplication, 14-13
 multipliers, single-phase, 14-13
 power transistor, breakdown sustaining, 14-5
 rating
 components, 5-4
 rectifiers, 14-3
 reference, IC, 20-43, 20-45
 reflection coefficient, transmission lines, 29-8-10
 regulators
 ac, 14-46
 ferresonant, 13-10-11
 gas tube, 16-43
 IC, 20-45
 linear transistor, 14-18-19
 saturation, power transistor, 14-5
 snubbers, 14-32
 standards, TTL, 20-34
 -tuned crossed-field oscillator, 16-36-37
 -tuned magnetron, 16-36-37
 variable capacitor, 18-12, 18-13
 Voltampere
 rating, transformer, 13-5-7
 Volume
 clutter equation, radar, 36-8-9
 Volume, gram atomic, 4-3, 4-4, 4-6
 VOR, 37-4
 VSAT, 27-4
 network configuration, 27-29-31
 V_{SS} , definition of, 20-11
 V_T (V_{BG}), definition of, 20-11
 V_{TD} , definition of, 20-11
 V_{TL} , definition of, 20-11

W

- W**, definition of, 20-11
WAAS, 37-8
Wafer flat, definition of, 20-11
Wafer stepper, definition of, 20-11
Wagner ground connection, 12-3
Walls, sound isolation, 40-19, 40-20-21
Water cooling, electron tube, 16-16
Watt, 3-3
 per meter kelvin, 3-3
 per square meter, 3-3
 per square meter steradian, 3-3
 per steradian, 3-3
Wattage rating, resistors, 5-9
Wave
 equation, 40-3
 propagation
 above 30 MHz, 33-13-29
 HF, 33-6-12
 LF, 33-3-6
 MF, 33-3-6
 VLF, 33-3
 Stokes, 41-19
 transmission, hops, 33-8
Waveform(s)
 channels, 25-19
 modulated, analytic signal
 representation of, 23-3
 periodic
 average power, 7-8
 and Fourier series, 7-11-13
Waveguide(s)
 attenuation in, beyond cutoff, 30-7
 bolt torque, 30-16-17
 cavity, simple, 30-24
 circuit elements, 30-18-19
 circular, 30-5-6
 current and voltage not normalized, 31-9
 dielectric-rod, 30-26-27
 dispersion, 22-22-24
 fiber, transmission, 22-19-26
 flange
 insertion loss, 30-17
 leakage, 30-16
 resistance, 30-16-17
 flexible, 30-8, 30-14-15
 hybrid junctions, 30-19-20
 losses, 30-14, 30-16-18
 and noise temperature, 30-17-18
 material and modes, 30-14, 30-16
 net power flow, 31-9
 normalized
 current and voltage, 31-9
 impedance and admittance, 31-9
 power flow, 31-2
 propagation of electromagnetic waves in, 30-3
 rectangular, 30-3-4
 reflection coefficient, 31-2
 ridged, 30-8-13
 slots, 32-20-21
 square, 30-7
 standard, 30-7, 30-10-11
Wavelength(s)
 division multiplexing (WDM), 26-35
 division optical network, 26-34
 -frequency, conversions, 1-2
 optical spectrum, 21-2
Waves
 optical, interaction with matter, 21-9-10
 plane, 40-3
 sound, theory of, 40-3-5
 spherical, 40-3
 TE, 30-3
 TEM, 30-3
Wayne-Kerr bridge, 12-10-11
Weather data, 49-3-5
Weather radar equation, 36-9
Weather radars, 36-16
Weather, space, 1-155
Weber, 3-3
Wedge frustrum, equations for, 47-6
Weighting factor, psophometric, 2-5
Weibull analysis, 45-22-23
Weibull distribution, 45-11
Wet etching, ICs, 20-19
Wheatstone bridge, 12-3
Wide Area Augmentation Service, 37-8
Wide area networks, high speed, 26-36-39
Wideband transformers, 13-15-17
 core material, 13-15
 design example, 13-15-16
 high frequency, 13-16-17
Wien bridge, 12-4
oscillator, 19-25
Wien's displacement law, 21-3, 21-6
Wind
 velocities and pressure, on steel towers, 49-22
 velocity extremes, 49-5
Winding
 methods, transformers, 13-9-10
 resistance loss, inductors, 13-27
Windowing, FFT, 28-11
Wire
 applications and insulations, NEC, 49-12
 color coding, 49-17
 identification, 49-16-17
 insulating coatings, comparison of, 13-12
 magnet, 6-3-4
 size, power transformer, determining, 13-6
 tables
 power transformer design, 13-8
 telephone, 38-11
Wireless access, 38-6
Wirewound resistors, 5-11-12, 5-15
Wiring
 diagrams, 49-26
 electronic equipment, 49-12, 49-14-15, 49-16
 power supply, 49-11-12
 temperature conditions, 49-16
Word, definition of, 20-11
Word error, probability of, 22-4
Workplace-industrial noise, assessment of, 40-12
World
 rain-rate climate regions, 33-20, 33-22
 temperatures, 49-4
 time chart, 49-5, 49-7
 time zones, 1-144, 1-145
World Administrative Radio Conferences (WARC), 35-26
Write
 enable, definition of, 20-11
 time, definition of, 20-11
WVW and WVVH, 1-148-157
 accuracy, 1-149
 audio signals by telephone, 1-157
 BCD time code, 1-152-154
 broadcast schedule, 1-149-150
 facilities and services, 1-149-151
 geophysical alerts, 1-155-156
 global positioning system, 1-156-157
 leap seconds, 1-154
 marine storm warnings, 1-156
 official announcements, 1-155-157
 services provided, 1-149-151
 silent periods, 1-152
 stability, 1-149
 standard audio frequencies, 1-151-152
 standard time intervals, 1-151
 time announcements, 1-151
 time code format, 1-152
 UT1 time corrections, 1-154
WVVB, 1-157-162
 accuracy, 1-157
 Global Positioning System (GPS), 1-160-162
 Internet time signals, 1-158, 1-160
 stability, 1-157
 station identification, 1-157
 telephone time signals, 1-158, 1-160
 time code, 1-157-158
Wye rectifier, 14-13

X
X-21, CCITT recommendation, 26-28-29
Xenon, 4-6-7, 4-10-11
XOR gate, definition of, 20-6

Y
Y- Δ transformation, network, 6-15-16
Yagi-Uda antenna, 32-24-28
Yard, 3-17
Year, 3-17
Yield, definition of, 20-11
Yield, IC fabrication process, 20-23, 20-26
Ytterbium, 4-6-7, 4-10-11
Yttrium, 4-6-7, 4-10-11

Z

- z-domain transfer functions, SC, 10-39-40
- z-transfer function, 15-16-17
- z-transform, 15-16, 15-17
- Z-transform, 28-3-5
- z-transformation, bilinear, 28-14
- Z_0 maximum available power, P_{iO} , definition of, 12-25
- Zener, definition of, 20-11
- Zener diodes, 14-4-5, 18-12-13
- Zener effect, 18-11
- Zero
 - padding, FFT, 28-11
 - relative level points, CCITT and CCIR recommendations, 2-3-4
 - transmission-level points, 2-3-4
- Zero placement formulas, 10-43
- Zinc, 4-6-7, 4-10-11
- Zirconium, 4-6-7, 4-10-11
- ZMR, 20-103
- Zone coverage, spacecraft antennas, 27-32
- Zones
 - Fresnel, 33-16-17
 - television broadcasting, 35-11
 - time, world, 1-144, 1-145

Reference Data for Engineers

Radio, Electronics, Computer, and Communications

NINTH EDITION

- ▶ For Engineers, Scientists, and Technologists
- ▶ Features contributions from 96 professionals
- ▶ Recommends additional sources of information

Reference Data for Engineers, Ninth Edition is the most respected, reliable, and indispensable information tool for technical professionals around the globe today. Written by professionals for professionals, this unique reference is the combined effort of ninety-six licensed engineers, scientists, educators and other recognized specialists in the fields of electronics, radio, computer, and communications technology.

By providing an abundance of information on essential, need-to-know topics without heavy emphasis on complicated mathematics, *Reference Data for Engineers* is an absolute "must-have" for every engineer who requires comprehensive electrical, electronics, and communications data at his or her fingertips. Featured in the Ninth Edition is updated coverage on intellectual property and patents, probability and design, antennas, power electronics, rectifiers, power supplies, and properties of materials. Useful information on units, constants and conversion factors, active filter design, antennas, integrated circuits, surface acoustic wave design, and digital signal processing is also included. The Ninth Edition also offers new knowledge in the fields of satellite technology, space communication, microwave science, telecommunication, global positioning systems, frequency data, and radar.

- Widely acclaimed as the most practical reference ever published for electronics and computer professionals, technicians and engineers.
- First revision in a decade brings the material completely up-to-date.
- Contributors represent some of the leading technology companies, organizations and universities throughout the world.

Related Titles:

Modern Dictionary of Electronics, 7th Edition

Rudolf F. Graf

0-7506-9866-7, Paperback, 869 pp

Telecommunications Engineer's Reference Book, 2nd Edition

Fraidoon Mazda

0-240-51491-2, Hardcover, 1224 pp

Electrical Engineer's Reference Book, 15th Edition

G R Jones, M A Laughton, M G Say

0-7506-1202-9, Hardcover, 1376 pp

 **Newnes**
An imprint of Butterworth-Heinemann

www.newnespress.com

