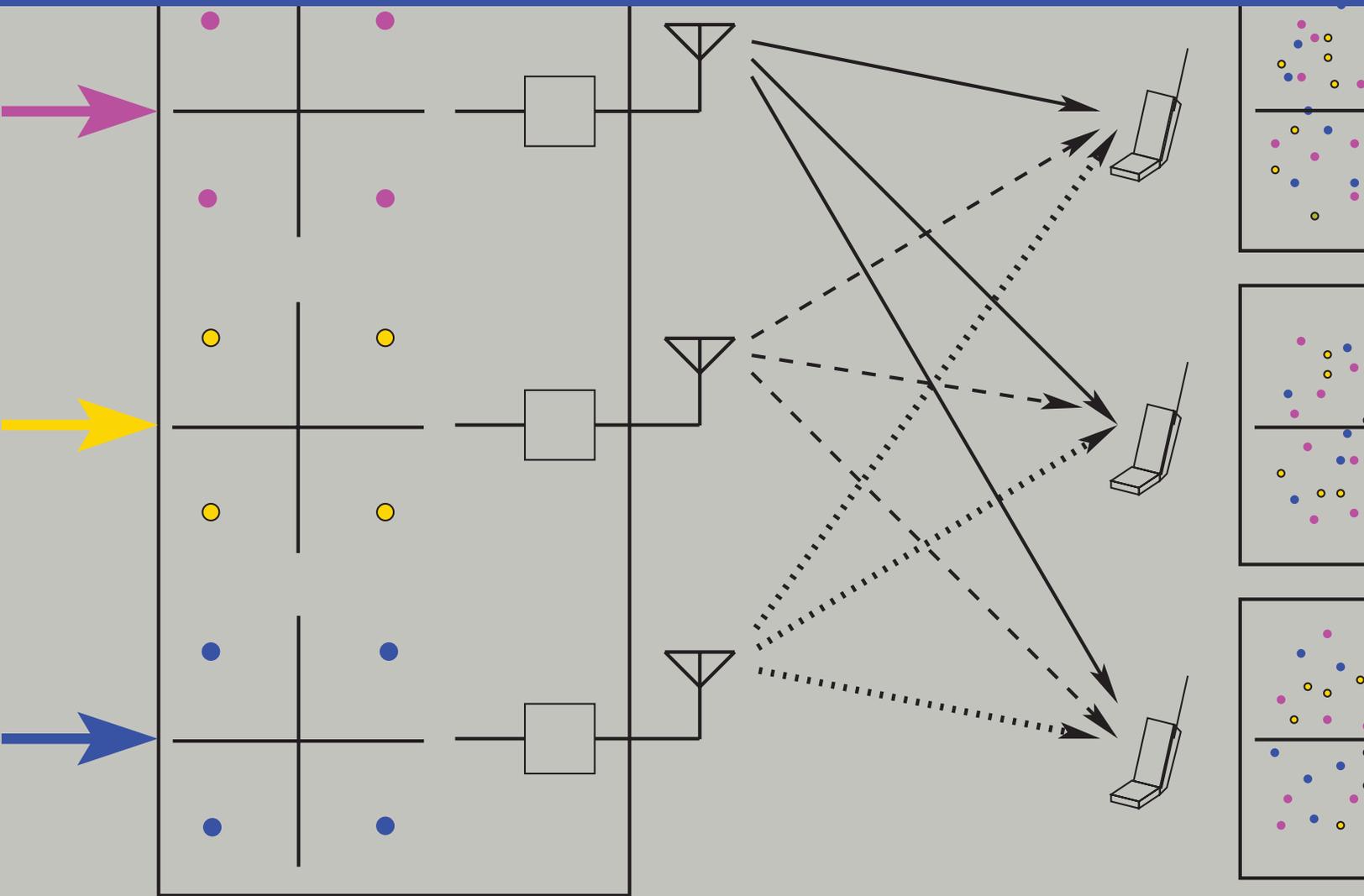


# MICROWAVE AND RF DESIGN RADIO SYSTEMS



# Microwave and RF Design

## *Radio Systems*

Volume 1

Third Edition

Michael Steer



# Microwave and RF Design

## *Radio Systems*

Volume 1

Third Edition

Michael Steer

Copyright © 2019 by M.B. Steer

Citation: Steer, Michael. *Microwave and RF Design: Radio Systems*. Volume 1. (Third Edition), NC State University, 2019. doi: [https://doi.org/10.5149/9781469656915\\_Steer](https://doi.org/10.5149/9781469656915_Steer)

This work is licensed under a Creative Commons Attribution-NonCommercial 4.0 International license (CC BY-NC 4.0). To view a copy of the license, visit <http://creativecommons.org/licenses>.

ISBN 978-1-4696-5690-8 (paperback)  
ISBN 978-1-4696-5691-5 (open access ebook)

Published by NC State University

**NC STATE UNIVERSITY**

Distributed by the University of North Carolina Press  
[www.uncpress.org](http://www.uncpress.org)

Printing: 1

To

Ross Lampe, in honor of his dedication to our profession

# Preface

The book series *Microwave and RF Design* is a comprehensive treatment of radio frequency (RF) and microwave design with a modern “systems-first” approach. A strong emphasis on design permeates the series with extensive case studies and design examples. Design is oriented towards cellular communications and microstrip design so that lessons learned can be applied to real-world design tasks. The books in the Microwave and RF Design series are:

- Microwave and RF Design: Radio Systems, Volume 1
- Microwave and RF Design: Transmission Lines, Volume 2
- Microwave and RF Design: Networks, Volume 3
- Microwave and RF Design: Modules, Volume 4
- Microwave and RF Design: Amplifiers and Oscillators, Volume 5

The length and format of each is suitable for automatic printing and binding.

## **Rationale**

The central philosophy behind this series’s popular approach is that the student or practicing engineer will develop a full appreciation for RF and microwave engineering and gain the practical skills to perform system-level design decisions. Now more than ever companies need engineers with an ingrained appreciation of systems and armed with the skills to make system decisions. One of the greatest challenges facing RF and microwave engineering is the increasing level of abstraction needed to create innovative microwave and RF systems. This book series is organized in such a way that the reader comes to understand the impact that system-level decisions have on component and subsystem design. At the same time, the capabilities of technologies, components, and subsystems impact system design. The book series is meticulously crafted to intertwine these themes.

## **Audience**

The book series was originally developed for three courses at North Carolina State University. One is a final-year undergraduate class, another an introductory graduate class, and the third an advanced graduate class. Books in the series are used as supplementary texts in two other classes. There are extensive case studies, examples, and end of chapter problems ranging from straight-forward to in-depth problems requiring hours to solve. A companion book, *Fundamentals of Microwave and RF Design*, is more suitable for an undergraduate class yet there is a direct linkage between the material in this book and the series which can then be used as a career-long reference text. I believe it is completely understandable for senior-level students where a microwave/RF engineering course is offered. The book series is a comprehensive RF and microwave text and reference, with detailed index, appendices, and cross-references throughout. Practicing engineers will find the book series a valuable systems primer, a refresher as needed, and a

reference tool in the field. Additionally, it can serve as a valuable, accessible resource for those outside RF circuit engineering who need to understand how they can work with RF hardware engineers.

### **Organization**

This book is a volume in a five volume series on RF and microwave design. The first volume in the series, *Microwave and RF Design: Radio Systems*, addresses radio systems mainly following the evolution of cellular radio. A central aspect of microwave engineering is distributed effects considered in the second volume of this book series, *Microwave and RF Design: Transmission Lines*. Here transmission lines are treated as supporting forward- and backward-traveling voltage and current waves and these are related to electromagnetic effects. The third volume, *Microwave and RF Design: Networks*, covers microwave network theory which is the theory that describes power flow and can be used with transmission line effects. Topics covered in *Microwave and RF Design: Modules*, focus on designing microwave circuits and systems using modules introducing a large number of different modules. Modules is just another term for a network but the implication is that it is packaged and often available off-the-shelf. Other topics that are important in system design using modules are considered including noise, distortion, and dynamic range. Most microwave and RF designers construct systems using modules developed by other engineers who specialize in developing the modules. Examples are filter and amplifier modules which once designed can be used in many different systems. Much of microwave design is about maximizing dynamic range, minimizing noise, and minimizing DC power consumption. The fifth volume in this series, *Microwave and RF Design: Amplifiers and Oscillators*, considers amplifier and oscillator design and develops the skills required to develop modules.

### **Volume 1: Microwave and RF Design: Radio Systems**

The first book of the series covers RF systems. It describes system concepts and provides comprehensive knowledge of RF and microwave systems. The emphasis is on understanding how systems are crafted from many different technologies and concepts. The reader gains valuable insight into how different technologies can be traded off in meeting system requirements. I do not believe this systems presentation is available anywhere else in such a compact form.

### **Volume 2: Microwave and RF Design: Transmission Lines**

This book begins with a chapter on transmission line theory and introduces the concepts of forward- and backward-traveling waves. Many examples are included of advanced techniques for analyzing and designing transmission line networks. This is followed by a chapter on planar transmission lines with microstrip lines primarily used in design examples. Design examples illustrate some of the less quantifiable design decisions that must be made. The next chapter describes frequency-dependent transmission line effects and describes the design choices that must be taken to avoid multimoding. The final chapter in this volume addresses coupled-lines. It is shown how to design coupled-line networks that exploit this distributed effect to realize novel circuit functionality and how to design networks that minimize negative effects. The modern treatment of transmission lines in this volume emphasizes planar circuit design and the practical aspects of designing

around unwanted effects. Detailed design of a directional coupler is used to illustrate the use of coupled lines. Network equivalents of coupled lines are introduced as fundamental building blocks that are used later in the synthesis of coupled-line filters. The text, examples, and problems introduce the often hidden design requirements of designing to mitigate parasitic effects and unwanted modes of operation.

### **Volume 3: Microwave and RF Design: Networks**

Volume 3 focuses on microwave networks with descriptions based on  $S$  parameters and  $ABCD$  matrices, and the representation of reflection and transmission information on polar plots called Smith charts. Microwave measurement and calibration technology are examined. A sampling of the wide variety of microwave elements based on transmission lines is presented. It is shown how many of these have lumped-element equivalents and how lumped elements and transmission lines can be combined as a compromise between the high performance of transmission line structures and the compactness of lumped elements. This volume concludes with an in-depth treatment of matching for maximum power transfer. Both lumped-element and distributed-element matching are presented.

### **Volume 4: Microwave and RF Design: Modules**

Volume 4 focuses on the design of systems based on microwave modules. The book considers the wide variety of RF modules including amplifiers, local oscillators, switches, circulators, isolators, phase detectors, frequency multipliers and dividers, phase-locked loops, and direct digital synthesizers. The use of modules has become increasingly important in RF and microwave engineering. A wide variety of passive and active modules are available and high-performance systems can be realized cost effectively and with stellar performance by using off-the-shelf modules interconnected using planar transmission lines. Module vendors are encouraged by the market to develop competitive modules that can be used in a wide variety of applications. The great majority of RF and microwave engineers either develop modules or use modules to realize RF systems. Systems must also be concerned with noise and distortion, including distortion that originates in supposedly linear elements. Something as simple as a termination can produce distortion called passive intermodulation distortion. Design techniques are presented for designing cascaded systems while managing noise and distortion. Filters are also modules and general filter theory is covered and the design of parallel coupled line filters is presented in detail. Filter design is presented as a mixture of art and science. This mix, and the thought processes involved, are emphasized through the design of a filter integrated throughout this chapter.

### **Volume 5: Microwave and RF Design: Amplifiers and Oscillators**

The fifth volume presents the design of amplifiers and oscillators in a way that enables state-of-the-art designs to be developed. Detailed strategies for amplifiers and voltage-controlled oscillators are presented. Design of competitive microwave amplifiers and oscillators are particularly challenging as many trade-offs are required in design, and the design decisions cannot be reduced to a formulaic flow. Very detailed case studies are presented and while some may seem quite complicated, they parallel the level of sophistication required to develop competitive designs.

### Case Studies

A key feature of this book series is the use of real world case studies of leading edge designs. Some of the case studies are designs done in my research group to demonstrate design techniques resulting in leading performance. The case studies and the persons responsible for helping to develop them are as follows.

1. Software defined radio transmitter.
2. High dynamic range down converter design. This case study was developed with Alan Victor.
3. Design of a third-order Chebyshev combline filter. This case study was developed with Wael Fathelbab.
4. Design of a bandstop filter. This case study was developed with Wael Fathelbab.
5. Tunable Resonator with a varactor diode stack. This case study was developed with Alan Victor.
6. Analysis of a 15 GHz Receiver. This case study was developed with Alan Victor.
7. Transceiver Architecture. This case study was developed with Alan Victor.
8. Narrowband linear amplifier design. This case study was developed with Dane Collins and National Instruments Corporation.
9. Wideband Amplifier Design. This case study was developed with Dane Collins and National Instruments Corporation.
10. Distributed biasing of differential amplifiers. This case study was developed with Wael Fathelbab.
11. Analysis of a distributed amplifier. This case study was developed with Ratan Bhatia, Jason Gerber, Tony Kwan, and Rowan Gilmore.
12. Design of a WiMAX power amplifier. This case study was developed with Dane Collins and National Instruments Corporation.
13. Reflection oscillator. This case study was developed with Dane Collins and National Instruments Corporation.
14. Design of a C-Band VCO. This case study was developed with Alan Victor.
15. Oscillator phase noise analysis. This case study was developed with Dane Collins and National Instruments Corporation.

Many of these case studies are available as captioned YouTube videos and qualified instructors can request higher resolution videos from the author.

### Course Structures

Based on the adoption of the first and second editions at universities, several different university courses have been developed using various parts of what was originally one very large book. The book supports teaching two or three classes with courses varying by the selection of volumes and chapters. A standard microwave class following the format of earlier microwave texts can be taught using the second and third volumes. Such a course will benefit from the strong practical design flavor and modern treatment of measurement technology, Smith charts, and matching networks. Transmission line propagation and design is presented in the context of microstrip technology providing an immediately useful skill. The subtleties of multimoding are also presented in the context of microstrip lines. In such

a class the first volume on microwave systems can be assigned for self-learning.

Another approach is to teach a course that focuses on transmission line effects including parallel coupled-line filters and module design. Such a class would focus on Volumes 2, 3 and 4. A filter design course would focus on using Volume 4 on module design. A course on amplifier and oscillator design would use Volume 5. This course is supported by a large number of case studies that present design concepts that would otherwise be difficult to put into the flow of the textbook.

Another option suited to an undergraduate or introductory graduate class is to teach a class that enables engineers to develop RF and microwave systems. This class uses portions of Volumes 2, 3 and 4. This class then omits detailed filter, amplifier, and oscillator design.

The fundamental philosophy behind the book series is that the broader impact of the material should be presented first. Systems should be discussed up front and not left as an afterthought for the final chapter of a textbook, the last lecture of the semester, or the last course of a curriculum.

The book series is written so that all electrical engineers can gain an appreciation of RF and microwave hardware engineering. The body of the text can be covered without strong reliance on this electromagnetic theory, but it is there for those who desire it for teaching or reader review. The book is rich with detailed information and also serves as a technical reference.

### **The Systems Engineer**

Systems are developed beginning with fuzzy requirements for components and subsystems. Just as system requirements provide impetus to develop new base technologies, the development of new technologies provides new capabilities that drive innovation and new systems. The new capabilities may arise from developments made in support of other systems. Sometimes serendipity leads to the new capabilities. Creating innovative microwave and RF systems that address market needs or provide for new opportunities is the most exciting challenge in RF design. The engineers who can conceptualize and architect new RF systems are in great demand. This book began as an effort to train RF systems engineers and as an RF systems resource for practicing engineers. Many RF systems engineers began their careers when systems were simple. Today, appreciating a system requires higher levels of abstraction than in the past, but it also requires detailed knowledge or the ability to access detailed knowledge and expertise. So what makes a systems engineer? There is not a simple answer, but many partial answers. We know that system engineers have great technical confidence and broad appreciation for technologies. They are both broad in their knowledge of a large swath of technologies and also deep in knowledge of a few areas, sometimes called the "T" model. One book or course will not make a systems engineer. It is clear that there must be a diverse set of experiences. This book series fulfills the role of fostering both high-level abstraction of RF engineering and also detailed design skills to realize effective RF and microwave modules. My hope is that this book will provide the necessary background for the next generation of RF systems engineers by stressing system principles immediately, followed by core RF technologies. Core technologies are thereby covered within the context of the systems in which they are used.

### Supplementary Materials

Supplementary materials available to qualified instructors adopting the book include PowerPoint slides and solutions to the end-of-chapter problems. Requests should be directed to the author. Access to downloads of the books, additional material and YouTube videos of many case studies are available at <https://www.lib.ncsu.edu/do/open-education>

### Acknowledgments

Writing this book has been a large task and I am indebted to the many people who helped along the way. First I want to thank the more than 1200 electrical engineering graduate students who used drafts and the first two editions at NC State. I thank the many instructors and students who have provided feedback. I particularly thank Dr. Wael Fathelbab, a filter expert, who co-wrote an early version of the filter chapter. Professor Andreas Cangellaris helped in developing the early structure of the book. Many people have reviewed the book and provided suggestions. I thank input on the structure of the manuscript: Professors Mark Wharton and Nuno Carvalho of Universidade de Aveiro, Professors Ed Delp and Saul Gelfand of Purdue University, Professor Lynn Carpenter of Pennsylvania State University, Professor Grant Ellis of the Universiti Teknologi Petronas, Professor Islam Eshrah of Cairo University, Professor Mohammad Essaaidi and Dr. Otman Aghzout of Abdelmalek Essaadi Univeristy, Professor Jianguo Ma of Guangdong University of Technology, Dr. Jayesh Nath of Apple, Mr. Sony Rowland of the U.S. Navy, and Dr. Jonathan Wilkerson of Lawrence Livermore National Laboratories, Dr. Josh Wetherington of Vadum, Dr. Glen Garner of Vadum, and Mr. Justin Lowry who graduated from North Carolina State University.

Many people helped in producing this book. In the first edition I was assisted by Ms. Claire Sideri, Ms. Susan Manning, and Mr. Robert Lawless who assisted in layout and production. The publisher, task master, and chief coordinator, Mr. Dudley Kay, provided focus and tremendous assistance in developing the first and second editions of the book, collecting feedback from many instructors and reviewers. I thank the Institution of Engineering and Technology, who acquired the original publisher, for returning the copyright to me. This open access book was facilitated by John McLeod and Samuel Dalzell of the University of North Carolina Press, and by Micah Vandergrift and William Cross of NC State University Libraries. The open access ebooks are host by NC State University Libraries.

The book was produced using LaTeX and open access fonts, line art was drawn using xfig and inkscape, and images were edited in gimp. So thanks to the many volunteers who developed these packages.

My family, Mary, Cormac, Fiona, and Killian, gracefully put up with my absence for innumerable nights and weekends, many more than I could have ever imagined. I truly thank them. I also thank my academic sponsor, Dr. Ross Lampe, Jr., whose support of the university and its mission enabled me to pursue high risk and high reward endeavors including this book.

**Michael Steer**  
**North Carolina State University**  
**Raleigh, North Carolina**  
**mbs@ncsu.edu**

## List of Trademarks

3GPP® is a registered trademark of the European Telecommunications Standards Institute.

802® is a registered trademark of the Institute of Electrical & Electronics Engineers .

APC-7® is a registered trademark of Amphenol Corporation.

AT&T® is a registered trademark of AT&T Intellectual Property II, L.P.

AWR® is a registered trademark of National Instruments Corporation.

AWRDE® is a trademark of National Instruments Corporation.

Bluetooth® is a registered trademark of the Bluetooth Special Interest Group.

GSM® is a registered trademark of the GSM MOU Association.

Mathcad® is a registered trademark of Parametric Technology Corporation.

MATLAB® is a registered trademark of The MathWorks, Inc.

NEC® is a registered trademark of NEC Corporation.

OFDMA® is a registered trademark of Runcom Technologies Ltd.

Qualcomm® is a registered trademark of Qualcomm Inc.

Teflon® is a registered trademark of E. I. du Pont de Nemours.

RFMD® is a registered trademark of RF Micro Devices, Inc.

SONNET® is a trademark of Sonnet Corporation.

Smith is a registered trademark of the Institute of Electrical and Electronics Engineers.

Touchstone® is a registered trademark of Agilent Corporation.

WiFi® is a registered trademark of the Wi-Fi Alliance.

WiMAX® is a registered trademark of the WiMAX Forum.

All other trademarks are the properties of their respective owners.



# Contents

Preface . . . . .	v
<b>1 Introduction to RF and Microwave Systems . . . . .</b>	<b>1</b>
1.1 Introduction . . . . .	1
1.2 RF and Microwave Engineering . . . . .	2
1.3 Communication Over Distance . . . . .	5
1.3.1 Electromagnetic Fields . . . . .	6
1.3.2 Biot-Savart Law . . . . .	7
1.3.3 Faraday's Law of Induction . . . . .	7
1.3.4 Ampere's Circuital Law . . . . .	7
1.3.5 Gauss's Law . . . . .	8
1.3.6 Gauss's Law of Magnetism . . . . .	8
1.3.7 Telegraph . . . . .	8
1.3.8 The Origins of Radio . . . . .	10
1.3.9 Maxwell's Equations . . . . .	10
1.3.10 Transmission of Radio Signals . . . . .	12
1.3.11 Early Radio . . . . .	13
1.4 Radio Architecture . . . . .	14
1.5 Conventional Wireless Communications . . . . .	15
1.6 RF Power Calculations . . . . .	17
1.6.1 RF Propagation . . . . .	17
1.6.2 Logarithm . . . . .	18
1.6.3 Decibels . . . . .	18
1.6.4 Decibels and Voltage Gain . . . . .	20
1.7 Photons and Electromagnetic Waves . . . . .	21
1.8 Summary . . . . .	22
1.9 References . . . . .	24
1.10 Exercises . . . . .	24
1.10.1 Exercises By Section . . . . .	26
1.10.2 Answers to Selected Exercises . . . . .	26
<b>2 Modulation . . . . .</b>	<b>27</b>
2.1 Introduction . . . . .	27
2.2 Radio Signal Metrics . . . . .	28
2.2.1 Crest Factor and Peak-to-Average Power Ratio . . . . .	29
2.2.2 Peak-to-Mean Envelope Power Ratio . . . . .	31
2.2.3 Two-Tone Signal . . . . .	34
2.3 Modulation Overview . . . . .	36
2.4 Analog Modulation . . . . .	37
2.4.1 Amplitude Modulation . . . . .	37
2.4.2 Phase Modulation . . . . .	40
2.4.3 Frequency Modulation . . . . .	41
2.4.4 Analog Modulation Summary . . . . .	44
2.5 Digital Modulation . . . . .	44

2.5.1	Modulation Efficiency . . . . .	46
2.6	Frequency Shift Keying, FSK . . . . .	47
2.6.1	Essentials of FSK Modulation . . . . .	47
2.6.2	Gaussian Minimum Shift Keying . . . . .	48
2.6.3	Doppler Effect . . . . .	49
2.6.4	Summary . . . . .	49
2.7	Carrier Recovery . . . . .	50
2.8	Phase Shift Keying Modulation . . . . .	50
2.8.1	Essentials of PSK . . . . .	51
2.8.2	Binary Phase Shift Keying . . . . .	53
2.8.3	Quadra-Phase Shift Keying, QPSK . . . . .	56
2.8.4	$\pi/4$ Quadrature Phase Shift Keying . . . . .	57
2.8.5	Differential Quadra Phase Shift Keying, DQPSK . . . . .	59
2.8.6	Offset Quadra Phase Shift Keying, OQPSK . . . . .	61
2.8.7	$3\pi/8$ -8PSK, Rotating Eight-State Phase Shift Keying . . . . .	63
2.8.8	Summary . . . . .	64
2.9	Quadrature Amplitude Modulation . . . . .	64
2.10	Digital Modulation Summary . . . . .	65
2.11	Interference and Distortion . . . . .	66
2.11.1	Cochannel Interference . . . . .	66
2.11.2	Adjacent Channel Interference . . . . .	67
2.11.3	Noise, Distortion, and Constellation Diagrams . . . . .	68
2.11.4	Comparison of GMSK and $\pi/4$ DQPSK Modulation . . . . .	68
2.11.5	Error Vector Magnitude . . . . .	69
2.12	Summary . . . . .	72
2.13	References . . . . .	73
2.14	Exercises . . . . .	74
2.14.1	Exercises By Section . . . . .	76
2.14.2	Answers to Selected Exercises . . . . .	76
<b>3</b>	<b>Transmitters and Receivers . . . . .</b>	<b>77</b>
3.1	Introduction . . . . .	77
3.2	Single-Sideband and Double-Sideband Modulation . . . . .	78
3.3	Early Modulation and Demodulation Technology . . . . .	80
3.3.1	Heterodyne Receiver . . . . .	80
3.3.2	Homodyne Receiver . . . . .	80
3.3.3	FM Modulator . . . . .	81
3.3.4	FM Demodulator . . . . .	82
3.3.5	Superheterodyne Receiver . . . . .	82
3.3.6	Summary . . . . .	82
3.4	Receiver and Transmitter Architectures . . . . .	82
3.4.1	Radio as a Cascade of Two-Ports . . . . .	83
3.4.2	Heterodyne Transmitter and Receiver . . . . .	83
3.4.3	Superheterodyne Receiver Architecture . . . . .	84
3.4.4	Single Heterodyne Receiver . . . . .	85
3.4.5	Transceiver . . . . .	86
3.4.6	Hartley Modulator . . . . .	87
3.4.7	The Hartley Modulator in Modern Radios . . . . .	87
3.5	Carrier Recovery . . . . .	88
3.6	Modern Transmitter Architectures . . . . .	88
3.6.1	Quadrature Modulator . . . . .	88

3.6.2	Quadrature Modulation . . . . .	89
3.6.3	Frequency Modulation . . . . .	90
3.6.4	Polar Modulation . . . . .	90
3.7	Modern Receiver Architectures . . . . .	91
3.7.1	Receiver Architectures . . . . .	91
3.7.2	Homodyne Frequency Conversion . . . . .	93
3.7.3	Heterodyne Frequency Conversion . . . . .	94
3.7.4	Direct Conversion Receiver . . . . .	94
3.7.5	Low-IF Receiver . . . . .	95
3.7.6	Subsampling Analog-to-Digital Conversion . . . . .	95
3.7.7	First IF-to-Baseband Conversion . . . . .	96
3.7.8	Bilateral Double-Conversion Receiver . . . . .	96
3.8	Introduction to Software Defined Radio . . . . .	97
3.9	SDR Quadrature Modulator . . . . .	98
3.9.1	Analog Quadrature Modulator . . . . .	98
3.9.2	Summary . . . . .	102
3.10	Case Study: SDR Transmitter . . . . .	103
3.10.1	Analog Quadrature Modulator . . . . .	103
3.10.2	Single-Sideband Suppressed-Carrier (SSB-SC) Modulation . . . . .	105
3.10.3	Digital Quadrature Modulation . . . . .	107
3.10.4	QAM Digital Modulation . . . . .	112
3.10.5	SDR Transmitter Using QAM Digital Modulation . . . . .	113
3.11	SDR Quadrature Demodulator . . . . .	120
3.12	SDR Receiver . . . . .	121
3.12.1	Demodulation of the I component . . . . .	122
3.12.2	Demodulation of the Q component . . . . .	124
3.13	SDR Summary . . . . .	125
3.14	Summary . . . . .	126
3.15	References . . . . .	126
3.16	Exercises . . . . .	126
<b>4</b>	<b>Antennas and the RF Link . . . . .</b>	<b>129</b>
4.1	Introduction . . . . .	129
4.2	RF Antennas . . . . .	130
4.3	Resonant Antennas . . . . .	131
4.3.1	Radiation from a Current Filament . . . . .	131
4.3.2	Finite-Length Wire Antennas . . . . .	133
4.4	Traveling-Wave Antennas . . . . .	136
4.5	Antenna Parameters . . . . .	137
4.5.1	Radiation Density and Radiation Intensity . . . . .	137
4.5.2	Directivity and Antenna Gain . . . . .	138
4.5.3	Effective Isotropic Radiated Power . . . . .	141
4.5.4	Effective Aperture Size . . . . .	141
4.5.5	Summary . . . . .	142
4.6	The RF Link . . . . .	143
4.6.1	Propagation Path . . . . .	143
4.6.2	Resonant Scattering . . . . .	145
4.6.3	Fading . . . . .	145
4.6.4	Link Loss and Path Loss . . . . .	150
4.6.5	Fresnel Zones . . . . .	153

---

4.6.6	Propagation Model in the Mobile Environment . . . . .	155
4.7	Multipath and Delay Spread . . . . .	156
4.7.1	Delay Spread . . . . .	156
4.7.2	Intersymbol Interference . . . . .	158
4.7.3	Summary . . . . .	160
4.8	Radio Link Interference . . . . .	160
4.8.1	Frequency Reuse Plan . . . . .	160
4.8.2	Summary . . . . .	163
4.9	Antenna array . . . . .	163
4.10	Summary . . . . .	165
4.11	References . . . . .	166
4.12	Exercises . . . . .	167
4.12.1	Exercises By Section . . . . .	172
4.12.2	Answers to Selected Exercises . . . . .	172
<b>5</b>	<b>RF Systems . . . . .</b>	<b>173</b>
5.1	Introduction . . . . .	173
5.2	Broadcast, Simplex, Duplex, Diplex, and Multiplex Operations	174
5.2.1	International Telecommunications Union Definitions .	175
5.2.2	Duplex Versus Diplex . . . . .	177
5.3	Cellular Communications . . . . .	178
5.3.1	Cellular Concept . . . . .	178
5.3.2	Personal Communication Services . . . . .	181
5.3.3	Call Flow and Handoff . . . . .	181
5.3.4	Cochannel Interference . . . . .	182
5.4	Multiple Access Schemes . . . . .	182
5.5	Spectrum Efficiency . . . . .	185
5.6	Processing Gain . . . . .	187
5.6.1	Energy of a Bit . . . . .	187
5.6.2	Coding Gain . . . . .	189
5.6.3	Spreading Gain . . . . .	190
5.6.4	Spreading Gain in Terms of Bandwidth . . . . .	191
5.6.5	Symbol Error Rate and Bit Error Rate . . . . .	192
5.6.6	Summary . . . . .	194
5.7	Early Generations of Cellular Phone Systems . . . . .	196
5.8	Early Generations of Radio . . . . .	197
5.8.1	1G, First Generation: Analog Radio . . . . .	197
5.8.2	2G, Second Generation: Digital Radio . . . . .	199
5.9	3G, Third Generation: Code Division Multiple Acces (CDMA)	201
5.9.1	Generation 2.5: Direct Sequence Code Division Multi- ple Access . . . . .	201
5.9.2	Multipath and Rake Receivers . . . . .	202
5.9.3	3G, Wideband CDMA . . . . .	204
5.9.4	Summary . . . . .	206
5.10	4G, Fourth Generation Radio . . . . .	208
5.10.1	Orthogonal Frequency Division Multiplexing . . . . .	209
5.10.2	Orthogonal Frequency Division Multiple Access . . . .	212
5.10.3	Cyclic Prefix . . . . .	212
5.10.4	FDD versus TDD . . . . .	213
5.10.5	Multiple Input, Multiple Output . . . . .	213
5.10.6	Carrier Aggregation . . . . .	215

- 5.10.7 IEEE 802.11n . . . . . 215
- 5.10.8 OFDM Modulator . . . . . 217
- 5.10.9 Summary of 4G . . . . . 217
- 5.11 5G, Fifth Generation Radio . . . . . 219
  - 5.11.1 Mesh Radio . . . . . 219
  - 5.11.2 Cognitive Radio . . . . . 220
  - 5.11.3 Massive MIMO . . . . . 220
  - 5.11.4 Active Antenna Systems . . . . . 221
  - 5.11.5 Microwave Frequency Operation . . . . . 222
  - 5.11.6 Millimeter-Wave Operation . . . . . 222
  - 5.11.7 Non Orthogonal Multiple Access . . . . . 222
  - 5.11.8 Summary . . . . . 223
- 5.12 6G, Sixth Generation Radio . . . . . 223
- 5.13 Radar Systems . . . . . 224
- 5.14 Summary . . . . . 228
- 5.15 References . . . . . 229
- 5.16 Exercises . . . . . 231
  - 5.16.1 Exercises By Section . . . . . 234
  - 5.16.2 Answers to Selected Exercises . . . . . 234
- 5.A Mathematics of Random Processes . . . . . 235
- Index** . . . . . 239



# Introduction to RF and Microwave Systems

1.1	Introduction .....	1
1.2	RF and Microwave Engineering .....	2
1.3	Communication Over Distance .....	5
1.4	Radio Architecture .....	14
1.5	Conventional Wireless Communications .....	15
1.6	RF Power Calculations .....	17
1.7	Photons and Electromagnetic Waves .....	21
1.8	Summary .....	22
1.9	References .....	24
1.10	Exercises .....	24

## 1.1 Introduction

Radio frequency (RF) systems drive the requirements of microwave and RF circuits, and the capabilities of RF and microwave circuits fuel the evolution of RF systems. This interdependence and the trade-offs required necessitate that the successful RF and microwave designer have an appreciation of systems. Today, communications is the main driver of RF system development, leading to RF technology evolution at an unprecedented pace. Similar relationships exist for national security including radar and sensors used in detection and ranging. Other radio systems have less immediate impact on RF technology but are very important for the smaller number of RF engineers working in the fields of navigation, astronomy, defense, and heating. No longer can many years be put aside for methodical trade-offs of circuit complexity, technology development, and architecture choices at the system level. As relationships have become more intertwined, RF communication, radar, and sensor engineers must develop a broad appreciation of technology, communication principles, and circuit design.

This book is the first volume in a series on microwave and RF design. A central aspect of microwave engineering is distributed effects considered in the second volume of his book series [1]. Here the transmission lines are treated as supporting forward and backward traveling voltage and current waves and these are related to electromagnetic effects. The third volume [2] covers microwave network theory which is the theory that describe power flow and can be used to describe transmission line effects. Topics covered in this volume include scattering parameters, Smith charts, and

**Table 1-1:** Broad electromagnetic spectrum divisions.

Name or band	Frequency	Wavelength
Radio frequency	3 Hz – 300 GHz	100 000 km – 1 mm
Microwave	300 MHz – 300 GHz	1 m – 1 mm
Millimeter (mm) band	110 – 300 GHz	2.7 mm – 1.0 mm
Infrared	300 GHz – 400 THz	1 mm – 750 nm
Far infrared	300 GHz – 20 THz	1 mm – 15 $\mu$ m
Long-wavelength infrared	20 THz – 37.5 THz	15–8 $\mu$ m
Mid-wavelength infrared	37.5 – 100 THz	8–3 $\mu$ m
Short-wavelength infrared	100 THz – 214 THz	3–1.4 $\mu$ m
Near infrared	214 THz – 400 THz	1.4 $\mu$ m – 750 nm
Visible	400 THz – 750 THz	750 – 400 nm
Ultraviolet	750 THz – 30 PHz	400 – 10 nm
X-Ray	30 PHz – 30 EHz	10 – 0.01 nm
Gamma Ray	> 15 EHz	< 0.02 nm

Gigahertz, GHz =  $10^9$  Hz; terahertz, THz =  $10^{12}$  Hz; pentahertz, PHz =  $10^{15}$  Hz; exahertz, EH =  $10^{18}$  Hz.

matching networks that enable maximum power transfer. The fourth volume [3] focuses on designing microwave circuits and systems using modules introducing a large number of different modules. Modules is just another term for a network but the implication is that is is packaged and often available off-the-shelf. Other topics in this chapter that are important in system design using modules are considered including noise, distortion, and dynamic range. Most microwave and RF designers construct systems using modules developed by other engineers who specialize in developing the modules. Examples are filter and amplifier chip modules which once designed can be used in many different systems. Much of microwave design is about maximizing dynamic range, minimizing noise, and minimizing DC power consumption. The fifth volume in this series [4] considers amplifier and oscillator design and develops the skills required to develop modules.

The books in the Microwave and RF Design series are:

- Microwave and RF Design: Radio Systems
- Microwave and RF Design: Transmission Lines
- Microwave and RF Design: Networks
- Microwave and RF Design: Modules
- Microwave and RF Design: Amplifiers and Oscillators

## 1.2 RF and Microwave Engineering

An RF signal is a signal that is coherently generated, radiated by a transmit antenna, propagated through air or space, collected by a receive antenna, and then amplified and information extracted. An RF circuit operates at the same frequency as the RF signal that is transmitted or received. That is, the frequency at which a circuit operates does not define that it is an RF circuit. The RF spectrum is part of the electromagnetic (EM) spectrum exploited by humans for communications. A broad categorization of the EM spectrum is shown in Table 1-1. Today radios operate from 3 Hz to 300 GHz, although the upper end will increase as technology progresses.

Microwaves refers to the frequencies where the size of a circuit or structure is comparable to or greater than the wavelength of the EM signal. The division is arbitrary but if a circuit structure is greater than  $1/20$  of a wavelength, then most engineers would regard the circuit as being a

microwave circuit. For now the microwave frequency range is generally taken as 300 MHz to 300 GHz. At these frequencies distributed effects, sometimes called transmission line effects, must be considered.

One of the key characteristics distinguishing RF signals from infrared and visible light is that an RF signal can be generated with coherent phase, and information can be transmitted in both amplitude and phase variations of the RF signal. Such signals can be easily generated up to 220 GHz. The necessary hardware becomes progressively more expensive as frequency increases. The upper limit of radio frequencies is about 300 GHz today, but the limit is extending slowly above this as technology progresses.

The RF bands are listed in Table 1-2 along with propagation modes and representative applications. The propagation of RF signals in free space follows one or more paths from a transmitter to a receiver at any frequency, with differences being in the size of the antennas needed to transmit and receive signals. The size of the necessary antenna is related to wavelength, with the typical dimensions ranging from a quarter of a wavelength to a few wavelengths if a reflector is used to focus the EM waves. On earth, and dependent on frequency, RF signals propagate through walls, diffract around objects, refract when the dielectric constant of the medium changes, and reflect from buildings and walls. The extent is dependent on frequency. Above ground the propagation at RF is affected by atmospheric loss, by charge layers induced by solar radiation in the upper atmosphere, and by density variations in the air caused by heating as well as the thinning of air with height above the earth. The ionosphere is the uppermost part of the atmosphere, at 60 to 90 km, and is ionized by solar radiation, producing a reflective surface, called the D layer, to radio signals up to 3 GHz. The D layer weakens at night and most radio signals can then pass through this weakened layer. The E layer extends from 90 to 120 km and is ionized by X-rays and extreme ultraviolet radiation, and the ionized regions, which reflect RF signals, form ionized clouds that last only a few hours. The F layer of the ionosphere extends from 200 to 500 km and ionization in this layer is due to extreme ultraviolet radiation. Refraction results from this charged layer rather than reflection, as the charges are widely separated. At night the F layer results in what is called the skywave, which is the refraction of radio waves around the earth. At low frequencies a radio wave penetrates the earth's surface and the wave can become trapped at the interface between two regions of different permittivity, the earth region and the air region. This radio wave is called the **surface wave** or **ground wave**.

Propagating RF signals in air are absorbed by molecules in the atmosphere primarily by molecular resonances such as the bending and stretching of bonds. This bending and stretching converts energy in the EM wave into vibrational energy of the molecules. The transmittance of radio signals versus frequency in dry air at an altitude of 4.2 km is shown in Figure 1-1. The first molecular resonance encountered in dry air as frequency increases is the oxygen resonance, centered at 60 GHz, but below that the absorption in dry air is very small. Attenuation increases with higher water vapor pressure (with a resonance at 22 GHz) and in rain. Within 2 GHz of 60 GHz a signal will not travel far, and this can be used to provide localized communication over a few meters as a local data link. Regions of low attenuation (i.e. high transmittance), are called windows and there are numerous low loss windows.

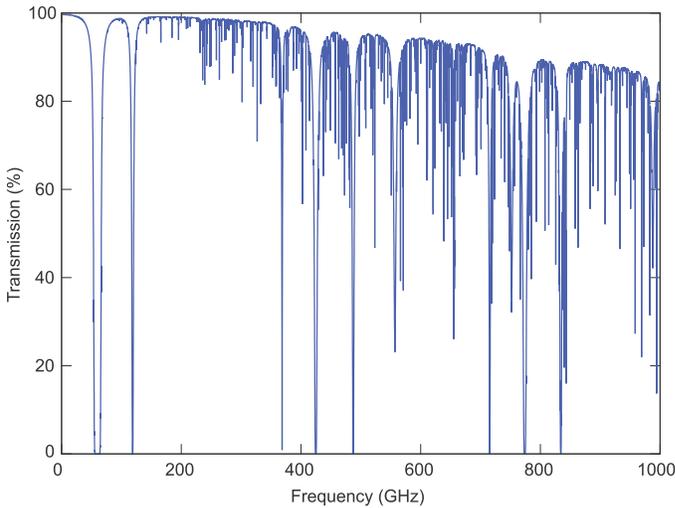
When a radio signal near 60 GHz passes through air an oxygen molecule of two bound oxygen atoms vibrates and EM energy is transferred to the mechanical energy of vibration and thus heat.

**Table 1-2:** Radio frequency bands, primary propagation mechanisms, and selected applications.

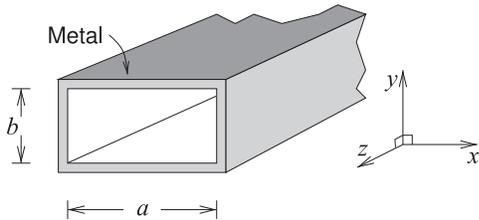
Band		Frequency wavelength	Propagation mode/applications
TLF	Tremendously low frequency	< 3 Hz > 100 000 km	Penetration of liquids and solids/Submarine communication
ELF	Extremely low frequency	3–30 Hz, 100 000– 10 000 km	Penetration of liquids and solids/Submarine communication
SLF	Super low frequency	30–300 Hz 10 000–1 000 km	Penetration of liquids and solids/Submarine communication
ULF	Ultra low frequency	300–3000 Hz 1 000–100 km	Penetration of liquids and solids/Submarine communication; communication within mines
VLF	Very low frequency	3–30 kHz 100–10 km	Guided wave trapped between the earth and the ionosphere/Navigation, geophysics
LF	Low frequency	30–300 kHz 10–1 km	Guided wave between the earth and the ionosphere's D layer; surface waves, building penetration/Navigation, AM broadcast, amateur radio, time signals, RFID
MF	Medium frequency	300–3000 kHz 1000–100 m	Surface wave, building penetration; day time: guided wave between the earth and the ionosphere's D layer; night time: sky wave/AM broadcast
HF	High frequency	3–30 MHz 100–10 m	Sky wave, building penetration/shortwave broadcast, over-the-horizon radar, RFID, amateur radio, marine and mobile telephony
VHF	Very high frequency	30–300 MHz 10–1 m	Line of sight, building penetration; up to 80 MHz, skywave during periods of high sunspot activity/FM and TV broadcast, weather radio, line-of-sight aircraft communications
UHF	Ultra high frequency	300–3000 MHz 10–1 cm	Line of sight, building penetration; sometimes tropospheric ducting/1G–4G cellular communications, RFID, microwave ovens, radio astronomy, satellite-based navigation
SHF	Super high frequency	3–30 GHz 10–1 cm	Line of sight/5G cellular communications, Radio astronomy, point-to-point communications, wireless local area networks, radar
EHF	Extremely high frequency	30–300 GHz 10–1 mm	Line of sight/5G cellular communications, Astronomy, remote sensing, point-to-point and satellite communications
THF	Terahertz or tremendously high frequency	300–3000 GHz 1000–100 $\mu$ m	Line of sight /Spectroscopy, imaging

RF signals diffract and so can bend around structures and penetrate into valleys. The ability to diffract reduces with increasing frequency. However, as frequency increases the size of antennas decreases and the capacity to carry information increases. A very good compromise for mobile communications is at UHF, 300 MHz to 4 GHz, where antennas are of convenient size and there is a good ability to diffract around objects and even penetrate walls. This choice can be seen with 1G–4G cellular communication systems operating in several bands from 450 MHz to 3.6 GHz where antennas do not dominate the size of the handset, and the ability to receive calls within buildings and without line of sight to the base station is well known.

RF bands have been further divided for particular applications. The



**Figure 1-1:** Atmospheric transmission at Mauna Kea, with a height of 4.2 km, on the Island of Hawaii where the atmospheric pressure is 60% of that at sea level and the air is dry with a precipitable water vapor level of 0.001 mm. After [5].



**Figure 1-2:** Rectangular waveguide with internal dimensions of  $a$  and  $b$ . Usually  $a \approx b$ . The EM waves are confined within the four metal walls and propagate in the  $\pm z$  direction. Little current flows in the waveguide walls and so resistive losses are small. Compared to coaxial lines rectangular waveguides have very low loss.

frequency bands for radar are shown in Table 1-3. The L, S, and C bands are referred to as having **octave bandwidths**, as the upper frequency of a band is twice the lower frequency. The other bands are half-octave bands, as the upper frequency limit is approximately 50% higher than the lower frequency limit. The same letter band designations are used by other standards. The most important alternative band designation is for the waveguide bands. These bands refer to the useful range of operation of a rectangular waveguide, which is a rectangular tube that confines a propagating signal within four conducting walls (see Figure 1-2). The waveguide bands are shown in Table 1-4 with the conventional letter designation of bands and standardized waveguide dimensions. Compared to coaxial lines rectangular waveguides have very low loss.<sup>1</sup>

### 1.3 Communication Over Distance

Communicating using EM signals has been an integral part of society since the transmission of the first telegraph signals over wires in the mid 19th century [7]. This development derived from an understanding of magnetic induction based on the experiments of Faraday in 1831 [8] in which he investigated the relationship of magnetic fields and currents. This work of Faraday is now known as Faraday's law, or Faraday's law of induction. It was one of four key laws developed between 1820 and 1835 that described the interaction of static fields and of static fields with currents. These four

<sup>1</sup> A semirigid coaxial line with an outer conductor diameter of 3.5 mm has a loss at 10 GHz of 0.5 dB/m while an X-band waveguide has a loss of 0.1 dB/m. At 100 GHz a 1 mm-diameter coaxial line has a loss of 12.5 dB/m compared to 2.5 dB/m loss for a W-band waveguide.

**Table 1-3:** IEEE radar bands [6]. The mm band designation is also used when the intent is to convey general information above 30 GHz.

Band	Frequency range
L "long"	1–2 GHz
S "short"	2–4 GHz
C "compromise"	4–8 GHz
X "extended"	8–12 GHz
K <sub>u</sub> "kurtz under"	12–18 GHz
K "kurtz" (short in German)	18–27 GHz
K <sub>a</sub> "kurtz above"	27–40 GHz
V	40–75 GHz
W	75–110 GHz
F	90–140 GHz
D	110–170 GHz
mm	110–300 GHz

In Table 1-4 the waveguide dimensions are specified in inches (use 25.4 mm/inch to convert to mm). The number in the WR designation is the long internal dimension of the waveguide in hundredths of an inch. The EIA is the U.S.-based Electronics Industry Association. Note that the radar band (see Table 1-3) and waveguide band designations do not necessarily coincide.

**Table 1-4:** Selected waveguide bands with operating frequencies and internal dimensions (refer to Figure 1-2).

Band	EIA waveguide band	Operating frequency (GHz)	Internal dimensions ( $a \times b$ , inches)
R	WR-430	1.70–2.60	4.300×2.150
D	WR-340	2.20–3.30	3.400×1.700
S	WR-284	2.60–3.95	2.840×1.340
E	WR-229	3.30–4.90	2.290×1.150
G	WR-187	3.95–5.85	1.872×0.872
F	WR-159	4.90–7.05	1.590×0.795
C	WR-137	5.85–8.20	1.372×0.622
H	WR-112	7.05–10.00	1.122×0.497
X	WR-90	8.2–12.4	0.900×0.400
Ku	WR-62	12.4–18.0	0.622×0.311
K	WR-51	15.0–22.0	0.510×0.255
K	WR-42	18.0–26.5	0.420×0.170
Ka	WR-28	26.5–40.0	0.280×0.140
Q	WR-22	33–50	0.224×0.112
U	WR-19	40–60	0.188×0.094
V	WR-15	50–75	0.148×0.074
E	WR-12	60–90	0.122×0.061
W	WR-10	75–110	0.100×0.050
F	WR-8	90–140	0.080×0.040
D	WR-6	110–170	0.0650×0.0325
G	WR-5	140–220	0.0510×0.0255

laws are the Biot–Savart law (developed around 1820), Ampere’s law (1826), Faraday’s law (1831), and Gauss’s law (1835). These are all static laws and do not describe propagating fields.

### 1.3.1 Electromagnetic Fields

We now know that there are two components of the EM field, the **electric field**,  $E$ , with units of volts per meter (V/m), and the **magnetic field**,  $H$ , with units of amperes per meter (A/m).  $E$  and  $H$  fields together describe the force between charges. There are also two flux quantities that are necessary to understand the interactions between these fields and vacuum or matter. The first is  $D$ , the **electric flux density**, with units of coulombs per square meter (C/m<sup>2</sup>), and the other is  $B$ , the **magnetic flux density**, with units of teslas (T).  $B$  and  $H$ , and  $D$  and  $E$ , are related to each other by the properties of the medium, which are embodied in the quantities  $\mu$  and  $\epsilon$  (with the calligraphic letter, e.g.  $\mathcal{B}$ , denoting a time-domain quantity):

$$\vec{D} = \mu \vec{H} \quad (1.1) \quad \vec{D} = \epsilon \vec{E}, \quad (1.2)$$

where the over bar denotes a vector quantity. The quantity  $\mu$  is called the **permeability** of the medium and describes the ability to store **magnetic energy** in a region. The permeability in free space (or vacuum) is denoted  $\mu_0 = 4\pi \times 10^{-7}$  H/m and the magnetic flux and magnetic field are related as

$$\vec{B} = \mu_0 \vec{H}. \quad (1.3)$$

The other material quantity is the **permittivity**,  $\varepsilon$ , which describes the ability to store energy in a volume and in a vacuum

$$\bar{D} = \varepsilon_0 \bar{E}, \quad (1.4)$$

where  $\varepsilon_0 = 8.854 \times 10^{-12}$  F/m is the permittivity of a vacuum. The **relative permittivity**,  $\varepsilon_r$ , is the ratio the permittivity of a material to that of vacuum:

$$\varepsilon_r = \varepsilon / \varepsilon_0. \quad (1.5)$$

Similarly, the **relative permeability**,  $\mu_r$ , refers to the ratio of permeability of a material to its value in a vacuum:

$$\mu_r = \mu / \mu_0. \quad (1.6)$$

### 1.3.2 Biot-Savart Law

The Biot-Savart law relates current to magnetic field as, see Figure 1-3,

$$d\bar{H} = \frac{I d\ell \times \hat{a}_R}{4\pi R^2}, \quad (1.7)$$

which has the units of amperes per meter in the SI system. In Equation (1.7)  $d\bar{H}$  is the incremental static  $H$  field,  $I$  is current,  $d\ell$  is the vector of the length of a filament of current  $I$ ,  $\hat{a}_R$  is the unit vector in the direction from the current filament to the magnetic field, and  $R$  is the distance between the filament and the magnetic field. The  $d\bar{H}$  field is directed at right angles to  $\hat{a}_R$  and the current filament. So Equation (1.7) says that a filament of current produces a magnetic field at a point. The total magnetic field from a current on a wire or surface can be found by modeling the wire or surface as a number of current filaments, and the total magnetic field at a point is obtained by integrating the contributions from each filament.

### 1.3.3 Faraday's Law of Induction

Faraday's law relates a time-varying magnetic field to an induced voltage drop,  $V$ , around a closed path, which is now understood to be  $\oint_{\ell} \bar{E} \cdot d\ell$ , that is, the closed contour integral of the electric field,

$$V = \oint_{\ell} \bar{E} \cdot d\ell = - \oint_s \frac{\partial \bar{B}}{\partial t} \cdot ds, \quad (1.8)$$

and this has the units of volts in the SI unit system. The operation described in Equation (1.8) is illustrated in Figure 1-4.

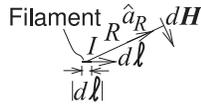
### 1.3.4 Ampere's Circuital Law

Ampere's circuital law, often called just Ampere's law, relates direct current and the static magnetic field  $\mathcal{H}$ . The relationship is based on Figure 1-5 and Ampere's circuital law is

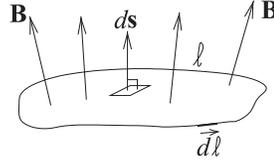
$$\oint_{\ell} \bar{H} \cdot d\ell = I_{\text{enclosed}}. \quad (1.9)$$

That is, the integral of the magnetic field around a loop is equal to the current enclosed by the loop. Using symmetry, the magnitude of the magnetic field at a distance  $r$  from the center of the wire shown in Figure 1-5 is

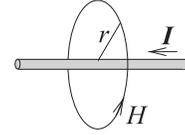
$$H = |I| / (2\pi r). \quad (1.10)$$



**Figure 1-3:** Diagram illustrating the Biot-Savart law. The law relates a static filament of current to the incremental  $H$  field at a distance.

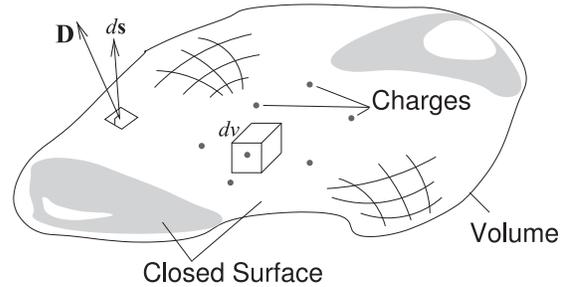


**Figure 1-4:** Diagram illustrating Faraday's law. The contour  $l$  encloses the surface.



**Figure 1-5:** Diagram illustrating Ampere's law. Ampere's law relates the current,  $I$ , on a wire to the magnetic field around it,  $H$ .

**Figure 1-6:** Diagram illustrating Gauss's law. Charges are distributed in the volume enclosed by the closed surface. An incremental area is described by the vector  $dS$ , which is normal to the surface and whose magnitude is the area of the incremental area.



### 1.3.5 Gauss's Law

The final static EM law is Gauss's law, which relates the static electric flux density vector,  $\bar{D}$ , to charge. With reference to Figure 1-6, Gauss's law in integral form is

$$\oint_s \bar{D} \cdot ds = \int_v \rho_v \cdot dv = Q_{\text{enclosed}}. \quad (1.11)$$

This states that the integral of the electric flux vector,  $\bar{D}$ , over a closed surface is equal to the total charge enclosed by the surface,  $Q_{\text{enclosed}}$ .

### 1.3.6 Gauss's Law of Magnetism

Gauss's law of magnetism parallels Gauss's law which now applies to magnetic fields. In integral form the law is

$$\oint_s \bar{B} \cdot ds = 0. \quad (1.12)$$

This states that the integral of the magnetic flux vector,  $\bar{D}$ , over a closed surface is zero reflecting the fact that magnetic charges do not exist.

### 1.3.7 Telegraph

With the static field laws established, the stage was set to begin the development of the transmission of EM signals over wires. While traveling by ship back to the United States from Europe in 1832, Samuel Morse learned of Faraday's experiments and conceived of an EM telegraph. He sought out partners in Leonard Gale, a professor of science at New York University, and Alfred Vail, "skilled in the mechanical arts," who constructed the telegraph models used in their experiments. In 1835 this collaboration led to an experimental version transmitting a signal over 16 km of wire. Morse was not

**Table 1-5:** International Morse code.

Symbol	Code	Symbol	Code	Symbol	Code	Symbol	Code	Symbol	Code
1	.-----	8	---..	E	.	L	.-..	S	...
2	..----	9	----.	F	..--	M	--	T	-
3	...---	0	-----	G	--.	N	-.	U	..-
4	....-	A	.-	H	....	O	---	V	...-
5	.....	B	-...	I	..	P	.-.-	W	.-.-
6	-....	C	-.-.	J	.----	Q	----	X	-.-.
7	--...	D	-.	K	-.-	R	.-.	Y	-.-.
								Z	--..

alone in imagining an EM telegraph, and in 1837 Charles Wheatstone opened the first commercial telegraph line between London and Camden Town, England, a distance of 2.4 km. Subsequently, in 1844, Morse designed and developed a line to connect Washington, DC, and Baltimore, Maryland. This culminated in the first public transmission on May 24, 1844, when Morse sent a telegraph message from the Capitol in Washington to Baltimore. This event is recognized as the birth of communication over distance using wires. This rapid pace of transition from basic research into electromagnetism (Faraday's experiment) to a fielded transmission system has been repeated many times in the evolution of wired and wireless communication technology.

The early telegraph systems used EM induction and multicell batteries that were switched in and out of circuit with the long telegraph wire and so created pulses of current. We now know that these current pulses created propagating magnetic fields that were guided by the wires and were accompanied by electric fields. In 1840 Morse applied for a U.S. patent for "Improvement in the Mode of Communicating Information by Signals by the Application of Electro-Magnetism Telegraph," which described "lightning wires" and "Morse code." By 1854, 37,000 km of telegraph wire crossed the United States, and this had a profound effect on the development of the country. Railroads made early extensive use of telegraph and a new industry was created. In the United States the telegraph industry was dominated by Western Union, which became one of the largest companies in the world. Just as with telegraph, the history of wired and wireless communication has been shaped by politics, business interests, market risk, entrepreneurship, patent ownership, and patent litigation as much as by the technology itself.

The first telegraph signals were just short bursts and slightly longer bursts of noise using Morse code in which sequences of dots, dashes, and pauses represent numbers and letters (see Table 1-5).<sup>2</sup> The speed of transmission was determined by an operator's ability to key and recognize the codes. Information transfer using EM signals in the late 19th century was therefore

<sup>2</sup> Morse code uses sequences of dots, dashes, and spaces. The duration of a dash (or "dah") is three times longer than that of a dot (or "dit"). Between letters there is a small gap. For example, the Morse code for PI is ".-.- .-.-". Between words there is a slightly longer pause and between sentences an even longer pause. Table 1-5 lists the international Morse code adopted in 1848. The original Morse code developed in the 1830s is now known as "American Morse code" or "railroad code." The "modern international Morse code" extends the international Morse code with sequences for non-English letters and special symbols.

about 5 **bits per second (bits/s)**. Morse achieved 10 words per minute.

### 1.3.8 *The Origins of Radio*

In the 1850s Morse began to experiment with wireless transmission, but this was still based on the principle of conduction. He used a flowing river, which as is now known is a medium rich with ions, to carry the charge. On one side of the river he set up a series connection of a metal plate, a battery, a Morse key, and a second metal plate. This formed the transmitter circuit. The metal plates were inserted into the water and separated by a distance considerably greater than the width of the river. On the other side of the river, metal plates were placed directly opposite the transmitter plates and this second set of plates was connected by a wire to a galvanometer in series. This formed the receive circuit, and electric pulses established by the transmitter resulted in the charge being transferred across the river by conduction and the pulses subsequently detected by the galvanometer. This was the first wireless transmission using electromagnetism, but it was not radio.

Morse relied entirely on conduction to achieve wireless transmission and it is now known that we need alternating electric and magnetic fields to propagate information over distance without charge carriers. The next steps in the progress to radio were experiments in induction. These culminated in an experiment by Loomis who in 1866 sent the first aerial wireless signals using kites flown by copper wires [9]. The transmitter kite had a Morse key at the ground end and an electric potential would have been developed between the ground and the kite itself. Closing the key resulted in current flow along the wire and this created a magnetic field that spread out and induced a current in the receive kite and this was detected by a galvanometer. However, not much of an electric field is produced and an EM wave is not transmitted. As such, the range of this system is very limited. Practical wireless communication requires an EM wave at a high-enough frequency that it can be efficiently generated by short wires.

### 1.3.9 *Maxwell's Equations*

The essential next step in the invention of radio was the development of Maxwell's equations in 1861. Before Maxwell's equations were postulated, several static EM laws were known. These are the Biot–Savart law, Ampere's circuital law, Gauss's law, and Faraday's law. Taken together they cannot describe the propagation of EM signals, but they can be derived from Maxwell's equations. Maxwell's equations cannot be derived from the static electric and magnetic field laws. Maxwell's equations embody additional insight relating spatial derivatives to time derivatives, which leads to a description of propagating fields. Maxwell's equations are

$$\nabla \times \bar{\mathcal{E}} = -\frac{\partial \bar{\mathcal{B}}}{\partial t} - \bar{\mathcal{M}} \quad (1.13) \quad \nabla \times \bar{\mathcal{H}} = \frac{\partial \bar{\mathcal{D}}}{\partial t} + \bar{\mathcal{J}} \quad (1.15)$$

$$\nabla \cdot \bar{\mathcal{D}} = \rho_V \quad (1.14) \quad \nabla \cdot \bar{\mathcal{B}} = \rho_m V. \quad (1.16)$$

Several of the quantities in Maxwell's equation have already been introduced, but now the electric and magnetic fields are in vector form. The other quantities in Equations (1.13)–(1.16) are

- $\vec{J}$ , the **electric current** density, with units of amperes per square meter (A/m<sup>2</sup>);
- $\rho_V$ , the **electric charge** density, with units of coulombs per cubic meter (C/m<sup>3</sup>);
- $\rho_{mV}$ , the magnetic charge density, with units of webers per cubic meter (Wb/m<sup>3</sup>); and
- $\vec{M}$ , the magnetic current density, with units of volts per square meter (V/m<sup>2</sup>).

Magnetic charges do not exist, but their introduction through the **magnetic charge** density,  $\rho_{mV}$ , and the **magnetic current** density,  $\vec{M}$ , introduce an aesthetically appealing symmetry to Maxwell's equations. Maxwell's equations are differential equations, and as with most differential equations, their solution is obtained with particular boundary conditions, which in radio engineering are imposed by conductors. Electric conductors (i.e., electric walls) support electric charges and hence electric current. By analogy, magnetic walls support magnetic charges and magnetic currents. Magnetic walls also provide boundary conditions to be used in the solution of Maxwell's equations. The notion of magnetic walls is important in RF and microwave engineering, as they are approximated by the boundary between two dielectrics of different permittivity. The greater the difference in permittivity, the more closely the boundary approximates a magnetic wall.

Maxwell's equations are fundamental properties and there is no underlying theory, so they must be accepted "as is," but they have been verified in countless experiments. Maxwell's equations have three types of derivatives. First, there is the time derivative,  $\partial/\partial t$ . Then there are two spatial derivatives,  $\nabla \times$ , called **curl**, capturing the way a field circulates spatially (or the amount that it curls up on itself), and  $\nabla \cdot$ , called the **div** operator, describing the spreading-out of a field. In rectangular coordinates, curl,  $\nabla \times$ , describes how much a field circles around the  $x$ ,  $y$ , and  $z$  axes. That is, the curl describes how a field circulates on itself. So Equation (1.13) relates the amount an electric field circulates on itself to changes of the  $B$  field in time. So a spatial derivative of electric fields is related to a time derivative of the magnetic field. Also in Equation (1.15) the spatial derivative of the magnetic field is related to the time derivative of the electric field. These are the key elements that result in self-sustaining propagation.

Div,  $\nabla \cdot$ , describes how a field spreads out from a point. So the presence of net electric charge (say, on a conductor) will result in the electric field spreading out from a point (see Equation (1.14)). In contrast, the magnetic field (Equation (1.16)) can never diverge from a point, which is a result of magnetic charges not existing (except when the magnetic wall approximation is used).

How fast a field varies with time,  $\partial \vec{B}/\partial t$  and  $\partial \vec{D}/\partial t$ , depends on frequency. The more interesting property is how fast a field can change spatially,  $\nabla \times \vec{E}$  and  $\nabla \times \vec{H}$ —this depends on wavelength relative to geometry. So if the cross-sectional dimensions of a transmission line are less than a wavelength ( $\lambda/2$  or  $\lambda/4$  in different circumstances), then it will be impossible for the fields to curl up on themselves and so there will be only one solution (with no or minimal spatial variation of the  $E$  and  $H$  fields) or, in some cases, no solution to Maxwell's equations.

### *1.3.10 Transmission of Radio Signals*

Now the discussion returns to the technological development of radio. About the same time as Loomis's induction experiments in 1864, James Maxwell [10] laid the foundations of modern EM theory in 1861 [11]. Maxwell theorized that electric and magnetic fields are different manifestations of the same phenomenon. The revolutionary conclusion was that if they are time varying, then they would travel through space as a wave. This insight was accepted almost immediately by many people and initiated a large number of endeavors. The period of 1875 to 1900 was a time of tremendous innovation in wireless communication.

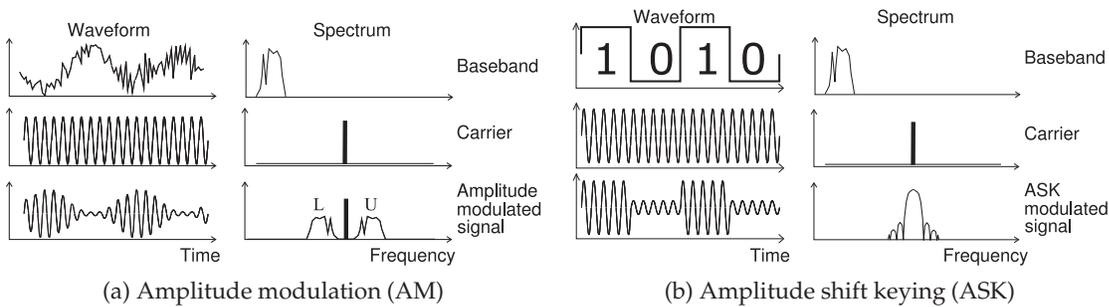
On November 22, 1875, Edison observed EM sparks. Previously sparks were considered to be an induction phenomenon, but Edison thought that he was producing a new kind of force, which he called the etheric force. He believed that this would enable communication without wires. To put this in context, the telegraph was invented in the 1830s and the telephone was invented in 1876.

The next stage leading to radio was orchestrated by D. E. Hughes beginning in 1879. Hughes experimented with a spark gap and reasoned that in the gap there was a rapidly alternating current and not a constant current as others of his time believed. The electric oscillator was born. The spark gap transmitter was augmented with a clockwork mechanism to interrupt the transmitter circuit and produce pulsed radio signals. He used a telephone as a receiver and walked around London and detected the transmitted signals over distance. Hughes noted that he had good reception at 180 feet. Hughes publicly demonstrated his "radio" in 1870 to the Royal Society, but the eminent scientists of the society determined that the effect was simply due to induction. This discouraged Hughes from continuing. However, Hughes has a legitimate claim to having invented radio, mobile digital radio at that, and probably was transmitting pulses on a 100 kHz carrier. In Hugeness's radio the RF carrier was produced by the spark gap oscillator and the information was coded as pulses. It was a small leap to a Morse key-based system.

The invention of practical radio can be attributed to many people, beginning with Heinrich Hertz, who in the period from 1885 to 1889 successfully verified the essential prediction of Maxwell's equations that EM energy could propagate through the atmosphere. Hertz was much more thorough than Hughes and his results were widely accepted. In 1891 Tesla developed what is now called the Tesla coil, which is a transformer with a primary and a secondary coil, one inside the other. When one of the coils was excited by an alternating signal, a large voltage was produced across the terminals of the other coil. Tesla pursued the application of his coils to radio and realized that the coils could be tuned so that the resulting resonance greatly amplified a radio signal.

The next milestone was the establishment of the first practical radio system by Marconi, with experiments beginning in 1894. Oscillations were produced in a spark gap, which were amplified by a Tesla coil. The work culminated in the transmission of telegraph signals across the Atlantic (from Ireland to Canada) by Marconi in 1901. In 1904, crystal radio kits to detect wireless telegraph signals could be readily purchased.

Spark gap transmitters could only send pulses of noise and not voice. One generator that could be amplitude modulated was an alternator. At the end



**Figure 1-7:** Waveform and spectra of simple modulation schemes. The modulating signal, at the top in (a) and (b), is also called the baseband signal.

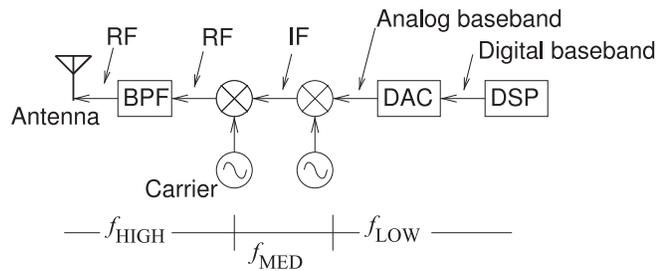
of the 19th century, readily available alternators produced a 60 Hz signal. Reginald Resplendent attempted to make a higher-frequency alternator and the best he achieved operated at 1 kHz. Resplendent realized that Maxwell's equations indicated that radiation increased dramatically with frequency and so he needed a much-higher-frequency signal source. Under contract, General Electric developed a 2-kW, 100-kHz alternator designed by Ernst Alexanderson. With this alternator, the first radio communication of voice occurred on December 23, 1900, in a transmission by Fessenden from an island in the Potomac River, near Washington, DC. Then on December 24, 1906, Fessenden transmitted voice from Massachusetts to ships hundreds of miles away in the Atlantic Ocean. This milestone is regarded as the beginning of the radio era.

Marconi subsequently purchased 50 and 200 kW Alexanderson alternators for his trans-Atlantic transmissions. Marconi was a great integrator of ideas, with particular achievements being the design of transmitting and receiving antennas that could be tuned to a particular frequency and the development of a coherer to improve detection of a signal.

### 1.3.11 Early Radio

Radio works by superimposing relatively slowly varying information, at what is called the **baseband** frequency, on a carrier sinusoid by varying the amplitude and/or phase of the sinusoid. Early radio systems were based on modulating an oscillating carrier either by pulsing the carrier (using for example Morse code)—this modulation scheme is called **amplitude shift keying (ASK)**—or by varying the amplitude of the carrier, i.e. **amplitude modulation (AM)**, in the case of analog, usually voice, transmission. The waveforms and spectra of these modulation schemes are shown in Figure 1-7. The information is contained in the baseband signal, which is also called the modulating signal. The spectrum of the baseband signal extends to DC or perhaps down to where it rolls off at a low frequency. The carrier is a single sinewave and contains no information. The amplitude of the carrier is varied by the baseband signal to produce the modulated signal. In general, there are many cycles of the carrier relative to variations of the baseband signal so that the bandwidth of the modulated signal is relatively small compared to the frequency of the carrier.

**Figure 1-8:** A simple transmitter with low,  $f_{LOW}$ , medium,  $f_{MED}$ , and high frequency,  $f_{HIGH}$ , sections. The mixers can be idealized as multipliers, shown as circles with crosses, that boost the frequency of the input baseband or IF signal by the frequency of the carrier.



AM and ASK radios are narrowband communication systems (they use a small portion of the EM spectrum), so to avoid interference with other radios it is necessary to search for an open part of the spectrum to place the carrier signal. In the decade of the 1900s there was little organization and a listener needed to search to find the desired transmission. The technology of the day necessitated this anyway, as the carrier would drift around by 10% or so since it was then not possible to build a stable oscillator. It was not until the *Titanic* sinking in 1912 that regulation was imposed on the wireless industry. Investigations of the *Titanic* sinking concluded that most of the lives lost would have been saved if a nearby ship had been monitoring its radio channels and if the frequency of the emergency channel was fixed. However, a second ship, but not close enough, did respond to *Titanic*'s "SOS" signal. A result of the investigations was the Service Regulations of the 1912 London International Radiotelegraph Convention. These early regulations were fairly liberal and radio stations were allowed to use radio wavelengths of their own choosing, but restricted to four broad bands: a single band at 1500 kHz for amateurs; 187.5 to 500 kHz, appropriated primarily for government use; below 187.5 kHz for commercial use, and 500 kHz to 1500 kHz, also a commercial band. Subsequent years saw more stringent assignment of narrow spectral bands and the assignment of channels. The standards and regulatory environment for radio were set—there would be assigned frequency bands for particular purposes. Very quickly strong government and commercial interests struggled for exclusive use of particular bands and thus the EM spectrum developed considerable value. Entities "owned" portions of the spectrum either through a license or through government allocation.

While most of the spectrum is allocated, there are several open bands where licenses are not required. The **instrumentation, scientific, and medical (ISM)** bands at 2.4 and 5.8 GHz are examples. Since these bands are loosely regulated, radios must cope with potentially high levels of interference.

## 1.4 Radio Architecture

A radio device is comprised of several key reasonably well-defined units. By frequency there are baseband, intermediate frequency (IF), and RF partitions. In a typical device, the information—either transmitted or received, bits or analog waveforms—is fully contained in the baseband unit. In the case of digital radios, the digital information originating in the baseband digital signal processor (DSP) is converted to an analog waveform typically using a digital-to-analog converter (DAC). This architecture is shown for a simple transmitter in Figure 1-8. When the basic information is analog,

say a voice signal in analog broadcast radios, the information is already a baseband analog waveform. This analog baseband signal can have frequency components that range from DC to many megahertz. However, the baseband signal can range from DC to gigahertz in the case of some radars and point-to-point links that operate at tens of gigahertz.

The RF hardware interfaces the external EM environment with the rest of the communication device. The information that is represented at baseband is translated to a higher-frequency signal that can more easily propagate over the air and for which antennas can be more easily built with manageable sizes. Thus the information content is generally contained in a narrow band of frequencies centered at the carrier frequency. The information content generally occupies a relatively small slice of the EM spectrum. The term “generally” is used as it is not strictly necessary that communication be confined to a narrow band: that is, narrow in percentage terms relative to the RF.

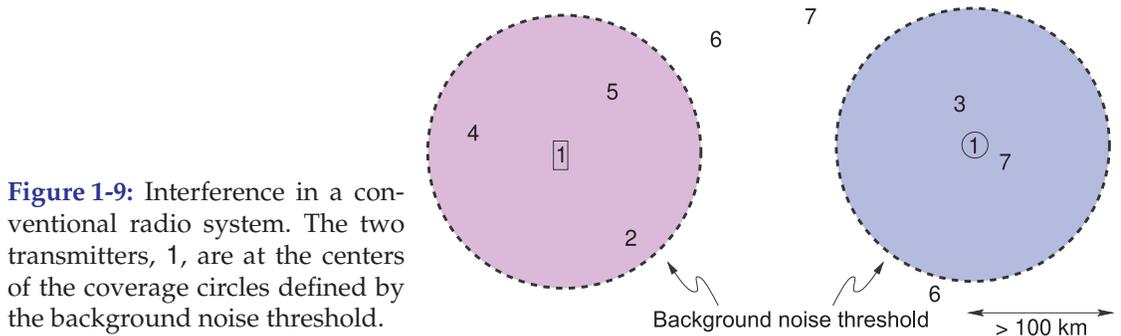
The trade-offs in the choice of carrier frequency are that lower-frequency EM signals require larger antennas, typically one-quarter to one-half wavelength long, but propagate over longer distances and tend to follow the curvature of the earth. AM broadcast radio stations operate around 1 MHz (where the wavelength,  $\lambda$ , is 300 m) using transmit antennas that are 100 m high or more, but good reception is possible at hundreds of kilometers from the transmitter. At higher frequencies, antennas can be smaller, a much larger amount of information can be transmitted with a fixed fractional bandwidth, and there is less congestion. An antenna at 2 GHz (where the free-space wavelength  $\lambda_0 = 15$  cm) is around 4 cm long (and smaller with a dielectric or when folded or coiled), which is a very convenient size for a hand-held communicator.

The concept of an IF is related to the almost universal architecture of transmitters in the 20th century when baseband signals were first translated, or heterodyned, to a band around an IF before a second translation to a higher RF. Initially the IF was just above the audible range and was known as the supersonic frequency. The same progression applies in reverse in a receiver where information carried at RF is first translated to IF before finally being converted to baseband. This architecture resulted in near-optimum noise performance and relatively simple hardware, particularly at RF, where components are much more expensive than at lower frequencies.

The above discussion is a broad description of how radios work. There are many qualifications, as there are many evolving architectures and significant rethinking of the way radios can operate. Architectures and basic properties of radios are trade-offs of the capabilities of technologies, signal processing capability, cost, market dynamics, and politics.

## 1.5 Conventional Wireless Communications

Up until the mid-1970s most wireless communications were based on centralized high-power transmitters, often operating in a wide-area broadcast mode, and reception (e.g., by a television or radio unit) was expected until the signal level fell below a noise-related threshold. These systems are particularly sensitive to interference, therefore systems transmitting at the same frequency were geographically separated so that a transmitted signal falls below the background noise threshold before there



**Figure 1-9:** Interference in a conventional radio system. The two transmitters, 1, are at the centers of the coverage circles defined by the background noise threshold.

is a chance of it interfering with a neighboring system operating at the same frequency. This situation is illustrated in Figure 1-9. Here there are a number of base stations, each operating at a frequency (or set of frequencies) designated by numerals referring to the frequency of operation, which are correspondingly designated as  $f_1$ ,  $f_2$ , etc. In Figure 1-9 the coverage by two base stations, 1, both operating at the frequency  $f_1$  are shown by the shaded regions. The shading indicates the geographical region over which the signals are above the minimum detectable signal threshold. The frequency reuse factor of these types of systems is low, as there is a large geographical area where there is no reception at a particular frequency. The coverage area will not be circular or constant because terrain is not flat, signals are blocked by and reflected from buildings, and background noise levels vary during the day and signal levels vary from season to season as vegetation coverage changes. Allowances must be made in the allocation of broadcast areas to account for the changing coverage level. At the same time, it is necessary, in conventional radio, for the coverage area to be large so that reception, particularly for mobile devices, is continuous over metropolitan-size areas.

The original mobile radio service in the United States is now called **0G** for zero-generation radio. Very few users could be supported in 0G mobile radio because there were very few channels. The first 0G mobile system, the Mobile Telephone Service introduced in 1946, had six channels. That is, only six calls could be made at any one time. Because of interference this was reduced to three channels. So a metropolitan area such as New York city could only support three calls at the same time. In the three-channel version, the channels were 60 kHz wide and with a little more than 60 kHz guard band between channels. More channels were eventually made available. However the maximum practical frequency at the time was 450 MHz and the spectrum from 1 MHz to 500 MHz was highly sought after. Other uses included AM and FM radio, TV broadcast, military communications, and radar. It was seen by regulatory authorities that it was not in the public interest to support more individual users if that meant that broadcast services that catered to many people had to be compromised. There were no cells, just one large coverage area. In every change in radio generation there have been multiple enhancements to improve capacity. So just providing more bandwidth so there could be more channels was not a viable option. Supporting the transition to 1G was more bandwidth, the concept of cells and handoff, narrow channels, and higher operating frequency (900 MHz to 1 GHz). The continued evolution to fifth generation (5G) radio and the concepts that supported it are described in Chapters 2–5.

## 1.6 RF Power Calculations

### 1.6.1 RF Propagation

As an RF signal propagates away from a transmitter the power density reduces conserving the power in the EM wave. In the absence of obstacles and without atmospheric attenuation the total power passing through the surface of a sphere centered on a transmitter is equal to the power transmitted. Since the area of the sphere of radius  $r$  is  $4\pi r^2$ , the power density, e.g. in  $\text{W}/\text{m}^2$ , at a distance  $r$  drops off as  $1/r^2$ . With obstacles the EM wave can diffract, reflect, and follow multiple paths to a receiver where it can combine destructively or constructively. It is the destructive interference that is of concern as this limits the reliable reception of a signal. There is a low probability of perfect cancellation occurring and instead it is found that the power density reduces as  $1/r^n$  where  $n$  ranges from 2 for free space to 5 for a dense urban environment with many obstacles, no line of sight, and multiple signal paths.

#### EXAMPLE 1.1 Signal Propagation

A signal is received at a distance  $r$  from a transmitter and the received power drops off as  $1/r^2$ . When  $r = 1$  km, 100 nW is received. What is  $r$  when the received power is 100 fW?

##### Solution:

The signal collected by the receiver is proportional to the power density of the EM signal. The received signal power  $P_r = k/r^2$  where  $k$  is a constant. This leads to

$$\frac{P_r(1 \text{ km})}{P_r(r)} = \frac{100 \text{ nW}}{100 \text{ fW}} = 10^6 = \frac{kr^2}{k(1 \text{ km})^2} = \frac{r^2}{(10^3 \text{ m})^2}; \quad r = \sqrt{10^{12} \text{ m}^2} = 1000 \text{ km} \quad (1.17)$$

#### EXAMPLE 1.2 Signal Propagation With Obstructions

A transmitter sends a signal to a receiver in a suburban environment that is a distance  $d$  away. When  $d = 5$  km the signal power received is 100 nW. At what distance from the transmitter is the reliably received signal 1 pW if the received signal power falls off as  $1/d^3$ .

##### Solution:

Note that the signal falls off faster than the  $1/d^2$  variation of free space. It is not sufficient to know the total power transmitted and instead the power density at a particular distance must be known. The power reliably received,  $P_R(5 \text{ km})$ , at 5 km is 100 nW and this is the power density,  $P_D(5 \text{ km})$ , multiplied by the effective area,  $A_r$ , of the receive antenna:

$$P_R(5 \text{ km}) = 100 \text{ nW} = P_D(5 \text{ km})A_r = \frac{k}{d^3} = \frac{k}{(5 \text{ km})^3} = \frac{k}{125 \text{ km}^3}.$$

Both  $A_r$  and  $k$  are constants and  $k = 12500 \text{ nW} \cdot \text{km}^3 = 1.25 \cdot 10^{-5} \text{ W} \cdot \text{km}^3$ . The power received at a distance  $d$  is 1 pW when

$$P_R(d) = 1 \text{ pW} = 10^{-12} \text{ W} = \frac{k}{d^3} = \frac{1.25 \cdot 10^{-5} \text{ W} \cdot \text{km}^3}{d^3}$$

$$d^3 = \frac{1.25 \cdot 10^{-5} \text{ W} \cdot \text{km}^3}{10^{-12} \text{ W}} = 1.25 \cdot 10^7 \text{ km}^3; \quad d = \sqrt[3]{1.25 \cdot 10^7} \text{ km} = 232.1 \text{ km}. \quad (1.18)$$

**Table 1-6:** Common logarithm formulas. In engineering  $\log x \equiv \log_{10} x$  and  $\ln x \equiv \log_2 x$ .

Description	Formula	Example
Equivalence	$y = \log_b(x) \leftrightarrow x = b^y$	$\log(1000) = 3$ and $10^3 = 1000$
Product	$\log_b(xy) = \log_b(x) + \log_b(y)$	$\log(0.13 \cdot 978) = \log(0.13) + \log(978)$ $= -0.8861 + 2.990 = 2.104$
Ratio	$\log_b(x/y) = \log_b(x) - \log_b(y)$	$\ln(8/2) = \ln(8) - \ln(2) = 3 - 1 = 2$
Power	$\log_b(x^p) = p \log_b(x)$	$\ln(3^2) = 2 \ln(3) = 2 \cdot 1.0986 = 2.197$
Root	$\log_b(\sqrt[p]{x}) = \frac{1}{p} \log_b(x)$	$\log(\sqrt[3]{20}) = \frac{1}{3} \log(20) = 0.4337$
Change of base	$\log_b(x) = \frac{\log_k(x)}{\log_k(b)}$	$\ln(100) = \frac{\log(100)}{\log(2)} = \frac{2}{0.30103} = 6.644$

### 1.6.2 Logarithm

A cellular phone can reliably receive a signal as small as 100 fW and the signal to be transmitted could be 1 W. So the same circuitry can encounter signals differing in power by a factor of  $10^{13}$ . To handle such a large range of signals a logarithmic scale is used.

Logarithms are used in RF engineering to express the ratio of powers using reasonable numbers. Logarithms are taken with respect to a base  $b$  such that if  $x = b^y$ , then  $y = \log_b(x)$ . In engineering,  $\log(x)$  is the same as  $\log_{10}(x)$ , and  $\ln(x)$  is the same as  $\log_e(x)$  and is called the natural logarithm ( $e = 2.71828 \dots$ ). Unfortunately in physics and mathematics (and in programs such as MATLAB),  $\log x$  means  $\ln x$ , so be careful. Common formulas involving logarithms are given in Table 1-6.

### 1.6.3 Decibels

RF signal levels are usually expressed in terms of the power of a signal. While power can be expressed in absolute terms such as watts (W) or milliwatts (mW), it is much more useful to use a logarithmic scale. The ratio of two power levels  $P$  and  $P_{\text{REF}}$  in bels<sup>3</sup> (B) is

$$P(B) = \log\left(\frac{P}{P_{\text{REF}}}\right), \quad (1.19)$$

where  $P_{\text{REF}}$  is a reference power. Here  $\log x$  is the same as  $\log_{10} x$ . Human senses have a logarithmic response and the minimum resolution tends to be about 0.1 B, so it is most common to use decibels (dB); 1 B = 10 dB. Common designations are shown in Table 1-7. Also, 1 mW = 0 dBm is a very common power level in RF and microwave power circuits where the m in dBm refers to the 1 mW reference. As well, dBW is used, and this is the power ratio with respect to 1 W with 1 W = 0 dBW = 30 dBm.

Working on the decibel scale enables convenient calculations using power numbers ranging from 10s of dBm to  $-110$  dBm to be used rather than numbers ranging from 100 W to 0.000000000000001 W.

<sup>3</sup> Named to honor Alexander Graham Bell, a prolific inventor and major contributor to RF communications.

**Table 1-7:** Common power designations: (a) reference power,  $P_{REF}$ ; (b) power ratio in decibels (dB); and (c) power in dBm and watts.

(a)			(c)	
$P_{REF}$	Bell units	Decibel units	Power	Absolute power
1 W	BW	dBW	-120 dBm	$10^{-12}$ mW = $10^{-15}$ W = 1 fW
$1 \text{ mW} = 10^{-3}$ W	Bm	dBm	0 dBm	1 mW
$1 \text{ fW} = 10^{-15}$ W	Bf	dBf	10 dBm	10 mW
			20 dBm	100 mW = 0.1 W
			30 dBm	1000 mW = 1 W
			40 dBm	$10^4$ mW = 10 W
			50 dBm	$10^5$ mW = 100 W
			-90 dBm	$10^{-9}$ mW = $10^{-12}$ W = 1 pW
			-60 dBm	$10^{-6}$ mW = $10^{-9}$ W = 1 nW
			-30 dBm	0.001 mW = 1 $\mu$ W
			-20 dBm	0.01 mW = 10 $\mu$ W
			-10 dBm	0.1 mW = 100 $\mu$ W

(b)	
Power ratio	in dB
$10^{-6}$	-60
0.001	-30
0.1	-20
1	0
10	10
1000	30
$10^6$	60

**EXAMPLE 1.3**      **Power Gain**

An amplifier has a power gain of 1200. What is the power gain in decibels? If the input power is 5 dBm, what is the output power in dBm?

**Solution:**

Power gain in decibels,  $G_{dB} = 10 \log 1200 = 30.79$  dB.

The output power is  $P_{out|dBm} = P_{dB} + P_{in|dBm} = 30.79 + 5 = 35.79$  dBm.

**EXAMPLE 1.4**      **Gain Calculations**

A signal with a power of 2 mW is applied to the input of an amplifier that increases the power of the signal by a factor of 20.

(a) What is the input power in dBm?

$$P_{in} = 2 \text{ mW} = 10 \cdot \log \left( \frac{2 \text{ mW}}{1 \text{ mW}} \right) = 10 \cdot \log(2) = 3.010 \text{ dBm} \approx 3.0 \text{ dBm}. \quad (1.20)$$

(a) What is the gain,  $G$ , of the amplifier in dB?  
The amplifier gain (by default this is power gain) is

$$G = 20 = 10 \cdot \log(20) \text{ dB} = 10 \cdot 1.301 \text{ dB} = 13.0 \text{ dB}. \quad (1.21)$$

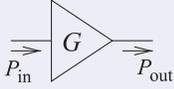
(b) What is the output power of the amplifier?

$$G = \frac{P_{out}}{P_{in}}, \quad \text{and in decibels } G|_{dB} = P_{out|dBm} - P_{in|dBm} \quad (1.22)$$

Thus the output power in dBm is

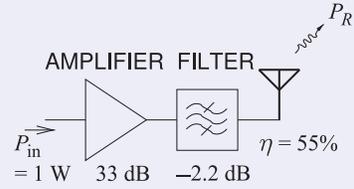
$$P_{out|dBm} = G|_{dB} + P_{in|dBm} = 13.0 \text{ dB} + 3.0 \text{ dBm} = 16.0 \text{ dBm}. \quad (1.23)$$

Note that dB and dBm are dimensionless but they do have meaning; dB indicates a power ratio but dBm refers to a power. Quantities in dB and one quantity in dBm can be added or subtracted to yield dBm, and the difference of two quantities in dBm yields a power ratio in dB.



**EXAMPLE 1.5** Power Calculations

The output stage of an RF front-end consists of an amplifier followed by a filter and then an antenna. The amplifier has a gain of 33 dB, the filter has a loss of 2.2 dB, and of the power input to the antenna, 45% is lost as heat due to resistive losses. If the power input to the amplifier is 1 W, then:



- (a) What is the power input to the amplifier expressed in dBm?  
 $P_{\text{in}} = 1 \text{ W} = 1000 \text{ mW}$ ,  $P_{\text{dBm}} = 10 \log(1000/1) = 30 \text{ dBm}$ .
- (b) Express the loss of the antenna in dB.  
 45% of the power input to the antenna is dissipated as heat.  
 The antenna has an efficiency,  $\eta$ , of 55% and so  $P_2 = 0.55P_1$ .  
 Loss =  $P_1/P_2 = 1/0.55 = 1.818 = 2.60 \text{ dB}$ .

- (c) What is the total gain of the RF front end (amplifier + filter + antenna)?

$$\begin{aligned} \text{Total gain} &= (\text{amplifier gain})_{\text{dB}} + (\text{filter gain})_{\text{dB}} - (\text{loss of antenna})_{\text{dB}} \\ &= (33 - 2.2 - 2.6) \text{ dB} = 28.2 \text{ dB} \end{aligned} \quad (1.24)$$

- (d) What is the total power radiated by the antenna in dBm?

$$\begin{aligned} P_R &= P_{\text{in}} |_{\text{dBm}} + (\text{amplifier gain})_{\text{dB}} + (\text{filter gain})_{\text{dB}} - (\text{loss of antenna})_{\text{dB}} \\ &= 30 \text{ dBm} + (33 - 2.2 - 2.6) \text{ dB} = 58.2 \text{ dBm}. \end{aligned} \quad (1.25)$$

- (e) What is the total power radiated by the antenna?

$$P_R = 10^{58.2/10} = (661 \times 10^3) \text{ mW} = 661 \text{ W}. \quad (1.26)$$

In Examples 1.3 and 1.4 two digits following the decimal point were used for the output power expressed in dBm. This corresponds to an implied accuracy of about 0.01% or 4 significant digits of the absolute number. This level of precision is typical for the result of an engineering calculation. See Section 2.A.1 of [1] for further discussion of precision and accuracy.

### 1.6.4 Decibels and Voltage Gain

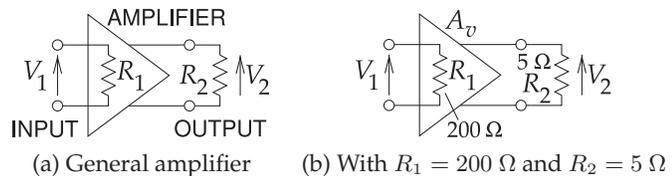
Figure 1-10(a) is an amplifier with input and output resistances that could be different. If  $A_v$  is the voltage gain of the RF amplifier, then

$$V_2 = A_v V_1 \quad (1.27)$$

and the input and output powers will be

$$P_{\text{in}} = \frac{V_1^2}{2R_1} \quad \text{and} \quad P_{\text{out}} = \frac{V_2^2}{2R_2}. \quad (1.28)$$

**Figure 1-10:** Amplifiers each with an input resistance  $R_1$  and output resistance  $R_2$ .



The '2' in the denominator arises because  $V_1$  and  $V_2$  are peak amplitudes of sinusoids in RF engineering. Thus the power gain is

$$G = \frac{P_{\text{out}}}{P_{\text{in}}} = \frac{V_2^2 2R_1}{V_1^2 2R_2} = \frac{R_1}{R_2} A_v^2. \quad (1.29)$$

The power gain depends on the input and output resistance ratio of the amplifiers and this is commonly used to realize significant power gain even if the voltage gain is quite small. If the input and output resistances of the amplifier are the same, then the power gain is just the voltage gain squared.

In handling this situation some authors have used the unit dBV (decibel as a voltage ratio). This should not be used, decibels should always refer to a power ratio, and it is needlessly confusing to use dBV in RF engineering.

#### EXAMPLE 1.6 Voltage Gain to Power Gain

Figure 1-10(b) is a differential amplifier with a  $200 \Omega$  input resistance and  $5 \Omega$  output resistance. If the voltage  $A_v$  is 0.6, what is the power gain of the amplifier in dB?

##### Solution:

The input and output powers are

$$P_{\text{in}} = \frac{1}{2} V_1^2 / R_1 \quad \text{and} \quad P_{\text{out}} = \frac{1}{2} V_2^2 / R_2 = \frac{1}{2} \frac{(A_v V_1)^2}{R_2}. \quad (1.30)$$

Thus the power gain is

$$G = \frac{P_{\text{out}}}{P_{\text{in}}} = \frac{(A_v V_1)^2}{R_2} \left( \frac{V_1^2}{R_1} \right)^{-1} = \frac{R_1}{R_2} A_v^2 = \frac{200}{5} 0.6^2 = 14.4 = 11.58 \text{ dB}. \quad (1.31)$$

The surprising result is that even with a voltage gain of less than 1, a significant power gain can be obtained if the input and output resistances are different. A result used in many RF amplifiers.

## 1.7 Photons and Electromagnetic Waves

Most of the time it is not necessary to consider the quantum nature of radio waves. Considering electric and magnetic fields, and Maxwell's equations, is sufficient to understand radio systems. Still the underlying physics is that currents in circuits derive from the movement of quantum particles, here electrons, and propagating EM signals are transported by photons, also quantum particles.

In a receiver system, information is translated by a metallic antenna from the propagating photons to moving charges in a conductor. This transfer of information is through the interaction of a photon with charge carriers. So the information transfer is based on the quantized energy of a photon and it is possible that quantum effects could be important. The relative level of the photon energy and the random kinetic energy of an electron is important in determining if quantum effects must be considered.

With all EM radiation, information is conveyed by photons and the energy of a photon is conventionally expressed in terms of electron volts (eV). An electron volt is the energy gained by an electron when it moves across an electric potential of 1 V and  $1 \text{ eV} = 1.602 \times 10^{-19} \text{ J}$  ( $1 \text{ J} = 6.241509 \times 10^{18} \text{ eV}$ ).

Physically  $kT$  is the amount of energy required to increase the entropy (corresponding to movement) of the electrons by a factor of  $e$  [12].

The energy of a photon is  $E = h\nu = hf$ , where  $h = 6.6260693 \times 10^{-34} \text{ J}\cdot\text{s}$  is the Planck constant and  $\nu$  and  $f$  are frequency, with the symbol  $\nu$  preferred by physicists and  $f$  preferred by engineers. So the energy of a photon is proportional to its frequency. At microwave frequencies the energy of a microwave photon ranges from  $1.24 \mu\text{eV}$  at 300 MHz to  $1.24 \text{ meV}$  at 300 GHz. This is a very small amount of energy.

The thermal energy of an electron in joules is  $E|_j = kT$  where  $k = 1.3806505 \times 10^{-23} \text{ J/K}$  is the Boltzmann constant and  $T$  is the temperature in kelvin. This applies to an electron that is moving in a group of electrons, say in a plasma, in thermodynamic equilibrium. Freely conducting electrons in a conductor are in a plasma. The thermal energy of such an electron is its random kinetic energy with  $kT = \frac{1}{2}mv^2$  where  $m$  is the mass of the electron and  $v$  is its velocity. (The thermal energy of an isolated electron, i.e. in a non-interacting electron gas, is  $\frac{3}{2}kT$ , but that is not the situation with microwave circuits.) At room temperature ( $T = 298 \text{ K}$ ) the thermal energy of a conducting electron in a metal is  $kT = 25.7 \text{ meV}$ . This is much more than the energy of a microwave photon ( $1.24 \mu\text{eV}$  to  $1.24 \text{ meV}$ ). Thus at room temperature discrete quantum effects are not apparent for microwave signals and so microwave radiation (at room temperature) can be treated as a continuum effect.

A photon has a dual nature, as a particle and as an EM wave. The break point as to which nature helps the most in understanding behavior depends on the energy of the photon versus the thermal energy of an electron. When a microwave photon is captured in a metal at room temperature it makes little sense to talk about the photon as increasing the energy state of an individual electron. Instead, it is best to think of the photon as an EM wave with an electric field that accelerates an ensemble of free electrons (in the conduction band). Thus the energy of one photon is transferred to a group of electrons as faster moving electrons but with the energy increase being so small relative to thermal energy that a quantized effect is not apparent. However even a very low power microwave signal has an enormous number of photons (a  $1 \text{ pW}$  300 GHz signal has  $5 \times 10^9$  photons per second) and each photon accelerates the free electrons a little bit with the combined effect being that the electric field of the microwave signal results in appreciable current. This is the view that we use with room temperature circuits and antennas at microwave frequencies; the EM signal (instead of discrete photons) interacts with the free electrons in a conductor and produces current. Quantum effects must be considered when temperature is very low (say below  $4 \text{ K}$ ) or frequency is very high, e.g. the photon energies for red light ( $400 \text{ THz}$ ) is  $1.7 \text{ eV}$ .

## 1.8 Summary

Today six billion individuals regularly transmit information wirelessly using transmitters and receivers that can be smaller than an infants hand, contain trillions of transistors, and provide connectivity in nearly any location.

The history of radio communications is remarkable and nearly every aspect of electrical engineering is involved. The most important factor of all is that consumers are prepared to part with a large portion of their income to have untethered connectivity. This overwhelming desire for ubiquitous communication surprised even the most optimistic proponents of personal wireless communications. Wireless communication is now a significant part

of every country's economy, and governments are very involved in setting standards and protecting the competitiveness of their own industries.

The history of radio communications has led to the current mode of operation, allocating a narrow slice of the EM spectrum to one or a few users. This choice dictated the need for very stable oscillators and high-rejection filters, for example. The stability of oscillators in consumer products is a few hertz around 1 GHz, a few parts in a billion, a level of precision that is unrivaled for any other physical quantity in a manufactured device.

The RF spectrum is used to support a tremendous range of applications, including voice and data communications, satellite-based navigation, radar, weather radar, mapping, environmental monitoring, air traffic control, police radar, perimeter surveillance, automobile collision avoidance, and many military applications. A big trend is the virtual disappearance of analog radio and now the almost complete use of digital radio. Digital radio is much more tolerant to interference and can use a much smaller slice of the EM spectrum. Another big trend is the tremendous demand for EM spectrum resulting in appreciable use of the spectrum up to 100 GHz and soon beyond that. Currently large parts of the spectrum are allocated for exclusive military use and there is pressure to reduce the spectrum allocated for government use of all kinds.

In RF and microwave engineering there are always considerable approximations made in design, partly because of necessary simplifications that must be made in modeling, but also because many of the material properties required in a detailed design can only be approximately known. Most RF and microwave design deals with frequency-selective circuits often relying on line lengths that have a length that is a particular fraction of a wavelength. Many designs can require frequency tolerances of as little as 0.1%, and filters can require even tighter tolerances. It is therefore impossible to design exactly. Measurements are required to validate and iterate designs. Conceptual understanding is essential; the designer must be able to relate measurements, which themselves have errors, with computer simulations. The ability to design circuits with good tolerance to manufacturing variations and perhaps circuits that can be tuned by automatic equipment are skills developed by experienced designers.

Chapter 2 describes modulation methods and the ideas that led to being able to transmit many bits of data per hertz of bandwidth. High orders of modulation send many bits of information and the higher the order of modulation the more sophisticated the modulation and demodulation schemes must be. Transmitters and receivers that implement the modulation and demodulation methods and up-convert and down-convert, respectively, between the low frequency baseband information and the high frequency radio signals are described in Chapter 3. Antennas reviewed in Chapter 4 are the interface between electronic circuits and freely propagating EM waves. Directional antennas are essential to supporting many users by enabling frequency reuse of the EM spectrum in the same cell. The 1G through to 5G cellular systems are described in Chapter 5. In addition other microwave systems such as radar and WiFi are briefly described and their development has closely followed or just preceded that of cellular radio. The tremendous advances are the result of the synergistic development of system concepts and of digital and microwave hardware.

## 1.9 References

- [1] M. Steer, *Microwave and RF Design, Transmission Lines*, 3rd ed. North Carolina State University, 2019.
- [2] —, *Microwave and RF Design, Networks*, 3rd ed. North Carolina State University, 2019.
- [3] —, *Microwave and RF Design, Modules*, 3rd ed. North Carolina State University, 2019.
- [4] —, *Microwave and RF Design, Amplifiers and Oscillators*, 3rd ed. North Carolina State University, 2019.
- [5] "Atmospheric microwave transmittance at mauna kea, wikipedia creative commons."
- [6] "IEEE standard 521-2002, *IEEE Standard Letter Designations for Radar-Frequency Bands*, 2002."
- [7] T. K. Sarkar, R. Mailloux, A. A. Oliner, M. Salazar-Palma, and D. L. Sengupta, *History of wireless*. John Wiley & Sons, 2006.
- [8] "IEEE Virtual Museum," at <http://www.ieee-virtual-museum.org> Search term: 'Faraday'.
- [9] M. Loomis, "Improvement in telegraphing," 1872, US Patent 129,971.
- [10] J. Rautio, "Maxwell's legacy," *IEEE Microwave Magazine*, vol. 6, no. 2, pp. 46–53, Jun. 2005.
- [11] J. Maxwell, *A Treatise on Electricity and Magnetism*. Clarendon Press (Reprinted, Oxford University Press, 1998), 1873.
- [12] P. Atkins and J. DePaula, *Physical Chemistry*, 9th ed. WH Freeman & Company, 2009.

## 1.10 Exercises

1. Consider a photon at 1 GHz.
  - (a) What is the energy of the photon in joules?
  - (b) Is this more or less than the random kinetic energy of an electron at room temperature?
2. Consider a photon at various frequencies.
  - (a) What is the photon's energy at 1 GHz in terms of electron-volts?
  - (b) What is the photon's energy at 10 GHz in terms of electron-volts?
  - (c) What is the photon's energy at 100 GHz in terms of electron-volts?
  - (d) What is the photon's energy at 1 THz in terms of electron-volts?
3. Consider a photon at 1 THz.
  - (a) What is the energy of the photon in terms of electron-volts?
  - (b) What is the energy of the photon in joules?
  - (c) Is this more or less than the random kinetic energy of an electron at room temperature (300 K)?
  - (d) Discuss if it is necessary to consider quantum effects of the 1 THz photon at room temperature.
4. Consider a photon at 10 GHz.
  - (a) What is the energy of the photon in terms of electron-volts?
  - (b) What is the energy of the photon in joules?
  - (c) What is the random kinetic energy of an electron at room temperature (300 K)?
  - (d) Calculate the temperature, in kelvins, at which the random kinetic energy of an electron is equal to the energy you calculated in (a).
5. A 10 GHz transmitter transmits a 1 W signal. How many photons are transmitted?
6. A receiver receives a 1 pW signal at 60 GHz. How many photons per second are received?
7. At what frequency is the photon energy equal to the thermal energy of an electron at 300 K?
8. What is the frequency at which the energy of a photon is equal to the thermal energy of an electron at 77 K?
9. What is the wavelength in free space of a signal at 4.5 GHz?
10. Consider a monopole antenna that is a quarter of a wavelength long. How long is the antenna if it operates at 3 kHz?
11. Consider a monopole antenna that is a quarter of a wavelength long. How long is the antenna if it operates at 500 MHz?
12. Consider a monopole antenna that is a quarter of a wavelength long. How long is the antenna if it operates at 2 GHz?
13. A dipole antenna is half of a wavelength long. How long is the antenna at 2 GHz?
14. A dipole antenna is half of a wavelength long. How long is the antenna at 1 THz?
15. Write your family name in Morse code (see Table 1-5).
16. A transmitter transmits an FM signal with a bandwidth of 100 kHz and the signal is received by a receiver at a distance  $r$  from the transmitter. When  $r = 1$  km the signal power received by the receiver is 100 nW. When the receiver moves

- further away from the transmitter the power received drops off as  $1/r^2$ . What is  $r$  in kilometers when the received power is 100 pW. [Parallels Example 1.1]
17. A transmitter transmits an AM signal with a bandwidth of 20 kHz and the signal is received by a receiver at a distance  $r$  from the transmitter. When  $r = 10$  km the signal power received is 10 nW. When the receiver moves further away from the transmitter the power received drops off as  $1/r^2$ . What is  $r$  in kilometers when the received power is equal to the received noise power of 1 pW? [Parallels Example 1.1]
  18. In a legacy, i.e. 0G, broadcast radio system a transmitter broadcasts an AM signal and the signal can be successfully received if the AM signal is 20 dB higher than the 10 fW noise power received. The received signal power when the transmitter and receiver are separated by  $r = 1$  km is 100 nW. The received signal power falls off as  $1/r^2$  as the receiver moves further away.
    - (a) What is the radius of the broadcast circle in which the broadcast signal is successfully received?
    - (b) At what distance does the power of the broadcast signal match the noise power?
    - (c) If two transmitters both transmit similar AM signals at the same frequency, how far should the transmitters be separated so that the interference received is 10 dB below the noise level?
  19. In a legacy radio system a transmitter broadcasts an FM signal and for noise-free reception the FM signal must be 30 dB higher than the received noise power of 10 fW. When the transmitter and receiver are separated by  $r = 1$  km the signal power received is 100 nW. The received signal power falls off as  $1/r^3$  with greater separation.
    - (a) What is the radius of the circle in which the broadcast signal is successfully received?
    - (b) At what distance does the power of the broadcast signal match the noise power?
    - (c) If two transmitters both transmit similar FM signals at the same frequency and power. One transmitter transmits the desired signal while the second transmits an interfering signal. How far should the transmitters be separated so that the interference received is 10 dB below the noise level?
  20. A transmitter broadcasts a signal to a receiver that is a distance  $d$  away. The noise power received is 1 pW and when  $d = 5$  km the signal power received is 100 nW. What is the radius of the noise threshold circle where the noise and signal powers are equal, when the received signal power falls off as:
    - (a)  $1/d^2$ ? [Parallels Example 1.1]
    - (b)  $1/d^{2.5}$ ? [Parallels Example 1.2]
  21. A signal is transmitted to a receiver that is a distance  $r$  away. The noise power received is 100 fW and when  $r$  is 1 km the received signal power is 500 nW. What is  $r$  when the noise and signal powers are equal when the received signal power falls off as:
    - (a)  $1/d^2$ ? [Parallels Example 1.1]
    - (b)  $1/d^3$ ? [Parallels Example 1.2]
  22. The logarithm to base 2 of a number  $x$  is 0.38 (i.e.,  $\log_2(x) = 0.38$ ). What is  $x$ ?
  23. The natural logarithm of a number  $x$  is 2.5 (i.e.,  $\ln(x) = 2.5$ ). What is  $x$ ?
  24. The logarithm to base 2 of a number  $x$  is 3 (i.e.,  $\log_2(x) = 3$ ). What is  $\log_2(\sqrt[3]{x})$ ?
  25. What is  $\log_3(10)$ ?
  26. What is  $\log_{4.5}(2)$ ?
  27. Without using a calculator evaluate  $\log\{[\log_3(3x) - \log_3(x)]\}$ .
  28. A 50  $\Omega$  resistor has a sinusoidal voltage across it with a peak voltage of 0.1 V. The RF voltage is  $0.1 \cos(\omega t)$ , where  $\omega$  is the radian frequency of the signal and  $t$  is time.
    - (a) What is the power dissipated in the resistor in watts?
    - (b) What is the power dissipated in the resistor in dBm?
  29. The power of an RF signal is 10 mW. What is the power of the signal in dBm?
  30. The power of an RF signal is 40 dBm. What is the power of the signal in watts?
  31. An amplifier has a power gain of 2100.
    - (a) What is the power gain in decibels?
    - (b) If the input power is  $-5$  dBm, what is the output power in dBm? [Parallels Example 1.3]
  32. An amplifier has a power gain of 6. What is the power gain in decibels? [Parallels Example 1.3]
  33. A filter has a loss factor of 100. [Parallels Example 1.3]
    - (a) What is the loss in decibels?
    - (b) What is the gain in decibels?
  34. An amplifier has a power gain of 1000. What is the power gain in dB? [Parallels Example 1.3]
  35. An amplifier has a gain of 14 dB. The input to the amplifier is a 1 mW signal, what is the output power in dBm?

36. An RF transmitter consists of an amplifier with a gain of 20 dB, a filter with a loss of 3 dB and then that is then followed by a lossless transmit antenna. If the power input to the amplifier is 1 mW, what is the total power radiated by the antenna in dBm? [Parallels Example 1.5]
37. The final stage of an RF transmitter consists of an amplifier with a gain of 30 dB and a filter with a loss of 2 dB that is then followed by a transmit antenna that loses half of the RF power as heat. [Parallels Example 1.5]
- (a) If the power input to the amplifier is 10 mW, what is the total power radiated by the antenna in dBm?
- (b) What is the radiated power in watts?
38. A 5 mW RF signal is applied to an amplifier that increases the power of the RF signal by a factor of 200. The amplifier is followed by a filter that loses half of the power as heat.
- (a) What is the output power of the filter in watts?
- (b) What is the output power of the filter in dBW?
39. The power of an RF signal at the output of a receive amplifier is 1  $\mu$ W and the noise power at the output is 1 nW. What is the output signal-to-noise ratio in dB?
40. The power of a received signal is 1 pW and the received noise power is 200 fW. In addition the level of the interfering signal is 100 fW. What is the signal-to-noise ratio in dB? Treat interference as if it is an additional noise signal. age gain of 1 has an input impedance of 100  $\Omega$ , a zero output impedance, and drives a 5  $\Omega$  load. What is the power gain of the amplifier?
41. A transmitter transmits an FM signal with a bandwidth of 100 kHz and the signal power received by a receiver is 100 nW. In the same bandwidth as that of the signal the receiver receives 100 pW of noise power. In decibels, what is the ratio of the signal power to the noise power, i.e. the signal-to-noise ratio (SNR), received?
43. An amplifier with a voltage gain of 20 has an input resistance of 100  $\Omega$  and an output resistance of 50  $\Omega$ . What is the power gain of the amplifier in decibels? [Parallels Example 1.6]
44. An amplifier with a voltage gain of 1 has an input resistance of 100  $\Omega$  and an output resistance of 5  $\Omega$ . What is the power gain of the amplifier in decibels? Explain why there is a power gain of more than 1 even though the voltage gain is 1. [Parallels Example 1.6]
45. An amplifier with a volt
46. An amplifier has a power gain of 1900.
- (a) What is the power gain in decibels?
- (b) If the input power is  $-8$  dBm, what is the output power in dBm? [Parallels Example 1.3]
47. An amplifier has a power gain of 20.
- (a) What is the power gain in decibels?
- (b) If the input power is  $-23$  dBm, what is the output power in dBm? [Parallels Example 1.3]
48. An amplifier has a voltage gain of 10 and a current gain of 100.
- (a) What is the power gain as a number?
- (b) What is the power gain in decibels?
- (c) If the input power is  $-30$  dBm, what is the output power in dBm?
- (c) What is the output power in mW?
49. An amplifier with 50  $\Omega$  input impedance and 50  $\Omega$  load impedance has a voltage gain of 100. What is the (power) gain in decibels?
50. An attenuator reduces the power level of a signal by 75%. What is the (power) gain of the attenuator in decibels?

### 1.10.1 Exercises By Section

†challenging, ‡very challenging

- |      |   |      |   |            |
|------|---|------|---|------------|
| §1.2 | 1, 2 <sup>†</sup> , 3 <sup>†</sup> , 4 <sup>†</sup> , 5, 6, 7 <sup>†</sup> , 8 <sup>†</sup> | §1.6 | 22, 23, 24, 25, 26, 27, 28, 29, 30  | 48, 49, 50 |
| §1.3 | 9, 10, 11, 12, 13, 14, 15   |      | 31 <sup>†</sup> , 32, 33, 34, 35, 36 <sup>†</sup> , 37 <sup>†</sup> , 38 <sup>†</sup> |            |
| §1.5 | 16, 17, 18 <sup>‡</sup> , 19 <sup>†</sup> , 20 <sup>†</sup> , 21 <sup>†</sup>               |      | 39, 40, 41, 42, 43, 44, 45, 46, 47  |            |

### 1.10.2 Answers to Selected Exercises

- |       |           |    |          |       |          |
|-------|-----------|----|----------|-------|----------|
| 2(d)  | 41.36 meV | 25 | 2.096    | 36    | 50.12 mW |
| 4(b)  | 662.6 fJ  | 29 | 10 dBm   | 37(b) | 3.162 W  |
| 11.12 | 3.25 cm   | 30 | 10 W     | 42(b) | $-6$ dB  |
| 35    | 1.301     | 32 | 7.782 dB |       |          |

# Modulation

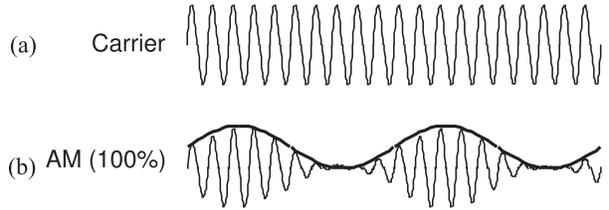
2.1	Introduction .....	27
2.2	Radio Signal Metrics .....	28
2.3	Modulation Overview .....	36
2.4	Analog Modulation .....	37
2.5	Digital Modulation .....	44
2.6	Frequency Shift Keying, FSK .....	47
2.7	Carrier Recovery .....	50
2.8	Phase Shift Keying Modulation .....	50
2.9	Quadrature Amplitude Modulation .....	64
2.10	Digital Modulation Summary .....	65
2.11	Interference and Distortion .....	66
2.12	Summary .....	72
2.13	References .....	73
2.14	Exercises .....	74

## 2.1 Introduction

Most radio communication systems superimpose slowly varying information on a sinusoidal carrier that is transmitted as a radio frequency (RF) signal. This modulated RF signal is sent through a medium, usually air, by a transmitter to a receiver. In the transmitter information is initially represented at what is called baseband. The process of transferring information from baseband to the much higher frequency carrier wave is called modulation. Most modulation schemes slowly vary the amplitude and/or phase of a sinusoidal carrier waveform. In the receiver the process is reversed using demodulation to extract the baseband information from the varying state, such as the amplitude and/or phase, of the modulated carrier.

Radio has evolved subject to constraints imposed by political, hardware, and compatibility considerations. New schemes generally must be compatible and co-exist with earlier schemes. This chapter discusses the many different modulation schemes that are used in radios. Nearly all modulation schemes are supported in modern radios such as 4G and 5G cellular radios, and many are supported in WiFi. Sometimes this is to provide support for legacy radios while in other situations they are used because simpler modulation formats tolerate higher levels of interference. Indeed the level of so-

**Figure 2-1:** AM showing the relationship between the carrier and modulation envelope: (a) carrier; and (b) 100% amplitude modulated carrier.



phistication of modulation methods may need to be frequently changed to accommodate varying interference environments. Legacy analog modulation schemes and the simpler digital modulation schemes were suitable for the relatively unsophisticated hardware of years past. High-order modulation schemes enable many digital bits to be sent in each hertz of bandwidth and are only possible because of the evolution of digital signal processing and because of advances in high-density, low-power digital electronics.

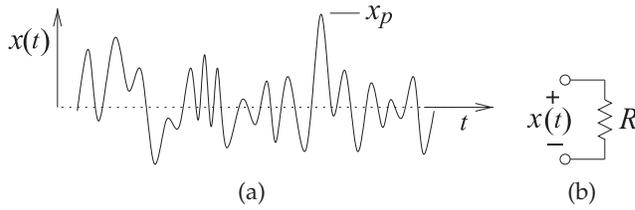
Section 2.2 introduces some of the metrics that are used to compare modulation schemes and Section 2.3 introduces modulation. Section 2.4 describes analog modulation. Then Section 2.5 describes digital modulation followed by sections that deal with the specifics of various digital modulation methods: frequency shift keying (FSK) in Section 2.6; phase shift keying (PSK) in Section 2.8; and quadrature amplitude modulation (QAM) in Section 2.9. Before the discussion of PSK a concept called carrier recovery is discussed in Section 2.7 as the necessity to do this was behind the development of a variety of PSK modulation schemes. This is followed by a discussion of the metrics that can be used to quantify interference and distortion of modulated signals.

Modulation, and the hardware architectures and circuits for modulating and demodulating radio signals, are presented largely in three chapters. There is an overlap of these topics but modulation itself is largely confined to this chapter although some architecture concepts must necessarily be introduced to understand the evolution of modulation schemes. The next chapter, Chapter 3, focuses on architectures and essential circuits for modulators and demodulators.

## 2.2 Radio Signal Metrics

Radio signals are engineered to trade-off efficient use of the EM spectrum with the complexity and performance of the required RF hardware. Ultimately the goal is to efficiently use spectrum through maximal packing of information, e.g. digital bits, in a given bandwidth while, for mobile radios especially, using as little prime power as possible. The choice of the type of modulation to use is at the core of the communication system design trade-off.

There are two families of modulation methods: analog and digital modulation. In analog modulation the RF signal has a continuous range of values; in digital modulation, the output has a number of discrete states at particular times called clock ticks, say every microsecond. There are just a few modulation schemes, all of which are digital, that achieve the optimum trade-offs of spectral efficiency and ease of use with acceptable hardware complexity. If hardware complexity is not a concern, which modulation scheme is used depends on noise and interference as well as the power required to transmit a signal, and the power required to process a received



**Figure 2-2:** Definition of crest factor: (a) arbitrary waveform; and (b) voltage across a resistor.

signal.

This section introduces several metrics that characterize the variability of the amplitude of a modulated signal, and this variability has a direct impact on how analog hardware performs are designed and how efficiently hardware can be used.

### 2.2.1 Crest Factor and Peak-to-Average Power Ratio

#### Introduction

In radio engineering crest factor (CF) is a metric that describes how the voltage of a modulated carrier signal varies with time, and peak-to-average power ratio (PAPR) describes how the instantaneous power of a carrier signal varies with time. Be aware that there is one metric, **peak-to-average ratio (PAR)**, that is defined differently in the power, communications theory, and microwave communities. In some communities CF is also called the **peak-to-average ratio (PAR)**. This can leads to problems. Consider, for example, the community that works on smart power metering which combines power measurement, communications theory, and microwave design. The solution to this inevitable confusion is to skip the use of PAR and use unambiguous metrics.

In standards PAR is defined as the ratio of the instantaneous peak value of a signal parameter to its time-averaged value. PAR is used with many signal parameters, e.g. voltage, current, power, and frequency [1].

#### Crest Factor

CF is the ratio of the maximum signal, such as a voltage, to its root-mean-square (rms) value. Referring to the arbitrary waveform shown in Figure 2-2(a),  $x_p$  is the absolute peak value of the waveform  $x(t)$ , if  $x_{\text{rms}}$  is its rms value, then the crest factor is [2]

$$\text{CF} = x_p / x_{\text{rms}}. \tag{2.1}$$

More formally, 
$$\text{CF} = \frac{\|x\|_{\infty}}{\|x\|_2}, \tag{2.2}$$

where  $\|x\|_{\infty}$  is the infinity norm, and here is the maximum value of  $x(t)$ ,  $\|x\|_{\infty} = \max[x(t)] = x_p$ , and  $\|x\|_2$  is just the rms value of  $x(t)$ :

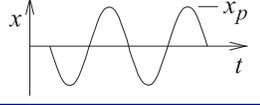
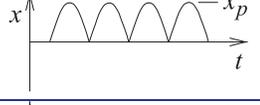
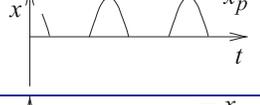
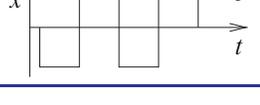
$$x_{\text{rms}} = \|x\|_2 = \lim_{T \rightarrow \infty} \sqrt{\frac{1}{T} \int_0^T x(t) \cdot dt}. \tag{2.3}$$

Note that CF is a voltage (or current) ratio rather than a power ratio. The CFs of several waveforms are given in Table 2-1.

#### Peak-to-Average Power Ratio (PAPR)

The peak-to-average power ratio (PAPR) is analogous to CF but for power. If  $x(t)$  is the voltage across a resistor, as shown in Figure 2-2(b), then the

**Table 2-1:** Crest factor (CF) and peak-to-average power ratio (PAPR) of several waveforms. ( $x_p$  is the peak value of the waveform.)

Waveform	$x(t)$	Max. value	rms ( $x_{\text{rms}}$ )	CF	PAPR
DC		$x_{\text{dc}}$	$x_{\text{dc}}$	1	0 dB
Sinewave		$x_p$	$\frac{x_p}{\sqrt{2}}$	1.414	3.01 dB
Full-wave rectified sinewave		$x_p$	$\frac{x_p}{\sqrt{2}} = 0.717 x_p$	1.414	3.01 dB
Half-wave rectified sinewave		$x_p$	$\frac{x_p}{2}$	2	6.02 dB
Triangle wave		$x_p$	$\frac{x_p}{\sqrt{3}} = 0.577 x_p$	1.732	4.77 dB
Square wave		$x_p$	$x_p$	1	0 dB

instantaneous peak power in the resistor is

$$P_p = |x_p|^2 / R, \quad (2.4)$$

where again  $x_p$  is the peak absolute value of the waveform.  $P_p$  is the power of the peak of a waveform treating it as though it was a DC signal. This is appropriate for a slowly varying signal such as a power frequency signal as it is this instantaneous power that determines thermal disruption of a power system. It is not the appropriate power to use with radio signals and a more suitable microwave signal metric is described in Section 2.2.2. The average power dissipated in the resistor is

$$P_{\text{avg}} = |x_{\text{rms}}|^2 / R. \quad (2.5)$$

Then 
$$\text{PAPR} = \frac{P_p}{P_{\text{avg}}} = \text{CF}^2 = (x_p / x_{\text{rms}})^2. \quad (2.6)$$

In decibels, 
$$\text{PAPR}|_{\text{dB}} = 10 \log(\text{PAPR}) \\ = 20 \log(\text{CF}) = 20 \log(x_p / x_{\text{rms}}). \quad (2.7)$$

The definition of PAPR above can be used with any waveform and can be used in all branches of electrical engineering. The PAPRs of several waveforms are given in Table 2-1.

## EXAMPLE 2.1

## Crest Factor and PAPR of an Offset Sinusoid.

What is the crest factor (CF) and peak-to-average power ratio (PAPR) of the signal  $x(t) = 0.1 + 0.5 \sin(\omega t)$ ?

**Solution:**

The signal is a sinusoid offset by a DC term. The peak value of  $x(t)$  is  $x_p = 0.6$ , and the rms value of the signal will be the square root of the rms values squared of the individual DC and sinusoidal components. This applies to any composite signal provided that the components are uncorrelated. So  $x_{\text{rms}} = \sqrt{0.1^2 + (0.5/\sqrt{2})^2} = 0.3674$ . The general solution for a signal  $x(t) = a + b \sin(\omega t)$  is, using Equation (2.3),

$$\begin{aligned} x_{\text{rms}} &= \sqrt{\lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T [x(t)]^2 .dt} = \sqrt{\lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T [a + b \sin(\omega t)]^2 .dt} \\ &= \sqrt{\lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T [a^2 + ab \sin(\omega t) + b^2 \sin^2(\omega t)] .dt} \\ &= \sqrt{\lim_{T \rightarrow \infty} \frac{1}{T} \left\{ \int_0^T a^2 .dt + \int_0^T ab \sin(\omega t) .dt + \int_0^T b^2 \frac{1}{2} [1 + \cos(2\omega t)] .dt \right\}} \\ &= \sqrt{\lim_{T \rightarrow \infty} \frac{1}{T} \{ a^2 T .dt + 0 + \frac{1}{2} b^2 T \}} \end{aligned} \quad (2.8)$$

since the integral of sin and cos over a period is zero. Thus

$$x_{\text{rms}} = \sqrt{a^2 + b^2/2} = \sqrt{0.1^2 + \frac{1}{2} 0.5^2} = 0.3674, \quad (2.9)$$

$$\text{the crest factor is } CF = \frac{x_p}{x_{\text{rms}}} = \frac{0.6}{0.3674} = 1.6331, \quad (2.10)$$

$$\text{and PAPR is } PAPR = 20 \log(1.6331) = 4.260 \text{ dB}. \quad (2.11)$$

There is a quicker way of calculating PAPR by dealing with the powers directly. The peak power of the waveform is  $P_p = x_p^2/R = 0.6^2/R = 0.36/R$ , where  $x$  is being treated as a voltage across a resistor  $R$ . The two parts of  $x(t)$ , i.e. the DC component and the sinewave, are uncorrelated, so the average power of the combined signal is the sum of the powers of the uncorrelated components, so

$$P_{\text{avg}} = \frac{1}{R} [0.1^2 + \frac{1}{2} 0.5^2] \frac{1}{R} = \frac{0.1350}{R}. \quad (2.12)$$

Thus, in decibels,

$$PAPR|_{\text{dB}} = 10 \log \left( \frac{P_p}{P_{\text{avg}}} \right) = \frac{x_p^2}{x_{\text{rms}}^2} = 10 \log \left( \frac{0.36}{0.135} \right) = 10 \log(2.667) = 4.260 \text{ dB}. \quad (2.13)$$

### 2.2.2 Peak-to-Mean Envelope Power Ratio

Another metric for characterizing signals is the peak-to-mean envelope power ratio (PMEPR) and this is particularly useful for modulated signals. The amount of information sent by a communication signal is proportional to its average power, however, RF hardware must be designed with enough margin to be able to handle peaks in the signal without producing appreciable distortion. The waveform of a narrowband modulated signal appears as a carrier that slowly changes in amplitude and phase. One sinewave of this modulated signal is called a **pseudo-carrier** and the power of one cycle of the pseudo-carrier when the amplitude of the modulated

signal is at its maximum (i.e. at the peak of the envelope) is called the **peak envelope power (PEP)** [1] ( $PEP = P_{PEP}$ ). The ratio of PEP to the average signal power (the power averaged over all time) is called the PMEPR.

Then if the average power of the modulated signal is  $P_{avg}$

$$PMEPR = \frac{PEP}{P_{avg}} = \frac{P_{PEP}}{P_{avg}}. \quad (2.14)$$

PMEPR is a good indicator of how sensitive a modulation format is to distortion introduced by the nonlinearity of RF hardware [3].

It is complex to determine the PMEPR for a general modulated signal. Below the mathematics is presented for an AM signal with a sinusoidal modulating signal. Determining the PMEPR otherwise requires numerical integration following the procedure outlined below.

### PMEPR of an AM Signal

A good estimate of the PMEPR of an AM signal can be obtained by considering a sinusoidal modulating signal (rather than an actual baseband signal). Let  $y(t) = \cos(2\pi f_m t)$  be a cosinusoidal modulating signal with frequency  $f_m$ . Then, for AM, the modulated carrier signal is

$$x(t) = A_c [1 + m \cos(2\pi f_m t)] \cos(2\pi f_c t) \quad (2.15)$$

where  $m$  is the modulation index (e.g. 100% AM has  $m = 1$ ). Thus if the power of just one quasi-period of  $x(t)$ , i.e. one cycle of the pseudo carrier, is considered then  $x(t)$  has a power that varies with time.

Consider a voltage  $v(t)$  across a resistor of conductance  $G$ . The power of the signal is determined by integrating over all time, which is work, and dividing by the time period. This yields the average power:

$$P_{avg} = \lim_{\tau \rightarrow \infty} \int_{-\tau}^{\tau} \frac{1}{2\tau} G v^2(t) dt. \quad (2.16)$$

Now, if  $v(t)$  is a sinusoidal,  $v(t) = A \cos \omega t$ , then

$$\begin{aligned} P_{avg} &= \lim_{\tau \rightarrow \infty} \frac{1}{2\tau} \int_{-\tau}^{\tau} A_c^2 G \cos^2(\omega t) dt \\ &= \lim_{\tau \rightarrow \infty} \frac{1}{2\tau} \int_{-\tau}^{\tau} A_c^2 G \frac{1}{2} [1 + \cos(2\omega t)] dt \\ &= \frac{1}{2} A_c^2 G \left\{ \lim_{\tau \rightarrow \infty} \frac{1}{2\tau} \int_{-\tau}^{\tau} 1 dt + \lim_{\tau \rightarrow \infty} \frac{1}{2\tau} \int_{-\tau}^{\tau} \cos(2\omega t) dt \right\} = \frac{1}{2} A_c^2 G. \end{aligned} \quad (2.17)$$

In the above equation, a useful equivalence has been employed by observing that the infinite integral of a cosinusoid can be simplified to just integrating over one period,  $T = 2\pi/\omega$ :

$$\lim_{\tau \rightarrow \infty} \frac{1}{2\tau} \int_{-\tau}^{\tau} \cos^n(\omega t) dt = \frac{1}{T} \int_{-T/2}^{T/2} \cos^n(\omega t) dt, \quad (2.18)$$

where  $n$  is a positive integer. In power calculations there are a number of other useful simplifying techniques based on trigonometric identities. Some

of the ones that will be used here are the following:

$$\begin{aligned}\cos A \cos B &= \frac{1}{2} [\cos(A - B) + \cos(A + B)] \\ \cos^2 A &= \frac{1}{2} [1 + \cos(2A)]\end{aligned}\quad (2.19)$$

$$\lim_{\tau \rightarrow \infty} \frac{1}{2\tau} \int_{-\tau}^{\tau} \cos \omega t \, dt = \frac{1}{T} \int_{-T/2}^{T/2} \cos(\omega t) \, dt = 0 \quad (2.20)$$

$$\frac{1}{T} \int_{-T/2}^{T/2} \cos^2(\omega t) \, dt = \frac{1}{T} \int_{-T/2}^{T/2} \frac{1}{2} [\cos(2\omega t) + \cos(0)] \, dt \quad (2.21)$$

$$\begin{aligned}&= \frac{1}{2T} \left[ \int_{-T/2}^{T/2} \cos(2\omega t) \, dt + \int_{-T/2}^{T/2} 1 \, dt \right] \\ &= \frac{1}{2T} (0 + T) = \frac{1}{2}.\end{aligned}\quad (2.22)$$

More trigonometric identities are given in Appendix 1.A.2 of [4]. Also, when cosinusoids  $\cos \omega_A t$  and  $\cos \omega_B t$ , having different frequencies ( $\omega_A \neq \omega_B$ ), are multiplied together, for large  $\tau$ ,

$$\int_{-\tau}^{\tau} \cos \omega_A t \cos \omega_B t \, dt = \int_{-\tau}^{\tau} \frac{1}{2} [\cos(\omega_A + \omega_B)t + \cos(\omega_A - \omega_B)t] \, dt = 0,$$

$$\text{and if } \omega_A \neq \omega_B \neq 0, \quad \int_{-\infty}^{\infty} \cos \omega_A t \cos^n \omega_B t \, dt = 0. \quad (2.23)$$

Now the discussion returns to characterizing an AM signal by considering the long-term average power and the maximum short-term power of the signal. The pseudo-carrier at its peak amplitude is, from Equation (2.15),

$$x_p(t) = A_c [1 + m] \cos(2\pi f_c t). \quad (2.24)$$

Then the power ( $P_{\text{PEP}}$ ) of the peak pseudo carrier is obtained by integrating over one period of the pseudo carrier:

$$\begin{aligned}P_{\text{PEP}} &= \frac{1}{T} \int_{-T/2}^{T/2} Gx^2(t) \, dt = \frac{1}{T} \int_{-T/2}^{T/2} A_c^2 G(1 + m)^2 \cos^2(\omega_c t) \, dt \\ &= A_c^2 G(1 + m)^2 \frac{1}{T} \int_{-T/2}^{T/2} \cos^2(\omega_c t) \, dt = \frac{1}{2} A_c^2 G(1 + m)^2.\end{aligned}\quad (2.25)$$

The **average power** ( $P_{\text{avg}}$ ) of the modulated signal is obtained by integrating over all time, so

$$\begin{aligned}P_{\text{avg}} &= \lim_{\tau \rightarrow \infty} \frac{1}{2\tau} \int_{-\tau}^{\tau} Gx^2(t) \, dt \\ &= A_c^2 G \lim_{\tau \rightarrow \infty} \frac{1}{2\tau} \int_{-\tau}^{\tau} \{[1 + m \cos(\omega_m t)] \cos(\omega_c t)\}^2 \, dt \\ &= A_c^2 G \lim_{\tau \rightarrow \infty} \frac{1}{2\tau} \int_{-\tau}^{\tau} \{[1 + 2m \cos(\omega_m t) + m^2 \cos^2(\omega_m t)] \cos^2(\omega_c t)\} \, dt \\ &= A_c^2 G \left[ \lim_{\tau \rightarrow \infty} \frac{1}{2\tau} \int_{-\tau}^{\tau} \cos^2(\omega_c t) \, dt + \lim_{\tau \rightarrow \infty} \frac{1}{2\tau} \int_{-\tau}^{\tau} 2m \cos(\omega_m t) \cos^2(\omega_c t) \, dt \right. \\ &\quad \left. + \lim_{\tau \rightarrow \infty} \frac{1}{2\tau} \int_{-\tau}^{\tau} m^2 \cos^2(\omega_m t) \cos^2(\omega_c t) \, dt \right] \\ &= A_c^2 G \left\{ \frac{1}{2} + 0 + m^2 \lim_{\tau \rightarrow \infty} \frac{1}{2\tau} \int_{-\tau}^{\tau} \frac{1}{4} [1 + \cos(2\omega_m t)] [1 + \cos(2\omega_c t)] \, dt \right\}\end{aligned}$$

$$\begin{aligned}
&= A_c^2 G \left\{ \frac{1}{2} + \frac{m^2}{4} \left[ \lim_{\tau \rightarrow \infty} \frac{1}{2\tau} \int_{-\tau}^{\tau} 1 dt + \lim_{\tau \rightarrow \infty} \frac{1}{2\tau} \int_{-\tau}^{\tau} \cos(2\omega_m t) dt \right. \right. \\
&\quad \left. \left. + \lim_{\tau \rightarrow \infty} \frac{1}{2\tau} \int_{-\tau}^{\tau} \cos(2\omega_c t) dt + \lim_{\tau \rightarrow \infty} \frac{1}{2\tau} \int_{-\tau}^{\tau} \cos(2\omega_m t) \cos(2\omega_c t) dt \right] \right\} \\
&= A_c^2 G \left[ \frac{1}{2} + m^2 \left( \frac{1}{4} + 0 + 0 + 0 \right) \right] = \frac{1}{2} A_c^2 G (1 + m^2/2). \tag{2.26}
\end{aligned}$$

Thus the rms voltage,  $x_{\text{rms}}$ , can be determined as  $P_{\text{avg}} = x_{\text{rms}}^2 G$ . So the PMEPR of an AM signal (i.e.,  $\text{PMEPR}_{\text{AM}}$ ) is

$$\text{PMEPR}_{\text{AM}} = \frac{P_{\text{PEP}}}{P_{\text{avg}}} = \frac{\frac{1}{2} A_c^2 G (1+m)^2}{\frac{1}{2} A_c^2 G (1+m^2/2)} = \frac{(1+m)^2}{1+m^2/2}.$$

For 100% AM described by  $m = 1$ , the PMEPR is

$$\text{PMEPR}_{100\% \text{AM}} = \frac{(1+1)^2}{1+1^2/2} = \frac{4}{1.5} = 2.667 = 4.26 \text{ dB}. \tag{2.27}$$

In expressing the PMEPR in decibels, the formula  $\text{PMEPR}_{\text{dB}} = 10 \log(\text{PMEPR})$  is used as PMEPR is a power ratio. As an example, for 50% AM, described by  $m = 0.5$ , the PMEPR is

$$\text{PMEPR}_{50\% \text{AM}} = \frac{(1+0.5)^2}{1+0.5^2/2} = \frac{2.25}{1.125} = 2 = 3 \text{ dB}. \tag{2.28}$$

### 2.2.3 Two-Tone Signal

In assessing, either through laboratory measurements or simulations, it is common and often necessary to use very simple representations of a baseband signal or even of a modulated signal. This greatly simplifies matters and there is a justified expectation that the performance with the test signal is a good indication of performance with an actual baseband or modulated signal. With simulation at the circuit level it is usually impossible to consider real baseband signals as simulation may not even be possible or simulation may take unacceptable times. Instead it is common to use single-tone, i.e. single sinusoid, or two-tone signals. A two-tone signal is a signal that is the sum of two cosinusoids:

$$y(t) = X_A \cos(\omega_A t) + X_B \cos(\omega_B t). \tag{2.29}$$

Generally the frequencies of the two tones are close ( $|\omega_A - \omega_B| \ll \omega_A$ ), with the concept being that both tones fit within the passband of a transmitter's or receiver's bandpass filters. A two-tone signal is not a form of modulation, but is commonly used to characterize the nonlinear performance of RF systems and has an envelope that is similar to that of many modulated signals. The composite signal,  $y(t)$ , looks like a pseudo-carrier with a slowly varying amplitude, not unlike an AM signal. The tones are uncorrelated so that the average power of the composite signal,  $y(t)$ , is the sum of the powers of each of the individual tones. The peak power of the composite signal is that of the peak pseudo-carrier, so  $y(t)$  has a peak amplitude of  $X_A + X_B$ . The peak pseudo carrier is the single RF sinusoid where the sinusoid of each sinusoid align as much as possible. Similar concepts apply to three-tone and  $n$ -tone signals.

**EXAMPLE 2.2** PMEPR of a Two-Tone Signal

What is the PMEPR of a two-tone signal with the tones having equal amplitude?

**Solution:**

Let the amplitudes of the two tones be  $X_A$  and  $X_B$ . Now  $X_A = X_B = X$ , and so the peak pseudo-carrier has amplitude  $2X$ , and the power of the peak RF carrier is proportional to  $\frac{1}{2}(2X)^2 = 2X^2$ . The average power is proportional to  $\frac{1}{2}(X_A^2 + X_B^2) = \frac{1}{2}(X^2 + X^2) = X^2$ , as each tone is independent of the other and so the powers can be added.

$$\text{PMEPR} = \frac{P_{\text{PEP}}}{P_{\text{avg}}} = \frac{2X^2}{X^2} = 2 = 3 \text{ dB.} \tag{2.30}$$

**EXAMPLE 2.3** PMEPR of Uncorrelated Signals

Consider the combination of two uncorrelated analog signals, e.g. a two-tone signal. One signal is denoted  $x(t)$  and the other  $y(t)$ , where  $x(t) = 0.1 \sin(10^9 t)$  and  $y(t) = 0.05 \sin(1.01 \cdot 10^9 t)$ . What is the PMEPR of this combined signal?

**Solution:**

These two signals are uncorrelated and this is key in determining the average power,  $P_{\text{avg}}$ , as the sum of the powers of each individual signal ( $k$  is a proportionality constant):

$$P_{\text{avg}} = \int_{-\infty}^{\infty} x^2(t) \cdot dt + \int_{-\infty}^{\infty} y^2(t) \cdot dt = \frac{k}{2}(0.1)^2 + \frac{k}{2}(0.05)^2 = \frac{k}{2}[0.01 + 0.0025] = 0.00625k.$$

The two carriers are close in frequency so that the sum signal  $z(t) = x(t) + y(t)$  looks like a slowly varying signal with a radian frequency near  $10^9$  rads per second. The peak amplitude of one pseudo-cycle of  $z(t)$  is  $0.1 + 0.05 = 0.15$ . Thus the power of the largest cycle is

$$P_{\text{PEP}} = \frac{1}{2}k(0.15)^2 = 0.01125k,$$

and so 
$$\text{PMEPR} = \frac{P_{\text{PEP}}}{P_{\text{avg}}} = \frac{0.01125}{0.00625} = 1.8 = 2.55 \text{ dB.} \tag{2.31}$$

**Summary**

The PMEPR is an important attribute of a modulation format and impacts the types of circuit designs that can be used. It is much more challenging to develop power-efficient hardware introducing only low levels of distortion when the PMEPR is high.

It is tempting to consider if the lengthy integrations can be circumvented. Powers can be added if the signal components (the tones making up the signal) are uncorrelated. If they are correlated, then the complete integrations are required. Consider two uncorrelated sinusoids of (average) powers  $P_1$  and  $P_2$ , respectively, then the average power of the composite signal is  $P_{\text{avg}} = P_1 + P_2$ . However, in determining the peak sinusoidal power, the RF cycle where the two largest pseudo-carrier sinusoids align is considered, and here the voltages add to produce a single cycle of a sinewave with a higher amplitude. So peak power applies to just one RF pseudo-cycle. Generally the voltage amplitude of the two sinewaves would be added and then the power calculated. If the uncorrelated carriers are modulated and the modulating signals (the baseband signals) are uncorrelated, then the average power can be determined in the same way, but the peak power calculation is much more complicated. The integrations are the only calculations that can always be relied on and can be used with all modulated signals.

Signals  $x(t)$  and  $y(t)$  are **uncorrelated** if the integral over all time and time offsets of their product is zero:  

$$C = \int_{-\infty}^{+\infty} x(t)y(t+\tau) dt = 0$$
 for all  $\tau$ .

The preferred usage of PAR, PAPR, or PMEPR in RF and microwave engineering is currently in a transition phase. The most common usage of PAR and PAPR in electrical engineering refers to the peak of a signal as being the instantaneous peak value, and in the case of PAPR, the instantaneous power of the signal is calculated as if the peak is a DC value. In the past, many RF and microwave publications have taken the peak as the peak power of a sinusoid having an amplitude equal to the peak voltage of the signal and used that to calculate PAR. This usage is inconsistent with the predominant usage in electrical engineering and is a particular problem when using wireless technology in other disciplines. PMEPR is the preferred usage for what RF and microwave engineers intend to refer to when using the term PAR. A reader of RF literature encountering PAR needs to determine how the term is being used. There is no confusion if PMEPR is used.

#### EXAMPLE 2.4 PAPR and PMEPR of an AM signal

What is the PAPR and PMEPR of a 100% AM signal?

##### Solution:

The signal is  $x(t) = A_c [1 + \cos 2\pi f_m t] \cos 2\pi f_c t$  and the PMEPR of this signal, from Equation (2.27), is 4.26 dB. Now PAPR uses the absolute maximum value of the signal rather than the maximum short-term power of the envelope. The peak value of  $x(t)$  is  $2A_c$  so the peak power (if the signal is a voltage across a conductance  $G$ ) is

$$P_{\text{peak,PAPR}} = (2A_c)^2 G. \quad (2.32)$$

$P_{\text{avg}}$  is the same for PAPR and PMEPR for the AM signal, see Equation (2.26), so that

$$\text{PAPR} = \frac{P_{\text{peak,PAPR}}}{P_{\text{avg}}} = \frac{(2A_c)^2 G}{\frac{1}{2}A_c^2 (1 + \frac{1}{2})} = \frac{4}{3/4} = \frac{16}{3} = 5.333 = 7.27 \text{ dB}. \quad (2.33)$$

So PAPR is 3 dB higher than PMEPR for a 100% modulated AM signal, see Equation (2.27). This is not always the case for other modulation schemes.

## 2.3 Modulation Overview

There are two families of modulation methods with analog modulation used in early radios including 1G cellular radio, and digital modulation used in modern radios starting with 2G cellular radio. While 1G cellular radio transmitted voice signals using analog modulation, 1G also used a simple type of digital modulation for signaling. With the exception of **ultra-wideband (UWB) pulse radio** [5], all modern radio modulation schemes slowly vary the amplitude, phase, or frequency of a sinusoidal signal called the carrier. This results in a narrow bandwidth modulated signal perhaps with fractional bandwidth typically in the range of 0.002% to 2%. The early spark-gap wireless telegraph systems were ultra-wideband but they were soon discontinued because they interfered with conventional radios which were soon developed and assigned specific parts, i.e. bands, of the spectrum. The initial pulse radio concept of the 1990s occupied most of the spectrum between 3.1 and 10.6 GHz but was never deployed mainly because capacity was relatively poor. The term ultra-wideband wireless is now widely taken to mean a wireless device such as a radar or radio with a bandwidth which is at least the lesser of 500 MHz or 20% of the carrier frequency [6]. So even the UWB millimeter-wave radios exploiting the high bandwidth available

at millimeter wave frequencies still employ a relatively slowly varying modulation of a carrier.

### 2.4 Analog Modulation

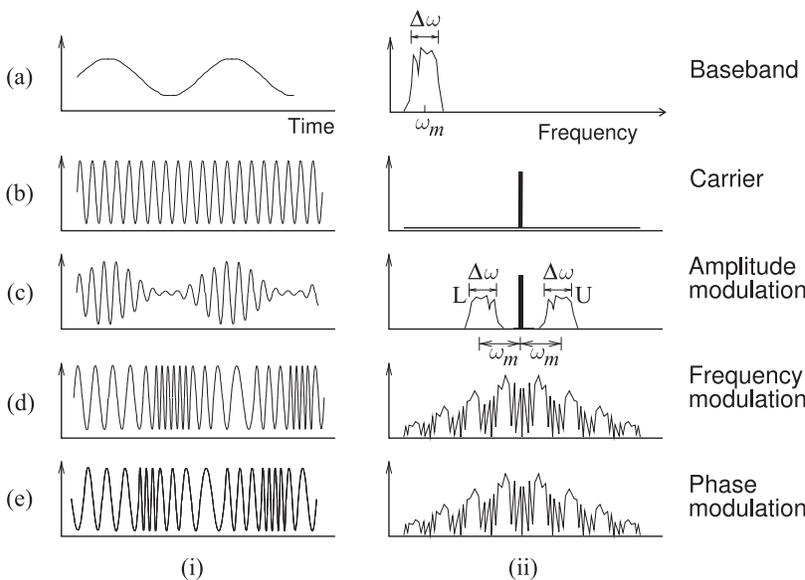
The waveforms and spectra of the signals with common analog modulation methods are shown in Figure 2-3. The modulating signal is generally referred to as the baseband signal and it contains all of the information to be transmitted and interpreted at the receiver. The waveforms in Figure 2-3 are stylized. They are presented this way so that the effects of modulation can be more easily seen. The baseband signal (Figure 2-3(a)) is shown as having a period that is not much greater than the period of the carrier (Figure 2-3(b)). In reality there would be hundreds or thousands of RF cycles for each cycle of the baseband signal so that the highest frequency component of the baseband signal is a tiny fraction of the carrier frequency. In this situation the spectra shown on the right in Figure 2-3(c-e) would be too narrow to enable any detail to be seen.

#### 2.4.1 Amplitude Modulation

Amplitude Modulation (AM) is the simplest analog modulation method to implement. Here a signal is used to slowly vary the amplitude of the carrier according to the level of the modulating signal. With AM (Figure 2-3(c)) the amplitude of the carrier is modulated, and this results in a broadening of the spectrum of the carrier, as shown in Figure 2-3(c)(ii). This spectrum contains the original carrier component and upper and lower sidebands, designated as U and L, respectively. In AM, the two sidebands contain identical information, so all the information contained in the baseband signal is conveyed if just one sideband is transmitted.

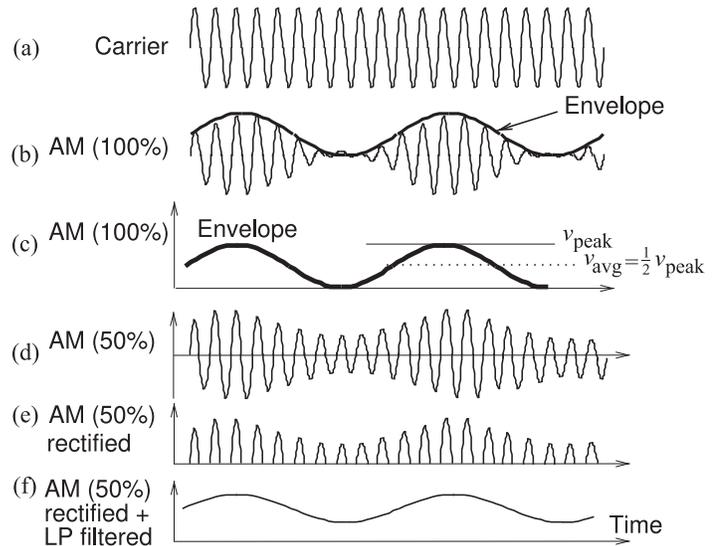
The basic AM signal  $x(t)$  has the form

$$x(t) = A_c [1 + my(t)] \cos(2\pi f_c t), \tag{2.34}$$



**Figure 2-3:** Basic analog modulation showing the (i) waveform and (ii) spectrum for (a) baseband signal; (b) carrier; (c) carrier modulated using amplitude modulation; (d) carrier modulated using frequency modulation; and (e) carrier modulated using phase modulation.

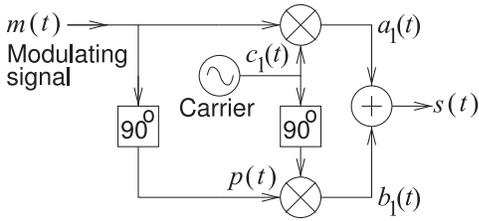
**Figure 2-4:** AM showing the relationship between the carrier and modulation envelope: (a) carrier; (b) 100% amplitude modulated carrier; (c) modulating or baseband signal; (d) 50% amplitude-modulated carrier; (e) rectified 50% AM modulated signal; and (f) rectified and lowpass (LP) filtered 50% modulated signal. The envelope contains only amplitude information and for AM the envelope is the same as the baseband signal.



where  $m$  is the modulation index,  $y(t)$  is the baseband information-bearing signal that has frequency components that are much lower than the carrier frequency  $f_c$ , and the maximum value of  $|y(t)|$  is one. Provided that  $y(t)$  varies slowly relative to the carrier,  $x(t)$  looks like a carrier whose amplitude varies slowly. To get an idea of how slowly the amplitude varies in an actual system, consider an AM radio that broadcasts at 1 MHz (which is in the middle of the AM broadcast band). The highest frequency component of the modulating signal corresponding to voice is about 4 kHz. Thus the amplitude of the carrier takes 250 carrier cycles to go through a complete amplitude variation. At all times a cycle of the modulated carrier, the pseudo-carrier, appears to be periodic, but in fact it is not quite.

The concept of the envelope of a modulated RF signal is introduced in Figure 2-4. The envelope is an important concept and is directly related to the distortion introduced by analog hardware and to the DC power requirements which determines the battery life for mobile radios. Figure 2-4(a) is the carrier and the amplitude-modulated carrier is shown in Figure 2-4(b). The outline of the modulated carrier is called the envelope, and for AM this is identical to the modulating, i.e. baseband, signal. The envelope is shown again in Figure 2-4(c). At the peak of the envelope, the RF signal has maximum short-term power (considering the power of a single RF cycle). With 100% AM,  $m = 1$  in Equation (2.34), there is no short-term RF power when the envelope is at its minimum. The modulated signal with 50% modulation,  $m = 0.5$ , is shown in Figure 2-4(d) and at all times there is an appreciable RF signal power.

Very simple analog hardware is required to demodulate the basic amplitude modulated signal, that is an AM signal with a carrier and both sidebands. The receiver requires bandpass filtering to select the channel from the incoming radio signal then rectifying the output of the bandpass filter. The waveform after rectification of a 50% AM signal is shown in Figure 2-4(e) and contains frequency components at baseband and sidebands around harmonics of the carrier, and the harmonics of the carrier itself. Lowpass filtering of the rectified waveform extracts the original baseband signal and completes demodulation, see Figure 2-4(f). The only electronics required is a



**Figure 2-5:** Hartley modulator implementing single-sideband suppressed-carrier (SSB-SC) modulation. The “90°” blocks shift the phase of the signal by +90°. The mixer indicated by the circle with a cross is an ideal multiplier, e.g.  $a_1(t) = m(t) \cdot c_1(t)$ .

single diode. The disadvantage is that more spectrum is used than required and the largest signal is the carrier that conveys no information but causes interference in other radios. Without the carrier and both sidebands being transmitted it is necessary to use DSP to demodulate the signal.

It is not possible to represent an actual baseband signal in a simple way and undertake the analytic derivations that illustrate the characteristics of modulation. Instead it is usual to use either a one-tone or two-tone signal, derive results, and then extrapolate the results for a finite bandwidth baseband signal. For a single-tone baseband signal  $y(t) = \cos(\omega_m t + \phi)$ , then the basic AM modulated signal, from Equation (2.34), is

$$\begin{aligned} x(t) &= [1 + m \cos(\omega_m t + \phi)] \cos(\omega_c t) \\ &= \cos(\omega_c t) + \frac{1}{2}m[\cos((\omega_c - \omega_m)t - \phi) + \cos((\omega_c + \omega_m)t + \phi)] \end{aligned} \quad (2.35)$$

which has three (radian) frequency components, one at the carrier frequency  $\omega_c$ , one just below the carrier at  $\omega_c - \omega_m$ , and one just above at  $\omega_c + \omega_m$  (since  $\omega_m \ll \omega_c$ ). The extension to a finite bandwidth baseband signal, see Figure 2-3(a)(ii), is to imagine that  $\omega_m$  ranges from a lower value  $\omega_m - \frac{1}{2}\Delta\omega$  to a higher value  $\omega_m + \frac{1}{2}\Delta\omega$ . The discrete tones in the modulated signal below and above the carrier then become finite bandwidth sidebands with a lower sideband L centered at  $\omega_c - \omega_m$  and an upper sideband U centered at  $(\omega_c + \omega_m)$  each having the same bandwidth,  $\Delta\omega$ , as the baseband signal, see Figure 2-3(c)(ii).

The AM modulator described so far produces a modulated signal with a carrier and two sidebands. This modulation is called double-sideband (DSB) modulation. There is identical information in each of the sidebands and so only one of the sidebands needs to be transmitted. The carrier contains no information so if only one sideband was transmitted then the received **single-sideband (SSB) suppressed-carrier SC** (together **SSB-SC**) signal has all of the information needed to recover the original baseband signal. However the simple demodulation process using rectification as described earlier in this section no longer works. The receiver needs to use DSP but the spectrum is used efficiently.

One circuit that implements SSB-SC AM is the **Hartley modulator** shown in Figure 2-5. As will be seen, this basic architecture is significant and used in all modern radios. In modern radios the Hartley modulator, or a variant, takes a modulated signal which is centered at an intermediate frequency and shifts it up in frequency so that it is centered at another frequency a little below or a little above the carrier of the Hartley modulator.

In a Hartley modulator both the modulating signal  $m(t)$  and the carrier are multiplied together in a mixer and then also 90° phase-shifted versions are mixed before being added together. The signal flow is as follows beginning with  $m(t) = \cos(\omega_m t + \phi)$ ,  $p(t) = \cos(\omega_m t + \phi - \pi/2) = \sin(\omega_m t + \phi)$  and

carrier signal  $c_1(t) = \cos(\omega_c t)$ :

$$\begin{aligned} a_1(t) &= \cos(\omega_m t + \phi) \cos(\omega_c t) = \frac{1}{2} [\cos((\omega_c - \omega_m)t - \phi) + \cos((\omega_c + \omega_m)t + \phi)] \\ b_1(t) &= \sin(\omega_m t + \phi) \sin(\omega_c t) = \frac{1}{2} [\cos((\omega_c - \omega_m)t - \phi) - \cos((\omega_c + \omega_m)t + \phi)] \\ s(t) &= a_1(t) + b_1(t) = \cos((\omega_c - \omega_m)t - \phi) \end{aligned} \quad (2.36)$$

and so the lower sideband (LSB) is selected. An interesting observation is that the phase,  $\phi$ , of the baseband signal is also translated up in frequency. A feature that is not exploited in AM but is in digital modulation.

### 2.4.2 Phase Modulation

In phase modulation (PM) the phase of the carrier depends on the instantaneous level of the baseband signal. The phase-modulated carrier is shown in Figure 2-3(e)(i) and it looks like the frequency of the modulated carrier is changing. What is actually happening is that when the phase is changing most quickly the apparent frequency of the RF waveform changes. Here, as the baseband signal is decreasing, the phase shift reduces and the effect is to increase the apparent frequency of the RF signal. As the baseband signal increases, the effect is to reduce the apparent frequency of the modulated RF signal. The result is that with PM is that the bandwidth of the time-varying signal is spread out, as seen in Figure 2-6. PM can be implemented using a phase-locked loop (PLL) but further details will be skipped here.

Consider a phase-modulated signal  $s(t) = \cos(\omega_c t + \phi(t))$  where  $\phi(t)$  is the baseband signal containing the information to be transmitted. The spectrum of  $s(t)$  can be determined by simplifying  $\phi(t)$  as a sinusoid with frequency  $f_m = 2\pi\omega_m$  so that  $\phi(t) = \beta \cos(\omega_m t)$  where  $\beta$  is the phase modulation index. (The maximum possible phase change is  $\pm\pi$  and then  $\beta = \pi$ .) The phase-modulated signal becomes

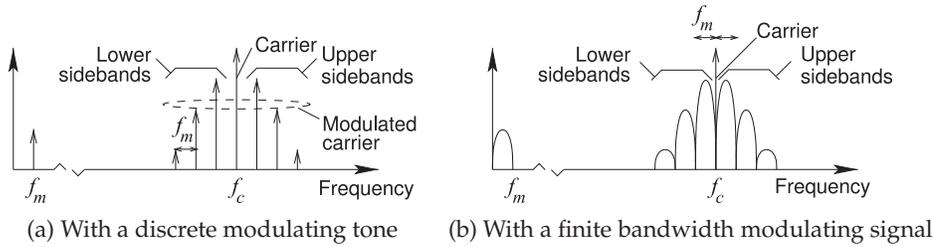
$$\begin{aligned} s(t) &= \cos(\omega_c t + \beta \cos(\omega_m t)) \\ &= \cos(\omega_c t) \cos(\cos(\beta\omega_m t)) - \sin(\omega_c t) \sin(\cos(\beta\omega_m t)) \end{aligned} \quad (2.37)$$

which has the Bessel function-based expansion

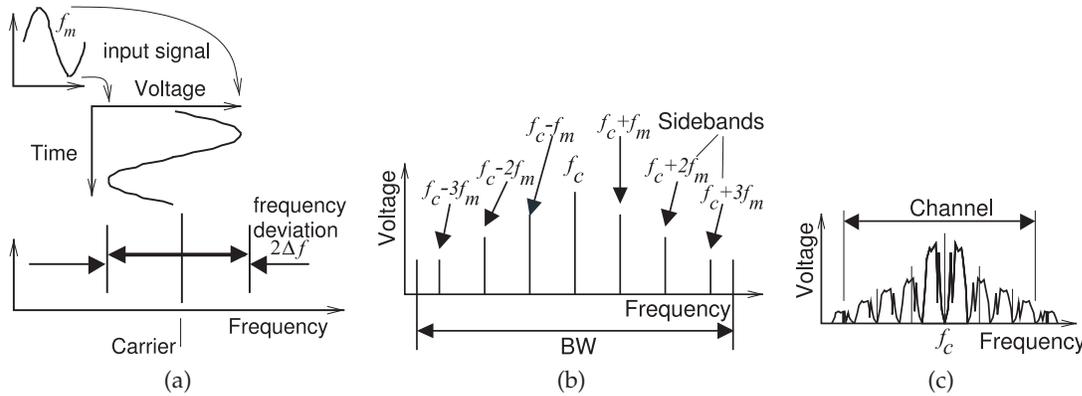
$$\begin{aligned} s(t) &= J_0(\beta) \cos(\omega_c t) \\ &+ J_1(\beta) \cos(\omega_c + \omega_m)t + \pi/2 + J_1(\beta) \cos(\omega_c - \omega_m)t + \pi/2 \\ &+ J_2(\beta) \cos(\omega_c + 2\omega_m)t + \pi + J_2(\beta) \cos(\omega_c - 2\omega_m)t + \pi \\ &+ J_3(\beta) \cos(\omega_c + 3\omega_m)t + 3\pi/2 + J_3(\beta) \cos(\omega_c - 3\omega_m)t + 3\pi/2 + \dots \end{aligned} \quad (2.38)$$

where  $J_n$  is the Bessel function of the first kind of order  $n$ . The spectrum of this signal is shown in Figure 2-6(a) which consists of discrete tones grouped as lower- and upper-sideband sets centered on the carrier at  $f_c$ . The discrete tones in the sidebands are separated from each other and from  $f_c$  by  $f_m$ . The sidebands have lower amplitude further away from the carrier.

If the modulating signal has a finite bandwidth, approximated by  $f_m$  varying from a minimum value,  $(f_m - \Delta f)$  up to the maximum frequency  $(f_m + \Delta f)$ , then the spectrum of the modulated signal becomes that shown in Figure 2-6(b), with the centers of adjacent sidebands separated by  $f_m$  and the first sidebands separated from the carrier by  $f_m$  as well. This is DSB



**Figure 2-6:** Spectrum of a phase-modulated carrier which includes the carrier at  $f_c$  and upper and lower sidebands with the spectrum of the discrete modulating signal at  $f_m$ .



**Figure 2-7:** Frequency modulation: (a) sinusoidal baseband signal shown varying the frequency of the carrier and so FM modulating the carrier; (b) the spectrum of the resulting FM-modulated waveform; and (c) spectrum of the modulated carrier when it is modulated by a broadband baseband signal such as voice.

modulation and there is a carrier (so it is not suppressed). The sidebands do not carry identical information and several, perhaps three below and three above the carrier, are required to enable demodulation of a PM signal. Thus a rather large bandwidth is required to transmit the modulated signal.

### 2.4.3 Frequency Modulation

The other analog modulation schemes commonly used is frequency modulation (FM), see Figure 2-3(d). The signals produced by FM and PM appear to be similar; the difference is in how the signals are generated. In FM, the amplitude of the baseband signal determines the frequency of the modulated carrier. Consider the FM waveform in Figure 2-3(d)(i). When the baseband signal is at its peak value the modulated carrier is at its minimum frequency, and when the signal is at its lowest value the modulated carrier is at its maximum frequency. (Depending on the hardware implementation it could be the other way around.) The result is that the bandwidth of the time-varying signal is spread out, as seen in Figure 2-7.

One way of implementing the FM modulator is to use a voltage-controlled

oscillator (VCO) with the baseband signal controlling the frequency of an oscillator. An FM receiver must compress, in frequency, the transmitted signal to re-create the original narrower bandwidth baseband signal. FM demodulation can be thought of as providing signal enhancement or equivalently noise suppression in a process that can be called analog processing gain. Only the components of the original FM signals are coherently collapsed to a narrower bandwidth baseband signal while noise, being uncorrelated, is still spread out (although rearranged). Thus the ratio of the signal to noise powers increases, as after demodulation only the power of the noise in the smaller bandwidth of the baseband signal is important. Thus compared to AM, FM significantly increases the tolerance to noise that may be added to the signal during transmission. PM has the same property, although the details of modulation and demodulation are different. For both FM and PM signals the peak amplitude of the RF **phasor** is equal to the average amplitude, and so the PMEPR is 1 or 0 dB.

Consider an FM signal  $s(t) = \cos([\omega_c + x(t)]t)$  where  $x(t)$  is the baseband signal containing the information to be transmitted. The spectrum of  $s(t)$  can be determined by simplifying  $x(t)$  as a sinusoid with frequency  $f_m = 2\pi\omega_m$  so that  $x(t) = \beta \cos(\omega_m t)$  where  $\beta$  is the frequency modulation index. The FM signal becomes

$$\begin{aligned} s(t) &= \cos([\omega_c + \beta \cos(\omega_m t)]t) \\ &= \cos(\omega_c t) \cos(\cos(\omega_m t)\beta t) - \sin(\omega_c t) \sin(\cos(\omega_m t)\beta t) \end{aligned} \quad (2.39)$$

which has the Bessel function-based expansion

$$\begin{aligned} s(t) &= J_0(\beta t) \cos(\omega_c t) \\ &\quad - J_1(\beta t) \sin(\omega_c + \omega_m)t + \pi/2) - J_1(\beta t) \sin(\omega_c - \omega_m)t + \pi/2) \\ &\quad - J_2(\beta t) \cos(\omega_c - 2\omega_m)t + \pi) + J_2(\beta t) \cos(\omega_c - 2\omega_m)t + \pi) \\ &\quad + J_3(\beta t) \sin(\omega_c + 3\omega_m)t + 3\pi/2) + J_3(\beta t) \sin(\omega_c - 3\omega_m)t + 3\pi/2) + \dots \end{aligned} \quad (2.40)$$

where  $J_n$  is the Bessel function of the first kind of order  $n$ . The spectrum of this signal is shown in Figure 2-7(b) which consists of discrete tones grouped as lower- and upper-sideband sets centered on the carrier at  $f_c$ . The discrete tones in the sidebands are separated from each other and from  $f_c$  by  $f_m$ . The sidebands have lower amplitude further away from the carrier.

If the modulating signal has a finite bandwidth, approximated by  $f_m = \omega_m/(2\pi)$  varying from a minimum value ( $f_m - \Delta f$ ) up to the maximum frequency ( $f_m + \Delta f$ ), then the spectrum of the modulated signal becomes that shown in Figure 2-7(c) with the centers of adjacent sidebands separated by  $f_m$  and the first sidebands from the carrier by  $f_m$  as well. This is DSB modulation and there is a carrier (so it is not suppressed but is smaller than with AM). The sidebands do not carry identical information and several, perhaps three on either side of the carrier, are required to enable demodulation of an FM signal. Thus a rather large bandwidth is required to transmit the modulated signal as it is not sufficient to transmit just one sideband to enable demodulation.

### Carson's Rule

Frequency- and phase-modulated signals have a very wide spectrum and the bandwidth required to reliably transmit a PM or FM signal is subjective. The best accepted criterion for determining the bandwidth requirement is called Carson's bandwidth rule or just Carson's rule [7, 8].

An FM signal is shown in Figure 2-7. In particular, Figure 2-7(a) shows the FM function. The level (typically voltage) of the baseband signal determines the frequency deviation of the carrier from its unmodulated value. The frequency shift when the modulating signal is a DC value  $x_m$  at its maximum amplitude is called the peak frequency deviation,  $\Delta f$ . So, if the modulating signal changes very slowly, the bandwidth of the modulated signal is  $2\Delta f$ .

A rapidly varying sinusoidal modulating signal produces a modulated signal with many discrete sidebands as seen in Figure 2-7(b). If the modulating baseband signal is broadband, then the sidebands have finite bandwidth as seen in Figure 2-7(c) and many are required to recover the original baseband signal. These sidebands continue indefinitely in frequency but rapidly reduce in power away from the frequency of the unmodulated carrier. Carson's rule provides an estimate of the bandwidth that contains 98% of the energy. If the maximum frequency of the modulating signal is  $f_m$ , and the maximum value of the modulating waveform is  $x_m$  (which would produce a frequency deviation of  $\Delta f$  if it is DC), then Carson's rule is that the

$$\text{bandwidth required} = 2 \times (f_m + \Delta f). \quad (2.41)$$

### Narrowband and Wideband FM

The most common type of FM signal, as used in FM broadcast radio, is called wideband FM, as the maximum frequency deviation is much greater than the highest frequency of the modulating or baseband signal, that is,  $\Delta f \gg f_m$ . In narrowband FM,  $\Delta f$  is close to  $f_m$ . Narrowband FM uses less bandwidth but requires a more sophisticated demodulation technique.

#### EXAMPLE 2.5 PAPR and PMEPR of FM Signals

Consider FM signals close in frequency but whose spectra do not overlap.

- What are PAPR and PMEPR of just one FM signal?
- What are PAPR and PMEPR of a signal comprised of two uncorrelated narrowband FM signals each having a small fractional bandwidth and having the same average power.

#### Solution:

- An FM signal has a constant envelope just like a single sinusoid, and so  $\text{PAPR} = 1.414 = 3.01 \text{ dB}$  and  $\text{PMEPR} = 1 = 0 \text{ dB}$ .
- Since the modulation is relatively slow, each of the FM signals will look like single tone signals and the combined signal will look like a two-tone signal. However this is not enough to solve the problem. A thought experiment is required to determine the largest pseudo-carrier when the FM signals combine. If the amplitude of each tone is  $X$ , then the amplitude when the FM signal waveforms align is  $2X$ . (This is the same as the peak of a two-tone signal but arrived at differently.) Then

$$P_{\text{avg}} = \text{sum of the powers of each FM signal} = 2k \frac{1}{2} X^2,$$

where  $k$  is a proportionality constant. For PAPR,

$$P_p = k(2X)^2 \quad \text{and} \quad \text{PMEPR} = \frac{P_p}{P_{\text{avg}}} = \frac{k(2X)^2}{2k\frac{1}{2}X^2} = 4 = 6.0 \text{ dB.} \quad (2.42)$$

For PMEPR,  $P_{\text{PEP}} = \text{power of the pseudo-carrier} = k\frac{1}{2}(2X)^2$

$$\text{and} \quad \text{PMEPR} = \frac{P_p''}{P_{\text{avg}}} = \frac{k\frac{1}{2}(2X)^2}{2k\frac{1}{2}X^2} = 2 = 3.0 \text{ dB.} \quad (2.43)$$

### 2.4.4 Analog Modulation Summary

Analog modulation was used in the first radios and in 1G cellular radios. Radio transmission using analog modulation, i.e. analog radio, has almost ceased as it does not use spectrum efficiently. Digital modulation along with error correction, can pack much more information in a limited bandwidth. A final comparison of the analog modulation techniques is given in Figure 2-8 emphasizing the PMEPR of AM and FM. The PMEPR of PM is the same as for FM.

One particular event in the development of radio is illustrative of the relationship of technology and business interests. Frequency modulation was invented by Edwin H. Armstrong and patented in 1933 [9, 10]. FM is virtually static free and clearly superior to AM radio. However, it was not immediately adopted largely because AM radio was established in the 1930s, and the adoption of FM would have resulted in the scrapping of a large installed infrastructure (seen as a commercial catastrophe) and so the introduction of FM was delayed by decades. The best technology does not always win immediately! Commercial interests and the large investment in an alternative technology have a great deal to do with the success of a technology [11].

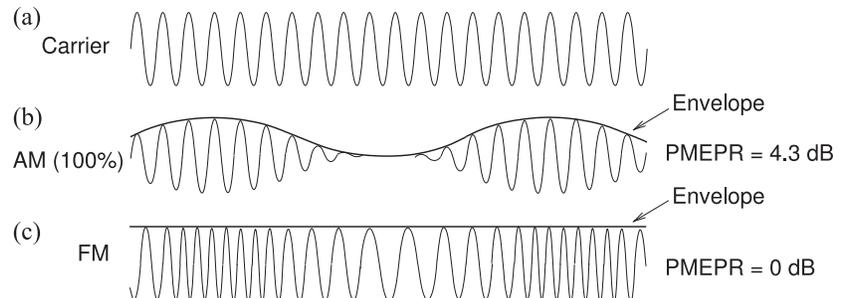
With FM and PM there are two sets of sidebands with one set above the carrier frequency and the other set below. The carrier itself is low-level but is not completely suppressed. Now SSB modulation refers to producing just one of the sideband sets. There is such as thing as SSB FM with just a few sidebands below (or above) the carrier but it is more like a combination of FM and AM [12], and was never deployed.

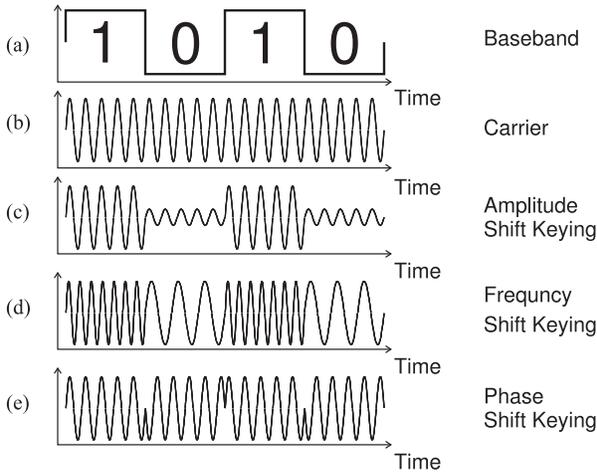
## 2.5 Digital Modulation

Digital radio transmits bits by creating discrete states, usually discrete amplitudes and phases of a carrier. The process of creating these discrete states from a digital bitstream is called digital modulation. A state is established at a particular time called a clock tick. What that means is that the

**Figure 2-8:**

Comparison of 100% AM and FM highlighting the envelopes of both: (a) carrier; (b) AM signal; and (c) FM signal with constant envelope.





**Figure 2-9:** Modes of digital modulation: (a) modulating bitstream; (b) **carrier**; (c) carrier modulated using amplitude shift keying (ASK); (d) carrier modulated using frequency shift keying (FSK); and (e) carrier modulated using binary phase shift keying (BPSK).

information in the signal is the state of the waveform, such as the amplitude and the phase of a phasor, at every clock tick such as every microsecond. The time it takes to go from one state to another (a clock tick interval) defines the bandwidth of the modulated signal. For example, if a clock tick is at every microsecond the bandwidth of the modulated signal is about one megahertz as it takes about one microsecond to go from one state to another. The inverse relationship of the interval between clock ticks to bandwidth is only approximate (as will be seen when software-defined radio is considered in a future chapter).

One important digital modulation method does not fit with the description above. This is Frequency Shift Keying (FSK) modulation where the carrier is set to a particular frequency at each clock tick.

The basic digital modulation formats are shown in Figure 2-9. The fundamental characteristic of digital modulation is that there are discrete states, each of which is also known as a symbol, with a symbol defining the value of one or more bits. For example, the states are different frequencies in FSK and different phases in phase shift keying (PSK). With the modulated waveforms shown in Figure 2-9 there are only two states, which is the same as saying that there are two symbols, each symbol having one bit of information (either 0 or 1). With multiple states groups of bits can be represented.

In this section many methods of digital modulation are described. The first few methods are binary modulation methods with just two symbols with one symbol indicating that a single bit is '0' and the other symbol indicating that it is a '1'. Then four-state modulation is introduced with four symbols with each symbol indicating the values of two bits. Higher-order modulation schemes can send more than more bits per symbol and thus more bits per second (bits/s) per hertz of bandwidth. There is a limit to the number of symbols as the "distance" between symbols becomes smaller and the effect of noise, interference, and circuit distortion can cause a symbol to be misinterpreted as another. A modulation method that sends more bits per symbol is said to have higher modulation efficiency. This and other metrics that enable modulation methods to be compared are defined in the next subsection.

### 2.5.1 Modulation Efficiency

With digital modulation, the information being sent is in the form of bits and it is possible to send more than one bit per second in one hertz of bandwidth. This is because in digital modulation there can be several bits per symbol, however the bandwidth of the modulated signal is determined by the rate of change from one state to another, whereas the number of bits per transition depends on the number of states. It is important for the transition to be no faster than required so as to minimize bandwidth.

The ratio of the **bit rate** in bits per second (bits/s) to the bandwidth (BW) in hertz is called the **modulation efficiency**,  $\eta_c$ , and has the units of bits per second per hertz (bits/s/Hz). The modulation efficiency is also called the channel efficiency, hence the subscript  $c$  on  $\eta_c$ . The bits here are the gross bits which includes the information bits and bits added for error correction and others added to aid in identifying the signal, and so  $\eta_c$  is a measure of the performance of the modulation method itself. Thus

$$\text{modulation efficiency} = \eta_c = \frac{\text{gross bit rate}}{\text{bandwidth}}. \quad (2.44)$$

The additional bits added to a bit stream are called **coding bits** and the process of adding the coding bits is called coding. If coding is used, then the information rate is lower than the gross bit rate transmitted. Thus gross bit rate refers to the bits actually transmitted and **information rate** (or **information bit rate**) refers to the bit rate of information transmission. The **link spectrum efficiency** is the information bit rate divided by the bandwidth. Often the term “link” is dropped and just **spectrum efficiency** is used (with units of bits/s/Hz). Thus

$$\text{link spectrum efficiency} = \frac{\text{information bit rate}}{\text{bandwidth}} \leq \eta_c. \quad (2.45)$$

#### EXAMPLE 2.6

#### Modulation Efficiency

A radio transmits a bit stream of 2 Mbits/s using a bandwidth of 1 MHz.

- What is the modulation efficiency?
- If 25% of the bits are used for error correction, what is the modulation efficiency?
- With error correction coding, what is the information rate?
- With error correction coding, what is the link spectrum efficiency?

#### Solution:

- The gross bit rate is 2 Mbits/s and the bandwidth is 1 MHz. So

$$\eta_c = \text{modulation efficiency} = \frac{\text{gross bit rate}}{\text{bandwidth}} = \frac{2 \text{ Mbits/s}}{1 \text{ MHz}} = 2 \text{ bits/s/Hz}.$$

- The modulation efficiency is unaffected by error correction coding. So the modulation efficiency is unchanged:

$$\eta_c = \text{modulation efficiency} = \frac{\text{gross bit rate}}{\text{bandwidth}} = \frac{2 \text{ Mbits/s}}{1 \text{ MHz}} = 2 \text{ bits/s/Hz}.$$

- With 25% of the bits in the gross bit stream being coding bits, the information rate is 75% of 2 Mbits/s or 1.5 Mbits/s.

- link spectrum efficiency =  $\frac{\text{information bit rate}}{\text{bandwidth}} = \frac{1.5 \text{ Mbits/s}}{1 \text{ MHz}} = 1.5 \text{ bits/s/Hz}.$

## 2.6 Frequency Shift Keying, FSK

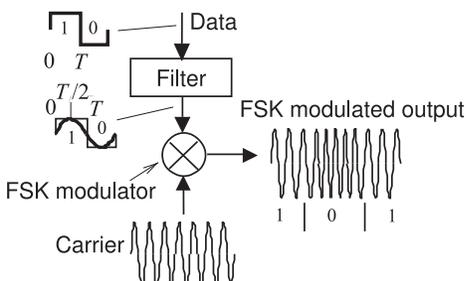
Frequency shift keying (FSK) is one of the simplest forms of digital modulation, with the frequency of the transmitted signal at a clock tick indicating a symbol, usually representing either one or two bits. Binary FSK (BFSK) is illustrated in Figure 2-9(d). It can be implemented by applying a discrete signal to the input of a voltage-controlled oscillator and so was ideally suited to early digital radio as simple high-performance FM modulators were available. Four-state FSK modulation is used in the GSM 2G cellular standard, a legacy standard still widely supported by modern cellular radios and sometimes the only modulation supported by the infrastructure (i.e. basestations) in some regions where it is not economically viable to retrofit old installations.

### 2.6.1 Essentials of FSK Modulation

The schematic of a binary FSK modulation system is shown in Figure 2-10. Here, a binary bitstream is lowpass filtered and used to drive an FSK modulator, one implementation of which shifts the frequency of an oscillator according to the voltage of the baseband signal. This function can be achieved using a VCO or a PLL circuit, and an FM demodulator can be used to receive the signal. Another characteristic feature of FSK is that the amplitude of the modulated signal is constant, so efficient saturating (and hence nonlinear) amplifiers can be used without the concern of frequency distortion. Not surprisingly, FSK was the first form of digital modulation used in mobile digital radio. A particular implementation of FSK is **Minimum Shift Keying (MSK)**, which uses a baseband lowpass filter so that the transitions from one state to another are smooth in time and limit the bandwidth of the modulated signal.

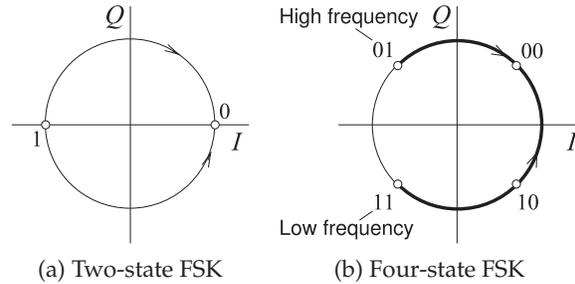
The **constellation diagram** is often thought of as being like a phasor diagram and this analogy works most of the time but it does not work for FSK modulation. A phasor diagram describes a phasor that is fixed in frequency. If the phasor is very slowly phase and/or amplitude modulated, then this approximation is good. FSK modulation cannot be represented on a phasor diagram, as the information is in the frequency at the clock ticks and not the than the phase and/or amplitude of a phasor. The symbols of two- and four-state FSK modulation are shown in Figure 2-11 which are called constellation diagrams

As an example consider an FSK-modulated signal with a bandwidth of 200 kHz and a carrier at 1 GHz (this approximately corresponds to the 2G GSM cellular system). This is a 0.02% bandwidth, so the phasor changes very slowly. Going from one FSK state to another takes about 1230 to 3692



**Figure 2-10:** The frequency shift keying (FSK) modulation system. In the GSM four-state cellular system adjacent constellation points differ in frequency by 33.25 kHz.

**Figure 2-11:** Constellation diagrams of FSK modulation. In two-state FSK a symbol indicates whether a bit is a '0' or a '1'. In four-state FSK there are four symbols and each symbol has a different frequency and indicates the values of two bits.



RF cycles depending on the frequency difference of the transition from one symbol to the next. With a 1 GHz carrier the frequencies of the four symbols are  $(1 \text{ GHz} - 33.25 \text{ kHz})$ ,  $(1 \text{ GHz} - 16.62 \text{ kHz})$ ,  $(1 \text{ GHz} + 16.62 \text{ kHz})$ , and  $(1 \text{ GHz} + 33.25 \text{ kHz})$ . This may seem like a very small frequency difference but hardware in the basestation and in the handset can easily achieve a frequency resolution of a few hertz at 1 GHz. In trying to represent FSK modulation on a pseudo-phaser diagram, the frequency is approximated as being fixed and the maximum real frequency shift is arbitrarily taken as being a significant shift of the pseudo-phaser.

In FSK, the states are on a circle in the constellation diagram (see Figure 2-11), with two-state FSK shown in Figure 2-11(a) and four-state FSK shown in Figure 2-11(b). Note that the constellation diagram indicates that the amplitude of the phasor is constant, as FSK modulation is a form of FM modulation. Consider four-state FSK more closely. There are four frequency states ranging from the low-frequency state to the high-frequency state as shown in Figure 2-11(b). In four-state FSK modulation a transition from the low-frequency state to the high-frequency state takes three times longer than a transition between adjacent states. While the '01' and '11' states appear to be adjacent, in reality the frequency transition must traverse through the other frequency states. Filtering of the baseband modulating signal is required to minimize the bandwidth of the modulated four-state FSK signal. This reduces modulation efficiency to less than the theoretical maximum of 2 bits/s/Hz.

In summary, there are slight inconsistencies and arbitrariness in using a phasor diagram for FSK, but FSK does have a defined constellation diagram that is closely related, but not identical, to a phasor diagram. Another difference is that a phasor diagram depends on the amplitude of the RF signal, while a constellation diagram is continuously being re-normalized to the average RF power level to maintain a constant size. With FSK modulation almost the entire modulation and demodulation paths can be implemented using analog circuitry and so was ideally suited to early cellular radios.

### 2.6.2 Gaussian Minimum Shift Keying

**Gaussian minimum shift keying (GMSK)** is the modulation scheme used in the **GSM** cellular wireless system and is a variant of **MSK** with waveform shaping coming from a Gaussian lowpass filter. It is a particular implementation of FSK modulation.

The modulation efficiency of GMSK as implemented in the GSM system (it depends slightly on the Gaussian filter parameters) is 1.35 bits/s/Hz. Unfiltered MSK has a constant RF envelope. However filtering is required

to limit the RF bandwidth and this results in amplitude variations of about 30%. This is still very little so one of the fundamental advantages of this modulation scheme is that nonlinear, power-efficient amplification can be used. GMSK is essentially a digital implementation of FM with discrete changes in the frequency of modulation with the input bitstream filtered so that the change in frequency from one state to the next is smooth. It is only at the clock ticks that the modulated signal must have the specified discrete frequency. The phase of the modulating signal is always continuous and there is no information in the phase of the modulated signal.

The ideal transitions in FSK follow a circle from one state to another as shown in Figure 2-11 so that the PMEPR of ideal FSK is 0 dB. With GMSK the transitions do not follow a circle because of the filtering and the transitions also overshoot. As such the amplitude of a GMSK modulated signal varies and the PMEPR of GMSK is 3.01 dB. This is the PMEPR for a single modulated carrier, combining multiple modulated carriers as done in a basestation increases the PMEPR. Statistically the envelopes are less likely to all align if there are multiple carriers. For example, with multi-carrier GMSK, PMEPR = 3.01 dB, 6.02 dB, 9.01 dB, 11.40 dB, 14.26 dB, and 17.39 dB for 1, 2, 4, 8, 16, and 32 carriers respectively. (These values were calculated numerically by simulating a multi-carrier system.)

GMSK and other FSK methods have the advantage that implementation of the baseband and RF hardware is relatively simple. A GMSK transmitter can use conventional frequency modulation. On receive, an FM discriminator, i.e. an FM receiver with sampling, can be used avoiding more complex  $I$  and  $Q$  demodulation.

### 2.6.3 Doppler Effect

Frequency is a physical parameter that can be established and measured with great accuracy, down to a few hertz at 1 GHz in a mobile handset for example. Thus if a receiver is stationary the frequency states at the clock ticks of an FSK modulated carrier can be measured with great accuracy. When a receiver and transmitter are moving relative to each other there will be a Doppler shift of the carrier frequency. If the relative velocity of the receiver and transmitter is  $v_s$  the Doppler shift will be

$$\Delta f = f v_s / c \quad (2.46)$$

where  $f$  is the frequency of the radio transmission and  $c$  is the speed of light. For a receiver moving at 100 km/hr receiving a 1 GHz signal from a fixed transmitter, the Doppler frequency shift is  $\Delta f = 92.6$  Hz which is much less than the 33 kHz frequency spacing of adjacent states in the FSK example above. Thus the Doppler shift is not of concern. This effective fixing of the constellation points is one of the advantages of GSM.

### 2.6.4 Summary

GSM was not the only 2G system. The 2G NADC (for North American Digital Cellular) system modulated the phase of a carrier using phase shift keying. The NADC cellular system had higher modulation efficiency than GSM yet MSK became the dominant 2G system and is still supported as a legacy modulation system in modern cellular radio. The main reason for

this is that GSM was more closely aligned with the business interests of the telephone operators of the day.

## 2.7 Carrier Recovery

Carrier recovery refers to establishing a local carrier reference signal which accurately reproduces the frequency and, with some modulation methods, the phase of the carrier of the modulated signal. All digital modulation methods require carrier recovery to establish a reference to determine the state of the carrier at the clock ticks. In addition digital modulation methods require that the timing of the clock ticks be established. Since radios using digital modulation all send packets of data, i.e. sequences of symbols, having a known sequence at the beginning of packet transmission enables the timing to be determined.

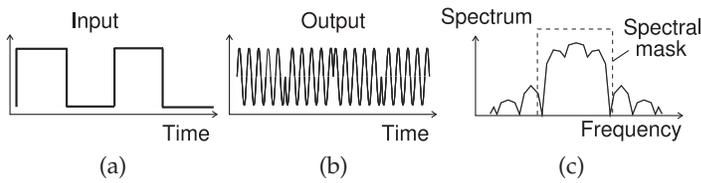
With FSK modulation the frequency at the clock ticks must be determined. This is relatively simple because the frequency at the clock ticks can be accurately measured as a local clock can be established within a few hertz because of the availability of accurate crystal references. The frequency of the received signal can still be shifted by the Doppler effect of the transmitter or receiver is moving but this is quite small compared to the frequency differences between the received states. With FSK it is not necessary to determine the phase of the carrier.

All digital modulation other than FSK modulates a carrier by shifting the carrier's phase and/or amplitude to a number of discrete states. Recovering the state of this modulated carrier requires that the phase of the carrier be recovered from the receive signal and to do this there must be a constant phase local version of the carrier. The circuits that implement the local version of the carrier are called carrier recovery circuits. These circuits modify a very stable internal oscillator in the receiver that after an initial setting of an approximate frequency, has a frequency and phase that can only change slowly. However, there must be a received signal at all times, because if the received signal falls below the noise level the carrier recovery circuit will try to track the noise. This requirement has led to a number of different modulation schemes that avoid the amplitude of the modulated signal from ever being small during a transition. This is important in 2G and 3G cellular radio but 4G and 5G cellular systems use pilot tones to achieve carrier recovery.

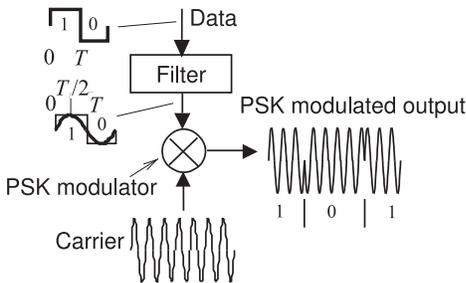
In early digital radios carrier recovery was implemented in analog circuitry and more modern radios implement carrier recovery by splitting the function between an analog oscillator signal that can be assigned to a large number of discrete states (providing coarse carrier recovery) and DSP of the (coarsely recovered) baseband signal to precisely recover the carrier signal. Thus in modern digital radios the carrier recovery circuit is implemented partially as an analog circuit and partially as a digital circuit.

## 2.8 Phase Shift Keying Modulation

There are many variations on phase shift keying (PSK) modulation with the methods differing by their spectral efficiencies, PMEPR, and suitability for carrier recovery. Compared to FSK more sophisticated digital signal processing is required to demodulate a PSK-modulated signal.



**Figure 2-12:** Binary PSK modulation: (a) modulating bitstream; (b) the modulated waveform; and (c) its spectrum after smoothing the transitions from one phase state to another.



**Figure 2-13:** A binary phase shift keying (PSK) modulation system.

### 2.8.1 Essentials of PSK

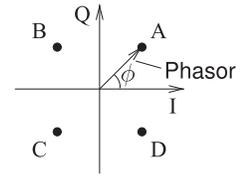
PSK is an efficient digital modulation scheme and can be simply implemented and demodulated using a phase-locked loop. The simplest scheme is binary PSK (BPSK) with two phase states. The waveform and spectrum of BPSK are shown in Figure 2-12. The incoming baseband bitstream shown in Figure 2-12(a) modulates the phase of a carrier producing the modulated signal shown in Figure 2-12(b). The spectrum of the modulated signal is shown in Figure 2-12(c). What is very interesting about this spectrum is that it approximately fills a square. So PSK modulation results in an efficient use of the spectrum. This can be contrasted with the spectrum of an FM signal shown in Figure 2-7(c), which does not fill the channel uniformly. A binary PSK modulation system is shown in Figure 2-13 where the binary input data causes 180° phase changes of the carrier. The abrupt changes in phase shown in the output modulated waveform result in more bandwidth than is necessary. However a practical PSK modulator first lowpass filters the binary data before the carrier is modulated. This filtering eliminates the abrupt changes in the phase of the modulated signal and so reduces the required bandwidth. It is the spectrum of this signal that is shown in Figure 2-12(c).

There are many variants of increasing complexity, called orders, of PSK, with the fundamental characteristics being the number of phase states (e.g. with  $2^n$  phase states,  $n$  bits of information can be transmitted) and how the phasor of the RF signal transitions from one phase state to another. PSK schemes are designed to shape the spectrum of the modulated signal to fit as much energy as possible within a spectral mask. This results in a modulated carrier whose amplitude varies (and thus has a time-varying envelope). Such schemes require highly linear amplifiers to preserve the amplitude variations of the modulated RF signal.

There are PSK methods that manage the phase transitions to achieve a constant envelope modulated RF signal but these have lower spectral efficiency. Military radios sometimes use this type of modulation scheme as it is much harder to detect and intercept communications if the amplitude of the modulated carrier is constant.

The communication limit of one symbol per hertz of bandwidth,

**Figure 2-14:** Phasor diagram of QPSK modulation. Here there are four discrete phase states of the phasor indicated by the points A, B, C, and D. The PSK modulator moves the phasor from one phase state to another. The task at the receiver is determining the phase of the phasor.



the **symbol rate**, comes from the **Nyquist signaling theorem**.<sup>1</sup> Nyquist determined that the number of independent pulses that could be put through a telegraph channel per unit of time is limited to twice the bandwidth of the channel. With a modulated RF carrier, this translates to the modulated carrier moving from one state to another in a unit of time equal to one over the bandwidth. The transition identifies a symbol, and hence one symbol can be sent per hertz of bandwidth. More accurately it could be said that the transition is a symbol rather than the end of the transition being a symbol. In PSK modulation the states are the phases of a phasor since the amplitude of the modulated signal is (ideally) constant.

The phase-shifted (i.e. phase-modulated) carrier of a PSK signal can be represented on a phasor diagram. Figure 2-14 is a phasor diagram with four phase states—A, B, C, D—and the phasor moves from one state to another under the control of the modulation circuit. What is shown here is 4-state PSK or quadrature phase shift keying (QPSK) and very often but less accurately called **quadrature phase shift keying**. The states, or symbols, are identified by their angle or equivalently by their rectangular coordinates, called I, for in-phase, and Q, for quadrature phase.

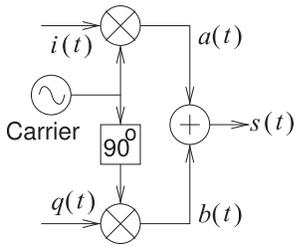
### PSK Modulation

In PSK modulation the phase of a carrier signal is set to one of a number of discrete values at the clock ticks. For example, in **QPSK** there are four discrete settings of the phase of the carrier, e.g.  $45^\circ$ ,  $135^\circ$ ,  $-135^\circ$ , and  $-45^\circ$ . Converting this to radians the discrete baseband signal is  $\phi(t) = \pi/4, 3\pi/4, 5\pi/4$ , and  $7\pi/4$ , at the clock ticks. Thus if the bandwidth of the baseband signal is 1 MHz what is shown as  $\phi(t)$  are the intended phases of the carrier every microsecond. Wave-shaping or filtering is used to provide a smooth variation of  $\phi(t)$  between the clock ticks and so the bandwidth of the modulated signal is constrained. High-order PSK modulation has more discrete states, e.g. 8-PSK has eight discrete phase states.

There are several ways to implement PSK modulation and one uses the quadrature modulator shown in Figure 2-15. The discrete baseband signal  $\phi(t)$  could be internal to a DSP which is then interpolated in time and output by the DSP's DAC as two smooth signals  $i(t) = \cos(\phi(t))$  and  $q(t) = \sin(\phi(t))$ . On a phasor diagram  $i(t)$  and  $q(t)$  at the clock ticks addresses one of QPSK's four states of the carrier's phasor, see Figure 2-14.

For PSK modulation the constellation diagram is very similar to a phasor diagram that is continuously being re-normalized to the average power of

<sup>1</sup> This theorem was discovered independently by several people and is also known as the Nyquist-Shannon sampling theorem, the Nyquist-Shannon-Kotelnikov, the Whittaker-Shannon-Kotelnikov, the Whittaker-Nyquist-Kotelnikov-Shannon (WKS), as well as the cardinal theorem of interpolation theory. The theorem states [13]: "If a function  $x(t)$  contains no frequencies higher than  $B$  hertz, it is completely determined by giving its ordinates at a series of points spaced  $1/(2B)$  seconds apart."



**Figure 2-15:** Quadrature modulator block diagram. In PSK modulation  $i(t)$  and  $q(t)$  have the same amplitude and indicate a phase  $\phi$  of the modulated carrier so that  $i(t) = \cos[\phi(t)]$  and  $q(t) = \sin[\phi(t)]$ . The particular example shows two possible values of  $I_k$  and  $Q_k$  and this indicates QPSK modulation.

the modulated signal. This a subtle but important distinction, for example, a PSK baseband signal has a constellation diagram even though the baseband signal does not have a phasor representation. PSK modulation using the block diagram shown in Figure 2-15 has a carrier that is directly input to the top multiplier and a  $90^\circ$  phase-shifted version input to the bottom multiplier. Let the carrier be  $\cos(\omega_c t)$  and so the version of the carrier input to the bottom multiplier is  $\cos(\omega_c t - \pi/2) = -\sin(\omega_c t)$ . So, with  $q(t)$  being a  $90^\circ$  phase-shifted version of  $i(t)$ , (using the identities in Section 1.A.2 of [4])

$$a(t) = \cos(\phi(t)) \cos(\omega_c t) = \frac{1}{2} [\cos(\omega_c t - \phi(t)) + \cos(\omega_c t + \phi(t))] \quad (2.47)$$

$$b(t) = \sin(\phi(t)) [-\sin(\omega_c t)] = -\frac{1}{2} [\cos(\omega_c t - \phi(t)) - \cos(\omega_c t + \phi(t))] \quad (2.48)$$

$$s(t) = a(t) + b(t) = \cos(\omega_c t + \phi(t)). \quad (2.49)$$

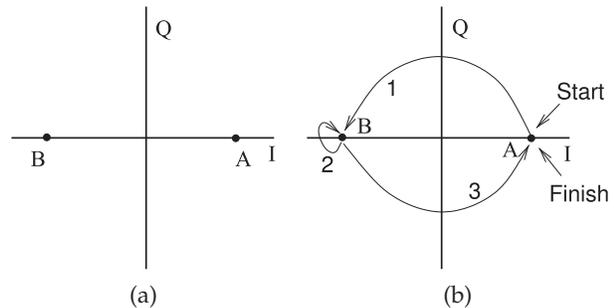
Thus  $s(t)$  is the single-sideband modulated carrier carrying information in the phase of the modulated carrier. The modulating signal  $\phi(t)$  is driven by a digital code that is designed so that  $\phi(t)$  changes at a minimum rate (it never has the same value for more than a few clock ticks). Thus there are no low frequency components of  $\phi(t)$  and thus there is no modulated signal at or very close to the carrier. Thus the carrier is suppressed but there is a sideband above and below the carrier frequency. This is SSB-SC modulation.

### 2.8.2 Binary Phase Shift Keying

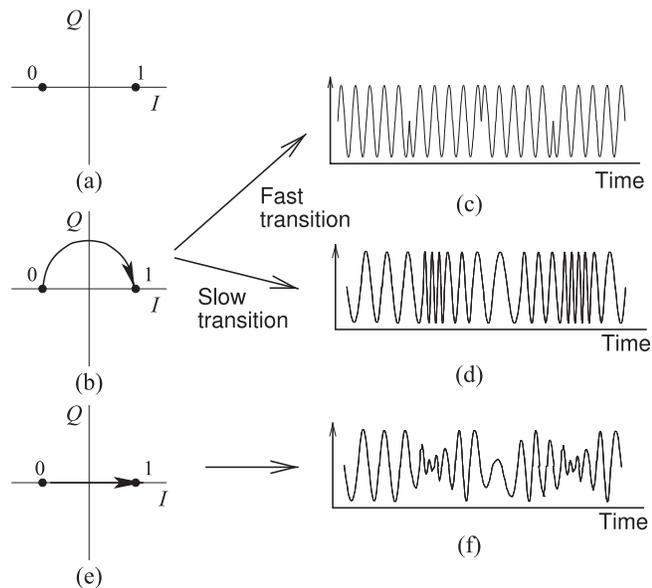
PSK uses prescribed phase shifts to define symbols, each of which can represent one, two, or more bits. **Binary Phase Shift Keying** (BPSK), illustrated in Figures 2-12 and 2-13, has two phase states and conveys one bit per symbol and is a relatively spectrally inefficient scheme, with a maximum (i.e. ideal) modulation efficiency of 1 bits/s/Hz. The reason why the practical modulation efficiency is less than this number is because the transition from one phase state to the other must be constrained to avoid the modulated signal becoming very small, and also because there are no ideal lowpass filters to filter the input binary data stream. Although it has low modulation efficiency, it is ideally suited to low-power applications. BPSK is commonly used in **Bluetooth**.

The operation of BPSK modulation can be described using the constellation diagram shown in Figure 2-16(a). The BPSK constellation diagram indicates that there are two states. These states can be interpreted as the rms values of  $i(t)$  and  $q(t)$  at the sampling times corresponding to the bit rate. The distance of a constellation point from the origin corresponds to (normalized) rms power of the pseudo-sinusoid of the modulated carrier at the sampling instant. (Normalization is with respect to the average power.) The curves in Figure 2-16(b) indicate three transitions. The states are at the ends of the transitions. If a 1, in Figure 2-16(b), is assigned to the positive  $I$  value and 0

**Figure 2-16:** BPSK modulation with constellation points A and B: (a) constellation diagram; and (b) constellation diagram with possible transitions from one phase state to the other, or possibly no change in the phase state. In practical systems the transition should not go through the origin, as then the RF signal would drop below the noise level.



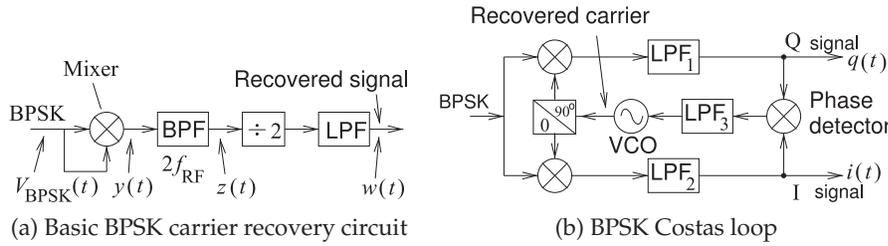
**Figure 2-17:** BPSK modulation: (a) constellation diagram; (b) constellation diagram with a constant amplitude transition; (c) time-domain waveform if the transition is fast; (d) time-domain waveform if the transition is slow; (e) constellation diagram with transition through the origin; and (f) time-domain waveform if the transition goes through the origin and is slow.



to a negative  $I$  value, then the bit sequence represented in Figure 2-16(b) is "1001."

Figure 2-17(a) is the constellation diagram of BPSK, with two symbols denoted as 0 and 1, and the trajectory of the transition from one constellation point to the other depending on the hardware used to implement the BPSK modulator. Figure 2-17(b) shows the transition from the '0' state to the '1' state (and back) while maintaining a constant amplitude. If this transition is very fast, then the waveform produced is as shown in Figure 2-17(c), where there are abrupt phase transitions and these have high spectral content. It is better to slow down the transitions, as then the waveform (shown in Figure 2-17(d)), has smooth transitions and the bandwidth of the modulated carrier is minimal. The preferred smooth transition is obtained by lowpass filtering the baseband signal. That is, the abrupt transitions in the modulated RF signal result in the modulated signal having a broad bandwidth. The graceful transition of BPSK modulation limits the bandwidth of the modulated carrier.

A simple implementation of BPSK modulation would result in direct transition from one state to the others causing the phasor to traverse the origin and the amplitude of the RF signal to become very small and less than the noise level (see Figure 2-17(e)). The resulting modulated RF waveform is



**Figure 2-18:** Block diagram of carrier recovery circuits for BPSK signals.

shown in Figure 2-17(f). This is a problem because the receiver would not be able to track the RF signal.

### Carrier Recovery

In a PSK demodulator, a local copy of the carrier must be produced to act as a reference in determining the phase of the modulated signal. The technique that produces the local copy of the unmodulated carrier is called carrier recovery. The circuit that directly implements carrier recovery of a BPSK signal is shown in Figure 2-18(a). At the clock ticks the waveform of a BPSK modulated signal is

$$v_{\text{BPSK}}(t) = A(t) \cos(\omega_{\text{RF}}t + n\pi), \quad (2.50)$$

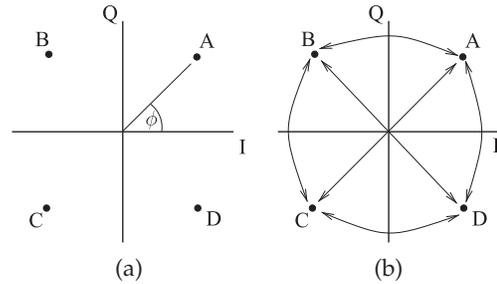
where the carrier frequency  $f_{\text{RF}} = \omega_{\text{RF}}/(2\pi)$  and  $n$  can have a value of 0 or 1. Squaring this produces a signal

$$\begin{aligned} y(t) = v_{\text{BPSK}}^2(t) &= A^2(t) \cos^2(\omega_{\text{RF}}t + n\pi) = \frac{1}{2}A^2(t) [1 + \cos(2\omega_{\text{RF}}t + n2\pi)] \\ &= \frac{1}{2}A^2(t) [1 + \cos(2\omega_{\text{RF}}t)]. \end{aligned} \quad (2.51)$$

This is a signal at twice the carrier frequency with no carrier modulation since  $n2\pi$  and 0 radians are indistinguishable. The squaring operation is performed by mixing  $v_{\text{BPSK}}(t)$  with itself. Bandpass filtering  $y(t)$  produces a signal  $z(t)$  at the second harmonic of the carrier. The divide-by-2 block is implemented using a phase-locked loop (PLL). The result is the recovered carrier,  $w(t)$ , that is used as the timing reference for sampling the demodulated I and Q components at precise times.

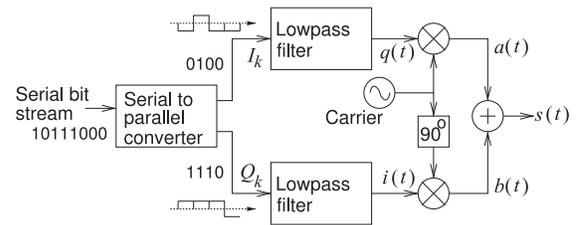
A better carrier recovery circuit than that in Figure 2-18(a) and described above is the Costas loop [14] shown in Figure 2-18(b). The BPSK Costas loop implements carrier recovery and I/Q demodulation simultaneously. In Figure 2-18(b)  $i(t)$  and  $q(t)$  are mixed to produce a signal applied at the input of the third lowpass filter, LPF<sub>3</sub>. The main function of this filter is to remove noise and to average the signal coming out of the phase detector. The output of LPF<sub>3</sub> drives a VCO in which the oscillation frequency is controlled by the applied voltage. The quadrature phase shifter then mixes the recovered carrier and a 90° shifted version with the BPSK signal.

It is critical that the signal-to-noise ratio (SNR), the ratio of the signal power to the noise power, of the BPSK signal be sufficiently large at all times or else the Costas loop will produce a noisy recovered carrier signal. If the modulated carrier becomes very small, for example when the trajectory on the constellation diagram goes through the origin (where the level of the carrier carrier falls below the noise level), the carrier will not be accurately recovered.



**Figure 2-19:** QPSK modulation: (a) constellation diagram; and (b) constellation diagram with possible transitions. Each constellation point indicates the phase,  $\phi$ , of the modulated carrier, i.e.  $\cos(\omega_c t + \phi)$  where  $\omega_c$  is the radian frequency of the carrier.

**Figure 2-20:** QPSK modulator block diagram.  $I_k$  and  $Q_k$  are similar to a stream of one-bit binary signals but are analog with a either a positive value or a negative voltage so that after lowpass filtering  $i(t)$  and  $q(t)$  each have either a positive or a negative voltage at each clock tick.



### 2.8.3 Quadra-Phase Shift Keying, QPSK

In QPSK wireless systems, modulation efficiency is obtained by sending more than one bit of information per hertz of bandwidth (i.e., more than one bit per symbol). In QPSK information is encoded in four phase states and two bits are required to identify a symbol (i.e., to identify a phase state). The constellation diagram of QPSK is shown in Figure 2-19(a) where the modulated RF carrier has four phase states identified as A, B, C, and D. So a QPSK modulator shifts the phase of the carrier to one of these phase states, and a QPSK demodulator must determine the phase of the received RF signal. The received RF signal is sampled with precise timing as determined by the recovered carrier signal. Thus two bits of information are transmitted per change of phase states. Each change of phase state requires at least 1 Hz of bandwidth with the minimum bandwidth obtained when the transition from one state to another is no faster than that required to reach the new phase state before the sampling instant. QPSK modulation is also referred to as quaternary PSK.

QPSK can be implemented using the modulator shown in Figure 2-20. In Figure 2-20, the input bitstream is first converted into two parallel bitstreams each containing half the number of bits of the original bit stream. Thus a two-bit sequence in the serial bitstream becomes one  $I_K$  bit and one  $Q_K$  bit. The  $(I_K, Q_K)$  pair constitutes the  $K$ th symbol. The bitstreams are converted into waveforms  $i(t)$  and  $q(t)$  by the wave-shaping circuit.

The constellation diagram of QPSK is the result of plotting  $I$  and  $Q$  on a rectangular graph as shown in Figure 2-19(a). All possible phase transitions are shown in Figure 2-19(b). In the absence of wave-shaping circuits,  $i(t)$  and  $q(t)$  have very sharp transitions, and the paths shown in Figure 2-16(b) occur almost instantaneously. This leads to large spectral spreads in the modulated waveform,  $s(t)$ . So to limit the spectrum of the RF signal  $s(t)$ , the shape of  $i(t)$  and  $q(t)$  is controlled; the waveform is shaped, usually by lowpass filtering. So a pulse-shaping circuit changes baseband binary information into a more smoothly varying signal. Each transition or path in Figure 2-16 represents the transfer of a symbol, with the best efficiency that can be obtained in wireless communication being one symbol per hertz

of bandwidth. Each symbol contains two bits so the maximum modulation efficiency of QPSK modulation is 2 bits/s/Hz of bandwidth. What is actually achieved depends on the wave-shaping circuits and on the criteria used to establish the bandwidth of  $s(t)$ .

### Carrier Recovery

Carrier recovery of a QPSK signal is similar to that for a BPSK signal. At the clock ticks an RF QPSK modulated signal

$$v_{\text{QPSK}}(t) = A(t) \cos(\omega_{\text{RF}}t + n\pi/2); \quad n = 0, 1, 2, 3, \quad (2.52)$$

where the carrier frequency  $f_{\text{RF}} = \omega_{\text{RF}}/(2\pi)$ . The fourth power of this produces

$$\begin{aligned} v_{\text{QPSK}}^4(t) &= A^4(t) \cos^4(\omega_{\text{RF}}t + n\pi/2) \\ &= \frac{1}{8}A^4(t) [3 + 4 \cos(2\omega_{\text{RF}}t + n\pi) + \cos(4\omega_{\text{RF}}t + n2\pi)]. \end{aligned} \quad (2.53)$$

Following bandpass filtering at  $4f_{\text{RF}}$  and then dividing the frequency by 4, the carrier is recovered. Circuits implementing this are similar to those for recovering the carrier of BPSK signals. This concept can be extended to carrier recovery for any  $M$ -PSK-modulated signal.

#### EXAMPLE 2.7

#### QPSK Modulation and Constellation

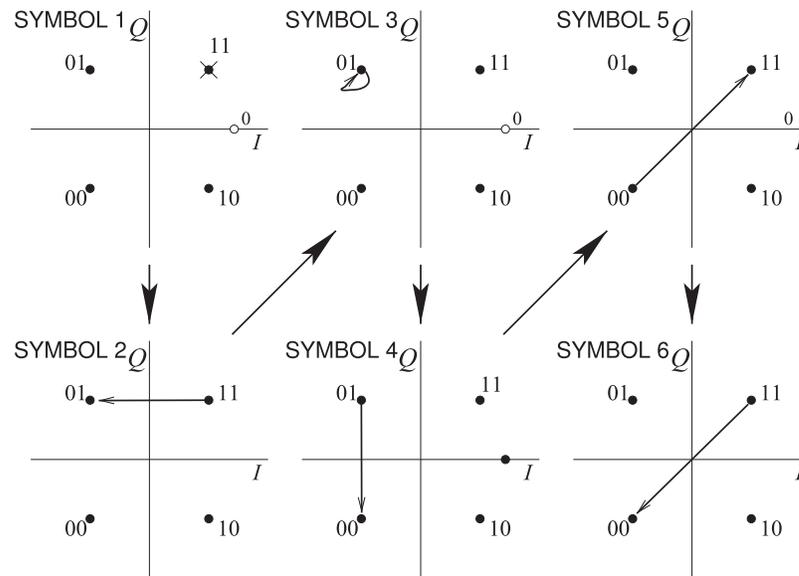
The bit sequence 110101001100 is to be transmitted using QPSK modulation. Show the transitions on a constellation diagram.

#### Solution:

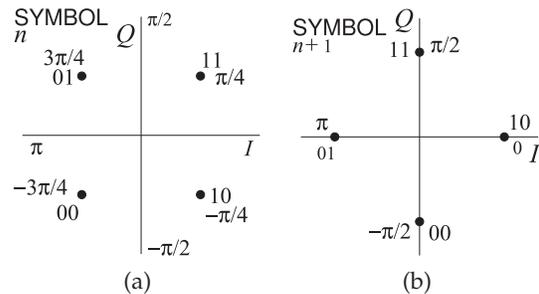
The bit sequence 110101001100 must be converted to a two-bit-wide parallel stream of symbols resulting in the sequence of symbols 11 01 01 00 11 00. The symbol 11 transitions to the symbol 01 and then to the symbol 01 and so on. The states (or symbols) and the transitions from one symbol to the next required to send the bitstream 110101001100 are shown in Figure 2-21. QPSK modulation results in the phasor of the carrier transitioning through the origin so that the average power is lower and the PMEPR is high. A more significant problem is that the phasor will fall below the noise floor, making carrier recovery almost impossible.

### 2.8.4 $\pi/4$ Quadrature Phase Shift Keying

A major objective in digital modulation is to ensure that the RF trajectory from one phase state to another does not go through the origin. The transition is slow, so that if the trajectory goes through the origin, the amplitude of the carrier will be below the noise floor for a considerable time and it will not be possible to recover the carrier reference. This is why the QPSK scheme is not used directly in 2G and 3G cellular radio. (The 4G and 5G cellular radio systems do use QPSK among other modulation schemes and use pilot tones to recover the carrier.) One of the solutions developed to address this problem is the  $\pi/4$  quadrature phase shift keying ( $\pi/4$ -QPSK) modulation scheme. In this scheme the constellation at each symbol is rotated  $\pi/4$  radians from the previous symbol, as shown in Figure 2-22. (In an alternative implementation of  $\pi/4$ -QPSK modulation the constellation diagram could



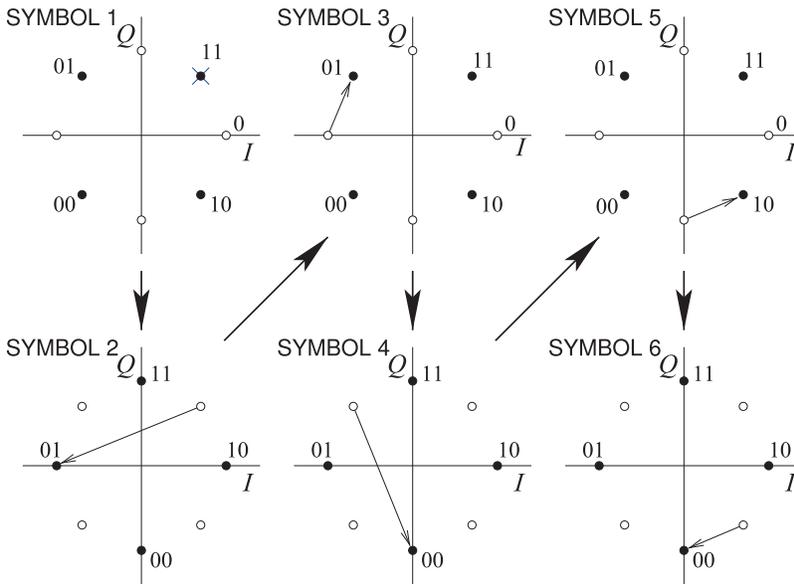
**Figure 2-21:** Constellation diagram and transitions for the bit sequence 110101001100 sent as the set of symbols 11 01 01 00 11 00 using QPSK. Note that symbols 2 and 3 are identical, so there is no transition. The SYMBOL numbers indicated reference the symbol at the end of the transition (end of the arrow). The assignment of bits to symbols (e.g., assigning the bits '11' to the symbol in the first quadrant) is arbitrary in general but the assignment of symbols is defined in a particular standard.



**Figure 2-22:** Constellation diagram of  $\pi/4$ -QPSK modulation: (a) initial constellation diagram at one symbol; and (b) the constellation diagram at the time of the next symbol.

rotate by  $\pi/4$  continuously rather than switching between conditions as described here.)

One of the unique characteristics of  $\pi/4$ -QPSK modulation is that there is always a change, even if a symbol is repeated. This helps with recovering the carrier frequency. If the binary bitstream itself (with sharp transitions in time) is the modulation signal, then the transition from one symbol to the next occurs instantaneously and hence the modulated signal has a broad spectrum around the carrier frequency. The transition, however, is slower if the bitstream is filtered, and so the bandwidth of the modulated signal will be less. Ideally the transmission of one symbol per hertz would be obtained. However, in  $\pi/4$ -QPSK modulation the change from one symbol to the next has a variable distance (and so a transition takes different times) so that the ideal modulation efficiency of one symbol per second per hertz (or 2



**Figure 2-23:** Constellation diagram states and transitions for the bit sequence 110101001000 sent as the set of symbols 11 01 01 00 10 00 using  $\pi/4$ QPSK modulation.

bits/s/Hz) is not obtained. In practice, with realistic filters and allowing for the longer transitions,  $\pi/4$ -QPSK modulation achieves 1.62 bits/s/Hz.

**EXAMPLE 2.8**  $\pi/4$ -QPSK Modulation and Constellation

The bit sequence 110101001000 is transmitted using  $\pi/4$ -QPSK modulation. Show the transitions on a constellation diagram.

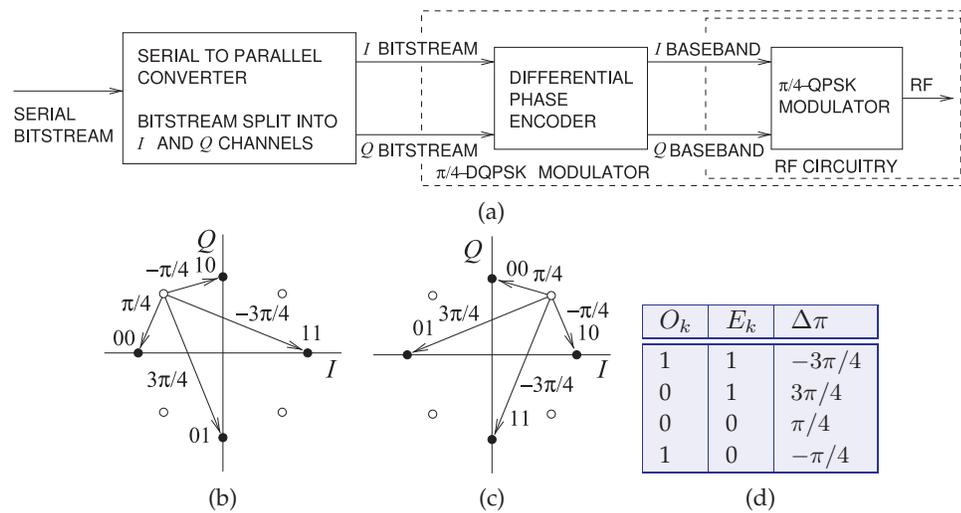
**Solution:**

The bit sequence 110101001000 must be converted to a two-bit-wide parallel stream of symbols, resulting in the sequence of symbols 11 01 01 00 10 00. The symbol 11 transitions to the symbol 01 and then to the symbol 01 and so on. The constellation diagram of  $\pi/4$ -QPSK modulation really consists of two QPSK constellation diagrams that are shifted by  $\pi/4$  radians, as shown in Figure 2-22. At one symbol (or time) the constellation diagram is that shown in Figure 2-22(a) and at the next symbol it is that shown in Figure 2-22(b). The next symbol uses the constellation diagram of Figure 2-22(a) and the process repeats. The states (or symbols) and the transitions from one symbol to the next that are required to send the bitstream 110101001000 are shown in Figure 2-23.

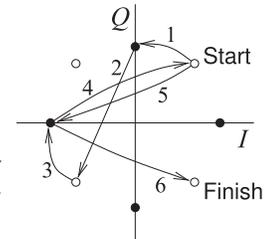
**2.8.5 Differential Quadra Phase Shift Keying, DQPSK**

Multiple transmission paths, or **multipaths**, due to reflections result in constructive and destructive interference and can result in rapid additional phase rotations. Thus relying on the phase of a phasor at the symbol sample time, at the clock ticks, to determine the symbol transmitted is prone to error. When an error results at one symbol, this error accumulates when subsequent symbols are extracted. The solution is to use encoding, and one of the simplest encoding schemes is differential phase encoding. In this scheme the information of the modulated signal is contained in changes in phase rather than in the absolute phase. That is, the transition defines the symbol rather than the end point of the transition.

The  $\pi/4$ -DQPSK modulation scheme is a differentially encoded form of



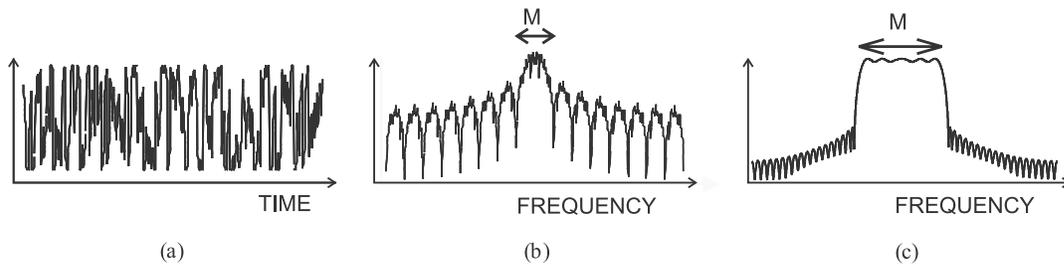
**Figure 2-24:** A  $\pi/4$ -DQPSK modulator: (a) differential phase encoder with a  $\pi/4$ -QPSK modulator; (b) constellation diagram of  $\pi/4$ -DQPSK; (c) a second constellation diagram; and (d) phase changes in a  $\pi/4$ -DQPSK modulation scheme. Note that the information is in the phase change rather than the phase state.



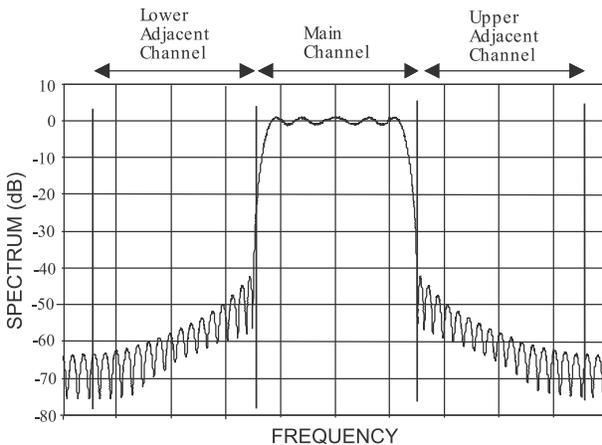
**Figure 2-25:** Constellation diagram of  $\pi/4$ -DQPSK modulation showing six symbol intervals coding the bit sequence 000110110101.

$\pi/4$ -QPSK. The  $\pi/4$ -DQPSK scheme incorporates the  $\pi/4$ -QPSK modulator and an encoding scheme, as shown in Figure 2-24(a). The scheme is defined with respect to its constellation diagram, shown in Figure 2-24(b) and repeated in Figure 2-24(c) for clarity. The D indicates **differential coding**, while the  $\pi/4$  denotes the rotation of the constellation by  $\pi/4$  radians from one interval to the next. This can be explained by considering Figure 2-24(a). A four-bit stream is divided into two quadrature **nibbles** of two bits each. These nibbles independently control the  $I$  and  $Q$  encoding, respectively, so that the allowable transitions rotate according to the last transition. The information or data is in the phase transitions rather than the constellation points themselves. The relationship between the symbol value and the transition is given in Figure 2-24(d). For example, the transitions shown in Figure 2-25 for six successive time intervals describes the input bit sequence 000110110101. Its waveform and spectrum are shown in Figure 2-26. More detail of the spectrum is shown in Figure 2-27. In practice with realistic filters and allowing for the longer transitions,  $\pi/4$ -DQPSK modulation achieves a modulation efficiency of 1.62 bits/s/Hz, the same as  $\pi/4$ -QPSK, but of course with greater resilience to changes in the transmission path.

In a differential scheme, the data transmitted are determined by



**Figure 2-26:** Details of digital modulation obtained using differential phase shift keying ( $\pi/4$ -DQPSK): (a) modulating waveform; (b) spectrum of the modulated carrier, with M denoting the main channel; and (c) details of the spectrum of the modulated carrier focusing on the main channel.

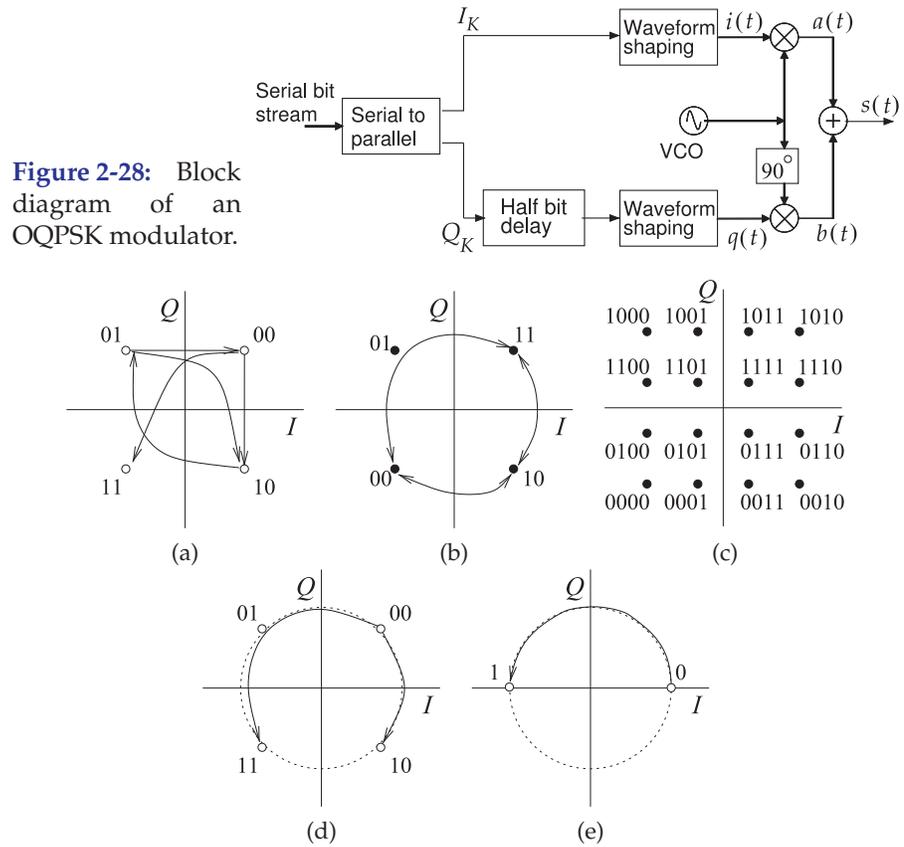


**Figure 2-27:** Detailed spectrum of a  $\pi/4$ -DQPSK signal showing the main channel and lower and upper adjacent channels.

comparing a symbol with the previously received symbol, so the data are determined from the change in phase of the carrier rather than the actual phase of the carrier. This process of inferring the data actually sent from the received symbols is called decoding. When  $\pi/4$ -DQPSK encoding was introduced in the early 1990s the DSP available for a mobile handset had only just reached sufficient complexity. Today, encoding is used with all digital radio systems and is more sophisticated than just the differential scheme of DQPSK. There are new ways to handle carrier phase ambiguity. The sophistication of modern coding schemes is beyond the scope of the hardware-centric theme of this book.

### 2.8.6 Offset Quadra Phase Shift Keying, OQPSK

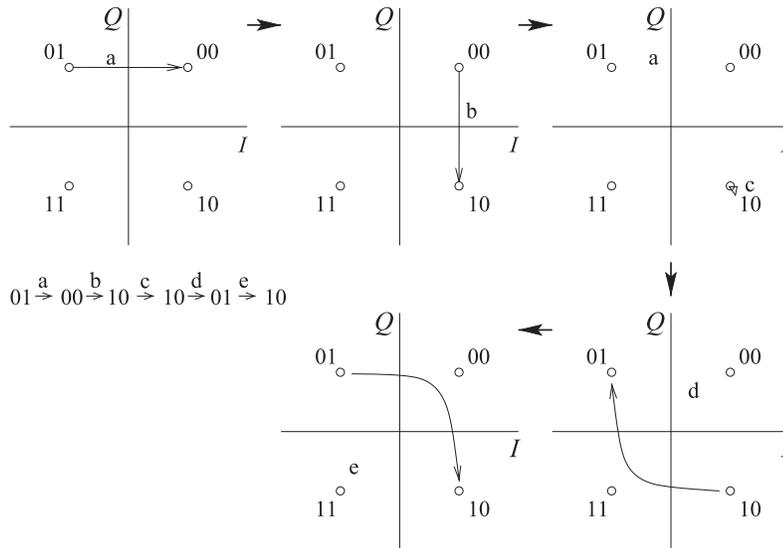
The **offset quadra phase shift keying (OQPSK)** modulation scheme avoids transitions passing through the origin on the constellation diagram (see Figure 2-29(a)). As in all QPSK schemes, there are two bits per symbol, but now one bit is used to immediately modulate the RF signal, whereas the other bit is delayed by half a symbol period, as shown in Figure 2-28. The maximum phase change for a bit transition is  $90^\circ$ , and as  $Q_K$  is delayed, a total phase change of approximately  $180^\circ$  is possible during one symbol. The



**Figure 2-29:** Constellation diagrams for various modulation formats: (a) OQPSK; (b) GMSK; (c) 16-QAM; (d) SOQPSK (also FOQPSK); and (e) SBPSK.

constellation diagram is shown in Figure 2-29(a).

The OQPSK modulator can be implemented using relatively simple electronics with a digital delay circuit delaying the  $Q$  bit by half a symbol period and lowpass filters shaping the  $I$  and  $Q$  bits. The OQPSK scheme is also called **staggered quadrature phase shift Keying (SQPSK)**. Better performance can be obtained by using DSP to shape the  $I$  and  $Q$  transitions so that they change smoothly and the phasor trajectory nearly follows a circle. Consequently  $I$  and  $Q$  change together, but in such a manner that the PMEPR is maintained close to 0 dB. Two modulation techniques that implement this are the **shaped offset QPSK (SOQPSK)** and the **Fehér QPSK (FQPSK)** schemes. The constellation diagrams for SOQPSK and FQPSK are shown in Figure 2-29(d). These are constant envelope digital modulation schemes. As with OQPSK, the  $Q$  bit is delayed by one-half of a symbol period and the  $I$  and  $Q$  baseband signals are shaped by a half-sine filter. The advantage is that high-efficiency saturating amplifier designs can be used and battery life extended. There is a similar modulation format called **shaped binary phase shift keying (SBPSK)** which, as expected, has two constellation points as shown in Figure 2-29(e). SOQPSK, FQPSK, and SBPSK are **continuous phase modulation (CPM)** schemes, as the phase



**Figure 2-30:** Constellation diagram of OQPSK modulation for the bit sequence 010010100110.

never changes abruptly. Instead, the phase changes smoothly, achieving high modulation efficiency and maintaining a constant envelope. Implementation of the receiver, however, is complex. CPM schemes have good immunity to interference and are commonly used in military systems.

**EXAMPLE 2.9** OQPSK Modulation

Draw the constellation diagrams for the bit sequence 010010100110 using OQPSK modulation.

**Solution:**

The bit sequence is first separated into the parallel stream 01-00-10-10-01-10. The *I* bit changes first, followed by the *Q* bit delayed by half of the time of a bit. Five constellation diagrams are shown in Figure 2-30 with the transitions sending the bit sequence.

**2.8.7  $3\pi/8$ -8PSK, Rotating Eight-State Phase Shift Keying**

The  $3\pi/8$ -8PSK modulation scheme is similar to  $\pi/4$ -DQPSK in the sense that rotation of the constellation occurs from one time interval to the next. This time, however, the rotation of the constellation from one symbol to the next is  $3\pi/8$ . This modulation scheme is used in the **enhanced data rates for GSM evolution (EDGE)** system, and provides three bits per symbol (ideally) compared to GMSK used in GSM which has two bits per symbol (ideally). With some other changes, GSM/EDGE provides data transmission of up to 128 kbits/s, faster than the 48 kbits/s possible with GSM.

Quadrature modulation schemes with four states, such as QPSK, have two *I* states and two *Q* states that can be established by lowpass filtering the *I* and *Q* bitstreams. For higher-order modulation schemes such as 8-PSK, this approach will not work. Instead,  $i(t)$  and  $q(t)$  are established in the DSP unit and then converted using a DAC to generate the analog signals applied

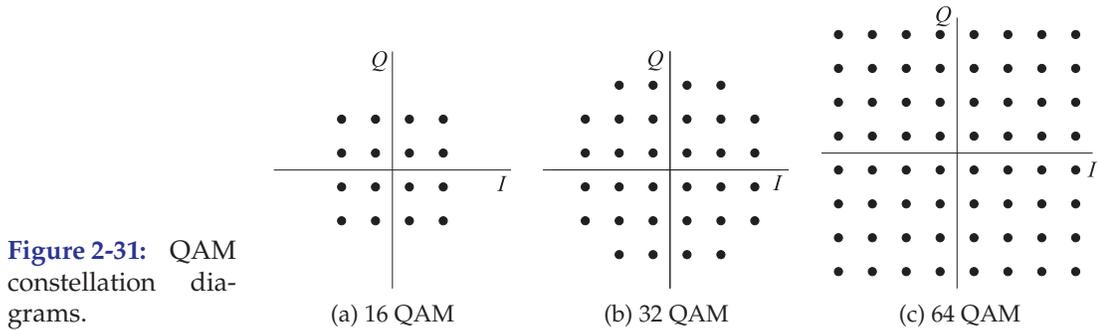


Figure 2-31: QAM constellation diagrams.

to the hardware modulator. Alternatively the modulated signal is created directly in the DSP and a DAC converts this to an IF and a hardware mixer up-converts this to RF. QAM

### 2.8.8 Summary

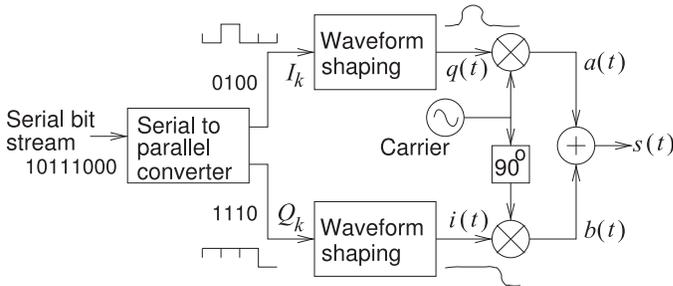
PSK modulation is implemented in many radio standards including all cellular standards after 2G. There was a 2G system that used  $\pi/4$ DQPSK but that is no longer supported. The modern radio standards support many modulation formats but in high interference situations BPSK, QPSK and 8-PSK have the best performance. While QPSK was dismissed in 2G and 3G because of difficulties with carrier recovery, 4G and 5G have another method for implementing carrier recovery which allows QPSK on its own to be used. GMSK is still supported by modern cellular phones but the infrastructure, i.e. basestations, are starting to be retired.

Most of the modulation schemes described in this section were introduced as optimum trade-offs of modulation efficiency, resistance to interference, and hardware complexity. Some, such as BPSK, draw very little power and are suited to the internet-of-things (IoT) applications which must have a battery lifetime of ten years.

## 2.9 Quadrature Amplitude Modulation

The digital modulation schemes described so far modulate the phase or frequency of a carrier to convey digital data and the constellation points lie on a circle of constant amplitude. The effect of this is to provide some immunity to amplitude changes to the signal. However, much more information can be transmitted if the amplitude is varied as well as the phase. With considerable signal processing it is possible to reliably use quadrature amplitude modulation (QAM) in which amplitude and phase are both changed.

A 16-state rectangular QAM, 16-QAM, constellation is shown in Figure 2-29(c). Since there are 16 ( $= 2^4$ ) symbols the values of 4 binary bits are uniquely specified by each symbol. In Figure 2-29(c) a gray-scale assignment of 4 bit values is shown. Several QAM schemes are shown in Figure 2-31. These constellations can be produced by separately amplitude modulating an  $I$  carrier and a  $Q$  carrier. Both carriers have the same frequency but are  $90^\circ$  out of phase. The two carriers are then combined so that the fixed carrier is suppressed. The most common form of QAM is square QAM, or rectangular QAM with an equal number of  $I$  and  $Q$  states. The most common forms are



**Figure 2-32:** QAM modulator block diagram. In QAM modulation  $i(t)$  and  $q(t)$  address the real and imaginary components of a phasor. The wave-shaping block ensures that the symbol has the correct amplitude and phase at each clock tick.

Modulation	bits/s/Hz
BPSK (ideal)	1
BFSK (actual)	1
QPSK (ideal)	2
GMSK (an actual FSK method)	1.354
$\pi/4$ -DQPSK (an actual QPSK method)	1.63
8-PSK (ideal)	3
$3\pi/8$ -8PSK (an actual 8PSK method)	2.7
16-QAM (ideal)	4
16-QAM (actual)	2.98
32-QAM (ideal)	4
32-QAM (actual)	3.35
64-QAM (ideal)	6
64-QAM (actual)	4.47
256-QAM (ideal)	8
256-QAM (actual, satellite & cable TV)	6.33
512-QAM (ideal)	9
1024-QAM (ideal)	10
2048-QAM (ideal)	11

**Table 2-2:** Modulation efficiencies of various modulation formats in bits/s/Hz (bits per second per hertz). The maximum (or ideal) modulation efficiencies obtained by modulation schemes (e.g., BPSK, BFSK, 64-QAM, 256-QAM) result in broad spectra. Actual modulation efficiencies achieved are less in an effort to manage bandwidth. For example, the values for  $\pi/4$ -DQPSK and  $3\pi/8$ -8PSK are actual. This reduction from ideal arises since symbol transitions are of different lengths and length corresponds to time durations. Since the symbol interval is fixed, it is the longest path that determines the bandwidth required.

16-QAM, 64-QAM, and 128-QAM, in 4G, and 256-QAM additionally in 5G. The constellation points are closer together with high-order QAM and so are more susceptible to noise and other interference. Thus high-order QAM can deliver more data, but less reliably than lower-order QAM.

The constellation in QAM can be constructed in many ways, and while rectangular QAM is the most common form, non rectangular schemes exist; for example, having two PSK schemes at two different amplitude levels. While there are sometimes minor advantages to such schemes, square QAM is generally preferred as it requires simpler modulation and demodulation.

One possible architecture of a QAM modulator is shown in Figure 2-32 and this can only be implemented in DSP since it is not sufficient to use analog lowpass filtering to implement the wave-shaping function as the  $i(t)$  and  $q(t)$  must be precisely the real and imaginary parts of the symbol at each clock tick.

### 2.10 Digital Modulation Summary

The modulation efficiencies of various digital modulation schemes are summarized in Table 2-2. For example, in 1 kHz of bandwidth the  $3\pi/8$ -8PSK scheme (supported in 3G cellular radio) transmits 2700 bits.

beta.69: It is critical to control interference in digital radio so that the error

in digital transmission is no more than one bit per symbol. Error correction can then be used to provide error-free digital communications.

The modulation efficiency of an actual modulation method is less than the ideal (see Table 2-2). With digital modulation wave-shaping at baseband is required to constrain the spectrum of the RF-modulated signal. Thus it will take different times for the phasor to make the transition from one symbol to another; to achieve longer transitions in the same time interval requires more bandwidth than that required for shorter transitions. As a result, the modulation efficiency of modulation methods other than binary methods will be less than the ideal. So in a QPSK-like scheme, 2 bits per symbol are achievable, but the longest transition takes the most time, so the bandwidth needs to be increased so that the transition is completed in time (i.e., in a fixed time equal to one over the bandwidth). Various modulation methods have relative merits in terms of modulation efficiency, tolerance to fading (due to destructive interference), carrier recovery, spectral spreading in nonlinear circuitry, and many other issues that are the purview of communication system theorists.

## 2.11 Interference and Distortion

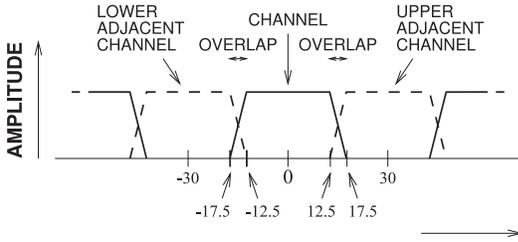
Demodulation of a received signal is equivalent to reconstructing the original constellation diagram of the modulation signal. Errors are caused by interference, noise, and distortion. Commonly all three effects are lumped together and called interference and treated as though they have noise-like Gaussian randomness. Ideally there is no interference so that a receiver correctly detects the correct symbol upon demodulation. Interference will result in perhaps an incorrect symbol choice and thus error. Errors can be reduced by increasing the signal level at the transmitter thus increasing the signal-to-interference ratio. This comes with a price as increasing the signal level results in higher levels of interference for other radios. The solution arrived at is to ensure that if there is an error, then the incorrectly detected symbol is no more than one symbol away from the actual symbol. This means that there is at most one bit in error while a symbol can carry multiple bits of information. Error correction coding can then be used to eliminate errors.

Émile Baudot used gray codes in telegraphy in 1878 [15]. The name derives from Frank Gray who used them in a pulse code modulation coding scheme [16].

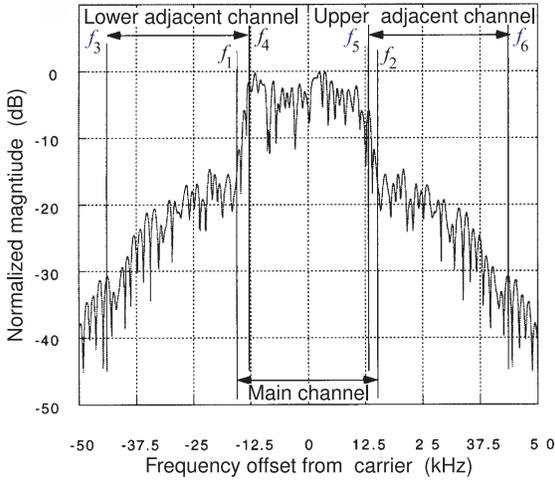
In QAM, symbols are assigned to constellation points using a **Gray code** in which nearest neighbor symbols change by only one bit [17], e.g. see Figure 2-29(c). Thus there is only one bit out of many that will be in error if there is noise and interference. If the error is greater then a lower-order of modulation is used so that if a symbol is incorrectly detected then the incorrect symbol is at most one symbol away from the actual transmitted symbol.

### 2.11.1 Cochannel Interference

The minimum signal detectable in conventional wireless systems is determined by the **signal-to-interference ratio (SIR)** at the input to a receiver, where interference refers to noise as well as interference from other radios. In cellular wireless systems the dominant interference is from other radios and is called cochannel interference. The degree to which cochannel interference can be controlled has a large effect on system capacity. Control of cochannel interference is largely achieved by controlling the power levels



**Figure 2-33:** Adjacent channels showing overlap in the AMPS and DAMPS cellular systems.



**Figure 2-34:** Spectrum defining adjacent channel and main channel integration limits using a  $\pi/4$ -DQPSK modulation scheme.

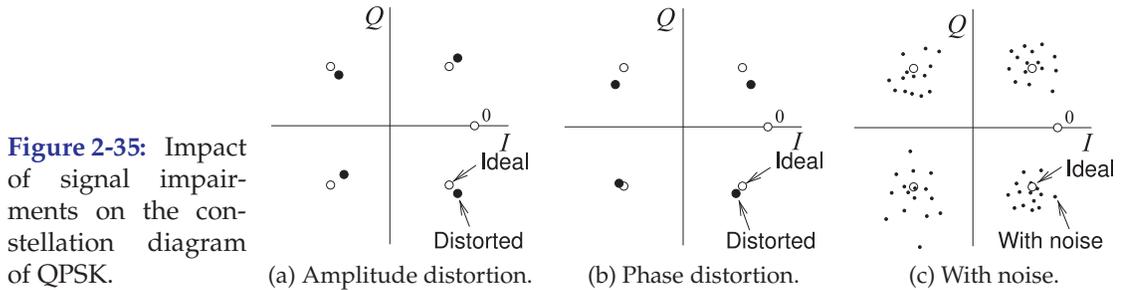
at the base station and at the mobile units.

### 2.11.2 Adjacent Channel Interference

Adjacent channel interference results from several factors. Since filtering is non ideal, there is inherent overlap of neighboring channels (Figure 2-33). For this reason, adjacent channels are assigned to different cells. The nonlinear behavior of transmitters also contributes to adjacent channel interference. Thus characterization of nonlinear phenomena is important in RF design. The spectrum of a signal modulated using a QPSK scheme is shown in Figure 2-34. The signal between frequencies  $f_1$  and  $f_2$  is due to the digital modulation scheme itself. Most of the signal outside this region is due to nonlinear effects and is called spectral regrowth, a process similar to distortion of a two-tone signal. Using the frequency limits defined in Figure 2-34, the lower channel ACPR is defined as

$$ACPR_{ADJ,LOWER} = \frac{\text{Power in lower adjacent channel}}{\text{Power in main channel}} = \frac{\int_{f_3}^{f_4} X(f)df}{\int_{f_1}^{f_2} X(f)df}, \quad (2.54)$$

where  $X(f)$  is the spectral power density of the RF signal. Upper channel ACPR,  $ACPR_{ADJ,UPPER}$ , is similarly defined. When ACPR is being referred to without indicating whether it is upper or lower ACPR, the larger (i.e., the worse) of  $ACPR_{ADJ,LOWER}$  and  $ACPR_{ADJ,UPPER}$  is used. ACPR is usually expressed in decibels, and while the definition is such that ACPR will be less than one, when expressed in decibels, a positive number is often used (e.g., 20 dB for an ACPR of 0.01 rather than the correct  $-20$  dB, so be careful).



### 2.11.3 Noise, Distortion, and Constellation Diagrams

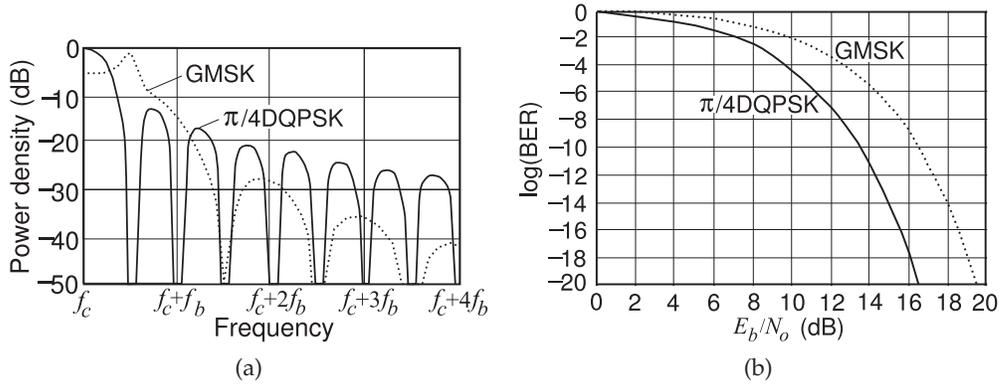
Noise and nonlinear distortion affect the ability to correctly demodulate signals and determine the transmitted symbols. These distortion effects can be described by their effect on the received constellation diagram, see Figure 2-35, which shows the state of the system at the sampling instant determined by the recovered carrier.

Figure 2-35(a) shows the effect of amplitude distortion errors on the demodulated signal. Sampling of the received signal will be a distorted constellation point that does not correspond to the ideal constellation point. A decision must be made by the DSP unit as to which ideal constellation point corresponds to the distorted constellation point. The effect of phase distortion is shown in Figure 2-35(b). Both amplitude and phase distortion could occur in the transmitter or receiver, or be the result of effects in the signal path. Figure 2-35(c) shows the effect of noise on signal impairment. Again the constellation point extracted from the RF signal is affected by noise and the sampled and ideal constellation points do not coincide. Associating the constellation point extracted from sampling the received RF waveform with the wrong constellation point creates a symbol error and thus a bit error. Errors in recovering the carrier further distort the constellation diagram. All mobile digital radio systems adjust the level of the transmitted RF signal, and additionally in 4G and 5G cellular radio change the order of modulation, so an acceptable BER is obtained. Using more power than necessary reduces battery life and causes additional interference in other radios.

### 2.11.4 Comparison of GMSK and $\pi/4$ DQPSK Modulation

This section presents the results of the type of analysis that is performed to characterize modulation methods. There is a large body of literature documenting the performance of modulation schemes and is usually the result of an assumed error model, that is the type of interference and the statistical description of that interference, and then numerical simulations. This section compares GMSK and  $\pi/4$ DQPSK modulation methods, the first widely-used cellular digital modulation methods.

The constellations of 4-state GMSK, see Figure 2-11(b), and  $\pi/4$ DQPSK, see Figure 2-19 are the same except that with  $\pi/4$  DQPSK the constellation rotates every clock tick. In GMSK the amplitude of the modulated carrier remains almost constant and the frequency of the carrier varies slowly. With GMSK the symbols correspond to different RF frequencies so it is not possible for symbols with frequencies at opposite ends of the frequency range to transition directly without traversing the other symbol points. This long transition results in reduction of GMSK's modulation efficiency. With



**Figure 2-36:** Comparison of GMSK and  $\pi/4$ DQPSK: (a) power spectral density as a function of frequency deviation from the carrier; and (b) BER versus signal-to-noise ratio (SNR) as  $E_b/N_o$  (or energy per bit divided by noise per bit).

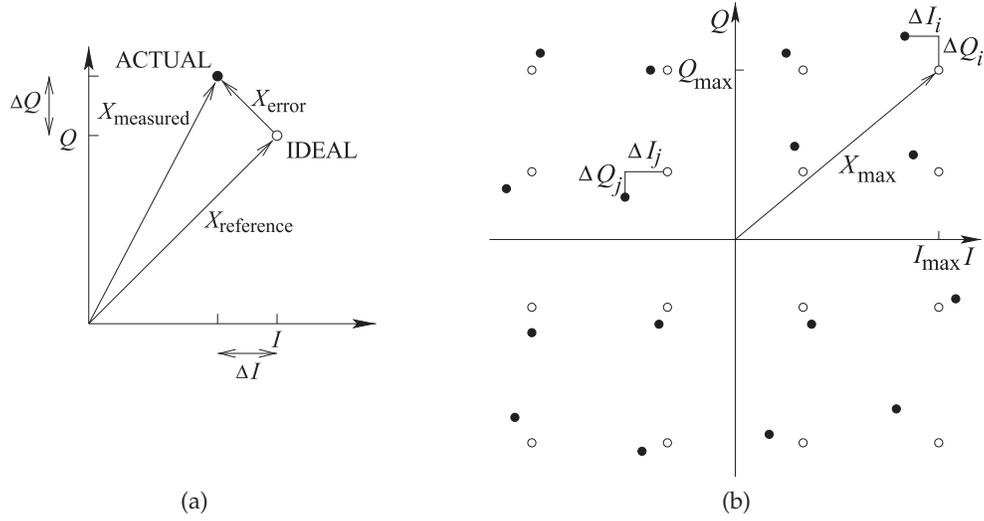
$\pi/4$ DQPSK there are direct transitions and the magnitude of the RF phasor does not stay constant. So while higher modulation efficiency is obtained compared to GMSK,  $\pi/4$ -DQPSK has a significantly time-varying envelope.

The modulation format used impacts the choice of circuitry, battery life, and the tolerance of the system to noise. Figure 2-36 contrasts the power density versus frequency and **bit error rate (BER)** of four-state GMSK and QPSK modulation. In Figure 2-36(a),  $f_c$  is the carrier frequency and  $f_b$  is the bit frequency, and it is seen that GMSK and QPSK have different spectral shapes. Most of the energy is contained within the bandwidth defined by half the bit frequency (this is the symbol frequency since there are two bits per symbol). At multiples of the bit frequency, the power density with GMSK is much lower than with QPSK, resulting in less interference (**adjacent channel interference [ACI]**) with radios in adjacent channels. This is an important metric with radios that is captured by the **adjacent channel power ratio (ACPR)**, the ratio of the power in the adjacent channel to the power in the main channel. Another important metric is the BER which is increased by noise in the main channel with different modulation formats differing in their susceptibility to interference. The level of noise is captured by the ratio of the power in a bit,  $E_b$ , to the noise power,  $N_o$ , in the time interval of a bit. This ratio,  $E_b/N_o$  (read as E B N O), is directly related to the **signal-to-noise ratio (SNR)**. In particular, consider the plot of the BER against  $E_b/N_o$  shown in Figure 2-36(b), where it can be seen that  $\pi/4$ DQPSK is less susceptible to noise than GMSK.

### 2.11.5 Error Vector Magnitude

The error vector magnitude (**EVM**) metric characterizes the accuracy of the waveform at the sampling instances and so is directly related to the BER in digital radio. EVM captures the combined effect of amplifier nonlinearities, amplitude and phase imbalances of separate  $I$  and  $Q$  signal paths, in-band amplitude ripple (e.g., due to filters), noise, non ideal mixing, non ideal carrier recovery, and DAC inaccuracies.

The EVM is a measure of the departure of a sampled phasor from the ideal



**Figure 2-37:** Partial constellation diagram showing quantities used in calculating EVM: (a) definition of error and reference signals; and (b) error quantities used when constellation points have different powers.

phasor located at the constellation point (see Figure 2-37(a)). Introducing an error vector,  $X_{\text{error}}$ , and a reference vector,  $X_{\text{reference}}$ , that points to the ideal constellation point, the EVM is defined as the ratio of the magnitude of the error vector to the reference vector so that

$$\text{EVM} = \frac{|X_{\text{error}}|}{|X_{\text{reference}}|}. \quad (2.55)$$

Expressing the error and reference in terms of the powers  $P_{\text{error}}$  and  $P_{\text{reference}}$ , respectively, enables EVM to be expressed as

$$\text{EVM} = \sqrt{\frac{P_{\text{error}}}{P_{\text{reference}}}}, \quad (2.56)$$

in decibels,  $\text{EVM}_{\text{dB}} = 10 \log \frac{P_{\text{error}}}{P_{\text{reference}}} = 20 \log \frac{|X_{\text{error}}|}{|X_{\text{reference}}|}$ ; (2.57)

or as a percentage,  $\text{EVM}(\%) = \frac{X_{\text{error}}}{X_{\text{reference}}} \cdot 100\%$ . (2.58)

If the modulation format results in constellation points having different powers (e.g., with 16-QAM), the constellation point with the highest power is used as the reference and the error at each constellation point is averaged. With reference to Figure 2-37, and with  $N$  constellation points,

$$\text{EVM} = \sqrt{\frac{\frac{1}{N} \sum_{i=1}^N (\Delta I_i^2 + \Delta Q_i^2)}{X_{\text{max}}^2}}, \quad (2.59)$$

where  $|X_{\text{max}}|$  is the magnitude of the reference vector to the most distant constellation point and  $\Delta I_i$  and  $\Delta Q_i$  are the  $I$  and  $Q$  offsets of the actual

constellation point and the ideal constellation point. Note that for QAM the constellation diagram corresponds to a phasor diagram that is being continuously normalized to the average received RF signal level. In the constellation diagram the  $I$  and  $Q$  coordinates correspond to RMS quantities. Thus  $|X_{\max}|$  is an RMS quantity. EVM is traditionally expressed as a percentage.

A similar measure of signal quality is the **modulation error ratio (MER)**, a measure of the average signal power to the average error power. In decibels it is defined as

$$\text{MER}_{\text{dB}} = 10 \log \frac{\frac{1}{N} \sum_{i=1}^N (I_i^2 + Q_i^2)}{\frac{1}{N} \sum_{i=1}^N (\Delta I_i^2 + \Delta Q_i^2)} = 10 \log \frac{\sum_{i=1}^N (I_i^2 + Q_i^2)}{\sum_{i=1}^N (\Delta I_i^2 + \Delta Q_i^2)}. \quad (2.60)$$

The advantage of the MER is that it relates directly to the SNR.

Another quantity related to both the EVM and MER concepts is the implementation margin,  $k$ . The implementation margin is a measure of the performance of particular hardware and is developed by design groups based on experience with similar designs. The required EVM can be estimated from the hardware implementation margin:

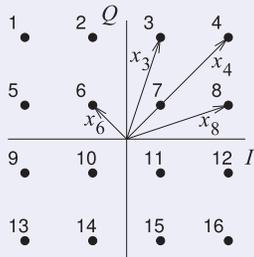
$$\text{EVM}_{\text{required}} = \sqrt{\frac{k}{\text{SNR} \cdot \text{PMEPR}}}. \quad (2.61)$$

In decibels,  $\text{EVM}_{\text{db, required}} = k_{\text{dB}} - \text{SNR}_{\text{dB}} - \text{PMEPR}_{\text{dB}}. \quad (2.62)$

**EXAMPLE 2.10** Modulation Error Ratio

A 16-QAM modulated signal has a maximum RF phasor rms value of 4 V. If the noise on the signal has an rms value of 0.1 V, what is the modulation error ratio of the modulated signal?

**Solution:**



The distance from the origin to each of the constellation points must be determined, but because of symmetry only the distances  $x_3, x_4, x_7 (= x_6)$  and  $x_8$  need to be calculated. The maximum RF phasor amplitude is 4 V, so the length  $x_4 = 4$ , with components

$$\begin{aligned} I_4 = Q_4 &= \sqrt{x_4^2/2} = 2.828 \\ I_6 &= \left(-\frac{1}{3}\right) \cdot I_4 = -0.943 = -Q_6, \text{ so } x_6 = 1.333 \\ I_3 &= 0.943; Q_3 = 2.828, \text{ so } x_3 = 2.981 = x_8 \end{aligned}$$

and 
$$\text{MER} = \frac{\sum_{i=1}^{16} (x_i^2)}{\sum_{i=1}^{16} (x_{\text{noise}}^2)}$$

The calculation is simplified by considering just one quadrant and the noise,  $x_{\text{noise}} = 0.1$ , is the same for each constellation point.

$$\begin{aligned} \text{MER} &= \frac{4(x_3^2 + x_4^2 + x_6^2 + x_8^2)}{16 \cdot x_{\text{noise}}^2} \\ &= \frac{4(2.981^2 + 4^2 + 1.333^2 + 2.981^2)}{16 \cdot 0.1^2} = \frac{142.2}{0.16} = 888.8. \\ \text{MER}_{\text{dB}} &= 10 \log(888.8) = 29.5 \text{ dB} \end{aligned}$$

Compare this to the SNR calculated at the individual constellation points:

Point 3,  $\text{SNR} = x_3^2/0.1^2 = 888.6 = 29.5 \text{ dB}$

Point 4,  $\text{SNR} = x_4^2/0.1^2 = 1600 = 32.0 \text{ dB}$

Point 7,  $\text{SNR} = x_7^2/0.1^2 = 177.7 = 22.5 \text{ dB}$ .

## 2.12 Summary

All modern modulation methods impress information on a sinusoidal carrier that is at a high enough frequency that it can be easily transmitted. There are many modulation techniques with the choice of which to use based on the technology available to implement the modulation scheme, the tolerance of the modulation scheme to interference, how efficiently the modulation scheme uses the EM spectrum, and the amount of DC power consumed. In military communications it is also important that a modulation scheme produce a noise-like signal that is difficult to detect and intercept.

The first widely adopted modulation schemes produced simple pulses as used in wireless telegraphy. The tolerance to interference was achieved through the relatively slow transmission of bits, and hence redundancy. More information was transmitted when amplitude modulation was introduced to superimpose voice on a carrier. With this scheme, interference was always a problem, and once interference appeared on a signal it could not be removed nor suppressed. The first significant advance in modulation techniques was the invention of frequency modulation. In this scheme a narrowband analog modulating signal (e.g. voice) became a relatively wideband frequency-modulated RF signal. When the modulated signal was received and demodulated, the wide bandwidth of the modulated signal was collapsed to the original relatively narrowband modulating signal. The demodulation process combined correlated components of the modulated signal and the uncorrelated components, noise and interference, were suppressed.

The introduction of digital modulation was a significant advance in the suppression of interference. Digital information was now being transmitted, and errors in the data caused by interference and noise would be completely unacceptable. The solution was to embed error-correcting codes in the data so that if a manageable number of bits were lost in the transmitted signal, the original data could still be fully recovered. As a result, digital radio (using digital modulation) could be used in situations with even more distortion than was acceptable in analog radio (using analog modulation). If interference is low, then today's wireless systems use high-order modulation switching to lower-order (and less spectrally efficient) modulation when necessary to cope with higher interference.

Several important metrics are used to provide a measure of the signal characteristics. The crest factor, peak-to-average ratio, and peak-to-mean envelope power ratio (PMEPR) are all indicators of how much care must be given to nonlinear circuit design, especially to amplifier and mixer design. Amplifiers must operate so that the peak signal is amplified with minimal distortion. It is the peak signal that determines the DC power drawn by an amplifier. However, the average RF output power of an RF front end is determined by the mean of the envelope. So a high PMEPR signal will result in lower amplifier efficiency.

Many of the techniques described in this chapter for modulating and demodulating RF signals were presented as circuit techniques. However many modern phones support multiple standards and hardware implementation would require multiple copies of similar versions of analog circuits. Today it is more cost effective to perform most of the operations in a DSP unit. Most of the time the DSP realization is close to the hardware implementation. An example is carrier recovery. For narrow-band communication signals in wireless communicators, carrier recovery can be performed using a digital implementation of the concepts described for the hardware carrier recovery circuits. While it is more power efficient to implement many of the techniques in hardware, the need to support multiple standards has necessitated the software reconfigurability available with a DSP unit. Which approach is used is the decision of the RF system designer—an experienced engineer with a rich background in wireless technologies. It is therefore important that the aspiring and practicing RF engineer have a broad perspective of RF circuits and of communications theory. Hence the emphasis of this book on a systems approach to RF and microwave design.

Frequency modulation, and the similar PM modulation method were used in the 1G analog cellular radio. With the addition of AM, the three schemes are the bases of all analog radio. Digital cellular radio began with 2G and there were two types of 2G cellular radios with the GSM system using GMSK modulation, a type of FSK modulation, and the NADC system using  $\pi/4$ -DQPSK modulation. The two 2G systems were incompatible. The 3G cellular radio used two types of QPSK modulation, one for the up-link from handset to basestation, and one from the basestation to a handset. The 1G–3G systems implemented most of the modulation and demodulation functions in analog hardware. With 4G and 5G cellular radio a large number of modulation schemes are supported choosing as high an order of modulation as allowed by the channel conditions. Most of the modulation and demodulation in 4G and 5G are implemented in DSP with just the translation to and from the radio frequency signal implemented in analog hardware.

## 2.13 References

- [1] “American National Standard T1.523-2001, Telecom Glossary 2011,” available on-line with revisions at <http://glossary.atis.org>, 2011, sponsored by Alliance for Telecommunications Industry Solutions.
- [2] S. Boyd, “Multitone signals with low crest factor,” *IEEE Trans. on Circuits and Systems*, vol. 33, no. 10, pp. 1018–1022, Oct. 1986.
- [3] A. Jones, T. Wilkinson, and S. Barton, “Block coding scheme for reduction of peak to mean envelope power ratio of multicarrier transmission schemes,” *Electronics Letters*, vol. 30, no. 25, pp. 2098–2099, Dec. 1994.
- [4] M. Steer, *Microwave and RF Design, Transmission Lines*, 3rd ed. North Carolina State University, 2019.
- [5] D. Porcino and W. Hirt, “Ultra-wideband radio technology: potential and challenges ahead,” *IEEE communications magazine*, vol. 41, no. 7, pp. 66–74, 2003.
- [6] “FCC (GPO) Title 47, Section 15 of the Code of Federal Regulations SubPart F: Ultra-wideband,” [http://www.access.gpo.gov/nara/cfr/waisidx\\_05/47cfr15.05.html](http://www.access.gpo.gov/nara/cfr/waisidx_05/47cfr15.05.html).
- [7] J. Carson, “Notes on the theory of modulation,” *Proc. of the Institute of Radio Engineers*, vol. 10, no. 1, pp. 57–64, Feb. 1922.
- [8] L. Couch III, *Digital and Analog Communication Systems*, 6th ed. Prentice-Hall, 2001.
- [9] E. Armstrong, “A method of reducing disturbances in radio signaling by a system of frequency modulation,” *Proc. of the Institute of Radio Engineers*, vol. 24, no. 5, pp. 689–740, May 1936.
- [10] —, “Radio telephone signaling,” US Patent US Patent 1 941 447, 12 26, 1933.
- [11] “Armstrong suit over fm settled. r.c.a. and n.b.c. to pay ‘\$1,000,000’ ending action begun

- by late inventor," New York Times, Dec. 31, 1954.
- [12] E. Bedrosian, "The analytic signal representation of modulated waveforms," *Proceedings of the IRE*, vol. 50, no. 10, pp. 2071–2076, 1962.
- [13] C. Shannon, "Communication in the presence of noise," *Proc. IRE*, vol. 37, no. 1, pp. 10–21, 1949.
- [14] J. Costas, "Synchronous communications," *Proc. of the IRE*, vol. 44, no. 12, pp. 1713–1718, Dec. 1956.
- [15] F. Hearth, "Origins of the binary code," *Scientific American*, pp. 76–83, Aug. 1972.
- [16] F. Gray, "Pulse code modulation," US Patent US Patent 11 111 111, 03 17, 1953.
- [17] C. Savage, "A survey of combinatorial gray codes," *SIAM Review*, vol. 39, no. 4, pp. 605–629, 1997.

## 2.14 Exercises

- Develop a formula for the average power of a signal  $x(t)$ . Consider  $x(t)$  to be a voltage across a  $1\ \Omega$  resistor.
- What is the PAPR of a 5-tone signal when the amplitude of each tone is the same?
- What is the PMEPR of a 10-tone signal when the amplitude of each tone is the same?
- Consider two uncorrelated analog signals combined together. One signal is denoted  $x(t)$  and the other  $y(t)$ , where  $x(t) = 0.1 \sin(10^9 t)$  and  $y(t) = 0.05 \sin(1.01 \cdot 10^9 t)$ . The combined signal is  $z(t) = x(t) + y(t)$ . [Parallels Example 2.3]
  - What is the PAPR of  $x(t)$  in decibels?
  - What is the PAPR of  $y(t)$  in decibels?
  - What is the PMEPR of  $x(t)$  in decibels?
  - Is it possible to calculate the PMEPR of  $z(t)$ ? If so, what is it?
- Consider two uncorrelated analog signals combined together. One signal is denoted  $x(t)$  and the other  $y(t)$ , where  $x(t) = 0.1 \sin(10^8 t)$  and  $y(t) = 0.05 \sin(1.01 \cdot 10^8 t)$ . What is the PMEPR of this combined signal? Express PMEPR in decibels. [Parallels Example 2.3]
- What is PMEPR of a three-tone signal when the amplitude of each tone is the same?
- What is PMEPR of a four-tone signal when the amplitude of each tone is the same?
- A tone  $x_1(t) = 0.12 \cos(\omega_1 t)$  is added to two other tones  $x_2(t) = 0.14 \cos(\omega_2 t)$  and  $x_3(t) = 0.1 \cos(\omega_3 t)$  to produce a signal  $y(t) = x_1(t) + x_2(t) + x_3(t)$ , where  $y(t)$ ,  $x_1(t)$ ,  $x_2(t)$  and  $x_3(t)$  are voltages across a  $100\ \Omega$  resistor. Consider that  $\omega_1$ ,  $\omega_2$ , and  $\omega_3$  are 10% apart and that the signals at these frequencies are uncorrelated.
  - What is the PMEPR of  $x_1(t)$ ? Express your answer in decibels.
  - Sketch  $y(t)$ .
  - The combined signal appears as a pseudo-carrier with a time-varying envelope. What is the power of the largest single cycle of the pseudo-carrier?
    - What is the average power of  $y(t)$ ?
    - What is the PMEPR of  $y(t)$ ? Express your answer in decibels.
- Consider two uncorrelated analog signals summed together. One signal is denoted  $x(t)$  and the other  $y(t)$ , where  $x(t) = \sin(10^9 t)$  and  $y(t) = 2 \sin(1.01 \cdot 10^9 t)$  so that the total signal is  $z(t) = x(t) + y(t)$ . What is the PMEPR of  $z(t)$  in decibels? [Parallels Example 2.3]
- What is the PMEPR of an FM signal at 1 GHz with a maximum modulated frequency deviation of  $\pm 10$  kHz?
- What is the PMEPR of a two-tone signal (consisting of two sinewaves at different frequencies that are, say, 1% apart)? First, use a symbolic expression, then consider the special case when the two amplitudes are equal. Consider that the two tones are close in frequency.
- What is the PMEPR of a three-tone signal (consisting of three equal-amplitude sinewaves, say, 1% apart in frequency)?
- A phase modulated tone  $x_1(t) = A_1 \cos[\omega_1 t + \phi_1(t)]$ . What is the PMEPR of  $x_1(t)$ ? Express your answer in decibels.
- What is the PMEPR of an AM signal with 75% amplitude modulation?
- Two FM voltage signals  $x_1(t)$  and  $x_2(t)$  are added together and then amplified by an ideal linear amplifier terminated in  $50\ \Omega$  with a gain of 10 dB and the output voltage of the amplifier is  $y(t) = \sqrt{10} [x_1(t) + x_2(t)]$ .
  - What is the PMEPR of  $x_1(t)$ ? Express your answer in decibels?
  - What effect does the amplifier have on the PMEPR of the signal?
  - If  $x_1(t) = A_1 \cos[\omega_1(t)t]$  and  $x_2(t) = A_2 \cos[\omega_2(t)t]$ , what is the PMEPR of the output of the amplifier,  $y(t)$ ? Express PMEPR in decibels. Consider that  $\omega_1(t)$  and  $\omega_2(t)$  are within 0.1% of each other.

16. An FM signal has a maximum frequency deviation of 20 kHz and a modulating signal between 300 Hz and 5 kHz. What is the bandwidth required to transmit the modulated RF signal when the carrier is 200 MHz? Is this considered to be narrowband FM or wideband FM?
17. A high-fidelity stereo audio signal has a frequency content ranging from 50 Hz to 20 kHz. If the signal is to be modulated on an FM carrier at 100 MHz, what is the bandwidth required for the modulated RF signal? The maximum frequency deviation is 5 kHz when the modulating signal is at its peak value.
18. Consider FM signals close in frequency but whose spectra do not overlap. [Parallels Example 2.5]
- What is the PMEPR of just one PM signal? Express your answer in decibels.
  - What is the PMEPR of a signal comprised of two uncorrelated narrowband PM signals with the same average power?
19. Consider two nonoverlapping equal amplitude FM signals having center frequencies within 1%.
- What is the PMEPR in dB of just one FM modulated signal?
  - What is the PMEPR in dB of a signal comprising two FM signals of the same power?
20. Consider a signal  $x(t)$  that is the sum of two uncorrelated signals, a narrowband AM signal with 50% modulation,  $y(t)$ , and a narrow-band FM signal,  $z(t)$ . The center frequencies of  $y(t)$  and  $z(t)$  are within 1%. The carriers have equal amplitude. Express answers in dB.
- What is the PAPR of the AM signal  $x(t)$ ?
  - What is the PAPR of the FM signal  $z(t)$ ?
  - What is the PAPR of  $x(t)$ ?
  - What is the PMEPR of the AM signal  $x(t)$ ?
  - What is the PMEPR of the FM signal  $z(t)$ ?
  - What is the PMEPR of  $x(t)$ ?
21. Two phase modulated tones  $x_1(t) = A_1 \cos[\omega_1 t + \phi_1(t)]$  and  $x_2(t) = A_2 \cos[\omega_2 t + \phi_2(t)]$  are added together as  $y(t) = x_1(t) + x_2(t)$ . What is the PMEPR of  $y(t)$  in decibels. Consider that  $\omega_1$  and  $\omega_2$  are within 0.1% of each other.
22. A radio uses a channel with a bandwidth of 25 kHz and a modulation scheme with a gross bit rate of 100 kbits/s that is made of an information bit rate of 60 kbits/s and a code bit rate of 40 kbits/s.
- What is the modulation efficiency in bits/s/Hz?
  - What is the spectral efficiency in bits/s/Hz?
23. A cellular communication system uses  $\pi/4$ -DQPSK modulation with a modulation efficiency of 1.63 bits/s/Hz to transmit data at the rate of 30 kbits/s. This would be the spectral efficiency in the absence of coding. However, 25% of the transmitted bits are used to implement a forward error correction code.
- What is the gross bit rate?
  - What is the information bit rate?
  - What is the bandwidth required to transmit the information and code bits?
  - What is the spectral efficiency in bits/s/Hz?
24. A radio uses a channel with a 5 MHz bandwidth and uses 256-QAM modulation with a modulation efficiency of 6.33 bits/s/Hz. The coding rate is  $3/4$  (i.e. of every 4 bits sent 3 are data bits and the other is an error correction bit).
- What is gross bit rate in Mbits/s?
  - What is information rate in Mbits/s?
  - What is the spectral efficiency in bits/s/Hz?
25. The following sequence of bits 0100110111 is to be transmitted using QPSK modulation. Take these data in pairs, that is, as 01 00 11 01 11. These pairs, one bit at a time, drive the  $I$  and  $Q$  channels. Show the transitions on a constellation diagram. [Parallels Example 2.7]
26. The following sequence of bits 0100110111 is to be transmitted using  $\pi/4$ -DQPSK modulation. Take these data in pairs, that is, as 01 00 11 01 11. These pairs, one bit at a time, drive the  $I$  and  $Q$  channels. Use five constellation diagrams, with each diagram showing one transition or symbol. [Parallels Example 2.8]
27. The following sequence of bits 0100110111 is transmitted using OQPSK modulation. Take these data in pairs, that is, as 01 00 11 01 11. These pairs, one bit at a time, drive the  $I$  and  $Q$  channels. Show the transitions on a constellation diagram.
28. Draw the constellation diagram of OQPSK.
29. Draw the constellation diagrams of  $3\pi/8$ -8DPSK and explain the operation of this system and describe its advantages.
30. How many bits per symbol can be sent using  $3\pi/8$ -8PSK?
31. How many bits per symbol can be sent using 8-PSK?
32. How many bits per symbol can be sent using 16-QAM?
33. Draw the constellation diagram of OQPSK modulation showing all possible transitions. You may want to use two diagrams.

34. What is the PMEPR of a 5-tone signal when the amplitude of each tone is the same? EVM of the modulated signal? [Parallels Example 2.10]
35. Draw the constellation diagram of 64QAM. 41. Consider a digitally modulated signal and describe the impact of a nonlinear amplifier on the signal. Include several negative effects.
36. How many bits per symbol can be sent using 32QAM?
37. How many bits per symbol can be sent using 16QAM? 42. A carrier with an amplitude of 3 V is modulated using 8-PSK modulation. If the noise on the modulated signal has an rms value of 0.1 V, what is the EVM of the modulated signal? [Parallels Example 2.10]
38. How many bits per symbol can be sent using 2048QAM?
39. Consider a two-tone signal and describe intermodulation distortion in a short paragraph and include a diagram. 43. Consider a 32-QAM modulated signal which has a maximum  $I$  component, and a maximum  $Q$  component, of the RF phasor of 5 V. If the noise on the signal has an RMS value of 0.1 V, what is the modulation error ratio of the modulated signal in decibels? Refer to Figure 2-31(b). [Parallels Example 2.10]
40. A 16-QAM modulated signal has a maximum RF phasor amplitude of 5 V. If the noise on the signal has an rms value of 0.2 V, what is the

### 2.14.1 Exercises By Section

§12.2 1, 2, 3, 4, 5, 6, 7, 8, 9	20, 21	§12.9 32, 33, 34, 35, 36, 37, 38
10, 11, 12	§12.5 22, 23, 24	§12.11 39, 40, 41, 42, 43
§12.4 13, 14, 15, 16, 17, 18, 19	§12.8 25, 26, 27, 28, 29, 30, 31	

### 2.14.2 Answers to Selected Exercises

5 2.55 dB	15 no effect	22(a) 4 bit/s/Hz
7(e) 3.78 dB	20(a) 6 dB	43 36.02 dB
8 0.00022 W	20(e) 0 dB	

# Transmitters and Receivers

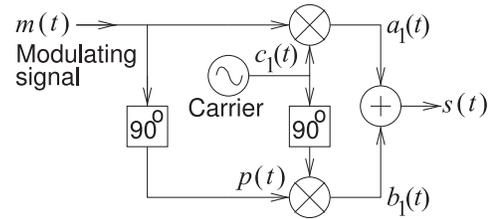
3.1	Introduction .....	77
3.2	Single-Sideband and Double-Sideband Modulation .....	78
3.3	Early Modulation and Demodulation Technology .....	80
3.4	Receiver and Transmitter Architectures .....	82
3.5	Carrier Recovery .....	88
3.6	Modern Transmitter Architectures .....	88
3.7	Modern Architectures .....	91
3.8	Introduction to Software Defined Radios .....	97
3.9	SDR Quadrature Modulators .....	98
3.10	Case Study: SDR Transmitter .....	103
3.11	SDR Quadrature Demodulator .....	120
3.12	SDR Receiver .....	121
3.13	SDR Summary .....	125
3.14	Summary .....	126
3.15	References .....	126
3.16	Exercises .....	126

## 3.1 Introduction

An essential function of a radio transmitter is modulation which is implemented by a modulator which converts the information at baseband to a finite bandwidth modulated radio signal centered at or near the carrier frequency. At the block level a modulator is described by its architecture and at a lower level by its circuit implementation. In a receiver a demodulator demodulates the radio signal extracting the baseband information.

The architectures of transmitters and receivers have changed significantly since the first radio signals were transmitted. Largely this is the result of increased integration of hardware but also due to advances in modulation concepts. In modern radios, as used with 4G and 5G cellular radio, many of the modulation functions once only implemented with analog hardware are now implemented digitally in a digital signal processing (**DSP**) unit. Only the final translation to the frequency of the radio signal is implemented in analog hardware. This concept is known as **software-defined radio** (SDR) which also has the additional implication that the analog hardware can be reconfigured under software control. SDR has two significant impacts.

**Figure 3-1:** Hartley modulator implementing single-sideband suppressed-carrier (SSB-SC) modulation. The “90°” blocks shift the phase of the signal by +90°. The mixer indicated by the circle with a cross is an ideal multiplier, e.g.  $a_1(t) = m(t) \cdot c_1(t)$ .



One is that it is feasible to support a large number of modulation schemes including legacy modulation schemes, without any changes to the analog hardware. The second is that by simplifying the functionality of the analog hardware it is possible to optimize analog hardware for maximum efficiency. Demodulation has also been significantly impacted by SDR with a similar transfer of functions from analog to digital hardware.

In some designs a receiver and a transmitter share components and then the RF front end is called a transceiver. This chapter is concerned with the architectures of receivers, transmitters, and transceivers. Design choices are made to maximize the tolerance to interference, manage hardware complexity, optimize power efficiency, and enable various radios to simultaneously operate together.

The presentation in this chapter closely follows the historical development of transmitters and receivers. As well, modulators and demodulators at different stages of evolution are considered together.

Section 3.2 introduces single-sideband and double-sideband modulation. Then early modulator and demodulator technology is considered in Section 3.3. These were circuits that could be implemented with just a few vacuum tubes or transistors and, in one case, a demodulator could be implemented with a single diode. With increasing circuit sophistication an architectural approach to modulator and demodulator design became possible. These are considered in Section 3.4. Demodulation requires local generation of a radio signal's carrier in a process called carrier recovery. This is discussed in Section 3.5. Following this is a discussion of more modern transmitters and receivers as used in 2G and 3G cellular radio in Sections 3.6 and 3.7 respectively. Current radios, e.g. 4G and 5G, are software defined radios (SDRs) where most of the modulation and demodulation is implemented in DSP and there is considerable software-based adaptation of the analog hardware that now undertakes the task of translating between and low intermediate frequency and the frequency of the radio signal. Sections 3.8–3.13 present the many aspects of SDR transmitters and receivers.

### 3.2 Single-Sideband and Double-Sideband Modulation

The simplest implementation of analog modulation results in a modulated carrier signal whose spectrum consists of the carrier and upper and lower sidebands. It is possible to eliminate one of the sidebands in AM modulation, or one of the sidebands sets in PM and FM modulation, producing single sideband modulation SSB. This however needs to be implemented in DSP. At the same time the carrier can be suppressed resulting in suppressed-carrier modulation or together SSB-SC modulation.

The simplest system that implements SSC-SC modulation is the **Hartley modulator** [1, 2], shown in Figure 3-1. This circuit results in **single-**

**sideband (SSB) modulation** or more precisely single-sideband modulation **suppressed-carrier (SSB-SC) modulation**. This circuit is used in all modern radios taking a modulated signal which is centered at an intermediate frequency and shifting it up in frequency so that its is centered at another frequency a little below or a little above the carrier of the Hartley modulator.

Both the modulating signal  $m(t)$  and the carrier are multiplied together in a mixer and then also  $90^\circ$  phase-shifted versions are also mixed before being added together. The signal flow is as follows beginning with  $m(t) = \cos(\omega_{m1}t)$ ,  $p(t) = \cos(\omega_{m1}t - \pi/2) = \sin(\omega_{m1}t)$  and carrier signal  $c_1(t) = \cos(\omega_c t)$ :

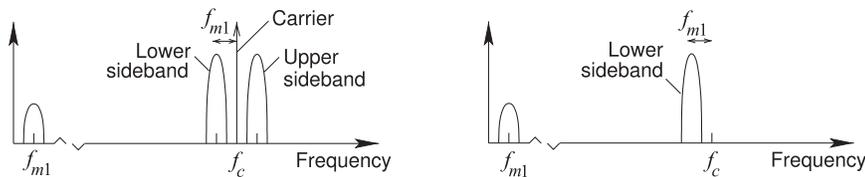
$$\begin{aligned} a_1(t) &= \cos(\omega_{m1}t) \cos(\omega_c t) = \frac{1}{2}[\cos((\omega_c - \omega_{m1})t) + \cos((\omega_c + \omega_{m1})t)] \\ b_1(t) &= \sin(\omega_{m1}t) \sin(\omega_c t) = \frac{1}{2}[\cos((\omega_c - \omega_{m1})t) - \cos((\omega_c + \omega_{m1})t)] \\ s_1(t) &= a_1(t) + b_1(t) = \cos((\omega_c - \omega_{m1})t) \end{aligned} \tag{3.1}$$

and so the lower sideband (USB) is selected.

That is, if a finite bandwidth modulating signal  $m(t)$  was mixed only once with the carrier  $c_1(t)$ , the spectrum of the output  $a_1(t)$  would include upper and lower sidebands as well as the carrier as shown in Figure 3-2(a). With the Hartley modulator the spectrum of Figure 3-2(b) is obtained.

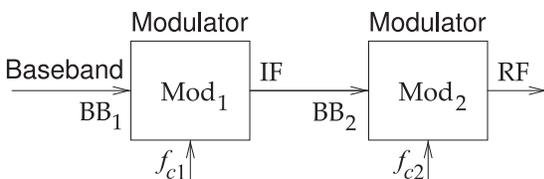
In digital modulation both the upper and lower sideband are retained but the carrier is suppressed. Both sidebands are required to recover the signal but the spectrum is used efficiently as the modulating signal is complex with two components, think of real/imaginary parts, or amplitude/phase information. Since a DSP unit is required the resulting modulated signal must be at a relatively low frequency and then a second frequency conversion stage is required to shift the modulated signal to the desired operating frequency, see Figure 3-3. This second stage is a SSB-SC modulator, however since the input to the second stage is a DSB-SC signal, the final RF signal is a DSB-SC signal.

All the concepts introduced in this section are still used in modern radios, just that now they are mostly implemented in a DSP unit rather than in analog hardware.



(a) Double sideband (DSB) (b) Single-sideband suppressed carrier (SSB-SC)

**Figure 3-2:** Spectrum of a modulated carrier with a modulating signal of finite bandwidth  $f_m$  with  $f_{m1}$  being the center frequency of the baseband signal.



**Figure 3-3:** Two-stage modulator with the baseband signal,  $BB_1$ , input to the first modulator,  $Mod_1$ , producing the intermediate frequency signal  $IF_1$ . This becomes the baseband signal,  $BB_2$ , for the second modulator,  $Mod_2$ , producing the radio frequency signal,  $RF$ .

### 3.3 Early Modulation and Demodulation Technology

Early modulators and demodulators are considered here in part because the terms associated with the historical transmitters and receivers are still used today, but also because the early trade-offs influenced the architectures used today. Today transmitters and receivers use DSP technology, very stable LOs, and sophisticated clock recovery schemes. This was not always so. One of the early problems was demodulating a signal when the frequency of transmitter oscillators, i.e. carriers, drifted by up to 10%. Radio at first used AM and the carrier was sent with the information-carrying sidebands. With this signal, a simple single-diode rectifier circuit connected to a bandpass filter could be used, but the reception was poor. To improve performance it was necessary to lock an oscillator in the receiver to the carrier and then amplify the received signal. Here some of the early schemes that addressed these problems are discussed. There were many more variants, but the discussion covers the essential ideas.

#### 3.3.1 Heterodyne Receiver

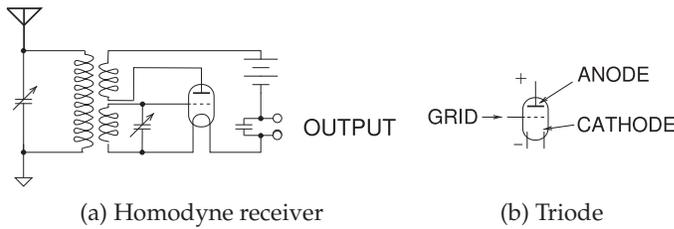
The heterodyning principle mixes a single-tone signal, the LO, with a finite bandwidth signal to produce a lower-frequency version of the information-bearing signal. With the LO frequency set appropriately, the low-frequency signal would be in the audio range. If the information-bearing signal is an AM signal, then the low-frequency version of the signal is the original audio signal, which is the envelope of the AM signal. This type of receiver is called a **tuned radio frequency (TRF) receiver**, and performance is critically dependent on the stability of the LO and the selectivity of the receive filter. The TRF receiver required the user to adjust a tunable capacitor so that, with a fixed inductor, a tunable bandpass filter was created. Such a filter has a limited  $Q$  and a bandwidth that is wider than the bandwidth of the radio channel.<sup>1</sup> Even worse, a user had to adjust both the frequency of the bandpass filter and the frequency of the LO. The initial radios based on this principle were called **audions**, used a triode vacuum tube as an amplifier, and were used beginning in 1906. They were an improvement on the **crystal detectors** (which used a single diode with filters), but there was a need for something better.

#### 3.3.2 Homodyne Receiver

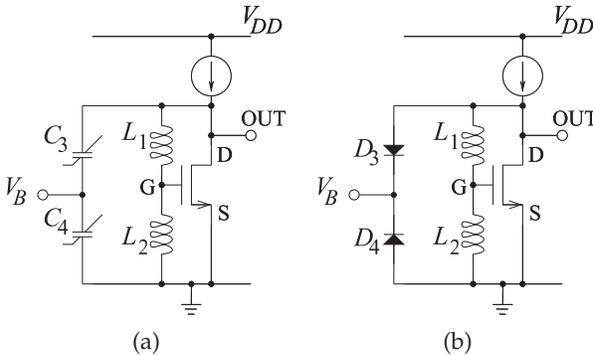
The homodyne [3], syncrodyne (for synchronous heterodyne) [4], and autodyne (for automatic heterodyne) circuits were needed improvements on the audion and are based on the regenerative circuit invented by Edwin Armstrong in 1912 while he was an electrical engineering student at New York City's Columbia University [5]. Armstrong's circuit fed the input signal into an amplifying circuit and a portion of this signal was coupled back into the input circuit so that the signal was amplified over and over again. This is a positive feedback amplifier. A small input RF signal was amplified to such a large extent that it resulted in the amplifying circuit becoming nonlinear and consequently it rectified the amplitude modulated RF signal.

---

<sup>1</sup>  $Q$  is the **quality factor** and is the ratio of the energy stored to the energy resistively lost in each cycle. Good frequency selectivity in a filter requires high- $Q$  components.



**Figure 3-4:** Colebrook’s original homodyne receiver: (a) circuit with an antenna, tunable band-pass filter, and triode amplifier; and (b) triode vacuum tube.



**Figure 3-5:** A common source Hartley voltage-controlled oscillator (VCO): (a) with nonlinear capacitors; and (b) with diodes which have a variable capacitance when reverse biased.

Colebrook used this principle and developed the original homodyne receiver shown in Figure 3-4(a). This serves to illustrate the operation of the family of regenerative receivers. The antenna shown on the left-hand side is part of a resonant circuit that is in the feedback path of a triode oscillator. The triode vacuum tube is shown in Figure 3-4(b). Here the grid coils (which control the flow of carriers between the bottom cathode<sup>2</sup> and top anode) are weakly coupled to the anode circuit. When an AC signal appears at the top anode, the part within the passband of the tuned circuit is fed back to the grid and the signal is reinforced. The radio signals of the day were AM and had a relatively large carrier, so the oscillator locked on to the carrier. The AM sidebands were then successfully heterodyned down to the desired audio frequencies.

The **autodyne** worked on a slightly different principle in that the oscillation frequency was tuned to a slightly different frequency from the carrier. Still, the autodyne combined the functions of an oscillator and detector in the same circuit.

### 3.3.3 FM Modulator

FM modulation can be implemented using a voltage-controlled oscillator (VCO) with the baseband signal controlling the frequency of an oscillator. A VCO can be very simple circuit and so was easily implemented in early radio. The circuits in Figure 3-5 are known as common source Hartley VCOs, where in Figure 3-5(a) the controllable elements are the nonlinear capacitors. These are generally implemented as reverse-biased diodes, as shown in Figure 3-5(b) where the bias,  $V_B$ , changes the capacitance of the reverse-biased diodes. Changing the capacitance changes the resonant frequency of the feedback loop formed by the inductors and the diode capacitances (called varactors).

<sup>2</sup> The **cathode** is heated (the heater circuit is not shown) and electrons are spontaneously emitted in a process called **thermionic emission**.

### 3.3.4 FM Demodulator

An FM demodulator is often implemented using a phase-locked loop with an error signal used to control the frequency of a voltage-controlled oscillator (VCO) with the loop arranged so that the VCO tracks the received signal. The desired baseband signal, i.e. the demodulated FM signal, being the loop's error signal.

### 3.3.5 Superheterodyne Receiver

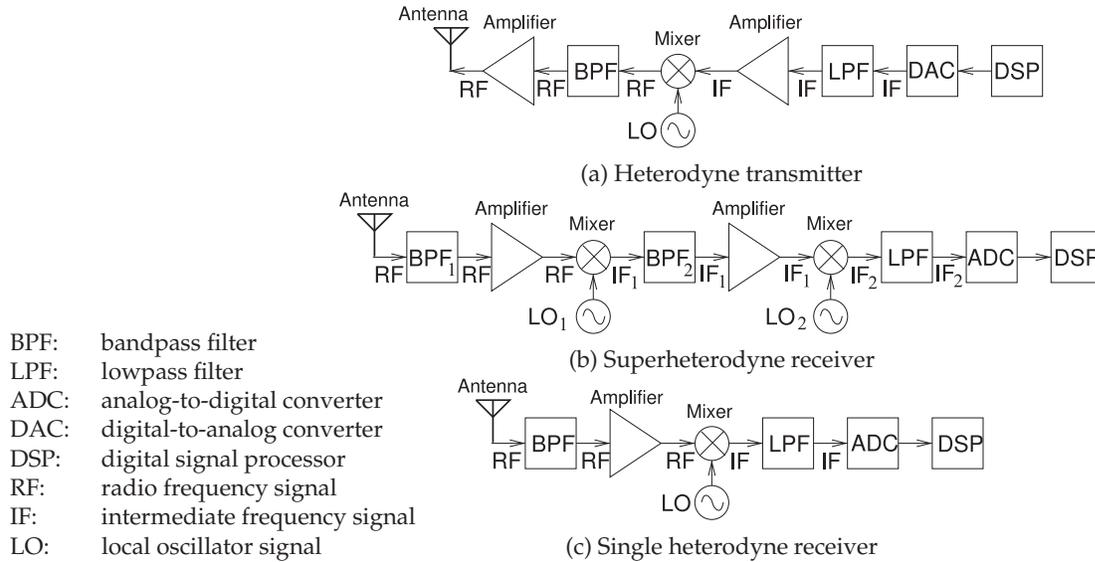
The superheterodyne receiver was invented by Edwin Armstrong in 1918 [6]. The key concept was to **heterodyne** down in two stages, use fixed filters, and use a tunable LO. The receiving antenna was connected to a bandpass filter that allowed several channels to pass. This relaxed the demands on the receive filter, but also filters with higher selectivity could be constructed if they did not need to be tuned. The filtered received signal is then mixed with an offset LO to produce what is called a **supersonic** signal—a signal above the audio range—and hence the name of this architecture. The performance of the superheterodyne (or superhet) receive architecture has only recently been achieved at cellular frequencies using direct conversion architectures.

### 3.3.6 Summary

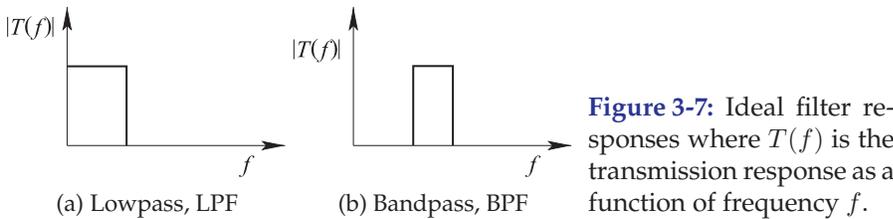
Early radio used AM and FM modulation and both modulation and demodulation could be performed with very simple circuits. However they used a lot of spectrum for the information that was transmitted.

## 3.4 Receiver and Transmitter Architectures

The essential function of a radio transmitter architecture is taking low-frequency information, the baseband signal, and transferring that information to much higher frequencies by superimposing the baseband signal on a high-frequency carrier, i.e. a sinewave. This could be done by slowly varying the amplitude, frequency, and/or phase of the carrier in what is called modulating the carrier to produce a modulated carrier (signal). A caveat here is that there are other less common ways of transferring the baseband information such as rapidly changing the frequency of the carrier at a faster rate than the maximum frequency of the baseband signal in a scheme called frequency hopping. Frequency hopping is particularly useful when operating in hostile situations, such as military communications, when the interference environment cannot be controlled. In reality there is a very large number of ways of producing a radio signal that carries the baseband information. The central theme of this book is to focusing on modulation schemes that slowly modulate characteristics of a carrier signal and other schemes are introduced as an exception. This section describes early architectural developments that led to the development of the regulations that still guide modern architectures; concepts that efficiently pack as much information as possible into a fixed RF bandwidth.



**Figure 3-6:** RF front ends: (a) a one-stage transmitter; (b) a receiver with two mixing (or heterodyning) stages; and (c) a receiver with one heterodyne stage.



**Figure 3-7:** Ideal filter responses where  $T(f)$  is the transmission response as a function of frequency  $f$ .

### 3.4.1 Radio as a Cascade of Two-Ports

The front end of an RF communication receiver or transmitter combines a number of subsystems in cascade. The design of the RF front end requires trade-offs of noise generated by the circuit, of frequency selectivity, and of power efficiency, which translates into battery life for a communication handset. There are only a few receiver and transmitter architectures that achieve the optimum trade-offs. The essentials of these architectures are shown in Figure 3-6. These architectures achieve frequency selectivity using bandpass filters (BPFs) and lowpass filters (LPFs) which ideally have the responses shown in Figure 3-7. The corner frequencies of these filters and, in the case of the BPFs, their center frequencies, are adjusted to minimize interfering signals and noise that passes through the system. The three architectures shown in Figure 3-6 have antennas that interface between circuits and the outside world. Antennas generally have a broad bandwidth, much greater than the bandwidth of an individual communication channel.

### 3.4.2 Heterodyne Transmitter and Receiver

First consider the transmitter architecture shown in Figure 3-6(a). In a transmitter, a low-frequency information-bearing signal is translated to a

frequency that can be more easily radiated. The information is contained in the baseband signal, which in many modern systems is generated as a digital signal within the digital signal processor (DSP). Then the **digital-to-analog converter, (DAC)**, converts the digital baseband to an analog signal called the analog baseband, identified in Figure 3-6(a) as the intermediate frequency (IF) signal at the output of the DAC. Older systems generate the IF signal using analog hardware as this reduces the power consumed by the digital electronics. The IF then passes through a lowpass filter (LPF) to remove the harmonics resulting from the DAC process, and the signal is then amplified before being applied to the mixer. There are several types of mixers, but the central concept is multiplying the IF signal by a much larger local oscillator (LO) signal at frequency  $f_{LO}$ . The LO will be a single cosinusoid and if its amplitude is  $A_{LO}$ , then the LO signal is  $x_{LO} = A_{LO} \cos(2\pi f_{LO})$ . While the IF signal will have a finite bandwidth, the operation of the mixer can be illustrated by considering that the IF is a single cosinusoid with amplitude  $A_{IF}$  and frequency  $f_{IF}$ , that is the RF signal is  $x_{IF} = A_{IF} \cos(2\pi f_{IF})$ . Then the output of the mixer is

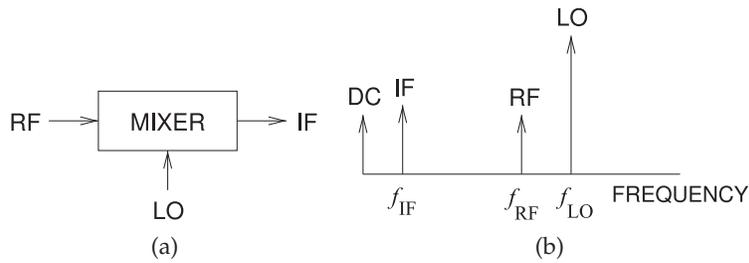
$$\begin{aligned} x_{RF} &= x_{IF} \times x_{LO} = A_{IF} A_{LO} \cos(2\pi f_{IF}) \cos(2\pi f_{LO}) \\ &= \frac{1}{2} A_{RF} A_{LO} \{ \cos[2\pi(f_{LO} - f_{IF})] + \cos[2\pi(f_{LO} + f_{IF})] \}. \end{aligned} \quad (3.2)$$

In Equation (3.2) a trigonometric expansion has been used. The LO is chosen so that its frequency is close to that of the desired RF so that multiplication by the mixer results in an output that has one component at  $f_{RF\Delta} = f_{LO} - f_{IF}$  and one at  $f_{RF\Sigma} = f_{LO} + f_{IF}$ . One of these is selected by the BPF, is further amplified, and then delivered to the antenna and radiated.

### 3.4.3 Superheterodyne Receiver Architecture

The first receiver architecture to be considered is the superheterodyne (or superhet) receiver architecture shown in Figure 3-6(b). Heterodyning refers to the use of a mixer and a superheterodyne circuit has two mixers. The antenna collects information from the environment at RF, and immediately this is bandpass filtered (by the BPF<sub>1</sub> block) to eliminate most of the interfering signals and noise. Thus the first BPF reduces the range of voltages presented to the first amplifier and so reduces the chance that the amplifier will distort the desired signal.

The RF at the output of the leftmost bandpass filter, BPF<sub>1</sub>, still has a spectrum that is much broader than that of the desired communication signal. For example, in 3G radio the communication channel is 5 MHz wide but the first BPF could be 50 MHz wide and the RF could be 1 GHz or 2 GHz. So it is still necessary to use additional frequency selectivity to isolate the single required channel. The optimum choice at this stage is to follow the first BPF with an amplifier to boost the level of the signal. This also boosts the level of the noise that was captured by the antenna along with the signal, but it means that the noise added by the circuitry after the first amplifier has much less importance. The next block in the receiver is the mixer that shifts the information down in frequency to the first intermediate frequency, IF<sub>1</sub>. The local oscillator, LO<sub>1</sub>, is chosen so that its frequency is close to that of the RF so that multiplication by the mixer results in a low frequency at the output (i.e., at  $f_{IF1} = f_{LO} - f_{RF}$ ), and one at a frequency nearly twice that of the LO (i.e., at the sum frequency,  $f_{\Sigma} = f_{LO} + f_{RF}$ ). The center frequency and



**Figure 3-8:** Simple mixer circuit: (a) block diagram; and (b) spectrum.

bandwidth of the second bandpass filter,  $BPF_2$ , is chosen so that only the signals around  $f_{IF1}$  pass through. Thus the main function of the first mixer stage and  $BPF_2$  in the superheterodyne receiver is to convert the information at the RF down to a lower frequency, here at  $IF_1$ . The operation of the mixer is shown in Figure 3-8.

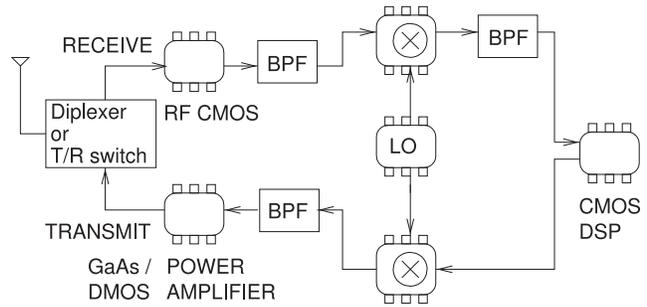
In the superheterodyne receiver architecture the output of the first mixer,  $f_{IF1}$ , is still at too high a frequency for the signal to be directly converted into digital form where it can be processed. So it is natural to ask why the frequency translation was not all the way down to baseband. The main reason for this is that there is substantial noise on an LO at frequencies very close to the oscillation frequency,  $f_{LO1}$ . This noise drops off quickly away from the oscillation frequency and its level at frequency  $f$  is proportional to  $(|f - f_{LO1}|)^n$ ,  $n = 1, 2, \dots$ . This noise will appear at the output of the mixer and will be substantial if the RF and LO are very close in frequency. So the optimum trade-off is to shift the frequency of the information-bearing signal in two stages. Following further amplification, the second stage of mixing converts the information-bearing component of the signal centered at the first intermediate frequency,  $IF_1$ , to the second intermediate frequency,  $IF_2$ , which is usually at the baseband frequency or slightly above it. The baseband signal is now analog and this is converted to digital form by the ADC, and then the signal, now the digital baseband signal, can be digitally processed by the DSP.

Since the second mixer operates at much lower frequencies than the first mixer, the phase noise on  $LO_2$  does not overlap the signal at  $IF_2$ . The reason why this architecture is called a superheterodyne receiver architecture is because when this architecture was first used,  $IF_2$  was an audio signal and  $IF_1$  was above the audible range and so was a supersonic signal. Thus the super in superheterodyne initially referred to the supersonic IF.

### 3.4.4 Single Heterodyne Receiver

The second receiver architecture shown in Figure 3-6(c) has a single heterodyne or mixing stage. If the LO has very low noise the IF will be at a lower frequency than in the superheterodyne architecture shown in Figure 3-6(b). A high-performance ADC is then required to convert the signal to digital form and deliver it to the DSP unit. Substantial digital processing power is required to translate the signal to a digital baseband signal. Alternatively a high-performance subsampling ADC could be used that in effect performs mixing during conversion. The advantage of this architecture is a simplified RF section and is of particular advantage when multiple RF communication standards are supported as the DSP and ADC can be common while the analog RF hardware for each band is considerably simplified.

**Figure 3-9:** RF front end organized as multiple chips. This corresponds to a combination of the receive architecture shown in Figure 3-6(c) and the transmitter architecture shown in Figure 3-6(a).

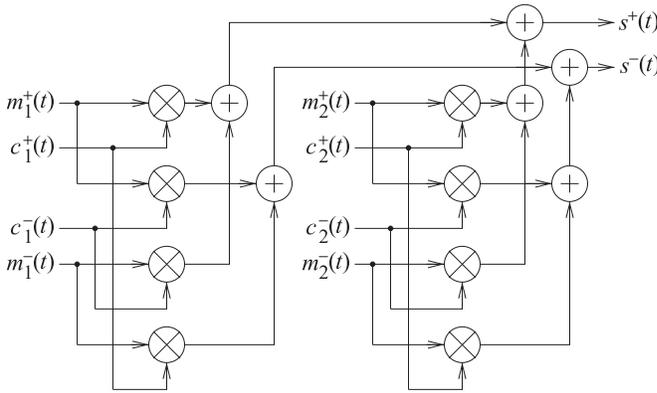


### 3.4.5 Transceiver

The major transistor- or diode-based **active elements** in the RF **front end** of both the transmitter and receiver are the **amplifiers, mixers, and oscillators**. These subsystems have much in common, using nonlinear devices to convert power at DC to power at RF. In the case of mixers, power at the LO is also converted to power at RF. The front end of a typical cellphone is shown in Figure 3-9. The components here are generally implemented in a module and use different technologies for the various elements, optimizing cost and performance.

Return now to the mixer-based **transceiver** (for transmitter and receiver) architecture shown in Figure 3-9. Here, a single antenna is used, and either a **diplexer**<sup>3</sup> (a combined lowpass and highpass filter) or a switch is used to separate the (frequency-spaced) transmit and receive paths. If the system protocol requires transmit and receive operations at the same time, a diplexer is required to separate the transmit and receive paths. This filter tends to be large, lossy, or costly (depending on the technology used). Consequently a transistor **switch** is preferred if the transmit and receive signals operate in different time slots. In the receive path the low-level received signal is amplified and the initial amplifier needs to have very good noise performance. Once the signal is larger the noise performance of the amplifiers is less critical. The initial amplifier is called a **low-noise amplifier (LNA)**. The amplified receive signal is then bandpass filtered and frequency down-converted by a mixer to IF that can be sampled by an ADC to produce a digital signal that is further processed by DSP. Variants of this architecture include one that has two mixing stages (as in the superheterodyne receiver shown in Figure 3-6(b)), and another with no mixing that relies instead on direct conversion of the receive signal using, as one possibility, a subsampling ADC. In the transmit path, the architecture is reversed, with a DAC driven by the DSP chip that produces an information-bearing signal at the IF which is then frequency up-converted by a mixer, bandpass filtered, and amplified by what is called a power amplifier to generate the tens to hundreds of milliwatts required.

<sup>3</sup> A diplexer separates transmitted and received signals and is often implemented as a filter called a diplexer. A **diplexer** separates two signals on different frequency bands so that they can use a common element such as an antenna. If the transmitted and received signals are in different time slots, the diplexer can be a switch.



**Figure 3-10:** Differential implementation of a Hartley SSB-SC modulator:  $s^+(t)$ ,  $s^-(t)$  are the positive and negative components of the differential modulated signal  $s(t)$ ;  $(c_1^+(t), c_1^-(t))$  is the differential carrier;  $(m_1^+(t), m_1^-(t))$  is the differential quadrature-modulated IF carrier;  $(c_2^+(t), c_2^-(t))$  is the differential carrier shifted  $90^\circ$ ; and  $m_2^+(t)$  and  $m_2^-(t)$  are the  $90^\circ$  phase-shifted versions of  $m_1^+(t)$  and  $m_1^-(t)$ , respectively.

### 3.4.6 Hartley Modulator

The Hartley modulator [1, 2], shown in Figure 3-1, results in SSB **single-sideband (SSB) modulation** or more precisely SSB **suppressed-carrier (SSB-SC) modulation**. This is one of the great inventions and variants of this circuit are used in all modern radios. In the Hartley modulator the modulating signal  $m(t)$  and the carrier are multiplied together in a mixer and then also  $90^\circ$  phase-shifted versions are also mixed before being added together. The signal flow is as follows beginning with  $m(t) = \cos(\omega_{m1}t)$ ,  $p(t) = \cos(\omega_{m1}t - \pi/2) = \sin(\omega_{m1}t)$  and carrier signal  $c_1(t) = \cos(\omega_c t)$ :

$$\begin{aligned} a_1(t) &= \cos(\omega_{m1}t) \cos(\omega_c t) = \frac{1}{2}[\cos((\omega_c - \omega_{m1})t) + \cos((\omega_c + \omega_{m1})t)] \\ b_1(t) &= \sin(\omega_{m1}t) \sin(\omega_c t) = \frac{1}{2}[\cos((\omega_c - \omega_{m1})t) - \cos((\omega_c + \omega_{m1})t)] \\ s_1(t) &= a_1(t) + b_1(t) = \cos((\omega_c - \omega_{m1})t) \end{aligned} \tag{3.3}$$

and so the lower sideband (LSB) is selected. That is, if a finite bandwidth modulating signal  $m(t)$  was mixed only once with the carrier  $c_1(t)$ , the spectrum of the output  $a_1(t)$  would include upper and lower sidebands as well as the carrier as shown in Figure 3-2(a). With the Hartley modulator the spectrum of Figure 3-2(b) is obtained.

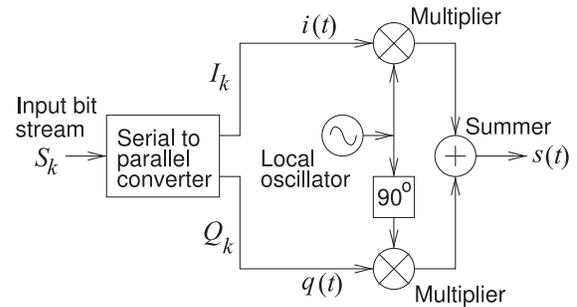
Every frequency component of  $m(t)$  needs to be shifted by  $90^\circ$  and this can be done using a polyphase circuit or digitally using a Hilbert transform. To select the upper sideband the summer in Figure 3-1 is replaced by a subtraction block and then the output becomes  $s'(t) = a_1(t) - b_1(t) = \cos((\omega_c + \omega_{m1})t)$ .

Another circuit that implements SSB-SC modulation is the **Weaver modulator** [7]. It uses only lowpass filters and mixers and is often used to implement SSB-SC in a digital-signal processor.

### 3.4.7 The Hartley Modulator in Modern Radios

In RF engineering, especially when RFICs are used, it is common to implement the phase shift of the modulating signal digitally using the Hilbert transform and to use differential signals. Then the Hartley SSB-SC modulator shown in Figure 3-1 is implemented as shown in Figure 3-10 where  $(m_1^+(t), m_1^-(t))$  is the differential form of the analog signal corresponding to the modulated signal  $s(t)$  at the output of the generic quadrature modulator in Figure 3-11. With the large number of modulation formats supported in modern cellular communication standards it is indeed

**Figure 3-11:** Quadrature modulator block diagram. An input bitstream,  $S_k$ , is divided into two bitstreams,  $I_k$  and  $Q_k$ , which are applied to the multipliers as (possibly filtered) waveforms  $i(t)$  and  $q(t)$ . The output of the multipliers (appropriately filtered) are summed to yield a modulated carrier signal,  $s(t)$ . The  $90^\circ$  block shifts the carrier by  $90^\circ$ .



fortunate that the quadrature modulators can be implemented digitally followed possibly by wave-shaping. The differential signal  $(m_2^+(t), m_2^-(t))$  is the phase-shifted form of  $(m_1^+(t), m_1^-(t))$  in which every frequency component is shifted by  $90^\circ$  usually implemented as the Hilbert transform of the digital forms of  $(m_1^+(t), m_1^-(t))$  (before wave-shaping).

### 3.5 Carrier Recovery

Demodulation requires the generation of a local version of the carrier of a received radio signal. This is simple with a DSB RF signal as the carrier is sent with the radio signal. DSB FM is a little different as strictly the carrier is not sent in the radio signal but the average frequency of the received signal is the carrier frequency. With all SSB schemes it is essential to create the carrier in a process called carrier recovery. This section considered carrier recovery for digitally modulated signals.

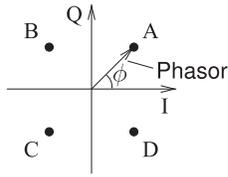
In modern radios carrier recovery is done in two stages. First a crude estimate of the carrier is generated and used to crudely demodulate the received signal. The crudely demodulated signal is then sampled to yield a digital demodulated signal. Then, for 2G and 3G radio, a fine carrier recovery procedure using a digital implementation of what has been described is used. With 4G and 5G radios the crude carrier recovery procedure is used then pilot tones transmitted along with the signal are used to develop a low-frequency version of the carrier

### 3.6 Modern Transmitter Architectures

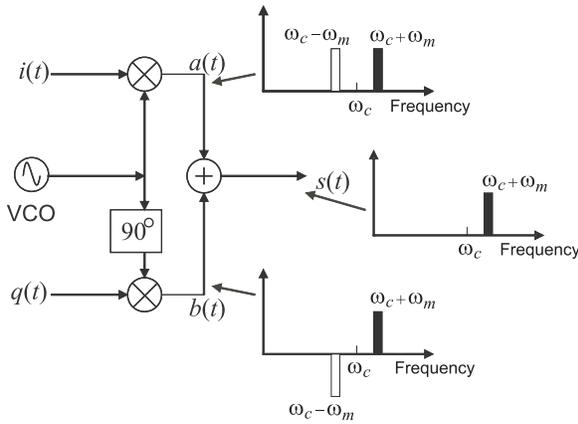
Modern transmitters maximize both spectral efficiency and electrical efficiency by using quadrature modulation and suppressing the carrier in the radio signal. Electrical efficiency must be achieved together with tight specifications on allowable distortion, and designs must achieve this with minimum manual adjustments. The discussion in this section focuses on narrowband communications when the modulated RF carrier can be considered as a slowly varying RF phasor. Modern radios must implement various order modulation schemes and also must implement legacy modulation methods.

#### 3.6.1 Quadrature Modulator

Most digital modulation schemes set the amplitude and phase of a carrier. The one exception is frequency shift keying (FSK) modulation, which changes the frequency of the carrier and has much in common with FM.



**Figure 3-12:** Phasor diagram with four discrete phase states that can be set by combining I and Q signals.



**Figure 3-13:** Quadrature modulator showing intermediate spectra.

So, in most digital modulation schemes, setting the amplitude and phase of a carrier addresses a point on a phasor diagram. A phasor diagram with four discrete states is shown in Figure 3-12. The circuit that implements this digital modulation is shown in Figure 3-11.

### 3.6.2 Quadrature Modulation

Quadrature modulation describes the frequency conversion process in which the real and imaginary parts of the RF phasor are varied separately. A subsystem that implements quadrature modulation is shown in Figure 3-13. This is quite an ingenious circuit. The operation of this subsystem is described by what is known as the generalized quadrature modulation equation:

$$s(t) = i(t) \cos [\omega_c t + \varphi_i(t)] + q(t) \sin [\omega_c t + \varphi_q(t)], \quad (3.4)$$

where,  $i(t)$  and  $q(t)$  embody the particular modulation rule for amplitude,  $\varphi_i(t)$  and  $\varphi_q(t)$  embody the particular modulation rule for phase, and  $\omega_c$  is the carrier radian frequency. In terms of the signals identified in Figure 3-13, the quadrature modulation equation can be written as

$$s(t) = a(t) + b(t) \quad (3.5)$$

$$a(t) = i(t) \cos [\omega_c t + \varphi_i(t)] \quad (3.6)$$

$$b(t) = q(t) \sin [\omega_c t + \varphi_q(t)]. \quad (3.7)$$

Figure 3-13 shows that both  $a(t)$  and  $b(t)$  have two bands, one above and one below the frequency of the carrier,  $\omega_c$ . However, there is a difference. The LO (here designated as the VCO) is shifted  $90^\circ$  (perhaps using an RC delay line) so that the frequency components of  $b(t)$  have a different phase relationship to the carrier than those of  $a(t)$ . When  $a(t)$  and  $b(t)$  are combined, the carrier content is canceled, as is one of the sidebands provided that  $q(t)$  is a  $90^\circ$  phase-shifted version of  $i(t)$ . This is exactly what is desired:

the carrier should not be transmitted, as it contains no information. Also, it is desirable to transmit only one sideband, as it contains all of the information in the modulating signal. This type of modulation is SSB-SC modulation. In the next section, frequency modulation is used to demonstrate SSB-SC operation.

### 3.6.3 Frequency Modulation

Frequency modulation is considered here to demonstrate SSB-SC operation. Let  $i(t)$  and  $q(t)$  be finite bandwidth signals centered at radian frequency  $\omega_m$  with their phases  $\phi_i(t)$  and  $\phi_q(t)$  chosen so that  $(\phi_q(t) - \phi_i(t))$  is  $90^\circ$  on average. This is shown in Figure 3-13, where  $\omega_m$  represents the frequency components of  $i(t)$  and  $q(t)$ . With reference to Figure 3-13,

$$i(t) = \cos(\omega_m t) \quad \text{and} \quad q(t) = -\sin(\omega_m t), \quad (3.8)$$

and the general quadrature modulation equation, Equation (3.4), becomes

$$s(t) = i(t) \cos(\omega_c t) + q(t) \sin(\omega_c t) = a(t) + b(t), \quad (3.9)$$

$$\text{where} \quad a(t) = \frac{1}{2} \{ \cos[(\omega_c - \omega_m)t] + \cos[(\omega_c + \omega_m)t] \} \quad (3.10)$$

$$\text{and} \quad b(t) = \frac{1}{2} \{ \cos[(\omega_c + \omega_m)t] - \cos[(\omega_c - \omega_m)t] \}. \quad (3.11)$$

Thus the combined frequency modulated signal at the output is

$$s(t) = a(t) + b(t) = \cos[(\omega_c + \omega_m)t], \quad (3.12)$$

and the carrier and lower sideband are both suppressed. The lower sideband,  $\cos[(\omega_c - \omega_m)t]$ , is also referred to as the image which may not be exactly zero because of circuit imperfections. In modulators it is important to suppress this image, and in demodulators it is important that undesired signals at the image frequency not be converted along with the desired signals.

### 3.6.4 Polar Modulation

In polar modulation, the  $i(t)$  and  $q(t)$  quadrature signals are converted to polar form as amplitude  $A(t)$  and phase  $\phi(t)$  components. This is either done in the DSP unit or, if a modulated RF carrier is all that is provided, using an envelope detector to extract  $A(t)$  and a limiter to extract the phase information corresponding to  $\phi(t)$ . Two polar modulator architectures are shown in Figure 3-14. In the first architecture, Figure 3-14(a),  $A(t)$  and  $\phi(t)$  are available and  $A(t)$  is used to amplitude modulate the RF carrier, which is then amplified by a power amplifier (PA). The phase signal,  $\phi(t)$ , is the input to a phase modulator implemented as a PLL. The output of the PLL is fed to an efficient amplifier operating near saturation (also called a saturating amplifier). The outputs of the two amplifiers are combined to obtain the large modulated RF signal to be transmitted.

In the second polar modulation architecture, Figure 3-14(b), a low-power modulated RF signal is decomposed into its amplitude- and phase-modulated components. The phase component,  $\phi(t)$ , is extracted using a limiter that produces a pulse-like waveform with the same zero crossings as the modulated RF signal. Thus the phase of the RF signal is captured. This

is then fed to a saturating amplifier whose gain is controlled by the carrier envelope, or  $A(t)$ . Specifically,  $A(t)$  is extracted using an envelope detector, with a simple implementation being a rectifier followed by a lowpass filter with a corner frequency equal to the bandwidth of the modulation.  $A(t)$  then drives a switching (and hence efficient) power supply that drives the **saturating power amplifier**.

### 3.7 Modern Receiver Architectures

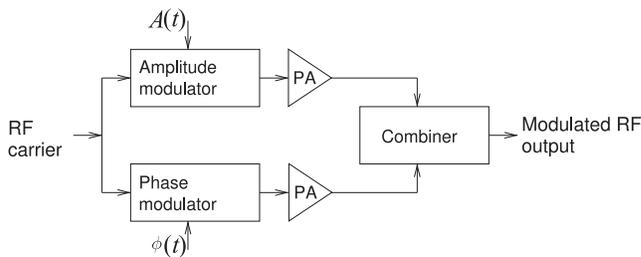
This section discusses transmitter and receiver architectures in the generations before software defined radio (as used in 4G and 5G). These are architectures that could be implemented in analog hardware.

#### 3.7.1 Receiver Architectures

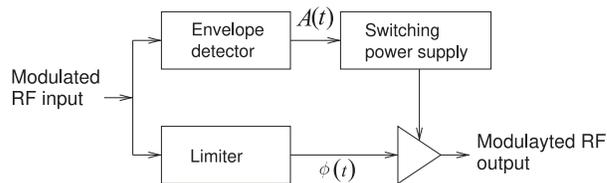
It is more challenging to achieve a high performance for a receiver than it is for a transmitter.

Communication receivers most commonly use mixing of the RF signal with a fixed signal called an LO to produce a lower-frequency replica of the modulated RF signal. Some receiver architectures use one stage of mixing, while others use two stages. In cellular systems, the receiver must be sensitive enough to detect signals of 100 fW or less.

Some of the architectures used in modern receivers are shown in Figure 3-15. Figure 3-15(a) is the superheterodyne architecture in much the same form that it has been used for a century. Key features of this architecture are that there are two stages of mixing, and filtering is required to suppress spurious mixing products. Each mixing stage has its own VCO. The receiver progressively reduces the frequency of the information-bearing signal. The image rejection mixer in the dashed box achieves rejection of the image frequency to produce an IF (or baseband frequency) that can be directly sampled. However, it is difficult to achieve the required amplitude and phase balance. Instead, the architecture shown in Figure 3-15(b) is sometimes used. The filter between the two mixers can be quite large. For example, if the

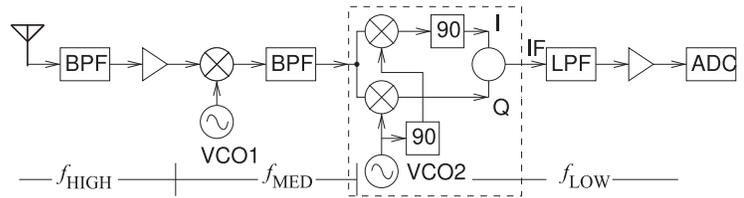


(a)

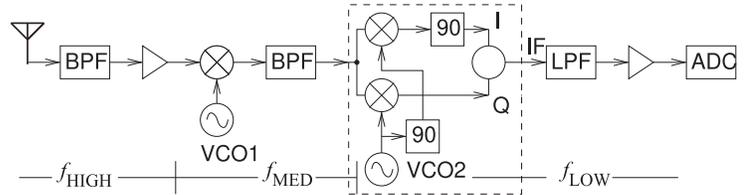


(b)

**Figure 3-14:** Polar modulator architectures: (a) amplitude- and phase-modulated components amplified separately and combined; and (b) the amplitude used to modulate a power supply driving a saturating amplifier with phase modulated input.

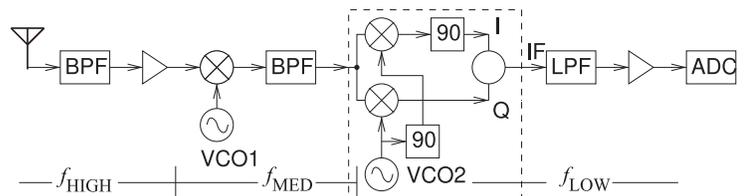


(a) Superheterodyne, Hartley image rejection receiver

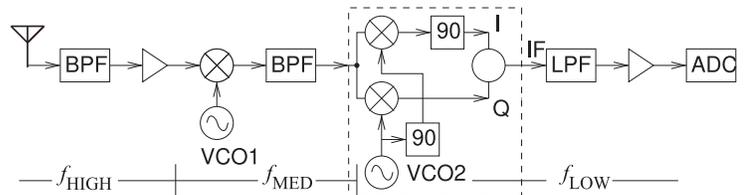


(b) Superheterodyne receiver

**Figure 3-15:** Architectures of modern receivers: (a) superheterodyne receiver using the **Hartley** architecture for image rejection; (b) superheterodyne receiver; (c) dual-conversion receiver; and (d) low-IF or zero-IF receiver. BPF, bandpass filter; LPF, lowpass filter; ADC, analog-to-digital converter; VCO, voltage-controlled oscillator; 90, 90° phase shifter;  $f_{HIGH}$ ,  $f_{MED}$ , and  $f_{LOW}$  indicate relatively high-, medium-, and low-frequency sections.



(c) Dual conversion receiver

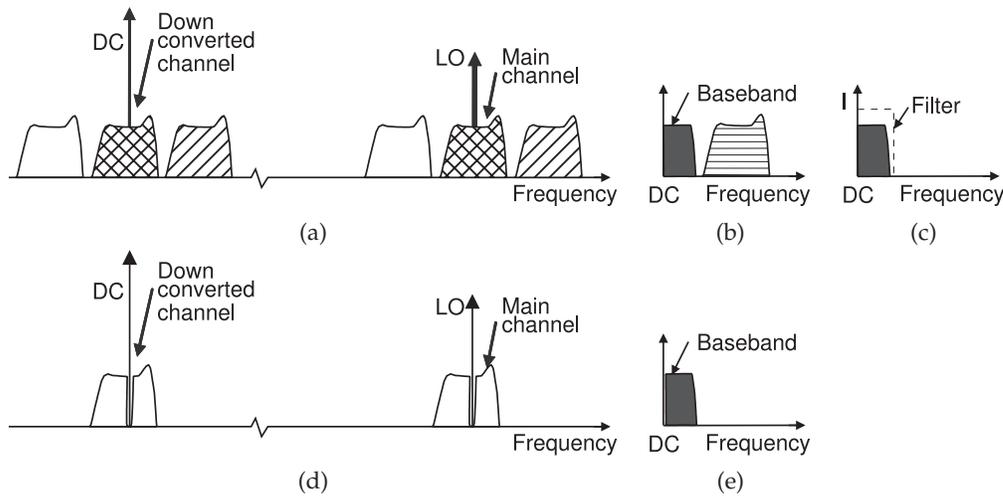


(d) Low-IF or zero-IF receiver

incoming signal is 1 GHz, the frequency of the signal after the first mixer could be 100 MHz.

Filters are smaller and have higher performance at higher frequencies. This is exploited in the dual-conversion receiver shown in Figure 3-15(c). This is similar to the traditional superheterodyne architecture except that the IF between the two mixers is high. For example, if the incoming signal is 1 GHz, the output of the first mixer could be 3 GHz. This architecture also enables broad radio operation with the band selected by choosing the frequencies of the two local oscillators.

The low-IF or zero-IF receiver shown in Figure 3-15(d) uses less hardware and is common in less demanding communication applications. In high-performance systems, such as the cellular phone system, this architecture requires more design time as well as calibration circuitry to trim the I and Q paths so that they are closely matched.



**Figure 3-16:** Frequency conversion using homodyne mixing: (a) the spectrum with a large LO and the low frequency products after down conversion; (b) the baseband spectrum showing only positive frequencies; (c) the baseband spectrum after mixing; (d) down conversion spectra when the radio signal has no spectral content at the carrier frequency; and (e) the lowpass filtered down-converted signal in (d).

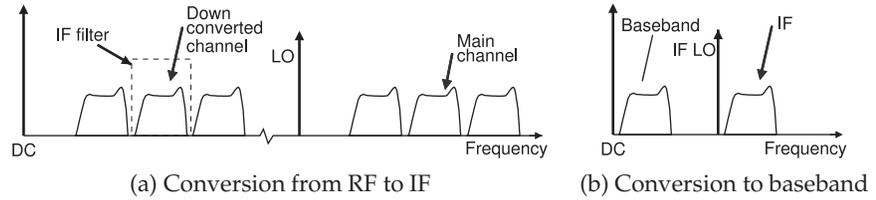
### 3.7.2 Homodyne Frequency Conversion

Homodyne mixing and detection is one of the earliest wireless receiver technologies and is used in AM radio. In homodyne mixing, the carrier of a modulated signal is regenerated and synchronized in phase with the incoming carrier frequency. Mixing the carrier with the RF signal results in an IF signal centered around zero frequency.

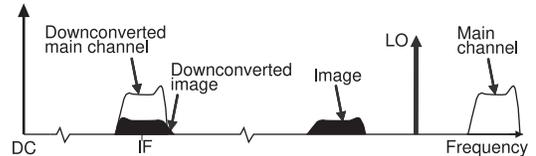
The signal spectra in homodyne mixing is shown in Figure 3-16. In Figure 3-16(a), the RF signals are shown on the right-hand side and the baseband signals are shown on the left-hand side. It is usual to show both positive and negative frequencies at the lower frequencies so that the conversion process is more easily illustrated. Of course negative frequencies do not exist. The characteristic of homodyne mixing is that the LO corresponds to the carrier and is in the middle of the desired RF channel. RF signal components mix with the LO, and it appears that the entire RF spectrum is down-shifted around DC. Of course, the actual baseband spectrum is only defined for positive frequencies, so the negative-frequency baseband signals and the positive-frequency baseband signals combine to yield the detected baseband spectrum shown in Figure 3-16(b). With other modulation schemes, this possible loss of information is avoided using quadrature demodulation. An amplitude modulated signal has identical modulation sidebands, so the collapsing of positive and what is shown as negative frequencies at baseband results in no loss of information. Then a simple amplitude detection circuitry, such as a rectifier, is used and the rectified signal was (typically) passed directly to a speaker.

**Figure 3-17:**

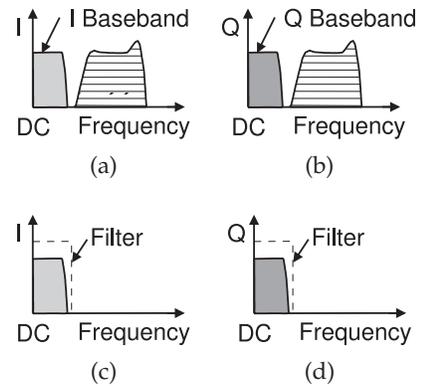
Frequency conversion using superheterodyne mixing.



**Figure 3-18:** Frequency conversion using heterodyne mixing showing the effect of image distortion with the down-converted image overlapping the down-converted main channel.



**Figure 3-19:** Frequency conversion using direct conversion quadrature mixing: (a) the baseband spectrum at the I output of the receiver; (b) the baseband spectrum at the Q output of the receiver; and (c and d) the spectrum of the I and Q channels following.

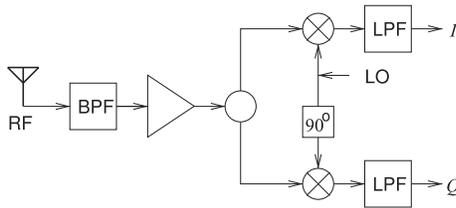


### 3.7.3 Heterodyne Frequency Conversion

In heterodyne mixing, the locally generated LO and the main RF channel are separated in frequency, as shown in Figure 3-17(a). In this figure the RF signals (shown as three discrete channels on the right-hand side of the spectrum) mix with the LO to produce signals at a lower frequency. This lower frequency is usually not the final baseband frequency, and so is called the intermediate frequency (IF). The IF of the main channel is at the difference frequency of the RF signal and the LO. There are several important refinements to this. The first of these is concerned with limiting the number of signals that can mix with the LO. This is done using an RF preselect filter. To see the difficulties introduced by the image channel, consider the frequency conversion to an intermediate frequency described in Figure 3-18. Filtering reduces the level of the image channel as shown in Figure 3-18(a). Note that the main channel and its image are equidistant from the LO, see Figure 3-18. Both down-convert to the same IF frequency. In the worst-case scenario, the IF image could be larger than that of the desired channel.

### 3.7.4 Direct Conversion Receiver

**Zero-IF** direct conversion receivers are similar to quadrature homodyne receivers in that the LO is placed near the center of the RF channel. The important characteristic of direct conversion receivers is that there is only one level of mixing. The conversion process is described in Figure 3-19. A particular advantage of direct conversion is that the relatively large IF filters are eliminated. They are invariably implemented as quadrature



**Figure 3-20:** Direct conversion quadrature demodulator.

demodulators, see Figure 3-20. used in cellular phones as it uses little DC power, and hence extends battery life, and is compatible with monolithic ICs. Direct conversion is now the preferred method of down conversion

The main nonideality of this design is the DC offset in the down-converted spectrum. DC offset results mostly from self-mixing, or rectification, of the LO. This DC offset can be much larger than the down-converted signal itself. One way of coping with the DC offset is to highpass filter the down-converted signal, but highpass filtering requires a large passive component (e.g., a series capacitor), at least to avoid the dynamic range problems of active filters. Highpass filtering the down-converted signal necessarily throws away information in the signal spectrum, and it is only satisfactory to do this if there is very little information around DC to begin with.

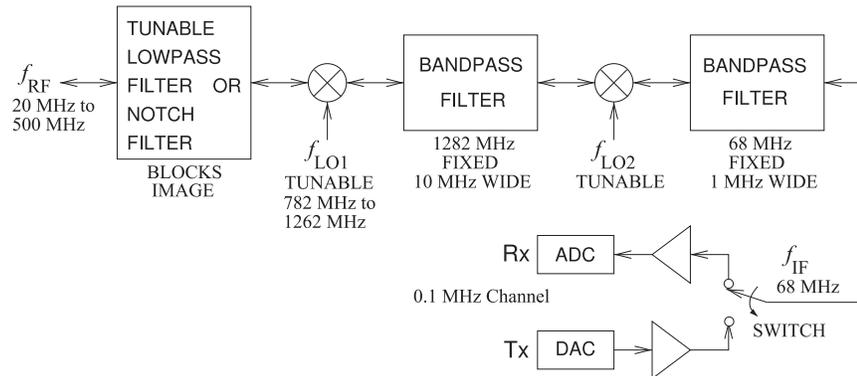
The primary design effort with zero-IF converters is overcoming the DC offset problem, and to a lesser extent coping with the jitter of the LO. However there is a scheme used in 4G and 5G that overcomes these limitations. This will be discussed in following chapters.

### 3.7.5 Low-IF Receiver

In a low-IF receiver, single-stage heterodyne mixing is used to down-convert the modulated RF carrier to a frequency just above DC, perhaps a few hundred kilohertz or a few megahertz, depending on the bandwidth of the RF channel. In doing this, the DC offset problem of a direct conversion receiver is avoided. This frequency offset can be just a few hundred hertz to be effective. Low IF conversion is now the preferred method of down-conversion in cellular phones as it requires very little batter usage, does not require large IF filters, and is compatible with monolithic ICs. Sometimes low-IF conversion is referred to as direct down conversion but there is a subtle difference.

### 3.7.6 Subsampling Analog-to-Digital Conversion

Subsampling receivers overcome the DC offset problem typical of other direct conversion receivers. The idea is to sample the modulated RF signal at a subharmonic of the carrier of the RF signal to be converted. The sampling rate must be at least twice the bandwidth of the baseband signal and the track-mode bandwidth must be greater than the carrier frequency. Thus the sampling aperture is the critical parameter and must be several times smaller than the period of the carrier. Fortunately the aperture times of CMOS tracking circuits are adequate. It is critical that an RF preselect filter be used to eliminate unwanted interferers and noise outside the communication band. Aliasing of signals outside the **Nyquist** bandwidth onto the baseband signal is a consequence of subsampling. Adjacent channel signals are converted without aliasing, but these will lie outside the bandwidth of the



**Figure 3-21:** Bilateral double-conversion transceiver for wideband operation of an emergency or military radio.

baseband signal. Flicker noise on the sampling clock is multiplied by the subsampling ratio and appears as additional noise at baseband. This was at one time a very attractive option but has been out-performed by the low-IF receiver.

### 3.7.7 First IF-to-Baseband Conversion

In a superheterodyne conversion architecture there are two heterodyne stages, with the IF of the first stage in modern systems in the range of 20 to 200 MHz. The assignment of frequencies is known as frequency planning, and this is treated as proprietary by the major radio vendors. This IF is then converted to a much lower IF, typically around 100 kHz to a few megahertz above the center frequency of the baseband signal. This frequency is generally called baseband, but strictly it is not because the signal is still offset in frequency from DC. Some direct conversion architectures leave the first heterodyne mixing stage in place and use direct conversion of the first IF to baseband (true baseband—around DC).

### 3.7.8 Bilateral Double-Conversion Receiver

The receivers considered so far are suitable for narrowband communications typical of point-to-point and consumer mobile radio. There are many situations where the range of received or transmitted RF signals covers a very wide bandwidth, such as with emergency radios, television, and military communications. Typically, however, the instantaneous bandwidth is small. If narrowband RF front-end architectures are used, a switchable filter bank would be required and this would result in an impractically large radio. One solution to covering very wide RF bandwidths is the double-conversion transceiver architecture shown in Figure 3-21. The frequency plan of a typical radio using 0.1 MHz channels between 20 MHz and 500 MHz is shown. The key feature of this radio is that bidirectional mixers are used. Following the RF chain from left to right, the RF is first mixed up in frequency, bandpass filtered using a high- $Q$  distributed filter, and then down-converted to a lower frequency that can be sampled directly by an ADC. A much higher performance passive (and hence bidirectional) filter can be realized

at gigahertz frequencies than at a few tens of megahertz. On transmit, the function is similar, with the mixers and LO reused. As a receiver, the notch filter or lowpass filter is used to block the image frequency of the first mixer so that only the upper sideband IF is presented to the first bandpass filter. The lowpass or notch filter may be fixed, although, with the plan shown, there must be at least two states of the filters. On transmit, the lowpass or notch filters prevent the image frequency from being radiated.

### 3.8 Introduction to Software Defined Radio

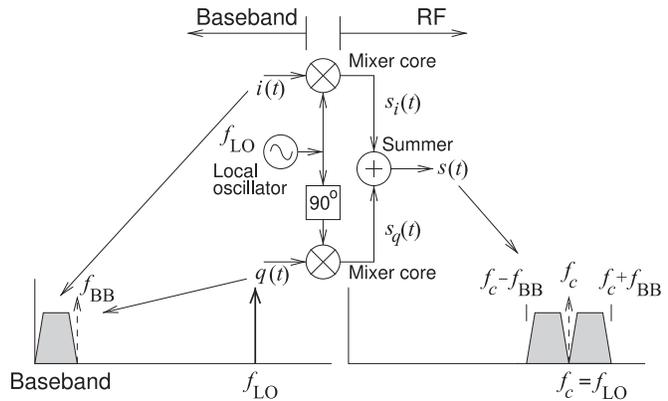
The 2G and 3G cellular radio schemes used very few modulation methods. With the introduction of 4G and now 5G many modulation methods must be supported with the highest-order modulation scheme used determined by the level of noise and interference on a channel. This support for multiple modulation methods is only possible if most of the demodulation process is performed in DSP where software controls demodulation. Such a radio is called a software defined radio (SDR). In an SDR many of the functions that traditionally would be performed using analog hardware are instead implemented digitally and only the RF functions close to the transmit/receive antenna are implemented in analog form. An extreme concept for the transmitter side of the radio is to implement all of the functions digitally with a final digital-to-analog converter (DAC) connected to an antenna. Then the performance is limited by the maximum frequency and output power of the DAC. This ultimate SDR provides maximum flexibility, for example, it is easy to change modulation schemes, but would also contribute to very high battery drain. So instead there is a trade-off with many aspects such as modulation and fine tuning of the RF carrier frequency done digitally while, with a transmitter, the final up-conversion done using analog hardware which only needs to be able to support the bandwidth of the modulated signal but otherwise is agnostic to the type of modulation used.

All radios today use quadrature modulation, a type of digital modulation, in which the real (in-phase) and imaginary (quadrature-phase) components of the phasor of the carrier are varied independently. While it is possible to vary the amplitude and phase of a carrier tone separately this is not common.

Quadrature demodulation recovers the original signals that varied the in-phase and quadrature-phase components of the carrier. Demodulation can proceed in stages with an initial analog separation of RF in-phase and quadrature signals which are sampled by an analog-to-digital converter (ADC) and demodulation completed in a digital signal processor.

The following sections discuss various aspects of an SDR radio. Section 3.9 begins with a description of quadrature modulation in a way that helps in the understanding of an SDR transmitter. In Section 3.10 a specific example of an SDR transmitter is presented in which the time-domain and frequency-domain signals are followed through first a digital signal processor (DSP) and then an analog up-converter to produce an RF signal. This discussion is followed by a description of an SDR receiver with the general SDR quadrature demodulator in Section 3.11 and then a specific example in Section 3.12.

**Figure 3-22:** Double-sideband suppressed-carrier (DSB-SC) modulation. Quadrature modulator with independent modulating baseband signals  $i(t)$ , the in-phase component input to the in-phase mixer core, and  $q(t)$ , the quadrature component driving the quadrature-phase mixer core. This results in a modulated signal with two sidebands around the local oscillator or carrier frequency. Each RF sideband has a mix of the information in  $i(t)$  and  $q(t)$ .



### 3.9 SDR Quadrature Modulator

The SDR transmitter uses two-stage modulation with DSB-SC modulation implemented in DSP to produce an IF signal which is output using DACs to produce analog I and Q channel IF signals. These IF signals are then input to an analog quadrature modulator implementing SSB-SC modulation with the resulting radio signal being a DSB-SC radio signal. As of the time of this writing the digital portion was implemented in what is called a baseband chip and the analog portion implemented in an RF modem chip. One can expect that eventually these would be combined into a single chip. As far as the RF modem chip is concerned, the IF signals input to the up-converter are baseband signals and this is how they are often referred to when the focus is on the RF modem chip.

Quadrature modulation, see Figure 3-22, comprises two mixer cores which are driven by a modulating in-phase component  $i(t)$  and a modulating quadrature-phase component  $q(t)$  where in-phase and quadrature phase refer to the phase of the local oscillator input to the mixer cores. Here  $i(t)$  and  $q(t)$  are baseband signals with a spectrum extending from (near) DC to  $f_{BB}$  and in today's radios they are produced internally in a DSP unit. The finite bandwidth  $i(t)$  and  $q(t)$  signals contain  $I$ -channel and  $Q$ -channel information respectively. The top mixer core is driven directly by the local oscillator and the other, the quadrature mixer core, is also driven by the local oscillator but now it is phase shifted by  $90^\circ$ , i.e. it is in quadrature. This scheme produces double-sideband suppressed carrier DSB-SC modulation and  $s(t)$  is the modulated output signal with each sideband having bandwidth  $f_{BB}$ . The block-level schematic illustrates the basic architecture of a quadrature modulator which is expanded if the signals are differential signals with additional variations according to whether the mixers are implemented as analog multipliers or as switches controlled by the LO. The whole structure shown is referred to as a mixer and each mixer core on its own is also often referred to as a mixer. This operation can be implemented without error in DSP.

#### 3.9.1 Analog Quadrature Modulator

The second stage of an SDR transmitter implements DSB-SC modulation using analog circuitry producing an RF signal.

An analog quadrature modulator using multipliers is shown in Figure

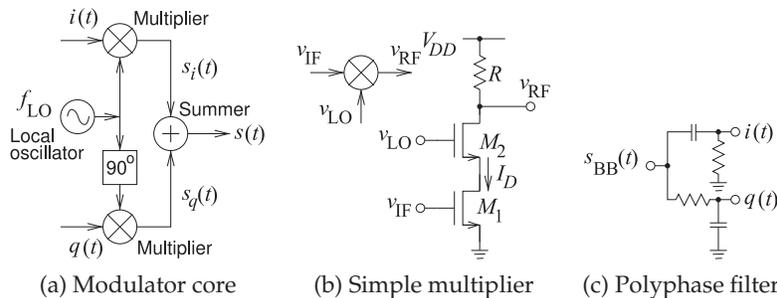
3-23(a) and consists of two multipliers each of which has two inputs and one output with the outputs summed yielding a modulated output signal  $s(t)$ . One particular characteristic of a quadrature modulator is that the LO at frequency  $f_{LO}$  is directly input to one of the multipliers but the second multiplier is driven by a version of the LO with a  $90^\circ$  phase lag, i.e. the LO input to the first mixer is in-phase and the LO input to the second multiplier has quadrature-phase (the phase is shifted by  $90^\circ$ ). This second LO is also called the quadrature LO. The  $90^\circ$  phase difference of the two LOs is where the quadrature in quadrature modulator comes from. So if the LO is  $\sin(2\pi f_{LO})$ ,  $i(t)$  is multiplied by  $\sin(2\pi f_{LO})$ . Then  $q(t)$  is multiplied by  $\sin(2\pi f_{LO} - \pi/2) = -\cos(2\pi f_{LO})$ . The second inputs of the multipliers in Figure 3-23(a) are the signals  $i(t)$  and  $q(t)$  with  $i$  indicating that the signal is driving the in-phase multiplier and  $q$  indicating that the signal is driving the quadrature-phase multiplier. The signals  $i(t)$  and  $q(t)$  may be independent, or the frequency components of  $q(t)$  may phase lag  $i(t)$  by  $90^\circ$  but otherwise be the same as  $i(t)$ . These two options yield modulated output signals with different bandwidths.

**Transistor-Based Multiplier**

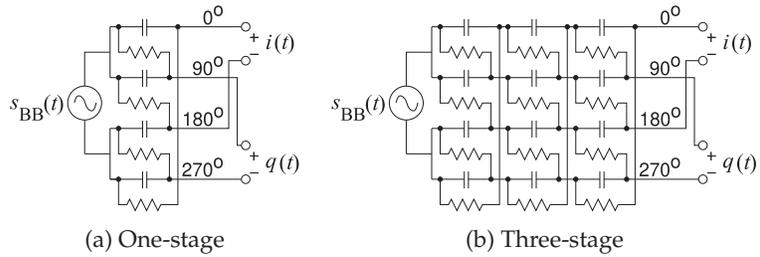
There are several ways to implement the mixer core in Figure 3-23(a) with the most common being as a multiplier or as a switch. Both can be conveniently implemented using transistors. The analog multiplier shown in Figure 3-23(b) is based on a cascode amplifier with one input applied to the gate of transistor  $M_1$ . Instead of the gate of  $M_2$  being held at a DC voltage as with a cascode amplifier, the gate of  $M_2$  is also an input. Approximately, the drain current,  $I_D$ , of  $M_1$  is proportional to the gate voltage  $v_{IF}$  and the voltage gain of  $M_2$ , i.e.  $v_{RF}/v_{LO}$  is proportional to  $I_D$ . Thus the RF output voltage  $v_{RF}$  is proportional to the product of  $v_{IF}$  (which in the modulator is either  $i(t)$  or  $q(t)$ ) and  $v_{LO}$ . So when  $v_{IF}$  and  $v_{LO}$  are sinewaves the output  $v_{RF}$  will be the trigonometric expansion of the product of two sinewaves and this product will also comprise two sinewaves at the sum and difference frequencies. Then circuit symmetry is used to select just one of these.

**Polyphase Filter**

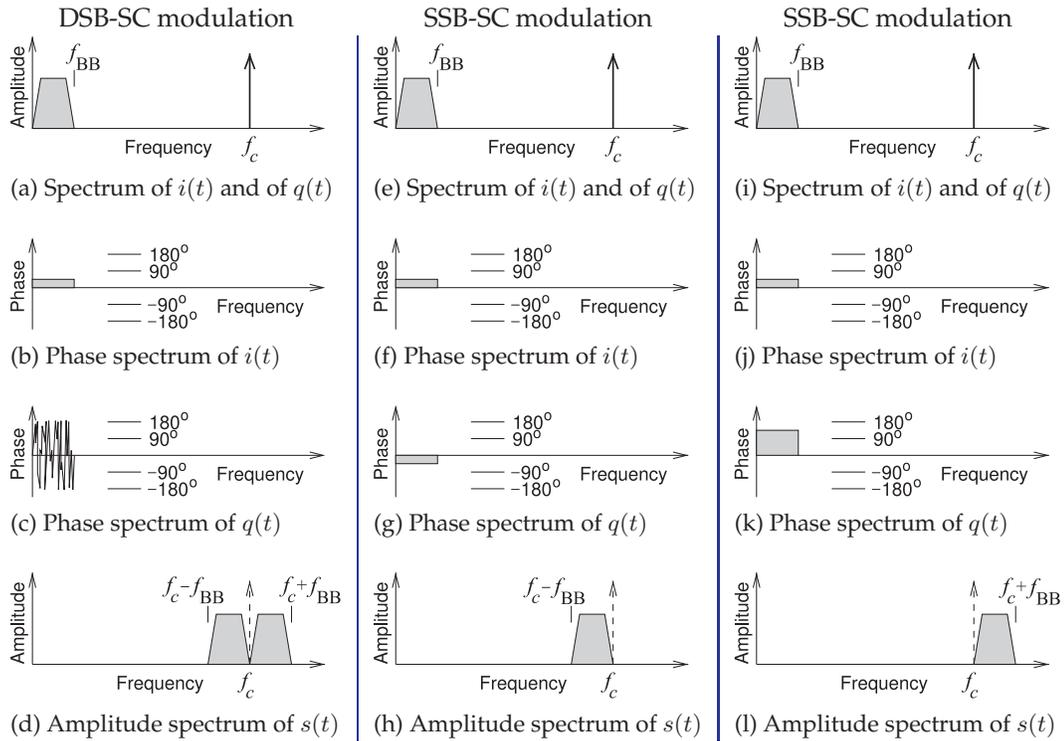
A polyphase filter, such as the one-stage polyphase filter in in Figure 3-23(c), takes an input analog input signal and outputs two signals that are the same except that the frequency components are shifted by  $90^\circ$ . This circuit can be used to produce the quadrature LO signal or to shift the frequency components of a baseband signal. More commonly a polyphase filter is



**Figure 3-23:** Quadrature modulator. The simple modulator in (b) is based on a FET cascode amplifier and the polyphase filter in (c) has a  $90^\circ$  phase difference between  $i(t)$  and  $q(t)$ .



**Figure 3-24:** Differential polyphase filters.



**Figure 3-25:** Spectra of the baseband  $i(t)$  and  $q(t)$  signals and the modulated  $s(t)$  signal, see Figure 3-23(a), for double sideband (DSB) and single sideband (SSB) suppressed-carrier (SC) modulation. Here the carrier frequency  $f_c = f_{LO}$ , the LO frequency in Figure 3-23(a).

realized in differential form as shown in Figure 3-24(a). The polyphase filters in Figures 3-23(c) and 3-24(a) are narrowband but the bandwidth can be increased using more stages, see Figure 3-24(b).

**Double Sideband Modulation**

When  $i(t)$  and  $q(t)$  are independent, effectively pseudo-random, signals the result is double sideband (DSB) modulation, see Figures 3-25(a–d) with each sideband having the bandwidth of the baseband signals. The amplitude spectra of  $i(t)$  and  $q(t)$  will be the same as shown in Figure 3-25(a) and each

has a bandwidth  $f_{\text{BB}}$ .<sup>4</sup> As seen in Figure 3-25(b) the frequency components of the  $i(t)$  spectrum are arbitrarily assigned a  $45^\circ$  phase. Since  $i(t)$  and  $q(t)$  are independent the phase of  $q(t)$  relative to the phase of  $i(t)$  is random, see Figure 3-25(c). The modulated signal  $s(t)$  then has a sideband below the carrier frequency  $f_c$  and a sideband above  $f_c$ , see Figure 3-25(d), for a total bandwidth  $2f_{\text{BB}}$ .

The multipliers in a quadrature modulator are implemented as mixer circuits and one type of mixer in particular is the multiplicative mixer shown in Figure 3-23(a). Ideally a multiplicative mixer multiplies two sinewaves together to produce the trigonometric expansion of the product of two sinewaves. For example,  $\sin(A) \cdot \sin(B) = \frac{1}{2}[\cos(A - B) - \cos(A + B)]$ . Following the signal paths in Figure 3-23(a) and considering the frequency component  $A_i(f_i) \sin(\omega_i t)$  of  $i(t)$  at radian frequency  $\omega_i = 2\pi f_i$  and a frequency component  $A_q(f_q) \sin(\omega_q t)$  of  $q(t)$  at radian  $\omega_q = 2\pi f_q$  the modulated signal with the LO frequency replaced by  $f_c$  (with radian carrier frequency  $\omega_c = 2\pi f_c$ ) is

$$\begin{aligned}
 s(t) &= s_i(t) + s_q(t) \\
 &= [A_i(f_i) \sin(\omega_i t) \sin(\omega_c t)] + [A_q(f_q) \sin(\omega_q t) \sin(\omega_c t - \pi/2)] \\
 &= [A_i(f_i) \sin(\omega_i t) \sin(\omega_c t)] - [A_q(f_q) \sin(\omega_q t) \cos(\omega_c t)] \\
 &= \frac{1}{2} [A_i(f_i) \cos(\omega_c - \omega_i)t - A_i(f_i) \cos(\omega_c + \omega_i)t \\
 &\quad - A_q(f_q) \sin(\omega_c + \omega_q)t + A_q(f_q) \sin(\omega_c - \omega_q)t] \\
 &= \frac{1}{2} [A_i(f_i) \cos(\omega_c - \omega_i)t + A_q(f_q) \sin(\omega_c - \omega_q)t \\
 &\quad - \frac{1}{2} [A_i(f_i) \cos(\omega_c + \omega_i)t + A_q(f_q) \sin(\omega_c + \omega_q)t]. \tag{3.13}
 \end{aligned}$$

The DSB-SC modulated signal is the signal in Equation (3.13) summed for all  $f_i$  and  $f_q$  components from DC to  $f_{\text{BB}}$ .

The expansion in Equation (3.13) can be repeated for all of the frequency components of  $i(t)$  and  $q(t)$ . So while the expansion is only performed for discrete frequencies, all that is necessary is that the multiplier be practically ideal, something that is typically achieved by an analog multiplier to better than 1%. If the DSB-SC signal was produced digitally then multiplication can be precisely implemented and the DSB-SC modulation is ideal although the maximum frequency is limited by the capabilities of the digital circuitry. A DSP-based DSB-SC modulation has a low carrier frequency as keeping the (digital) carrier frequency low reduces DC power requirements. If  $i(t)$  and  $q(t)$  are independent, Equation (3.13) indicates a lower modulated sideband at the range of frequencies  $(f_{\text{LO}} - f_i)$  and  $(f_{\text{LO}} - f_q)$  and an upper modulated sideband at the range of frequencies  $(f_{\text{LO}} + f_i)$  and  $(f_{\text{LO}} + f_q)$  for all  $f_i$  and  $f_q$  from 0 (DC) to  $f_{\text{BB}}$ . That is, this is DSB-SC modulation, as seen in Figure 3-25(d). In demodulation both sidebands are needed to recover  $i(t)$  and  $q(t)$ .

### Single Sideband Modulation

When  $i(t)$  and  $q(t)$  are the same signal except that every frequency component of  $q(t)$  is shifted by  $90^\circ$  the result is single sideband (SSB) modulation and the modulated output signal has a bandwidth  $f_{\text{BB}}$ . The

<sup>4</sup> The short-term spectra will be different because  $i(t)$  and  $q(t)$  are different signals but over a long time interval the envelope of the amplitude spectra will become similar.

carrier itself does not exist in the output with a quadrature modulator using multipliers so this modulator implements SSB suppressed-carrier (SSB-SC) modulation. The modulated output signal is obtained with  $q(t) = A_i(f_i) \sin(\omega_i - \pi/2) = -A_i(f_i) \cos(\omega_i)$ . Then Equation (3.13) becomes (but now  $f_{LO}$  is used to distinguish it from the carrier frequency which is defined by the characteristics of the modulating signal)

$$\begin{aligned}
 s(t) &= s_i(t) + s_q(t) \\
 &= [A_i(f_i) \sin(\omega_i t) \sin(\omega_{LO} t)] - [A_i(f_i) \cos(\omega_i t) \sin(\omega_{LO} t - \pi/2)] \\
 &= A_i(f_i) [\sin(\omega_i t) \sin(\omega_{LO} t) + \cos(\omega_i t) \cos(\omega_{LO} t)] \\
 &= \frac{1}{2} A_i(f_i) \{ \cos[(\omega_{LO} - \omega_i)t] - \cos[(\omega_{LO} + \omega_i)t] \\
 &\quad + \cos[(\omega_{LO} + \omega_i)t] + \cos[(\omega_{LO} - \omega_i)t] \} \\
 &= A_i(f_i) \cos[(\omega_{LO} - \omega_i)t].
 \end{aligned} \tag{3.14}$$

Equation (3.14) indicates that just the lower sideband is present and this is SSB-SC modulation as seen in Figure 3-25(h) and the bandwidth of the modulated output signal is  $f_{BB}$ . The original  $i(t)$  signal can be recovered from this one sideband but that is because  $q(t)$  contains exactly the same information as  $i(t)$  (although phase shifted).

If instead the phase of each frequency component of  $q(t)$  led the same frequency component of  $i(t)$  by  $+90^\circ$  then  $s(t)$  would comprise the upper sideband and this is still this would be SSB-SC modulation, see Figures 3-25(i-l). For SSB-SC modulation each frequency component of  $q(t)$  must have a phase that differs from the corresponding component of  $i(t)$  by  $90^\circ$ . A lumped-element circuit that realizes this is the polyphase filter, see Figure 3-23(c), but the phase shift can also be realized in DSP.

Earlier, just before Equation (3.14), it was stated that the frequency of the carrier was defined by the characteristics of the modulated signal which in turn depends on the characteristics of the modulating signal. This modulating signal, the  $i(t)$  input to the SSB-SC modulator, could also be modulated as is usually the case in SDR where DSB-SC modulation is done in a DSP and this is followed by SSB-SC modulation done at RF using analog hardware. Identifying the correct RF carrier is required for demodulation. In identifying the carrier there are two situations to consider. If the input,  $i(t)$ , of the SSB-SC modulator is not modulated, e.g. it is just a baseband signal, then the carrier frequency is just the frequency of the LO of the SSB-SC modulator as shown in Figure 3-26(a), i.e.  $f_c = f_{LO}$ . If the input signal to the SSB-SC modulator is itself a DSB-SC signal (produced by a DSB-SC modulator) so that it has its own intermediate carrier frequency  $f_{c,IF}$ , then the carrier frequency  $f_c = f_{LO} - f_{c,IF}$ . This situation is shown in Figure 3-26(b). (Note that the carrier frequency would be above  $f_{LO}$  if the frequency components of  $q(t)$  were advanced in phase by  $90^\circ$  relative to the phase of the frequency components of  $i(t)$ .)

### 3.9.2 Summary

This section discussed quadrature modulation and showed how the same circuit can be used for DSB and for SSB modulation. The difference is in whether or not  $i(t)$  and  $q(t)$  are related. In modern radios DSB is implemented in DSP to produce an IF modulated signal with the spectra shown in Figure 3-25(d) and  $f_c$  is very low, perhaps even  $f_c = f_{BB}$ . Then

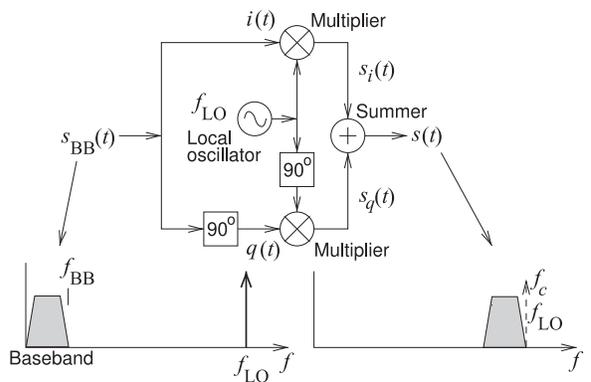
this DSB-SC signal becomes the baseband of an analog SSB modulator that produces the RF modulated signal. This RF modulated signal is a DSB signal with a (suppressed) carrier in the middle of the spectrum.

### 3.10 Case Study: SDR Transmitter

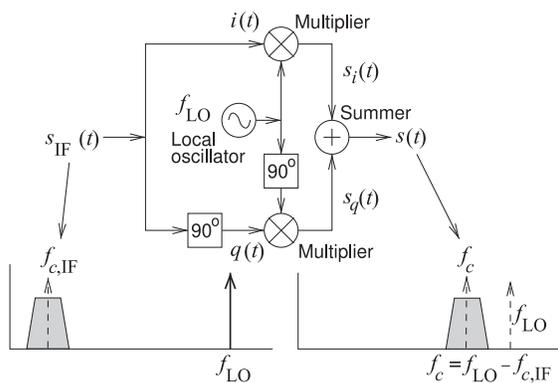
An SDR transmitter combines DSB-SC modulation having a relatively low intermediate frequency carrier with a broadband analog SSB-SC modulation to produce a DSB-SC RF signal. There are many possible implementations. This case study presents system simulations of an SDR transmitter with specific parameters. First a DSB-SC quadrature modulator is studied with sinusoidal in-phase and quadrature baseband signals and then a SSB-SC modulator is examined. This is followed by the study of a direct analog modulation of a digital signal. The final study corresponds to a typical SDR transmitter with DSP-based intermediate frequency DSB-SC modulation followed by SSB-SC analog modulation producing a DSB-SC modulated RF signal.

#### 3.10.1 Analog Quadrature Modulator

The quadrature modulator of Figure 3-23(a) is simulated here with a 10 MHz sinewave for  $i(t)$ , a 15 MHz sinewave for  $q(t)$ , and a 1 GHz sinewave LO. The local oscillator is a 1 GHz sinewave. The resulting RF waveform at the



(a) SSB-SC modulator with baseband input signal  $s_{BB}(t)$



(b) SSB-SC modulator with IF input signal  $s_{IF}(t)$

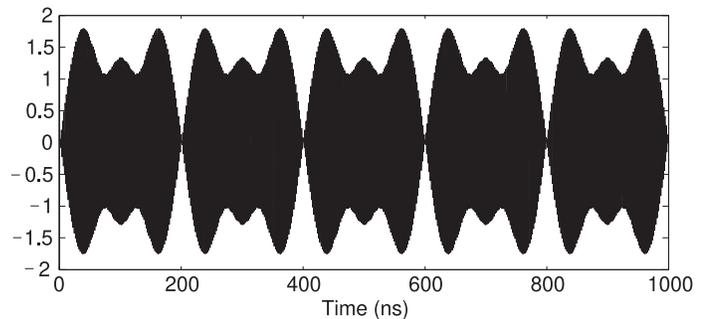
**Figure 3-26:** Quadrature modulator with the  $i(t)$  component derived directly from the input signal and  $q(t)$  derived from a  $90^\circ$  negatively phase-shifted input signal. The IF signal in (b) is typically a DSB-SC signal with intermediate carrier frequency  $f_{c,IF}$  set in DSP.

output,  $s(t)$ , is shown in Figure 3-27. The waveform is plotted on a 1  $\mu\text{s}$  scale so that there are 10 cycles of  $i(t)$  and 15 cycles of  $q(t)$  which are modulated on 1,000 cycles of the 1 GHz RF carrier.

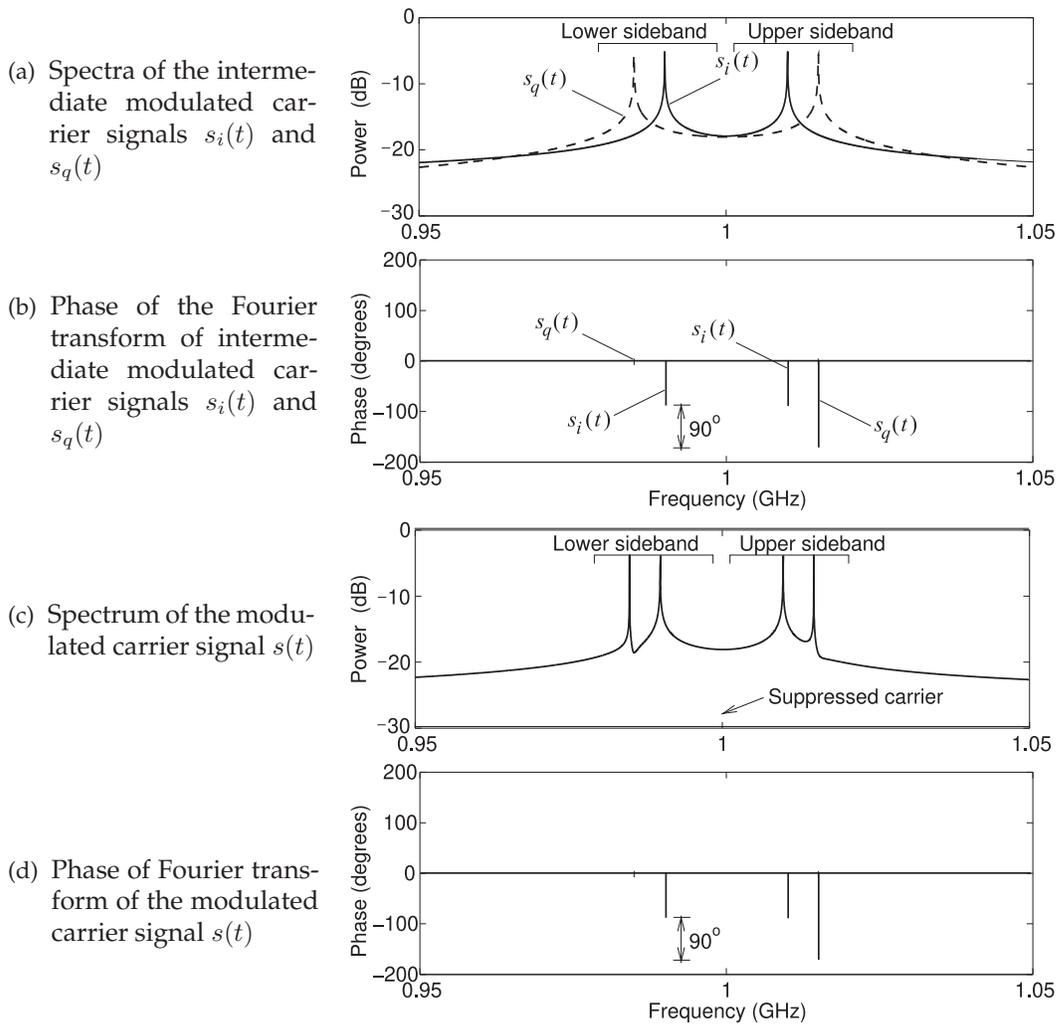
The spectra of the signals at the outputs of the two multipliers, i.e. of  $s_i(t)$  and  $s_q(t)$ , are shown in Figure 3-28(a). Each spectrum has two peaks offset from the 1 GHz LO frequency by 10 MHz for the  $I$  signal,  $s_i(t)$ , and by 15 MHz for the  $Q$  signal,  $s_q(t)$ . The amplitudes of each spectra are symmetrical around the LO frequency but there is a difference in the phase of the Fourier transforms of  $s_i(t)$  and  $s_q(t)$ . The phases of the Fourier transforms of the signals are shown in Figure 3-28(b). The upper and lower sideband phases of the 10 MHz  $I$  channel signal are equal but the phases of the upper and lower sidebands of the 15 MHz  $Q$  channel signal differ by  $180^\circ$ , i.e. the upper sideband of  $s_q(t)$  is the negative of the lower sideband of  $s_q(t)$ . In contrast the phases of the upper and lower realized sidebands of  $s_i(t)$  are equal.

The amplitude and phase spectra of the combined signal,  $s(t) = s_i(t) + s_q(t)$ , are shown in Figures 3-28(c and d) and are as expected from combining the spectra of the components. The spectrum of  $s(t)$ , consists of two pairs of peaks with the peaks of one pair being above and below the LO frequency by 10 MHz, and the peaks of the other pair being offset by 15 MHz. The components are in a lower sideband and an upper sideband and hence this modulation is DSB modulation. Also there is not a component of the output signal at the local oscillator frequency,  $f_{LO}$ . The middle of the output signal bandwidth is generally the carrier frequency  $f_c$  which here is the same as  $f_{LO}$ . Thus the carrier does not exist in the output and so this is called suppressed carrier (SC) modulation. Together this is double sideband suppressed carrier (DSB-SC) modulation. Suppression of the carrier is a property of the multiplicative mixer but some types of modulators, e.g. amplitude modulators, have the carrier in the RF output signal and hence the suppressed-carrier distinction being used with quadrature modulation.

Note that each of the sidebands of the DSB-SC modulated signal have both  $I$  and  $Q$  information. With finite and equal bandwidth  $i(t)$  and  $q(t)$  signals, the  $I$  and  $Q$  information in the RF modulated signal will have overlapping lower sideband and upper sideband. It is only possible to recover and separate the original  $i(t)$  and  $q(t)$  signals if both sidebands are used in demodulation.



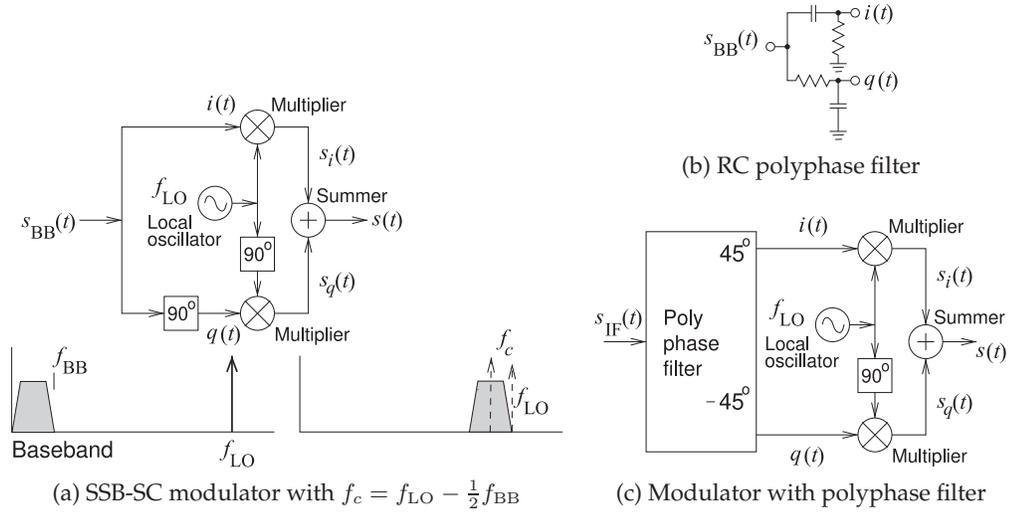
**Figure 3-27:** A 1 GHz carrier modulated by a 10 MHz sinusoidal  $i(t)$  and a 15 MHz sinusoidal  $q(t)$ .



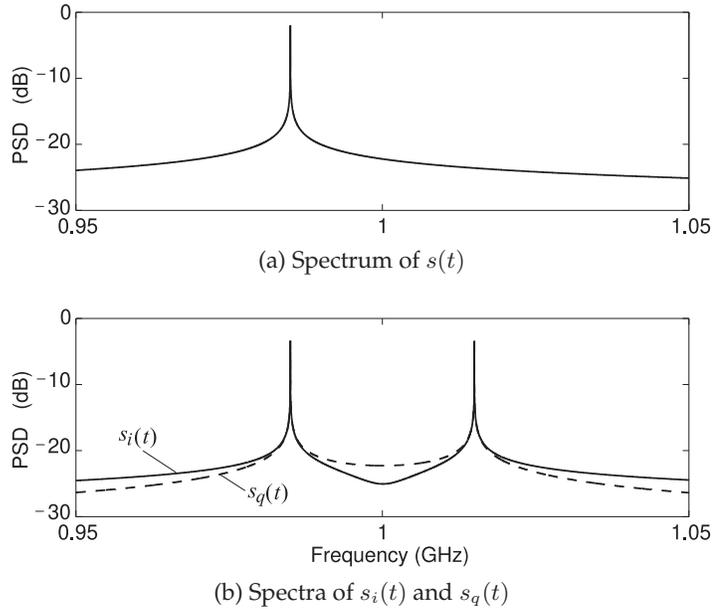
**Figure 3-28:** Spectrum of signals in the quadrature modulator of Figure 3-23(a) with 10 MHz in-phase,  $i(t)$ , and 15 MHz quadrature-phase,  $q(t)$ , modulating signals. (Simulated time = 65.536  $\mu$ s, using a  $2^{23}$  point FFT.)

### 3.10.2 Single-Sideband Suppressed-Carrier (SSB-SC) Modulation

Examination of the phase differences of the lower and upper sideband components with the DSB-SC modulator considered in the previous section leads to the design of a SSB-SC modulator. This is obtained when  $i(t)$  and  $q(t)$  are the same signal except that the frequency components of  $q(t)$  lag those of  $i(t)$  by  $90^\circ$ . Thus the phase components of  $s_q(t)$  shown in Figure 3-28(b) are shifted by  $-90^\circ$  so that the lower sideband components of  $s_i(t)$  and  $s_q(t)$  (now having the same frequency) combine constructively but the upper sideband components cancel. The variation of the quadrature modulator that implements this is shown in Figure 3-29(a) with the baseband signal



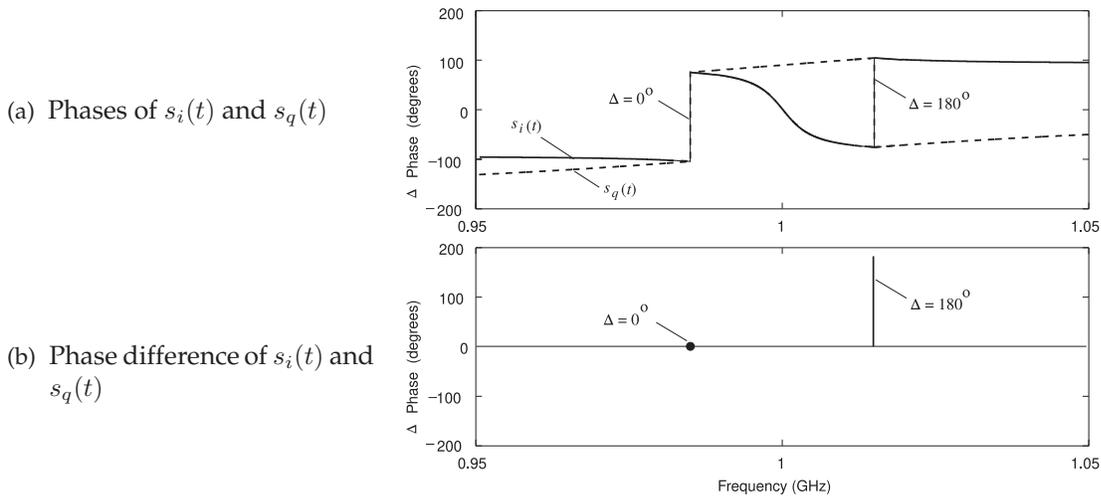
**Figure 3-29:** Quadrature modulator as a single-sideband, suppressed-carrier (SSB-SC) modulator. At the design frequency the polyphase filter in (b) the phase of  $i(t)$  is advanced by  $45^\circ$  and the phase of  $q(t)$  is retarded by  $45^\circ$ .



**Figure 3-30:** Spectra of the modulated signal with  $f_{LO} = 1$  GHz. (Simulated time =  $65.536 \mu s$ ,  $2^{23}$  point FFT.)

$s_{BB}(t) = i(t)$ , and  $q(t)$  is the same signal except that  $q(t)$  lags  $i(t)$  by  $90^\circ$ . The signals in this modulator will now be examined.

With  $s_{BB} = i(t)$  being a 15 MHz sinusoidal signal, the  $q(t)$  is a 15 MHz sinewave that lags  $i(t)$  by  $90^\circ$ . The spectrum at the output of the quadrature modulator of Figure 3-29(a and c) is as shown in Figure 3-30(a). Now there is only one sideband. Insight into the operation of SSB modulation is obtained by examining the spectra of the signals at the output of the multipliers, see Figure 3-30(b). It is seen that the amplitude spectra of  $s_i(t)$  and of  $s_q(t)$  are



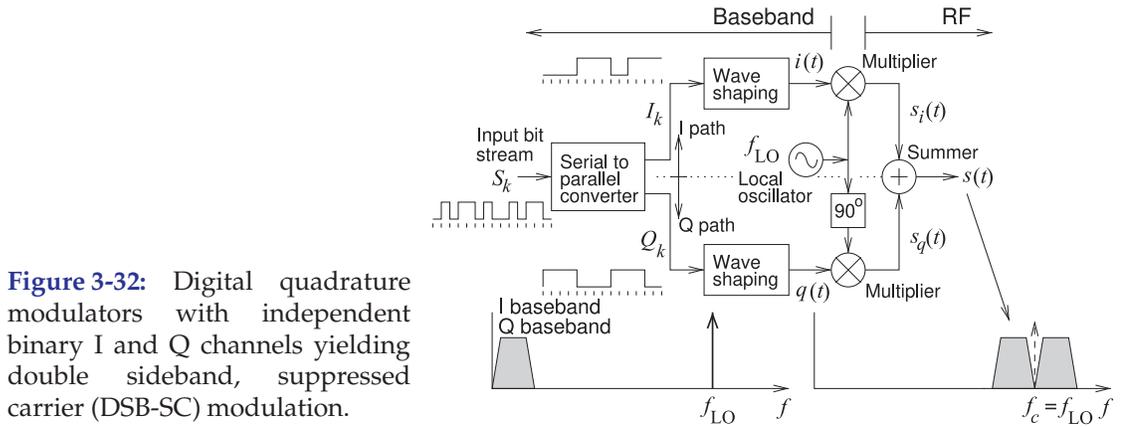
**Figure 3-31:** Phase of the modulated signal. (Simulated time = 65.536  $\mu$ s,  $2^{23}$  point FFT.)

essentially the same and both have upper and lower sidebands. However there is a difference in the phases of their Fourier transforms, see Figure 3-31. It is seen that at 15 MHz offset from the 1 GHz LO there is no differences in the phases of  $s_i(t)$  and  $s_q(t)$  components in the lower sideband but there is a  $180^\circ$  difference in the upper sideband, see Figure 3-31(b). Thus when  $s_i(t)$  and  $s_q(t)$  are summed their lower sideband components add but their upper sideband components cancel thus suppressing the upper sideband in the combined signal  $s(t)$ . This illustrates the most important metric of a SSB modulator as having I/Q paths that are matched as any imbalance in implementation that results in an amplitude or phase difference of the I and Q paths will result in a spurious (upper) sideband. Provided that the balance is good a filter is not required to suppress an upper sideband.

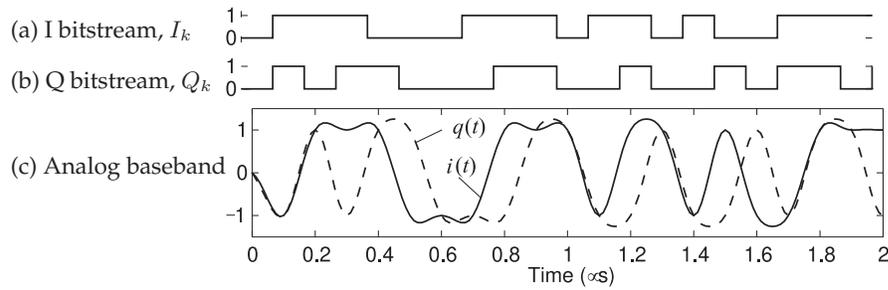
The SSB-SC modulator was modeled here with a sinusoidal input signal,  $s_{BB}(t)$ , but upper sideband suppression will also be obtained with a finite bandwidth baseband signal. This is achieved if every frequency component of  $s_{BB}(t)$  is  $90^\circ$  phase-shifted to become  $q(t)$  so that  $q(t)$  is the quadrature-phase version of the in-phase  $i(t)$ . A circuit that implements this is the polyphase filter and one type is the RC circuit shown in Figure 3-29(b). This has a finite bandwidth but this can be broadened by using multiple stages, e.g. see Figure 3-24(b). Incorporating the polyphase filter in the quadrature modulator leads to the implementation in Figure 3-29(c).

### 3.10.3 Digital Quadrature Modulation

A digital quadrature modulator is shown in Figure 3-32. The input bitstream  $S_k$  is converted into two independent binary bitstreams  $I_k$  and  $Q_k$  which are then lowpass filtered, or in general wave-shaped, to provide two independent analog signals  $i(t)$  and  $q(t)$ . That is, each pair of bits in the  $S_k$  bitstream becomes one  $I_k$  bit and one  $Q_k$  bit. The binary waveforms of  $I_k$  and  $Q_k$  are then filtered to obtain the analog signals  $i(t)$  and  $q(t)$  each of which has the same baseband bandwidth indicated in the spectrum on the lower



**Figure 3-32:** Digital quadrature modulators with independent binary I and Q channels yielding double sideband, suppressed carrier (DSB-SC) modulation.

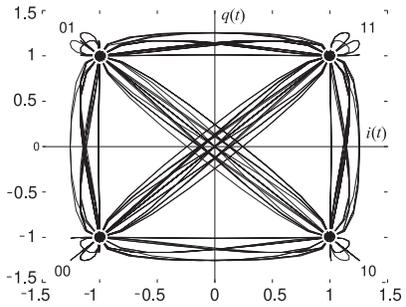


**Figure 3-33:** Baseband signals for a four state binary quadrature modulator. The  $I_k$  and  $Q_k$  bitstreams are derived from the random bitstream  $S_k = 00\ 11\ 10\ 11\ 01\ 00\ 10\ 11\ 11\ 00\ 10\ 11\ 00\ 10\ 01\ 00\ 11\ 11\ 10$  and taken as  $(I_k, Q_k)$  pairs. The initial setting is  $i(0) = 0 = q(0)$ .

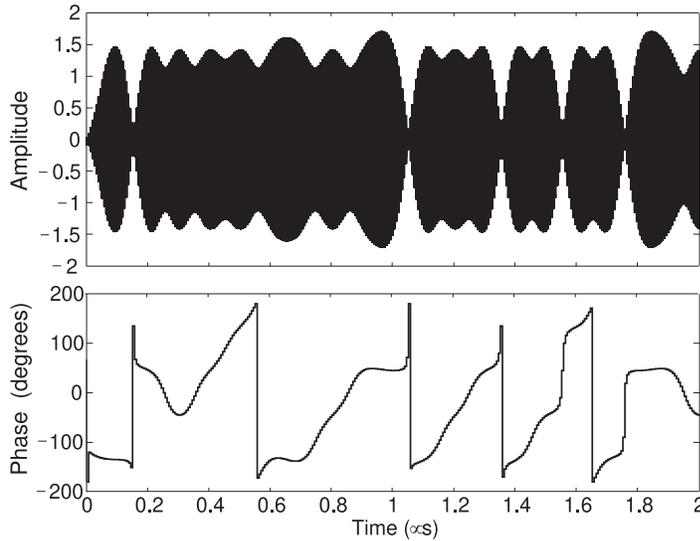
left in Figure 3-32. The in-phase baseband signal,  $i(t)$ , is multiplied by the local oscillator signal having frequency  $f_{LO}$ . Also, the quadrature baseband signal,  $q(t)$ , is multiplied by the local oscillator signal having frequency  $f_{LO}$  but now the LO is delayed by  $90^\circ$ . Summing  $s_i(t)$  and  $s_q(t)$ , the outputs of the multipliers, yields the double sideband suppressed carrier RF signal  $s(t)$  which has twice the bandwidth of each of the analog baseband signals as indicated in the lower right in Figure 3-32.

The modulator of Figure 3-32 results in four states of the modulated carrier signal. That is, if the modulated carrier signal  $s(t)$  is sampled at a time indicated by the common clock of  $I_k$  and  $Q_k$ , the sample of the carrier will have a particular amplitude and phase. Since this is QPSK modulation the amplitudes at these multiple sampling instances will be the same but the phases will have four values addressed by the coordinate  $(I_k, Q_k)$ .

Now consider a particular quadrature multiplier where  $S_k$  is a 20 Mbit/s random bitstream so that  $I_k$  and  $Q_k$  are the 10 Mbit/s bit streams shown in Figure 3-33(a and b). The  $I_k$  and  $Q_k$  bitstreams are then level shifted so that the transitions are between  $-1$  and  $1$  instead of between  $0$  and  $1$ . The level-shifted bitstreams are then lowpass filtered to obtain the analog baseband signals  $i(t)$  and  $q(t)$  shown in Figure 3-33(c). The lowpass filter used here is a raised cosine filter which is commonly used in digital modulation rather than a lumped-element lowpass filter and here has an effective corner frequency



**Figure 3-34:** Constellation diagram for the binary quadrature modulator. This is also the phasor diagram (with appropriate scaling) of the modulated carrier signal with constellation points sampled at  $0.1 \mu\text{s}$  intervals shown as large dots. There are four constellation points corresponding to four symbols and each symbol represents two bits of information.

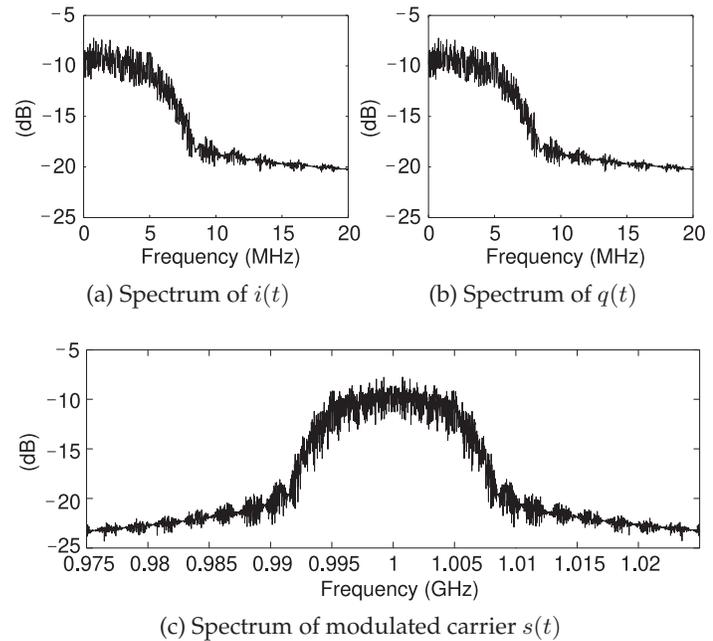


**Figure 3-35:** Amplitude and phase of the modulated carrier signal  $s(t)$  as the output of a binary modulator with a 20 Mbit/s digital baseband signal. Note that the rapid phase transitions occur when the amplitude of the modulated carrier goes to zero. (The rapid switch between  $180^\circ$  and  $-180^\circ$  when the amplitude is not zero is a continuous smooth phase change.)

of 7 MHz. In practice the raised cosine filter is implemented in DSP so that the signals  $i(t)$  and  $q(t)$  follow DACs.

One feature of the raised cosine filter is that the filtered response at the clock-derived sampling times is exact. That is, if the filtered analog baseband signals are sampled every  $0.1 \mu\text{s}$  then the sampled values of  $i(t)$  and  $q(t)$  will be exactly either  $+1.00$  or  $-1.00$  and (after level-shifting by adding one and multiplying by half) the original bitstream is exactly recovered as 0 or 1. If an analog lowpass filter was used the sampled values would not be exactly right. The raised cosine filter introduces no sampling distortion and also the transitions are minimal, i.e. they have the minimum bandwidth compared to what would be obtained if an analog filter was used for wave-shaping. The bandwidth of the filtered signal obtained using a lumped-element filter would be 10 MHz or more and even then the carrier samples would never correspond exactly with the constellation points. Plotting  $q(t)$  against  $i(t)$  yields the transitions shown in Figure 3-34. Samples of the baseband signal every  $0.1 \mu\text{s}$  coincides with one of the four states and enables two bits of information to be recovered.

The next stage of the DSB-SC modulator multiplies the  $i(t)$  signal by a 1 GHz LO and the  $q(t)$  signal is multiplied by a  $90^\circ$  phase-shifted LO. Then the outputs of the multipliers are summed to produce the modulated RF signal  $s(t)$  shown in Figure 3-35. The LO here is in the middle of the RF bandwidth and so here the carrier frequency is the same as the LO frequency.



**Figure 3-36:** Spectra of signals in the digital modulator with a 20 Mbit/s input stream,  $S_k$ , and a 1 GHz LO. (1024 symbols, a time duration of 102.4  $\mu$ s, and using a 524,288 point FFT)

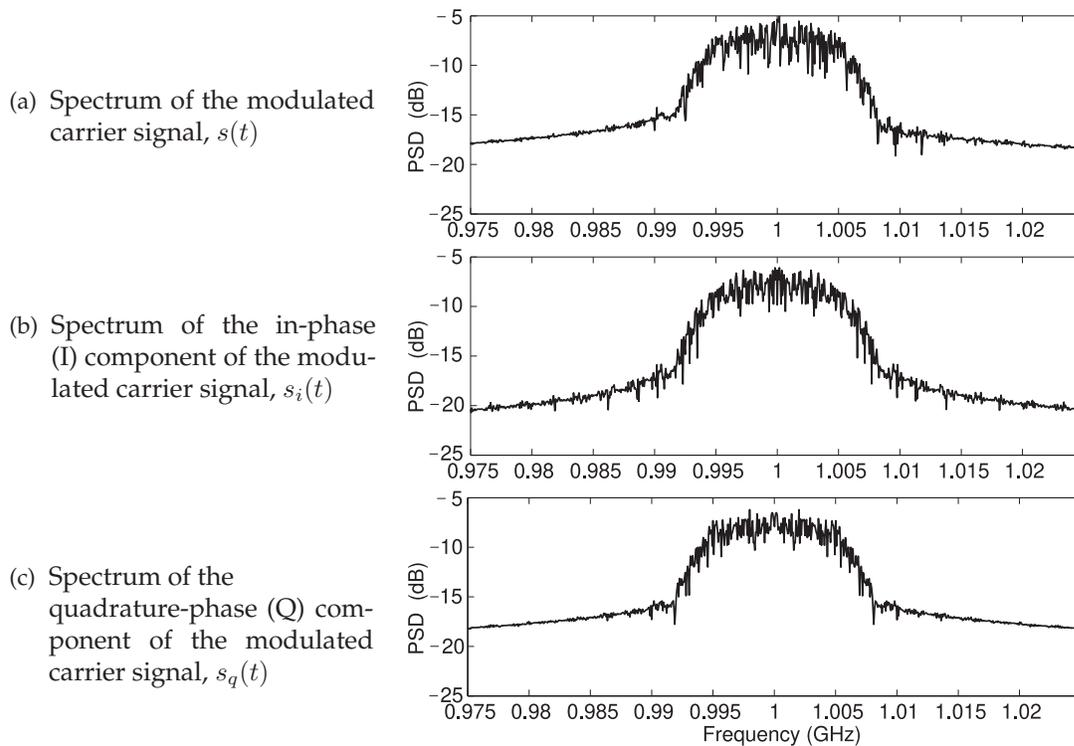
Sampling the RF signal  $s(t)$  every 0.1  $\mu$ s provides the amplitudes and phases of the carrier with each sample coinciding with one of the constellation points. This modulation scheme is generally called QPSK modulation for quadrature phase shift keying but it is sometimes, but less accurately, called quadrature phase shift keying. In QPSK modulation each constellation point corresponds to one of four phases of the RF signal:  $45^\circ$ ,  $135^\circ$ ,  $-45^\circ$ , or  $-135^\circ$ .

QPSK modulation was originally implemented with simpler hardware so the QPSK term is used rather than 4-QAM for four-state quadrature amplitude modulation which more accurately reflects the modulation process described in this section, it does not have to be done this way if only the phase of the carrier is to be adjusted.

Plotting the phasor of the RF carrier signal will result in a phasor diagram identical to the transitions shown in the constellation diagram of Figure 3-34 although it may be necessary to scale the amplitude of the phasor signal to match the values in the constellation diagram. Note that the constellation diagram does not change when the average signal level changes. With the sampling clock aligned to the original clock for the  $I_k$  and  $Q_k$  bitstreams, samples of the RF phasor will precisely coincide with one of the constellation points in Figure 3-34 (and hence  $s(t)$  is on the way to being demodulated).

Digital radio sends data in finite length packets and here the packet is 2048 bits long. With 2 bits per symbol in QPSK modulation there are 1024 symbols and with a symbol interval of 0.1  $\mu$ s, the duration of the packet is 102.4  $\mu$ s.

The spectra of the baseband and the output RF signals in the modulator are shown in Figure 3-36. The spectra of  $i(t)$  and  $q(t)$ , Figures 3-36(a and b) respectively, are not identical because each bitstream is independent. For each bitstream the spectra extends down to DC. It would be possible to use a coding scheme for the bitstreams that ensures that there is not a DC component and this aids automatic carrier recovery in pre-4G cellular systems. Instead 4G and 5G systems use separate mechanisms enabling carrier recovery. The bandwidth of the in-phase and quadrature-phase

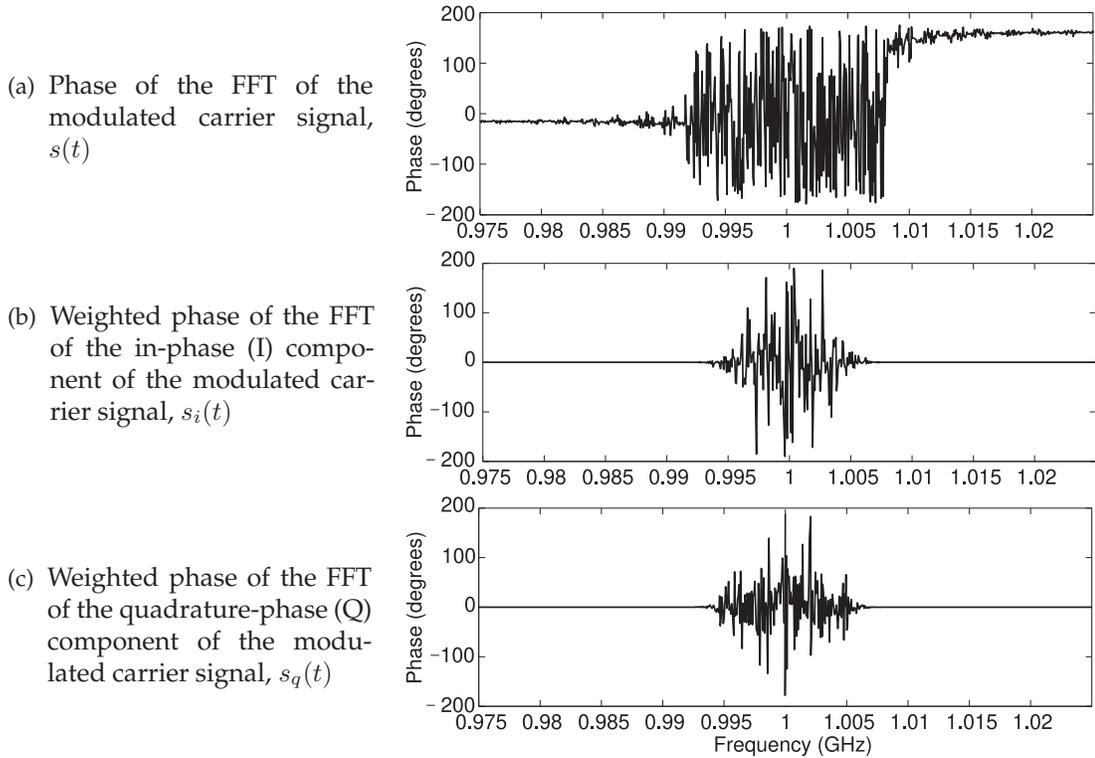


**Figure 3-37:** Spectra of the digitally modulated RF signals for 1024 symbols (4096 bits), a time duration of  $102.4 \mu\text{s}$ , and using a 524,288 point FFT.

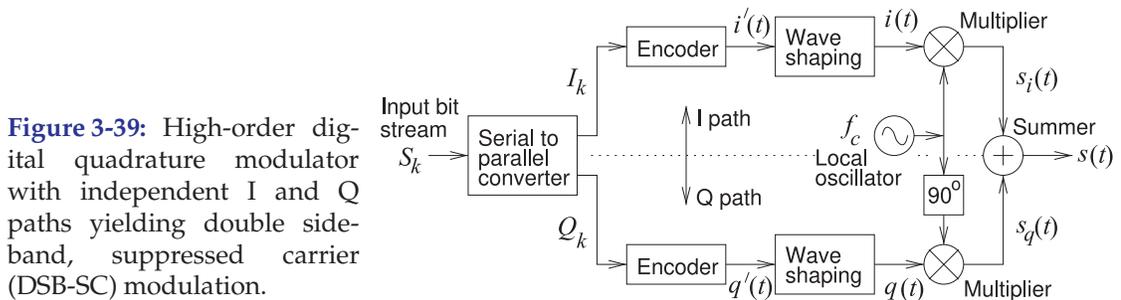
analog baseband signals are slightly less than 10 MHz. The spectrum of the RF output signal is shown in Figure 3-36(c) and the bandwidth of this double-sideband modulated signal is slightly less than 20 MHz. The RF signal has two sidebands, one below 1 GHz and one above even though there is no clear demarcation such as a dip at 1 GHz. The continuous spectrum through 1 GHz is a consequence of the baseband spectra extending to DC.

As was noted previously, and in the absence of noise, the constellation points are faithfully replicated if the RF signal is sampled at  $0.1 \mu\text{s}$  intervals (provided that the phase and the frequency of the carrier has been replicated accurately). The replication of the constellation points is a property of the raised cosine filter which also results in reduced bandwidth baseband analog signals than would be obtained if analog lowpass filtering was used. Also note that (with the DSB-SC quadrature modulation described here) the carrier frequency, the center of the RF signal spectrum, is 1 GHz, the same as the LO frequency.

The spectra of all of the RF signals are shown in Figure 3-37. The spectrum of the RF output, Figure 3-37(a), is centered at 1 GHz and this is the carrier frequency and is also the LO frequency for this DSB-SC transmitter. A curiosity is examining the phases of the RF signals, however little insight is obtained from the phase plots in Figure 3-38.

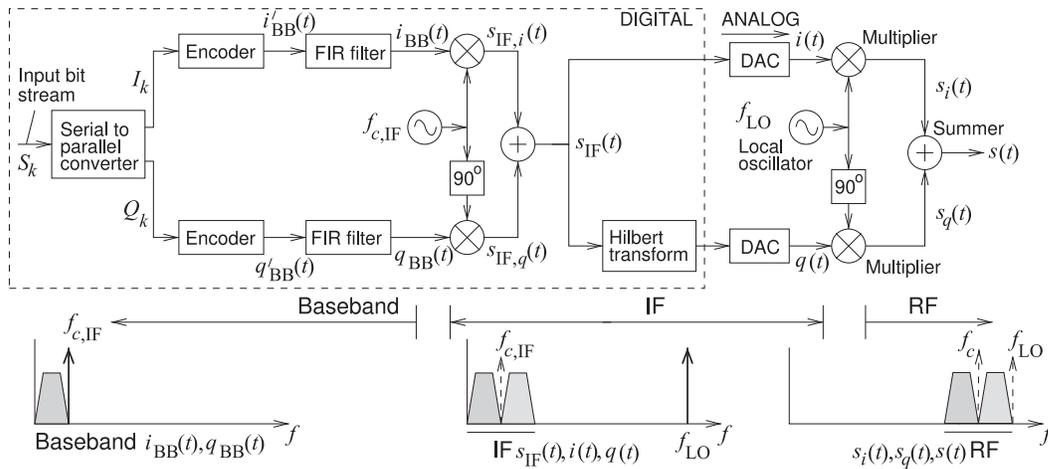


**Figure 3-38:** Phase of the FFT of the digitally modulated RF signal and its components. The phase of the  $s_i(t)$  and  $s_q(t)$  has been weighted by their spectrum amplitudes as otherwise meaningless numerical noise dominates the plots outside the passband. (Using 1024 symbols, 102.4  $\mu$ s)



### 3.10.4 QAM Digital Modulation

The QPSK digital modulator considered in the previous section had only four states. This is a consequence of the digital baseband bitstreams,  $I_k$  and  $Q_k$ , being used one bit at a time. If two or more bits are used at a time to address a constellation point, more data can be sent in the same RF bandwidth. Modulation that does this is called higher-order digital modulation. A higher-order quadrature modulator is shown in Figure 3-39 where now groups of two or more  $I_k$  bits and  $Q_k$  bits can be encoded as an analog signal. That is, if there are four possible values of  $i(t)$  at the clock ticks, then pairs of  $I_k$  bits are encoded as one of four analog values. Such as the  $I_k$



**Figure 3-40:** Complete SDR transmitter with bistream  $S_k$  input to an IF quadrature modulator producing a DSB-SC signal,  $s_{IF}$ , with IF carrier of frequency  $f_{c,IF}$ . Each finite-impulse response (FIR) filter implements raised cosine filtering. Then an analog SSB-SC quadrature modulator up-converts the IF as a DSB-SC RF signal centered at the RF carrier frequency  $f_c = f_{LO} - f_{c,IF}$ . The Hilbert transform implements a  $90^\circ$  phase shift.

bit pair 00 resulting in  $i(t)$  before filtering, i.e.  $i'(t)$ , being encoded as  $-1$  V, 01 being encoded as  $-\frac{1}{3}$  V, 10 being encoded as  $+\frac{1}{3}$  V, and 11 being encoded as  $+1$  V.  $Q_k$  would be encoded the same way. This results in 16 possible states and this is called 16-QAM modulation. Operation of this quadrature modulator will be analyzed in the next section in the context of a complete SDR transmitter.

### 3.10.5 SDR Transmitter Using QAM Digital Modulation

This section describes a 16-QAM modulator as used in modern communications systems (4G and 5G). This reflects the capabilities of baseband DSPs which can perform many of the modulation operations. The architecture of a QAM quadrature-modulator is depicted in block diagram form in Figure 3-40. This modulator consists of two quadrature modulators, a DSP-based digital IF quadrature modulator on the left, and an analog RF quadrature modulator on the right. The digital IF quadrature modulator on the left is a DSB-SC quadrature modulator taking an input bitstream and partitioning it by converting from a serial input bit stream,  $S_k$ , into two parallel but independent bitstreams,  $I_k$  and  $Q_k$ . The output of the IF quadrature modulator,  $s_{IF}(t)$ , is a DSB-SC IF signal centered at the IF carrier frequency  $f_{c,IF}$ . This is then modulated by the analog RF quadrature modulator to create a SSB-SC RF signal which has the IF signal,  $s_{IF}(t)$ , as its modulating signal. The RF signal however retains the double sidebands of the IF signal so that  $s(t)$  is a DSB-SC signal with carrier frequency  $f_c = f_{LO} - f_{c,IF}$ .

The 16-QAM system has the parameters given in Table 3-1. Since pairs of  $I_k$  bits and pairs of  $Q_k$  bits are used to encode the baseband analog signals  $i'_{BB}$  and  $q'_{BB}$  there are 16 possible states and 16-QAM modulation is the result. The constellation diagram of 16-QAM modulation is shown in

Input bitstream $S_k$	40 Mbit/s		
$I_k$ bitstream	20 Mbit/s		
$Q_k$ bitstream	20 Mbit/s		
Nominal bandwidth $i'_{\text{BB}}$	10 MHz (discrete)	Bandwidth $i(t)$	20 MHz (analog)
Nominal bandwidth $q'_{\text{BB}}$	10 MHz (discrete)	Bandwidth $q(t)$	20 MHz (analog)
Intermediate carrier, $f_{c,\text{IF}}$	10 MHz (discrete)	Local oscillator, $f_{\text{LO}}$	1 GHz
Nominal bandwidth $s_{\text{IF},i}$	20 MHz (discrete)	Bandwidth $s_i(t)$	20 MHz (analog)
Nominal bandwidth $s_{\text{IF},q}$	20 MHz (discrete)	Bandwidth $s_q(t)$	20 MHz (analog)
Nominal bandwidth $s_{\text{IF}}$	20 MHz (discrete)	Bandwidth $s(t)$	20 MHz (analog)
IF carrier, $f_{c,\text{RF}}$	10 MHz	RF carrier, $f_{c,\text{RF}}$	990 MHz

**Table 3-1:** Parameters of the SDR transmitter considered in Section 3.10.4.

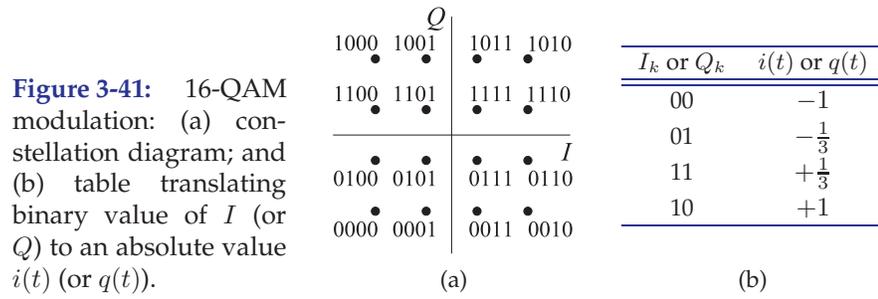
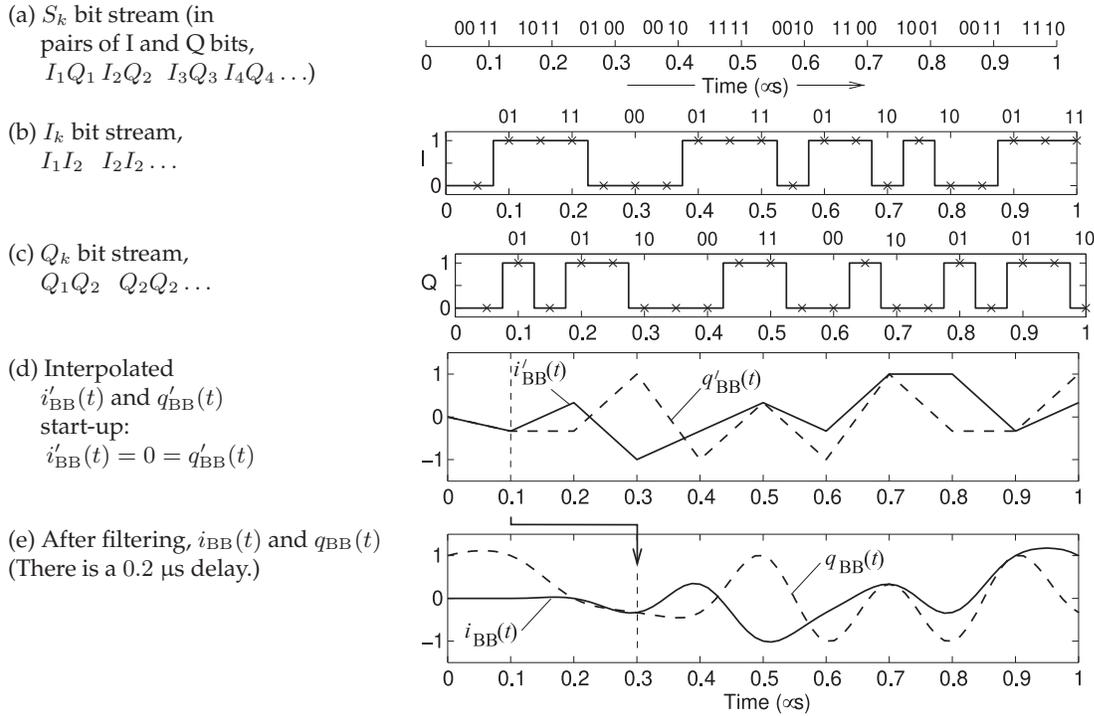


Figure 3-41(a) with the bit-wise addresses of the constellation points. There are 16 constellation points each of which indicates a symbol which provides four bits of information. In the case of QAM modulation the constellation diagram corresponds to a phasor diagram with each constellation point indicating the amplitude and phase of the RF signal at a point in time. The difference between the constellation diagram and a phasor diagram is that the constellation diagram does not change when the average power of the RF signal changes. Thus the constellation diagram of QAM corresponds to a phasor diagram that is continuously being re-normalized to the average RF power.

The table in Figure 3-41(b) shows the encoding that enables pairs of  $I_k$  bits and pairs of  $Q_k$  bits to address a symbol in the constellation diagram. That is, each of these bitstreams is encoded so that two bits are encoded as one analog value. Figure 3-42(a) shows the bits in the input bitstream,  $S_k$ , which are in groups of four bits with the first and third bits of the group becoming  $I$  bits and the second and fourth bits in the group becoming  $Q$  bits. The  $I_k$  and  $Q_k$  bitstreams are shown in Figures 3-42(b) and 3-42(c) respectively. A pair of  $I_k$  ( $Q_k$ ) bits yields an encoded  $i'_{\text{BB}}$  ( $q'_{\text{BB}}$ ) value every  $0.1 \mu\text{s}$ . Here  $I_k = I_1 I_2 I_3 I_4 I_5 I_6 = 011100$  so  $i'(0.1 \mu\text{s}) = -\frac{1}{3}$ ,  $i'(0.2 \mu\text{s}) = +\frac{1}{3}$ , and  $i'(0.3 \mu\text{s}) = -1$ . Linearly interpolating these yields the  $i'_{\text{BB}}$  and  $q'_{\text{BB}}$  waveforms shown in Figure 3-42(d). In DSP of course these are discrete values but have been plotted as continuous waveforms which is more easily visualized.

FIR filtering of  $i'_{\text{BB}}$  and  $q'_{\text{BB}}$  yields the band-limited waveforms  $i(t)$  and  $q(t)$  shown in Figure 3-42(e). Filtering introduces a delay and here the delay is  $0.2 \mu\text{s}$ . Thus the first encoded value before filtering is established at  $0.1 \mu\text{s}$ , see

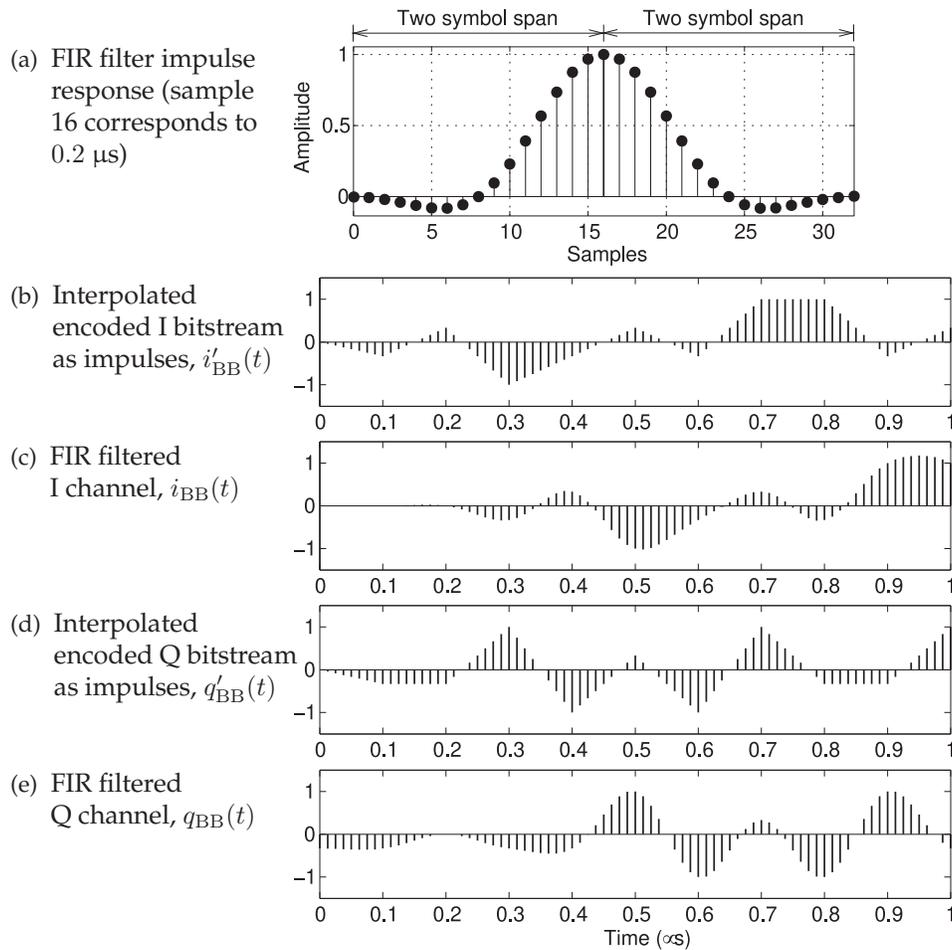


**Figure 3-42:** Baseband signals for 16-QAM modulation with markers showing the baseband bit impulse times. The bit pairs shown above each bit waveform 0.1  $\mu\text{s}$  intervals.

the vertical dashed line in Figure 3-42(d), and the first symbol (represented by  $i_{BB}$  and  $q_{BB}$ ) is at 0.3  $\mu\text{s}$ , see the vertical dashed line in Figure 3-42(e). The baseband values are actually discrete samples. With an oversampling factor of 8, the  $i'_{BB}$  samples, the interpolated encoded values, are as shown in Figure 3-43(b) and the interpolated  $q'_{BB}$  samples in Figure 3-43(d). In this SDR transmitter baseband filtering is a raised cosine filter implemented digitally as a finite impulse response (FIR) filter.

The interpolated encoded samples  $i'_{BB}$  and  $q'_{BB}$  are filtered digitally using a finite impulse response (FIR) filter yielding the filtered samples  $i_{BB}$  and  $q_{BB}$  shown in Figures 3-43(c and e) respectively. The raised-cosine FIR filter has a span of 4 symbols (i.e. 0.4  $\mu\text{s}$ ) and since the oversampling factor is 8 the FIR filter's response is 33 samples long, see Figure 3-43(a). (With the 0th and 33rd bins being zero, it could be said to be 31 or 32 samples long.) (Specification of the raised cosine filter is completed with a roll-off factor of 0.7 which describes how quickly the impulse response rolls-off around its peak response.) The peak response of the filter is at the 16th sample, after 0.2  $\mu\text{s}$ , and the impulse responses at 0 ( $-0.2 \mu\text{s}$ ), 8 ( $-0.1 \mu\text{s}$ ), 24 (0.1  $\mu\text{s}$ ), and 32 (0.2  $\mu\text{s}$ ) samples are zero. Thus the FIR filter introduces a 0.2  $\mu\text{s}$  delay. There are other types of FIR filters including two-dimensional FIR filters that operate on the I and Q channels together and can provide even better management of the bandwidth of the modulated IF signal.

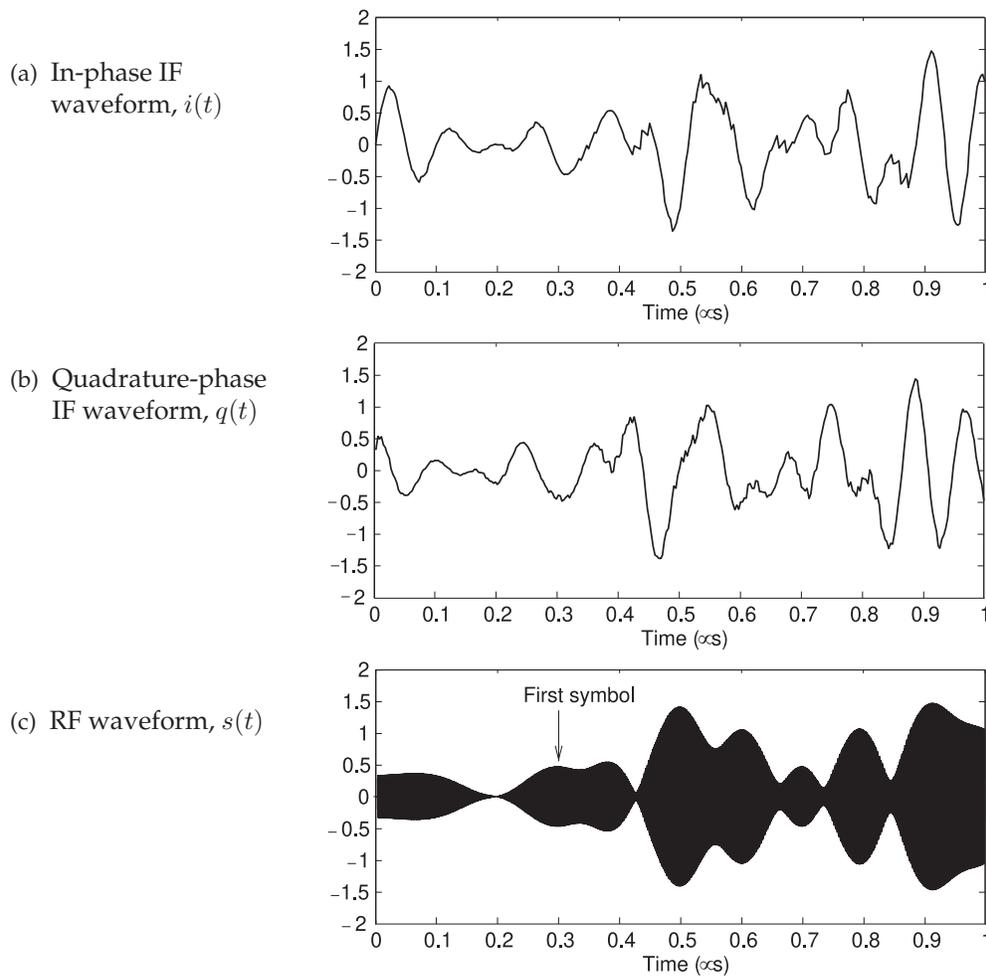
The next stage in the transmitter is multiplying  $i_{BB}$  by the IF LO having



**Figure 3-43:** Sampled baseband signals for 16-QAM modulation corresponding to the analog baseband signals in Figure 3-42.

frequency  $f_{c,IF}$  and multiplying  $q_{BB}$  by a  $90^\circ$  phase-shifted IF LO. This results in the modulated IF signal,  $s_{IF}$ , shown in Figure 3-44(a).

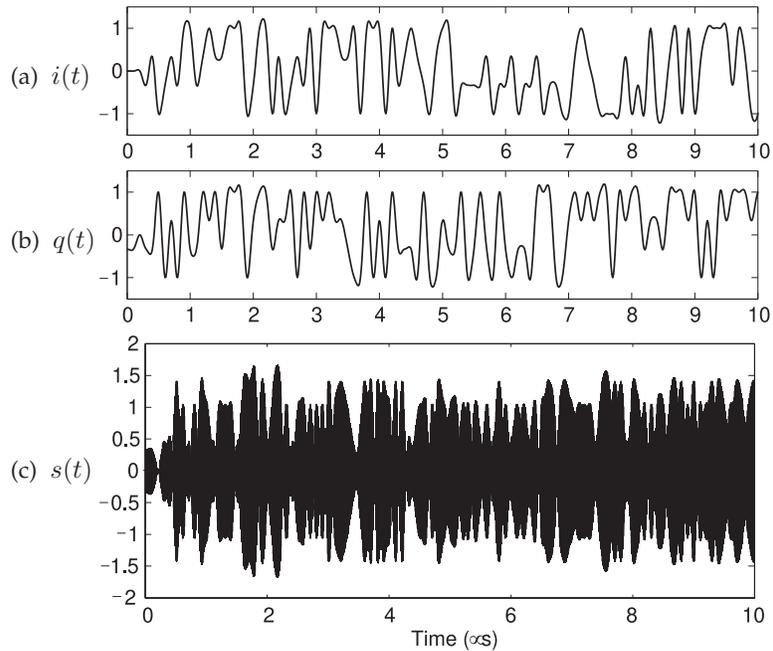
The procedure of producing a SSB-SC modulated RF signal from the IF signal is to drive the multipliers in the RF quadrature modulator of the second mixing stage with an in-phase IF signal and its quadrature, i.e.  $90^\circ$  phase-shifted version. In this case every frequency component of  $s_{IF}$  must be phase shifted by  $90^\circ$ . The mathematical procedure that does this is the Hilbert transform operating on a finite number of samples of  $s_{IF}$  and here 8192 symbols with an oversampling factor of 8 were considered so that the Hilbert transform operates on 65,536 samples. In DSP an FFT can be used so that an FFT of  $s_{IF}$  yields the amplitude and phase of discrete frequency components of  $s_{IF}$  and then the individual frequency samples are phase-shifted by  $90^\circ$ . Then an inverse FFT yields the required Hilbert transform in the time domain. A DAC then converts the digital samples into the analog signals  $i(t)$  and  $q(t)$ . Their waveforms are shown in Figures 3-44(a and b).



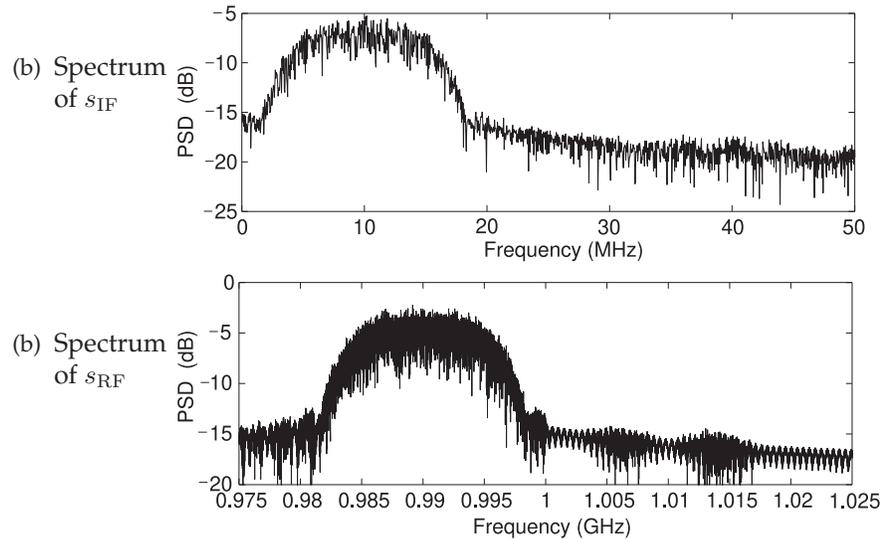
**Figure 3-44:** Waveforms in the RF quadrature modulator where the in-phase IF signal  $i(t)$  is just the output,  $s_{\text{IF}}(t)$  of the IF quadrature multiplier. The first 0.2  $\mu\text{s}$  can be ignored as the FIR filter output is settling. (The IF sampling interval is 3.125 ns. The Hilbert transform operates on  $2^{16} = 65,537$  IF samples or 204.8  $\mu\text{s}$  of data corresponding to 2048 symbols or 8192 bits. The RF time step is 31.25 ps. The RF signal spectrum was calculated using a  $2^{22}$  point FFT.

Then these signals are quadrature modulated using an RF LO of 1 GHz yielding the RF waveform shown in Figure 3-44(c). The SSB-SC RF analog modulator has translated the DSB-SC IF signal to RF to produce a DSB-SC RF signal. Plots of  $i(t)$ ,  $q(t)$  and  $s(t)$  over a longer time are shown in Figure 3-45.

The spectra of the IF and RF signals are shown in Figure 3-46. It is seen that the bandwidth of the IF modulated signal is just under 20 MHz and the bandwidth of the RF signal is the same. The carrier of the IF signal is 10 MHz, the center of the DSB-SC modulated IF. The RF signal is centered around 990 MHz and this is the carrier frequency of the radio signal. That is the 10 MHz carrier frequency,  $f_{c,\text{IF}}$  has been up-converted to the lower



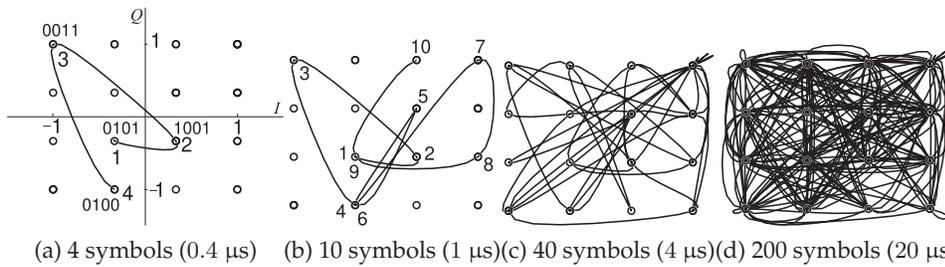
**Figure 3-45:** In-phase and quadrature-phase waveforms and RF waveforms in the RF quadrature modulator.



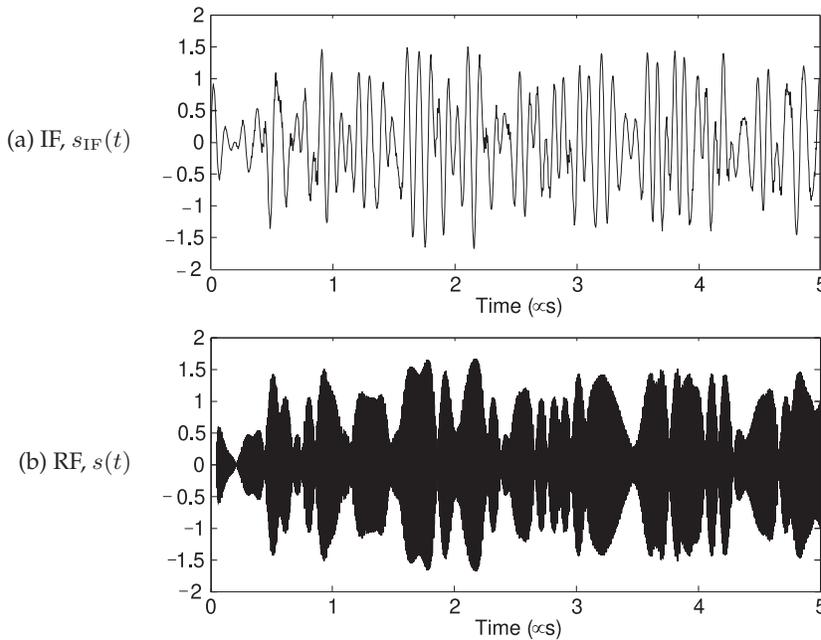
**Figure 3-46:** Spectrum of the IF and RF modulated carriers.

sideband of the 1 GHz LO.

The constellation diagram together with trajectories of the modulated carrier signal are shown in Figure 3-47 for various numbers of symbols. This also corresponds to the phasor diagram of the QAM modulated signal except that the constellation diagram is effectively a phasor diagram that is re-normalized to the average power level so that the constellation diagram does not change with the average RF power changes. A phasor diagram would change of course as the average RF power changed. The trajectories in Figure 3-47 go through a constellation point every 0.1  $\mu$ s. So provided that



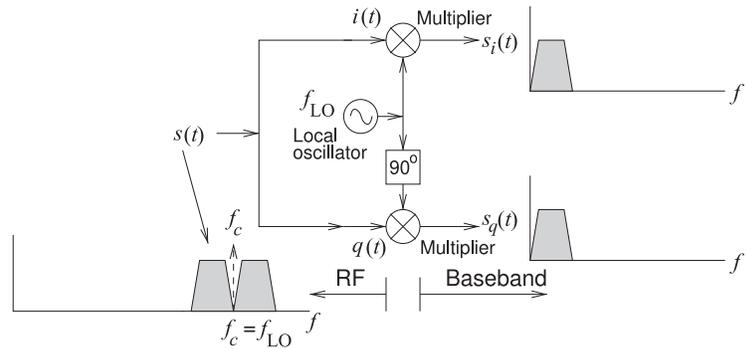
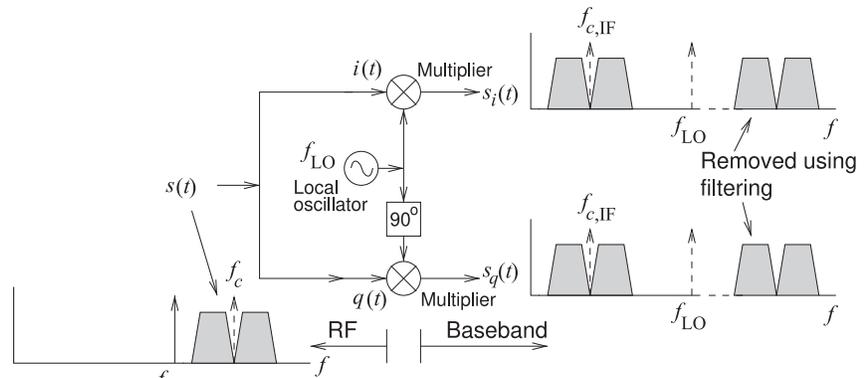
**Figure 3-47:** Constellation diagram with trajectories of the phasor of  $f_c$  for various symbol lengths with the first symbol at  $0.3 \mu s$ . Each symbol has 4 bits. In (a) and (b) the first symbol is labeled '1', etc.



**Figure 3-48:** Comparison of the IF and RF modulated signals.

the clock in a receiver is correctly aligned and samples the RF signal every  $0.1 \mu s$ , the transmitted symbols will be precisely recovered. This of course is in the absence of noise and circuit non-idealities.

The transmitter described here has two modulation stages. The IF quadrature modulator modulates the baseband signal as a DSB-SC IF signal with a bandwidth of around 20 MHz on a 10 MHz carrier. The second stage, the RF quadrature modulator, takes the IF waveform and SSB-SC modulates it on a 990 MHz carrier and maintains the IF bandwidth. The relationship between the IF and RF waveforms is easier to see by viewing the waveforms over a  $5 \mu s$  span corresponding to 50 cycles of the 10 MHz IF carrier and approximately 5000 cycles of the RF LO (and 5051 cycles of the RF carrier), see Figure 3-48.

(a) Quadrature demodulator with  $f_{LO} = f_c$ (b) Quadrature demodulator with  $f_{LO} < f_c$ ,  $f_{c,IF} = f_c - f_{LO}$ .

**Figure 3-49:** Quadrature demodulators used to down-convert an I/Q modulated RF signal.

### Summary

The SDR transmitter involves many technologies. This interplay was best demonstrated using the specific parameters of the SDR transmitter considered in this section. The main limitation is on the bandwidth of the IF modulated signal as this relates directly to the performance of the DSP unit which in turn impacts battery drain. The advantages of implementing a first stage using DSP are tremendous and nearly any modulation scheme can be supported simply by changing the parameters in the DSP unit. The analog RF unit can support very wide bandwidths as fixed frequency filters are not needed. The most critical performance parameter is balancing the I and Q paths through the entire transmitter chain. With most of the modulation occurring in the DSP unit, balancing of the digital portions of the I and Q paths is achieved precisely. The challenge then falls to balancing the I and Q paths in the analog RF modulator.

### 3.11 SDR Quadrature Demodulator

The SDR circuit used to demodulate an I/Q modulated RF signal is very similar to the standard quadrature demodulator. Complete demodulation and separation of the I channel and the Q channel is obtained when the LO of the demodulator coincides with the RF carrier frequency, see Figure 3-49(a). If the LO is not equal to the carrier frequency the RF modulated signal will be down-converted to an intermediate frequency signal which corresponds to a

band around the difference frequency of the carrier and the LO, and a band around the sum frequency of the carrier and the LO, see Figure 3-49(a). Then lowpass filtering would be used to eliminate the higher frequency band. The signal centered at the intermediate frequency can then be I/Q demodulated by another quadrature demodulator.

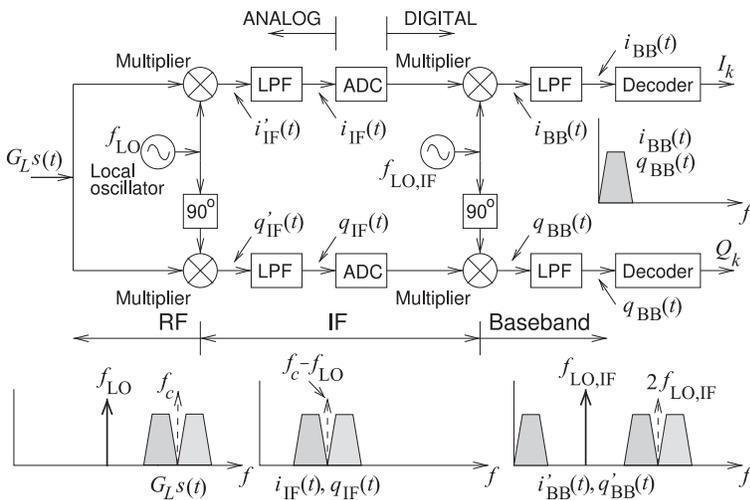
### 3.12 SDR Receiver

An SDR receiver implements the demodulation shown in Figure 3-49(a) in two stages with analog RF quadrature demodulation in the first stage similar to that shown in Figure 3-49(b). There is a lot of detail in this section but this is required to understand the implementation of an SDR receiver at the level of design.

The block diagram of a particular implementation of an SDR receiver is shown in Figure 3-50 which has a first stage in the analog domain implemented in an RF modem chip, and a second digital stage that implements DSP in the baseband chip. While these are separate chips at the time of this writing they will probably be combined in a single chip someday.

The analog portion of the SDR receiver is shown in Figure 3-50. This first stage separates the I and Q channels and outputs baseband frequency signals that are sampled and input to a second stage of quadrature demodulation to produce the final baseband signals. This second stage is implemented digitally. The signals at the output of the first analog stage are also called intermediate frequencies. Digital demodulation, on the right in Figure 3-50, is performed in a DSP unit commonly referred to as a baseband modem chip. A key attribute is that the LO frequency of the first demodulation stage,  $f_{LO}$  can be set at relatively few discrete values while the LO in the second stage,  $f_{LO,IF}$  is set finely for full carrier recovery. Together  $f_{LO} + f_{LO,IF} = f_c$ , the frequency of the carrier of the received signal  $s(t)$ . The DSP unit can implement decoding such as CDMA- or OFDM-decoding.

The rest of this section traces the signal flow through the SDR receiver. The signal presented to the demodulator is  $G_L s(t)$  where  $s(t)$  is the transmitted DSB-SC signal in Equation (3.13).  $G_L$ , which will be very small, is the link gain accounting for loss in transmission and receiver gain and, in the absence



**Figure 3-50:** SDR receiver with RF input signal  $G_L s(t)$  with the first LO frequency,  $f_{LO}$ , less than the carrier signal  $f_c$ . Note that  $2f_{LO,IF} = f_c - f_{LO}$ .

of interference,  $s(t)$  is the transmitted radio signal.

### 3.12.1 Demodulation of the I component

If the transmitted signal is the DSB-SC signal  $s(t)$  with carrier  $f_c$  in Equation (3.13) (replacing  $\omega_{LO}$  by  $f_c$ ) then the signal presented to the demodulator in the receiver is  $G_L s(t)$  where  $G_L$ , which will be very small, is the link gain accounting for loss in transmission and receiver gain. The spectrum of the received DSB-SC signal,  $G_L s(t)$ , is shown on the bottom left of Figure 3-50 with the carrier frequency  $f_c$  identified by the dash arrow. The LO frequency of the first demodulator at  $f_{LO}$  is also shown. The in-phase component of the demodulated signal after the first quadrature demodulator is, see Figure 3-50,

$$\begin{aligned}
 i'_{IF}(t) &= G_L s(t) \sin(\omega_{LO} t) \\
 &= \frac{1}{2} [A_i(f_i) \cos(\omega_c - \omega_i)t + A_q(f_q) \sin(\omega_c - \omega_q)t \\
 &\quad - A_i(f_i) \cos(\omega_c + \omega_i)t - A_q(f_q) \sin(\omega_c + \omega_q)t] \sin(\omega_{LO} t) \\
 &= \frac{1}{2} [A_i(f_i) \sin(\omega_c + \omega_{LO} - \omega_i)t - A_i(f_i) \sin(\omega_c - \omega_{LO} - \omega_i)t \\
 &\quad + A_q(f_q) \cos(\omega_c - \omega_{LO} - \omega_q)t - A_q(f_q) \cos(\omega_c + \omega_{LO} - \omega_q)t \\
 &\quad - A_i(f_i) \sin(\omega_c + \omega_{LO} + \omega_i)t + A_i(f_i) \sin(\omega_c - \omega_{LO} + \omega_i)t \\
 &\quad - A_q(f_q) \cos(\omega_c - \omega_{LO} + \omega_q)t + A_q(f_q) \cos(\omega_c + \omega_{LO} + \omega_q)t]. \quad (3.15)
 \end{aligned}$$

Equation (3.15) includes intermediate frequency terms centered at  $\omega_c - \omega_{LO}$  and high frequency terms centered at  $\omega_c + \omega_{LO}$ . These high frequency terms can be eliminated through lowpass filtering leaving the lowpass filtered received signal

$$\begin{aligned}
 i_{IF}(t) &= \frac{1}{2} [-A_i(f_i) \sin(\omega_c - \omega_{LO} - \omega_i)t + A_q(f_q) \cos(\omega_c - \omega_{LO} - \omega_q)t \\
 &\quad + A_i(f_i) \sin(\omega_c - \omega_{LO} + \omega_i)t - A_q(f_q) \cos(\omega_c - \omega_{LO} + \omega_q)t] \\
 &= \frac{1}{2} A_i(f_i) [-\sin(\omega_c - \omega_{LO} - \omega_i)t + \sin(\omega_c - \omega_{LO} + \omega_i)t] \\
 &\quad + \frac{1}{2} A_q(f_q) [\cos(\omega_c - \omega_{LO} - \omega_q)t - \cos(\omega_c - \omega_{LO} + \omega_q)t] \\
 &= \frac{1}{2} A_i(f_i) [\sin(\omega_i - \omega_c + \omega_{LO})t + \sin(\omega_i + \omega_c - \omega_{LO})t] \\
 &\quad + \frac{1}{2} A_q(f_q) [\cos(\omega_q - \omega_c + \omega_{LO})t - \cos(\omega_q + \omega_c - \omega_{LO})t] \quad (3.16)
 \end{aligned}$$

and this is a DSB-SC signal with a carrier frequency  $f_c - f_{LO}$ . The location of  $i_{IF}$  is shown in Figure 3-50 and its spectrum is the middle spectrum at the bottom of the figure. Of course Equation (3.16) is a discrete signal and the spectrum shows a finite bandwidth signal describing  $i_{IF}(t)$  over the range of  $f_i$  components. (That is, the linear sum of  $i_{IF}(t)$  in Equation (3.16) for all  $f_i$  components.)

### Single Stage Demodulation with $f_{LO} = f_c$

The IF signal,  $i_{IF}(t)$  is then mixed with the IF LO with radian frequency  $\omega_{LO,IF}t = \omega_c t - \omega_{LO}t + \phi_i$  producing the recovered I channel baseband signal

$$\begin{aligned}
 i_{BB}(t) &= \frac{1}{2} A_i(f_i) [\sin(\omega_i t + \phi_i) + \sin(\omega_i t - \phi_i)] \\
 &\quad + \frac{1}{2} A_q(f_q) [\cos(\omega_q t + \phi_i) - \cos(\omega_q t - \phi_i)]. \quad (3.17)
 \end{aligned}$$

Then if  $\phi_i = 0$ , which is when the phase of the signal at  $f_{LO,IF}$  has been recovered correctly, the demodulated baseband signal for the I channel is

$$i_{BB}(t) = A_i(f_i) \sin(\omega_i t). \quad (3.18)$$

### Two stage demodulation $f_{LO} < f_c$

In an SDR demodulator there is two-stage demodulation with an analog stage followed by a digital stage and this is what is shown in Figure 3-50. The analog stage, producing  $i_{IF}(t)$  in Equation (3.16), is in a chip commonly referred to as the RF modem chip. The second stage is usually in a separate chip called the baseband processing chip. One implementation of an SDR receiver uses an analog LO adjusted in discrete frequency steps and, with  $f_{LO} < f_c$  the resulting down-converted signal  $i_{IF}(t)$  is a DSB-SC signal with a carrier at  $f_c - f_{LO}$ , see the middle spectrum at the bottom of Figure 3-50. A second stage of quadrature demodulation operates on  $i_{IF}$ . The  $i_{IF}$  signal is sometimes called the baseband signal if the focus is on the RF modem chip.

The intermediate signal  $i_{IF}(t)$  in Equation (3.17) is mixed with an IF carrier signal at frequency  $f_{LO,IF}$  as follows

$$\begin{aligned} i'_{BB}(t) = & \left\{ \frac{1}{2} A_i(f_i) [\sin(\omega_i - \omega_c + \omega_{LO})t + \sin(\omega_i + \omega_c - \omega_{LO})t] \right. \\ & \left. + \frac{1}{2} A_q(f_q) [\cos(\omega_q - \omega_c + \omega_{LO})t - \cos(\omega_q + \omega_c - \omega_{LO})t] \right\} \\ & \times \cos(\omega_{LO,IF})t. \end{aligned} \quad (3.19)$$

The spectrum of this signal is shown in the bottom right of Figure 3-50 as the spectrum immediately adjacent to DC.

The intermediate LO frequency  $f_{LO,IF}$  can be set with fine adjustment so that  $f_{LO,IF} = f_c - f_{LO}$  which produces a baseband signal  $i'_{BB}(t)$ . Carrier recovery then becomes the process of determining the phase-synchronized  $f_{LO,IF}$  which is just  $f_c$  frequency offset by the numerical value of  $f_{LO}$ . (Note that  $\omega_c - \omega_{LO} - \omega_{LO,IF} = 0$ ,  $-\omega_c + \omega_{LO} + \omega_{LO,IF} = 0$ ,  $\omega_c - \omega_{LO} + \omega_{LO,IF} = 2\omega_{LO,IF}$ , and  $-\omega_c + \omega_{LO} - \omega_{LO,IF} = -2\omega_{c,IF}$ . The I channel baseband signal is then

$$\begin{aligned} i'_{BB}(t) = & i'_{IF}(t) \cos(\omega_{LO,IF}t) \\ = & G_L \left\{ \frac{1}{2} A_i(f_i) [\sin(\omega_i - \omega_c + \omega_{LO})t + \sin(\omega_i + \omega_c - \omega_{LO})t] \right. \\ & \left. + \frac{1}{2} A_q(f_q) [\cos(\omega_q - \omega_c + \omega_{LO})t - \cos(\omega_q + \omega_c - \omega_{LO})t] \right\} \\ & \times \cos(\omega_{LO,IF}t) \\ = & \frac{1}{2} G_L A_i(f_i) [\sin(\omega_i + 0)t + \sin(\omega_i - 2\omega_{LO,IF})t \\ & + \sin(\omega_i + 2\omega_{LO,IF})t + \sin(\omega_i + 0)t] \\ & + \frac{1}{2} G_L A_q(f_q) [+ \cos(\omega_q - 2\omega_{LO,IF})t + \cos(\omega_q + 0)t \\ & - \cos(\omega_q + 0)t - \cos(\omega_q + 2\omega_{LO,IF})t]. \end{aligned} \quad (3.20)$$

After lowpass filtering (in the DSP) to eliminate components centered at  $\pm 2\omega_{LO,IF}$

$$i_{BB}(t) = G_L A_i(f_i) \sin(\omega_i t) \quad (3.21)$$

This is the final desired signal and its spectrum is shown as the inset on the middle right of Figure 3-50 (between the two decoders). This analysis has

been undertaken with a single tone for the I channel but by extension this holds for the actual I channel signal. Sampling  $i_{\text{BB}}(t)$  at the clock ticks yields the sequence of symbols that were transmitted which after decoding yields the bitstream  $I_k$ .

### 3.12.2 Demodulation of the Q component

Demodulation of the Q channel proceeds similarly. The quadrature-phase component is demodulated in a similar way except that the phase of the LO signals differ by  $90^\circ$ . If the transmitted signal is the DSB-SC signal  $s(t)$  with carrier  $f_c$  in Equation (3.13) (and replacing  $\omega_{\text{LO}}$  by  $f_c$ ) and with a receiver LO frequency  $f_{\text{LO}}$ , the in-phase demodulated signal is

$$\begin{aligned}
 q'_{\text{IF}}(t) &= G_L s(t) \sin(\omega_{\text{LO}}t - \pi/2) \\
 &= -s(t) \cos(\omega_{\text{LO}}t) \\
 &= -\frac{1}{2} [A_i(f_i) \cos(\omega_c - \omega_i)t + A_q(f_q) \sin(\omega_c - \omega_q)t \\
 &\quad - A_i(f_i) \cos(\omega_c + \omega_i)t - A_q(f_q) \sin(\omega_c + \omega_q)t] \cos(\omega_{\text{LO}}t) \\
 &= -\frac{1}{2} \{A_i(f_i) \cos(\omega_c + \omega_{\text{LO}} - \omega_i)t + A_i(f_i) \cos(\omega_c - \omega_{\text{LO}} - \omega_i)t \\
 &\quad + A_q(f_q) \sin(\omega_c - \omega_{\text{LO}} - \omega_q)t + A_q(f_q) \sin(\omega_c + \omega_{\text{LO}} - \omega_q)t \\
 &\quad - A_i(f_i) \cos(\omega_c + \omega_{\text{LO}} + \omega_i)t - A_i(f_i) \cos(\omega_c - \omega_{\text{LO}} + \omega_i)t \\
 &\quad - A_q(f_q) \sin(\omega_c - \omega_{\text{LO}} + \omega_q)t + A_q(f_q) \sin(\omega_c + \omega_{\text{LO}} + \omega_q)t\}. \quad (3.22)
 \end{aligned}$$

This includes intermediate frequency terms centered at  $\omega_c - \omega_{\text{LO}}$  and high frequency terms centered at  $\omega_c + \omega_{\text{LO}}$ . These high frequency terms can be eliminated through lowpass filtering leaving the lowpass filtered receiver signal

$$\begin{aligned}
 q_{\text{IF}}(t) &= -\frac{1}{2} [A_i(f_i) \cos(\omega_c - \omega_{\text{LO}} - \omega_i)t + A_q(f_q) \sin(\omega_c - \omega_{\text{LO}} - \omega_q)t \\
 &\quad - A_i(f_i) \cos(\omega_c - \omega_{\text{LO}} + \omega_i)t - A_q(f_q) \sin(\omega_c - \omega_{\text{LO}} + \omega_q)t] \\
 &= \frac{1}{2} A_i(f_i) [-\cos(\omega_c - \omega_{\text{LO}} - \omega_i)t + \cos(\omega_c - \omega_{\text{LO}} + \omega_i)t] \\
 &\quad + \frac{1}{2} A_q(f_q) [-\sin(\omega_c - \omega_{\text{LO}} - \omega_q)t + \sin(\omega_c - \omega_{\text{LO}} + \omega_q)t] \\
 &= \frac{1}{2} A_i(f_i) [-\cos(\omega_i - \omega_c + \omega_{\text{LO}})t + \cos(\omega_i + \omega_c - \omega_{\text{LO}})t] \\
 &\quad + \frac{1}{2} A_q(f_q) [\sin(\omega_q - \omega_c + \omega_{\text{LO}})t + \sin(\omega_q + \omega_c - \omega_{\text{LO}})t]. \quad (3.23)
 \end{aligned}$$

#### Single-Stage Demodulation With $f_{\text{LO}} = f_c$

If the LO frequency equals the carrier frequency but there is an offset of  $\phi_i$  so that  $\omega_{\text{LO}}t = \omega_c t + \phi_i$ , then the recovered baseband signal is

$$\begin{aligned}
 q_{\text{BB}}(t) &= \frac{1}{2} A_i(f_i) [-\cos(\omega_i t + \phi_i) + \cos(\omega_i t - \phi_i)] \\
 &\quad + \frac{1}{2} A_q(f_q) \{\sin(\omega_q t + \phi_i) + \sin(\omega_q t - \phi_i)\}. \quad (3.24)
 \end{aligned}$$

Then if  $\phi_i = 0$  the demodulated baseband signal for the I channel is

$$q'_{\text{BB}}(t) = A_q(f_q) \sin(\omega_q t) \quad (3.25)$$

### Two stage demodulation $f_{LO} < f_c$

When  $q_{IF}(t)$  in Equation (3.17) is mixed with an IF LO signal at frequency  $f_{LO,IF}$  the resulting signal is

$$\begin{aligned} q'_{BB}(t) &= \left\{ \frac{1}{2} A_i(f_i) [\sin(\omega_i - \omega_c + \omega_{LO})t + \sin(\omega_i + \omega_c - \omega_{LO})t] \right. \\ &\quad \left. + \frac{1}{2} A_q(f_q) [\cos(\omega_q - \omega_c + \omega_{LO})t - \cos(\omega_q + \omega_c - \omega_{LO})t] \right\} \\ &\quad \times \cos(\omega_{LO,IF}t). \end{aligned} \quad (3.26)$$

The intermediate LO frequency  $f_{LO,IF}$  can be set with fine adjustment so that  $f_{LO,IF} = f_c - f_{LO}$  which produces a baseband signal  $q_{BB}(t)$ . (Note that  $\omega_c - \omega_{LO} - \omega_{LO,IF} = 0$ ,  $-\omega_c + \omega_{LO} + \omega_{LO,IF} = 0$ ,  $\omega_c - \omega_{LO} + \omega_{LO,IF} = 2\omega_{LO,IF}$ , and  $-\omega_c + \omega_{LO} - \omega_{LO,IF} = -2\omega_{LO,IF}$ . The Q channel baseband signal is then

$$\begin{aligned} q'_{BB}(t) &= q_{IF}(t) \cos(\omega_{LO,IF}t) \\ &= G_L \left\{ \frac{1}{2} A_i(f_i) [-\cos(\omega_i - \omega_c + \omega_{LO})t + \cos(\omega_i + \omega_c - \omega_{LO})t] \right. \\ &\quad \left. + \frac{1}{2} A_q(f_q) [\sin(\omega_q - \omega_c + \omega_{LO})t + \sin(\omega_q + \omega_c - \omega_{LO})t] \right\} \\ &\quad \times \cos(\omega_{LO,IF}t) \\ &= \frac{1}{2} G_L A_i(f_i) [-\cos(\omega_i + 0)t + \cos(\omega_i - 2\omega_{LO,IF})t \\ &\quad + \cos(\omega_i + 2\omega_{LO,IF})t + \cos(\omega_i + 0)t] \\ &\quad + \frac{1}{2} G_L A_q(f_q) [\sin(\omega_q - 2\omega_{LO,IF})t + \sin(\omega_q + 0)t \\ &\quad - \sin(\omega_q + 2\omega_{LO,IF})t - \cos(\omega_q + 0)t]. \end{aligned} \quad (3.27)$$

After lowpass filtering (in the DSP) to eliminate components centered at  $\pm 2\omega_{LO,IF}$

$$q_{BB}(t) = G_L A_q(f_q) \cos(\omega_q t). \quad (3.28)$$

This can be compared to the original quadrature signal which led to the DSB-SC modulated signal  $s(t)$ . The  $f_q$  frequency component of  $q(t)$  was  $A_q f_q \sin(\omega_q t)$ . Thus it is necessary to change the phase of  $q_{BB}(t)$  to obtain the original  $q(t)$  signal:

$$\begin{aligned} q_{BB}(t) \text{ (phase shifted by } 90^\circ) &= G_L A_q(f_q) \cos(\omega_q t - \pi/2) q(t) \\ &= G_L A_q(f_q) \sin(\omega_q t) \end{aligned} \quad (3.29)$$

and this can be implemented using a Hilbert transform just as was done in the SDR transmitter.

### 3.13 SDR Summary

The software defined radio described in this chapter used two stages of quadrature modulation in the SDR transmitter and two stages of quadrature demodulation in the SDR receiver. Each of the quadrature modulators and demodulators used ideal multipliers to multiply signals by an LO signal. There are other ways of implementing quadrature de/modulators but the principles of one stage of the quadrature modulator being implemented in a baseband DSP unit and the final analog RF modulator being analog; and one of the quadrature demodulator being analog and operating on the RF signal and the second stage quadrature demodulator being implemented in DSP operating on an intermediate frequency signal is essential to implementations of SDR.

### 3.14 Summary

This chapter considered the architectures of transmitters and receivers and in particular the implementation of modulators and demodulators. Since the beginning of radio there has been a tremendous increase in the sophistication of modulators and demodulators. The earliest AM demodulators required just a single diode and now the demodulators in software defined radio use hundreds of millions of transistors. Following early radios the discussion tracked the development of modulators and demodulators for cellular systems but the increasingly advanced techniques are used in other radio and radar systems. The 1G systems used FM modulation to transmit voice and simple FSK modulation to transmit limited digital data such as phone such as channel information. The 2G and 3G systems were confined to using just one or two types of modulation such four-state modulation. The level of complexity increased tremendously with 4G where now it is essential to use software defined radio to support a very large number of different modulation formats including mostly various orders of QAM but also supporting legacy modulation formats. The 5G modulation and demodulation schemes are the same as those used in 4G with the 5G advances being in the move to millimeter-waves and beam steering. There is another layer of modulation used in 3G (i.e. wideband CDMA) which can layer multiple users in the same bandwidth. Strictly this is an access scheme and will be considered in Chapter 5. The 4G systems uses orthogonal frequency division modulation (OFDM) modulation. That is not considered here as it is a hybrid of an access scheme and a modulation scheme. It too will be considered in Chapter 5.

### 3.15 References

- [1] S. Darlington, "On digital single-sideband modulators," *IEEE Trans. on Circuit Theory*, vol. 17, no. 3, pp. 409–414, 1970.
- [2] R. V. Hartley, "Modulating system." 1918, uS Patent 1,287,982.
- [3] F. Colebrook, "Homodyne," *Wireless World and Radio. Rev.*, 1924.
- [4] D. Tucker, "The synchrodyne," *Electronic Engineering*, vol. 19, pp. 75–76, Mar. 1947.
- [5] E. Armstrong, "Some recent developments in the audion receiver," *Proc. of the Institute of Radio Engineers*, vol. 3, no. 3, pp. 215–238, Sep. 1915.
- [6] L. Lessing, *Man of High Fidelity: Edwin Howard Armstrong, a Biography*. J.B. Lippincott, 1956.
- [7] B. A. Weaver, "A new, high efficiency, digital, modulation technique for am or ssb sound broadcasting applications," *IEEE Trans. on Broadcasting*, vol. 38, no. 1, pp. 38–42, 1992.

### 3.16 Exercises

1. A superheterodyne receiver has, in order, an antenna, a low-noise amplifier, a bandpass filter, a mixer, a second bandpass filter, a second mixer, a lowpass filter, an ADC, and a DSP that will implement quadrature demodulation. Develop the frequency plan of the receiver if the RF input is at 2 GHz and has a 200 kHz single-channel bandwidth. The final signal applied to the ADC must be between DC and 400 kHz so that  $I/Q$  demodulation can be done in the DSP unit. Noise considerations mandate that the LO of the first mixer must be more than 10 MHz away from the input RF. Also, for a bandpass filter to have minimum size, the center frequency of the filter should be as high as possible. It has been determined that the appropriate trade-off of physical size and cost is to have a 100 MHz bandpass filter between the two mixers. (Note: 100 MHz is the center frequency of the bandpass filter.)
  - (a) Draw a block diagram of the receiver and annotate it with symbols for the frequencies

- of the LOs and the RF and IF signals.
- (b) What is the LO frequency  $f_{LO1}$  of the first mixer?
  - (c) What is the LO frequency  $f_{LO2}$  of the second mixer?
  - (d) Specify the cutoff frequency of the lowpass filter following the second mixer.
  - (e) Discuss in less than  $\frac{1}{2}$  a page other design considerations relating to frequency plan, filter size, and filter specification.
2. Short answer questions. Each part requires a short paragraph of about five lines and a figure where appropriate to illustrate your understanding.
    - (a) Explain the operation of a superheterodyne receiver.
    - (b) Compare zero-IF and low-IF receivers.



# Antennas and the RF Link

4.1	Introduction .....	129
4.2	RF Antennas .....	130
4.3	Resonant Antennas .....	131
4.4	Traveling-Wave Antennas .....	136
4.5	Antenna Parameters .....	137
4.6	The RF link .....	143
4.7	Multipaths and Delay Spread .....	156
4.8	Radio Link Interference .....	160
4.9	Antenna array .....	163
4.10	Summary .....	165
4.11	References .....	166
4.12	Exercises .....	167

## 4.1 Introduction

An antenna interfaces circuits with free-space with a transmit antenna converting a guided wave signal on a transmission line to an electromagnetic (EM) wave propagating in free space, while a receive antenna is a transducer that converts a free-space EM wave to a guided wave on a transmission line and eventually to a receiver circuit. Together the transmit and receive antennas are part of the RF link. The RF link is the path between the output of the transmitter circuit and the input of the receiver circuit (see Figure 4-1). Usually this path includes the cable from the transmitter to the transmit antenna, the transmit antenna itself, the propagation path, the receive antenna, and the transmission line connecting the receive antenna to the receiver circuit. The received signal is much smaller than the transmitted signal. The overwhelming majority of the loss is from the propagation path as the EM signal spreads out, and usually diffracts, reflects, and is partially blocked by objects such as hills and buildings.

The first half of this chapter is concerned with the properties of antennas. One of the characteristics of antennas is that the energy can be focused in a particular direction, a phenomenon captured by the concept of antenna gain, which can partially compensate for path loss. The second half of this chapter considers modeling the RF link and the geographical arrangement of antennas that manage interference from other radios while providing support for as many users as possible.

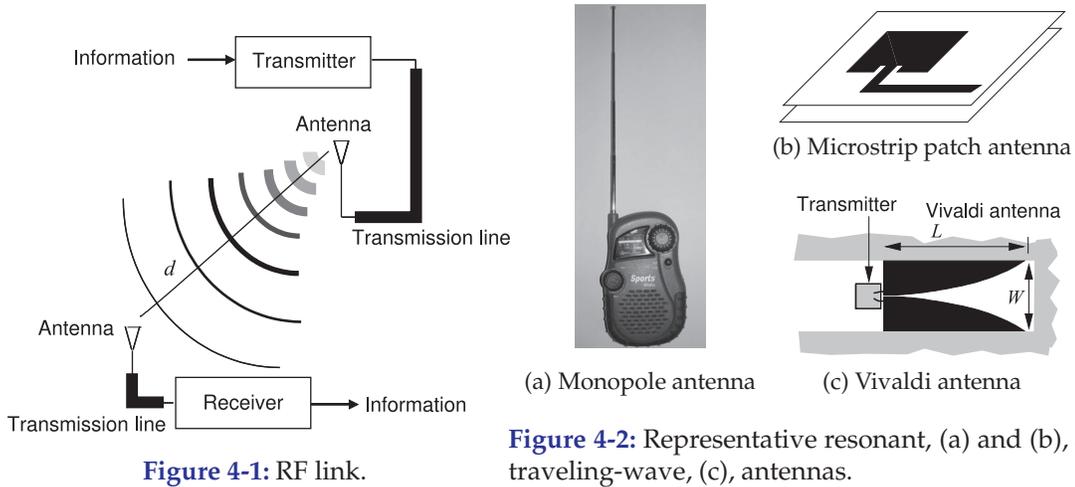


Figure 4-1: RF link.

Figure 4-2: Representative resonant, (a) and (b), and traveling-wave, (c), antennas.

## EXAMPLE 4.1

## Interference

In the figure there are two transmitters,  $T_{X1}$  and  $T_{X2}$ , operating at the same power level, and one receiver, Rx.  $T_{X1}$  is an intentional transmitter and its signal is intended to be received at Rx.  $T_{X1}$  is separated from Rx by  $D_1 = 2$  km.  $T_{X2}$  uses the same frequency channel as  $T_{X1}$ , and as far as Rx is concerned it transmits an interfering signal. Assume that the antennas are omnidirectional (i.e., they transmit and receive signals equally in all directions) and that the transmitted power density drops off as  $1/d^2$ , where  $d$  is the distance from the transmitter. Calculate the signal-to-interference ratio (SIR) at Rx.

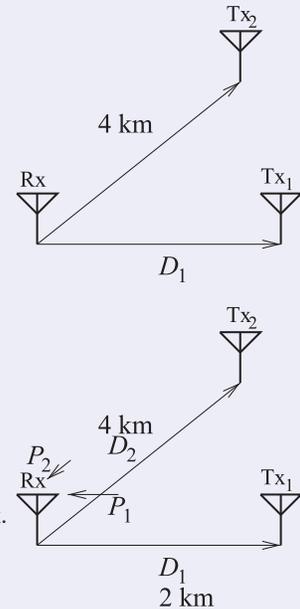
**Solution:**

$D_1 = 2$  km and  $D_2 = 4$  km.

$P_1$  is the signal power transmitted by  $T_{X1}$  and received at Rx.

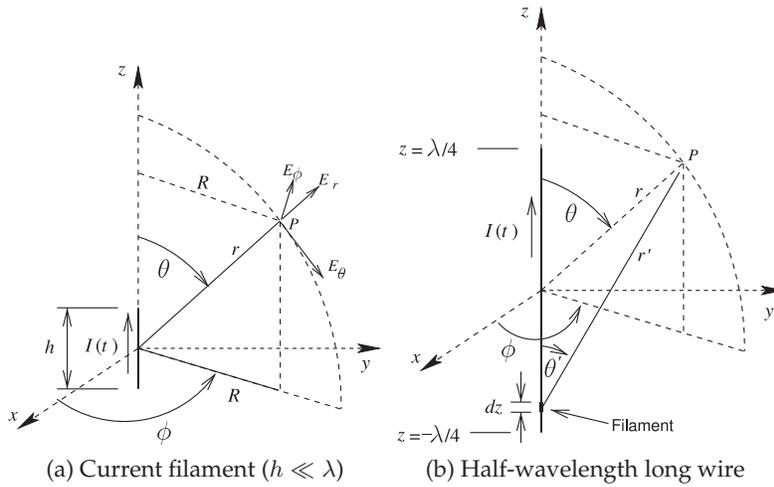
$P_2$  is the interference power transmitted by  $T_{X2}$  and received at Rx.

$$\text{So SIR} = \frac{P_1}{P_2} = \left(\frac{D_2}{D_1}\right)^2 = 4 = 6.02 \text{ dB.}$$



## 4.2 RF Antennas

Antennas are of two fundamental types: **resonant** and **traveling-wave antennas**, see Figure 4-2. Resonant antennas establish a standing wave of current with required resonance usually established when the antenna section is either a quarter- or half-wavelength long. These antennas are also known as standing-wave antennas. Resonant antennas are inherently narrowband because of the resonance required to establish a large standing wave of current. Figures 4-2(a and b) show two representative resonant antennas. The physics of the operation of a resonant antenna is understood by considering the time domain. First consider the physical operation of a transmit antenna. When a sinusoidal voltage is applied to the conductor



**Figure 4-3:** Wire antennas. The distance from the center of an antenna to the field point  $P$  is  $r$ .  $R = r \sin \theta$ .

of an antenna, charges, i.e. free electrons, accelerate or decelerate under the influence of an applied voltage source which typically arrives at the antenna from a traveling wave voltage on a transmission line. When the charges accelerate (or decelerate) they produce an EM field which radiates away from the antenna. At a point on the antenna there is a current sinusoidally varying in time, and the net acceleration of the charges (in charge per second per second) is also sinusoidal with an amplitude that is directly proportional to the amplitude of the current sinewave. Hence having a large standing wave of current, when the antenna resonates sinusoidally, results in a large charge acceleration and hence large radiated fields.

When an EM field impinges upon a conductor the field causes charges to accelerate and hence induce a voltage that propagates on a transmission line connected to the receive antenna. Traveling-wave antennas, an example of one is shown in Figures 4-2(c), operate as extended lines that gradually flare out so that a traveling wave on the original transmission line transitions into free space. Traveling-wave antennas tend to be two or more wavelengths long at the lowest frequency of operation. While relatively long, they are broadband, many 3 or more octaves wide (e.g., extending from 500 MHz to 4 GHz or more). These antennas are sometimes referred to as **aperture antennas**.

### 4.3 Resonant Antennas

With a resonant antenna the current on the antenna is directly related to the amplitude of the radiated EM field. Resonance ensures that the standing wave current on the antenna is high.

#### 4.3.1 Radiation from a Current Filament

The fields radiated by a resonant antenna are most conveniently calculated by considering the distribution of current on the antenna. The analysis begins by considering a short filament of current, see Figure 4-3(a). Considering the sinusoidal steady state at radian frequency  $\omega$ , the current on the filament with phase  $\chi$  is  $I(t) = |I_0| \cos(\omega t + \chi)$ , so that  $I_0 = |I_0| e^{-j\chi}$  is the phasor of the current on the filament. The length of the filament is  $h$ , but it has no other

dimensions, that is, it is considered to be infinitely thin.

Resonant antennas are conveniently modeled as being made up of an array of current filaments with spacings and lengths being a tiny fraction of a wavelength. Wire antennas are even simpler and can be considered to be a line of current filaments. Ramo, Whinnery, and Van Duzer [1] calculated the spherical EM fields at the point  $P$  with the spherical coordinates  $(\phi, \theta, r)$  generated by the  $z$ -directed current filament centered at the origin in Figure 4-3. The total EM field components in phasor form are

$$H_\phi = \frac{I_0 h}{4\pi} e^{-jk r} \left( \frac{jk}{r} + \frac{1}{r^2} \right) \sin \theta, \quad \bar{H}_\phi = H_\phi \hat{\phi} \quad (4.1)$$

$$E_r = \frac{I_0 h}{4\pi} e^{-jk r} \left( \frac{2\eta}{r^2} + \frac{2}{j\omega\epsilon_0 r^3} \right) \cos \theta, \quad \bar{E}_r = E_r \hat{r} \quad (4.2)$$

$$E_\theta = \frac{I_0 h}{4\pi} e^{-jk r} \left( \frac{j\omega\mu_0}{r} + \frac{1}{j\omega\epsilon_0 r^3} + \frac{\eta}{r^2} \right) \sin \theta, \quad \bar{E}_\theta = E_\theta \hat{\theta}, \quad (4.3)$$

where  $\eta$  is the free-space characteristic impedance,  $\epsilon_0$  is the permittivity of free space, and  $\mu_0$  is the permeability of free space. The variable  $k$  is called the **wavenumber** and  $k = 2\pi/\lambda = \omega\sqrt{\mu_0\epsilon_0}$ . The  $e^{-jk r}$  terms describe the variation of the phase of the field as the field propagates away from the filament. Equations (4.1)–(4.3) are the complete fields with the  $1/r^2$  and  $1/r^3$  dependence describing the near-field components. In the far field, i.e. large  $r \gg \lambda$ , the components with  $1/r^2$  and  $1/r^3$  dependence become negligible and the field components left are the propagating components  $H_\phi$  and  $E_\theta$ :

$$H_\phi = \frac{I_0 h}{4\pi} e^{-jk r} \left( \frac{jk}{r} \right) \sin \theta, \quad E_r = 0, \quad \text{and} \quad E_\theta = \frac{I_0 h}{4\pi} e^{-jk r} \left( \frac{j\omega\mu_0}{r} \right) \sin \theta. \quad (4.4)$$

Now consider the fields in the plane normal to the filament, that is, with  $\theta = \pi/2$  radians so that  $\sin \theta = 1$ . The fields are now

$$H_\phi = \frac{I_0 h}{4\pi} e^{-jk r} \left( \frac{jk}{r} \right) \quad \text{and} \quad E_\theta = \frac{I_0 h}{4\pi} e^{-jk r} \left( \frac{j\omega\mu_0}{r} \right) \quad (4.5)$$

and the wave impedance is

$$\eta = \frac{E_\theta}{H_\phi} = \frac{I_0 h}{4\pi} e^{-jk r} \frac{j\omega\mu_0}{r} \left( \frac{I_0 h}{4\pi} e^{-jk r} \frac{jk}{r} \right)^{-1} = \frac{\omega\mu_0}{k}. \quad (4.6)$$

Note that the strength of the fields is directly proportional to the magnitude of the current. This proves to be very useful in understanding spurious radiation from microwave structures. Now  $k = \omega\sqrt{\mu_0\epsilon_0}$ , so

$$\eta = \frac{\omega\mu_0}{\omega\sqrt{\mu_0\epsilon_0}} = \sqrt{\frac{\mu_0}{\epsilon_0}} = 377 \, \Omega, \quad (4.7)$$

as expected. Thus an antenna can be viewed as having the inherent function of an impedance transformer converting from the lower characteristic impedance of a transmission line (often 50  $\Omega$ ) to the 377  $\Omega$  characteristic impedance of free space. Further comments can be made about the propagating fields (Equation (4.4)). The EM field propagates in all directions

except not directly in line with the filament. For fixed  $r$ , the amplitude of the propagating field increases sinusoidally with respect to  $\theta$  until it is maximum in the direction normal to the filament.

The power radiated is obtained using the **Poynting vector**, which is the cross-product of the propagating electric and magnetic fields. From this the time-average propagating power density is (with the SI units of  $\text{W}/\text{m}^2$ )

$$P_R = \frac{1}{2} \Re(E_\theta H_\phi^*) = \frac{\eta k^2 |I_0|^2 h^2}{32\pi^2 r^2} \sin^2 \theta, \quad (4.8)$$

and the power density is proportional to  $1/r^2$ . In Equation (4.8)  $\Re(\dots)$  indicates that the real part is taken.

### 4.3.2 Finite-Length Wire Antennas

The EM wave launched by a wire antenna of finite length is obtained by considering the wire as being made up of many filaments and the field is then the superposition of the fields from each filament. As an example, consider the antenna in Figure 4-3(b) where the wire is half a wavelength long. As a good approximation the current on the wire is a standing wave and the current on the wire is in phase so that the current phasor is

$$I(z) = I_0 \cos(kz). \quad (4.9)$$

From Equation (4.4) and referring to Figure 4-3 the fields in the far field are

$$H_\phi = \int_{-\lambda/4}^{\lambda/4} \frac{I_0 \cos(kz)}{4\pi} e^{-jkr'} \left( \frac{jk}{r'} \right) \sin \theta' dz \quad (4.10)$$

$$E_\theta = \int_{-\lambda/4}^{\lambda/4} \frac{I_0 \cos(kz)}{4\pi} e^{-jkr'} \left( \frac{j\omega\mu_0}{r'} \right) \sin \theta' dz, \quad (4.11)$$

where  $\theta'$  is the angle from the filament to the point  $P$ . Now  $k = 2\pi/\lambda$  and at the ends of the wire  $z = \pm\lambda/4$  where  $\cos(kz) = \cos(\pm\pi/2) = 0$ . Evaluating the equations is analytically involved and will not be done here. The net result is that the fields are further concentrated in the plane normal to the wire. At large  $r$ , of at least several wavelengths distant from the antenna, only the field components decreasing as  $1/r$  are significant. At large  $r$  the phase differences of the contributions from the filaments is significant and results in shaping of the fields. The geometry to be used in calculating the far field is shown in Figure 4-4(a). The phase contribution of each filament, relative to that at  $z = 0$ , is  $(kz \sin \theta)/\lambda$  and Equations (4.10) and (4.11) become

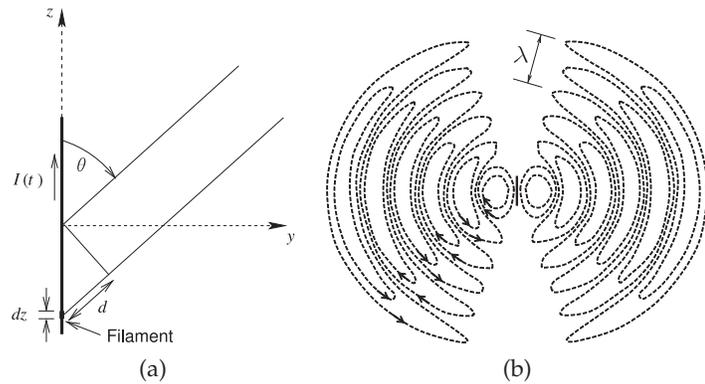
$$H_\phi = I_0 \left( \frac{jk}{4\pi r} \right) \sin(\theta) e^{-jkr} \int_{-\lambda/4}^{\lambda/4} \frac{\sin(kz)}{4\pi} \sin(z \sin(\theta)) dz \quad (4.12)$$

$$E_\theta = I_0 \left( \frac{j\omega\mu_0}{4\pi r} \right) \sin(\theta) e^{-jkr} \int_{-\lambda/4}^{\lambda/4} \sin(kz) \sin(z \sin(\theta)) dz. \quad (4.13)$$

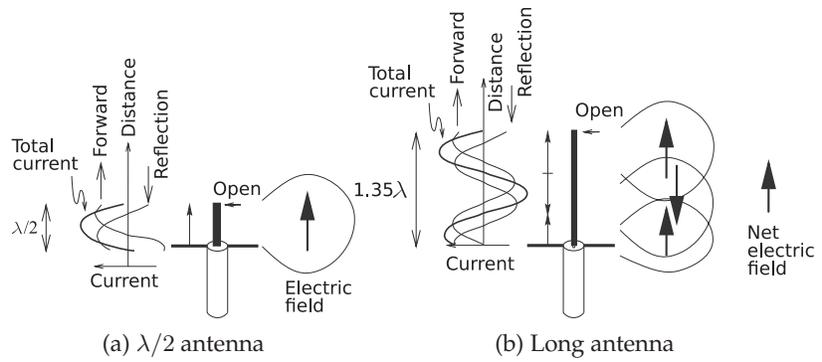
Figure 4-4(b) is a plot of the near-field electric field in the  $y$ - $z$  plane calculating  $E_r$  and  $E_\theta$  (recall that  $E_\phi = 0$ ) every  $90^\circ$ . Further from the antenna the  $E_r$  component rapidly reduces in size, and  $E_\theta$  dominates.

A summary of the implications of the above equations are, first, that the strength of the radiated electric and magnetic fields are proportional to the

**Figure 4-4:** Wire antenna: (a) geometry for calculating contributions from current filaments of length  $dz$  with coordinate  $z$  ( $d = -z \sin \theta$ ); and (b) instantaneous electric field in the  $y$ - $z$  plane due to a  $\lambda/2$  long current element. There is also a magnetic field.



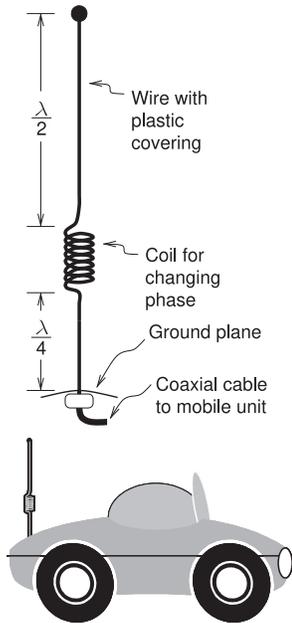
**Figure 4-5:** Monopole antenna showing total current and forward- and backward-traveling currents: (a) a  $\frac{1}{2}\lambda$ -long antenna; and (b) a relatively long antenna.



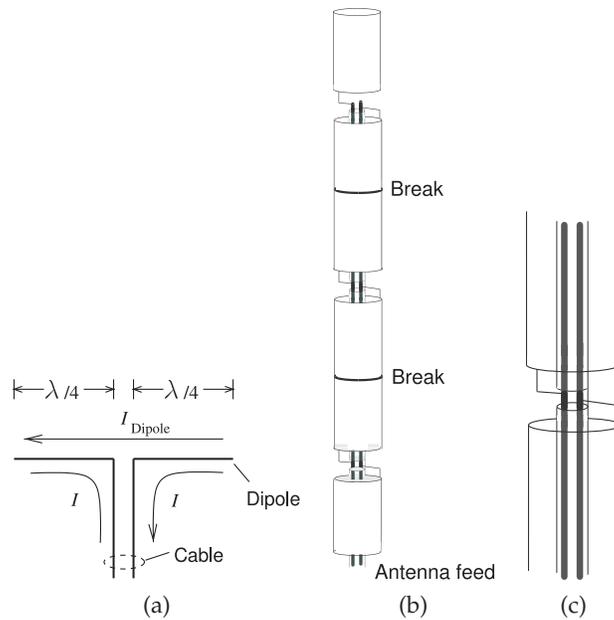
current on the wire antenna. So establishing a standing current wave and hence magnifying the current is important to the efficiency of a wire antenna. A second result is that the power density of freely propagating EM fields in the far field is proportional to  $1/r^2$ , where  $r$  is the distance from the antenna. A third interpretation is that the longer the antenna, the flatter the radiated transmission profile; that is, the radiated energy is more tightly confined to the  $x$ - $y$  (i.e.  $\Theta = 0$ ) plane. For the wire antenna the peak radiated field is in the plane normal to the antenna, and thus the wire antenna is generally oriented vertically so that transmission is in the plane of the earth and power is not radiated unnecessarily into the ground or into the sky.

To obtain an efficient resonant antenna, all of the current should be pointed in the same direction at a particular time. One way of achieving this is to establish a standing wave, as shown in Figure 4-5(a). At the open-circuited end, the current reflects so that the total current at the end of the wire is zero. The initial and reflected current waves combine to create a standing wave. Provided that the antenna is sufficiently short, all of the total current—the standing wave—is pointed in the same direction. The optimum length is about a half wavelength. If the wire is longer, the contributions to the field from the oppositely directed current segments cancel (see Figure 4-5(b)).

In Figure 4-5(a) a coaxial cable is attached to the monopole antenna below the ground plane and often a series capacitor between the cable and the antenna provides a low level of coupling leading to a larger standing wave. The capacitor also approximately matches the characteristic impedance of the cable to the input impedance  $Z_{in}$  of the antenna. If the length of the monopole is reduced to one-quarter wavelength long it is again resonant, and the input impedance,  $Z_{in}$ , is found to be  $36 \Omega$ . Then a  $50 \Omega$  cable can



**Figure 4-6:** Mobile antenna with phasing coil extending the effective length of the antenna.

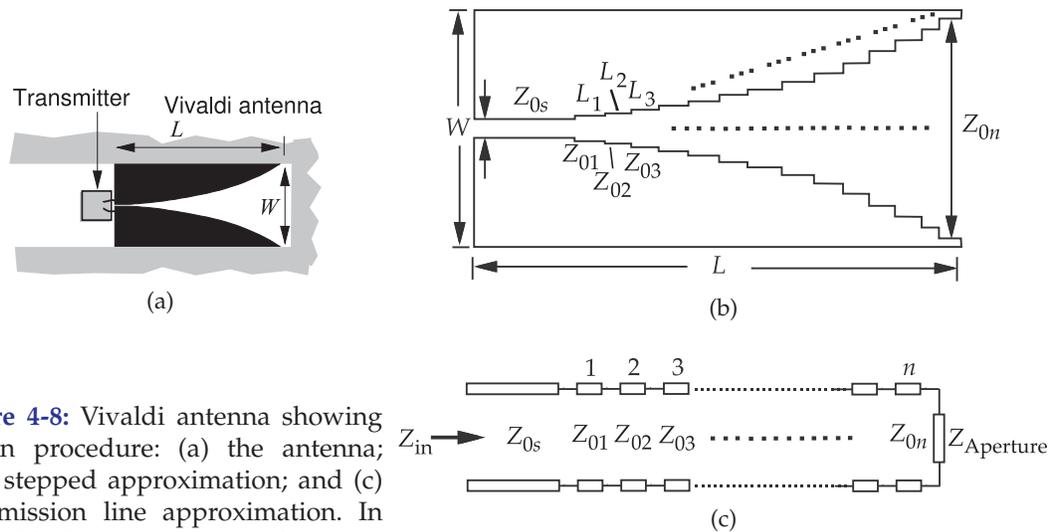


**Figure 4-7:** Dipole antenna: (a) current distribution; (b) stacked dipole antenna; and (c) detail of the connection in a stacked dipole antenna.

be directly connected to the antenna and there is only a small mismatch and nearly all the power is transferred to the antenna and then radiated.

Another variation on the monopole is shown in Figure 4-6, where the key component is the phasing coil. The phasing coil (with a wire length of  $\lambda/2$ ) rotates the electrical angle of the current phasor on the line so that the current on the  $\lambda/4$  segment is in the same direction as on the  $\lambda/2$  segment. The result is that the two straight segments of the loaded monopole radiate a more tightly confined EM field. The phasing coil itself does not radiate (much).

Another ingenious solution to obtaining a longer effective wire antenna with same-directed current (and hence a more tightly confined RF beam) is the stacked dipole antenna (Figure 4-7). The basis of the antenna is a dipole as shown in Figure 4-7(a). The cable has two conductors that have equal amplitude currents,  $I$ , but flowing as shown. The wire section is coupled to the cable so that the currents on the two conductors realize a single effective current  $I_{\text{dipole}}$  on the dipole antenna. The stacked dipole shown in Figure 4-7(b) takes this geometric arrangement further. Now the radiating element is hollow and a coaxial cable is passed through the antenna elements and the half-wavelength sections are fed separately to effectively create a wire antenna that is several wavelengths long with the current pointing in one direction. Most cellular antennas using wire antennas are stacked dipole antennas.



**Figure 4-8:** Vivaldi antenna showing design procedure: (a) the antenna; (b) a stepped approximation; and (c) transmission line approximation. In (a) the black regions are metal.

### Summary

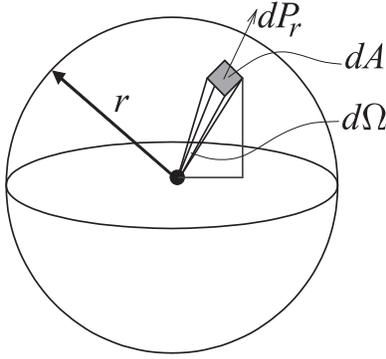
Standing waves of current can be realized by resonant structures other than wires. A microstrip patch antenna, see Figure 4-2(b), is an example, but the underlying principle is that an array of current filaments generates EM components that combine to create a propagating field. Resonant antennas are inherently narrowband because of the reliance on the establishment of a standing wave. A relative bandwidth of 5%–10% is typical.

## 4.4 Traveling-Wave Antennas

Traveling-wave antennas have the characteristics of broad bandwidth and large size. These antennas begin as a transmission line structure that flares out slowly, providing a low reflection transition from a transmission line to free space. The bandwidth can be very large and is primarily dependent on how gradual the transition is.

One of the more interesting traveling-wave antennas is the **Vivaldi antenna** of Figure 4-8(a). The Vivaldi antenna is an extension of a slotline in which the fields are confined in the space between two metal sheets in the same plane. The slotline spacing increases gradually in an exponential manner, much like that of a Vivaldi violin (from which it gets its name), over a distance of a wavelength or more. A circuit model is shown in Figures 4-8(b and c) where the antenna is modeled as a cascade of many transmission lines of slowly increasing characteristic impedance,  $Z_0$ . Since the  $Z_0$  progression is gradual there are low-level reflections at the transmission line interfaces. The forward-traveling wave on the antenna continues to propagate with a negligible reflected field. Eventually the slot opens sufficiently that the effective impedance of the slot is that of free space and the traveling wave continues to propagate in air.

The other traveling-wave antennas work similarly and all are at least a wavelength long, with the central concept being a gradual taper from the



**Figure 4-9:** Free-space spreading loss. The incremental power,  $dP_r$ , intercepted by the shaded region of incremental area  $dA$  is proportional to  $1/r^2$ . The solid angle subtended by the shaded area is the incremental solid angle  $d\Omega$ . The integral of  $dA$  over the surface of the sphere, i.e. the area of the sphere is  $4\pi r^2$ . The total solid angle subtended by the sphere is the integral of  $d\Omega$  over the sphere and is  $4\pi$  steradians (or  $4\pi$  sr).

characteristic impedance of the originating transmission line to free space. The final aperture is at least one-half wavelength across so that the fields can curl on themselves (i.e., loop back on themselves) and are self-supporting as they leave the antenna.

## 4.5 Antenna Parameters

This section introduces a number of antenna metrics that are used to characterize antenna performance.

### 4.5.1 Radiation Density and Radiation Intensity

Antennas do not radiate equally in all directions concentrating radiated power in (usually) one direction called the **main** (or **major**) **lobe** of the antenna. This focusing effect is called **directivity**. The power in a particular direction is characterized by the radiation density and the radiation intensity metrics. The radiation density,  $S_r$ , is the power per unit area with the SI units of  $\text{W}/\text{m}^2$ , and will be maximum in the main lobe. Referring to Figure 4-9 with an antenna located at the center of the sphere of radius  $r$  and radiating a total power  $P_r$ ,  $S_r$  is the incremental radiated power  $dP_r$  passing through the incremental shaded region of area,  $dA$ :

$$S_r = \frac{dP_r}{dA}. \quad (4.14)$$

$S_r$  reduces with distance falling of as  $1/r^2$  in free space. For a practical antenna  $S_r$  will vary across the surface of the sphere. The total power radiated is the closed integral over the surface  $S$  of the sphere:

$$P_r = \oint_S dP_r = \oint_S S_r dA. \quad (4.15)$$

An alternative measure of power concentration is the radiation intensity  $U$  which is in terms of the incremental solid angle  $d\Omega$  subtended by  $dA$  so that  $d\Omega = dA/r^2$  and (with the SI units of  $\text{W}/\text{steradian}$  or  $\text{W}/\text{sr}$ )

$$U = \frac{dP_r}{d\Omega} = \frac{dP_r}{dA} r^2 = r^2 S_r. \quad (4.16)$$

### Isotropic Antenna

It is useful to reference the directivity of an antenna with respect to a fictitious isotropic antenna that has no loss and radiates equally in all directions so

that  $S_r$  is only a function of  $r$ . Then integrating over the surface of the sphere yields the total radiated power

$$P_r|_{\text{Isotropic}} = \oint_S dP_r = \oint_S S_r dA = S_r \oint_S dA = S_r 4\pi r^2 = 4\pi U. \quad (4.17)$$

Since the isotropic antenna has no loss the power input to the antenna  $P_{\text{IN}}$  is equal to the power radiated  $P_r = P_{\text{IN}}$ . Thus for an isotropic antenna

$$S_r = \frac{P_r}{4\pi r^2} = \frac{P_{\text{IN}}}{4\pi r^2} \quad (4.18) \quad \text{and} \quad U|_{\text{Isotropic}} = r^2 S_r = \frac{P_{\text{IN}}}{4\pi r^2}. \quad (4.19)$$

### Antenna Efficiency

Antenna efficiency, sometimes called the **radiation efficiency**, describes losses in an antenna principally due to resistive ( $I^2R$ ) losses. Resonant antennas work by creating a large current that is maximized through the generation of a standing wave at resonance. There is a lot of current, and even just a little resistance results in substantial resistive loss. The power that is reflected from the input of the antenna is usually small. The total radiated power (in all directions),  $P_r$ , is the power input to the antenna less losses. The antenna efficiency,  $\eta_A$  is therefore defined as

$$\eta_A = P_r/P_{\text{IN}}, \quad (4.20)$$

where  $P_{\text{IN}}$  is the power input to the antenna and  $\eta_A < 1$  and usually expressed as a percentage. Antenna efficiency is very close to one for many antennas, but can be 50% for microstrip patch antennas.

**Antenna loss** refers to the same mechanism that gives rise to antenna efficiency. Thus an antenna with an antenna efficiency of 50% has an antenna loss of 3 dB. Generally losses are resistive due to  $I^2R$  loss and mismatch loss of the antenna that occurs when the input impedance is not matched to the impedance of the cable connected to the antenna. Because of confusion with antenna gain (they are not the opposite of each other) the use of the term 'antenna loss' is discouraged and instead 'antenna efficiency' preferred.

### 4.5.2 Directivity and Antenna Gain

The directivity of an antenna,  $D$ , is the ratio of the radiated power density to that of an isotropic antenna with the same total radiated power,  $P_r$ :

$$D = \frac{S_r}{S_r|_{\text{Isotropic}}} = \frac{U}{U|_{\text{Isotropic}}} \quad (4.21)$$

where  $S_r$  and  $U$  refer to the actual antenna and the power densities and intensities are measured at the same distance from the antennas. For an actual antenna  $D$  is dependent on the direction from the antenna, see Figure 4-10. The maximum value of  $D$  will be in the direction of the main lobe of the antenna and this is called **directivity gain**.

The focusing property of an antenna is characterized by comparing the radiated power density to that of an isotropic antenna with the same input power. The antenna gain,  $G_A$ , is the maximum value of  $D$  when the power input  $P_{\text{IN}} = P_r/\eta_A$  to the antenna and the isotropic antenna are the same:

$$G_A = \eta_A \max(D). \quad (4.22)$$

Antenna	Type	Figure	Gain (dBi)	Notes
Lossless isotropic antenna			0	
$\lambda/2$ dipole	Resonant	4-7(a)	2	$R_{in} = 73 \Omega$
$3\lambda$ diameter parabolic dish	Traveling	—	38	$R_{in} = \text{match}$
Patch	Resonant	4-2(b)	9	$R_{in} = \text{match}$
Vivaldi	Traveling	4-2(c)	10	$R_{in} = \text{match}$
$\lambda/4$ monopole on ground	Resonant	4-5(a)	2	$R_{in} = 36 \Omega$
$5/8\lambda$ monopole on ground	Resonant	4-5(a)	3	Matching required

**Table 4-1:** Several antenna systems.  $R_{in} = \text{match}$  for resonant antennas indicates that the antenna can be designed to have an input impedance matching that of a feed cable. Traveling-wave antennas are intrinsically matched.

Losses in the antenna are accounted for by the efficiency term  $\eta_A$ .

In Equation (4.22)  $G_A$  is a gain factor and is often expressed in terms of decibels (taking 10 times the log of  $G_A$ ) but dBi (with ‘i’ standing for with-respect-to isotropic) is used to indicate that it is not a power gain in the same sense as amplifier gain.  $G_A$  instead is the ratio of power densities for two different antennas. For example, an antenna that focuses power in one direction increasing the peak radiated power density by a factor of 20 relative to that of an isotropic antenna has an antenna gain of 13 dBi. With care  $G_A$  can be often used in calculations of power as as with amplifier gain.

Since it is almost impossible to calculate internal antenna losses, antenna gain is invariably only measured. The input power to an antenna can be measured and the peak radiated power density,  $P_D|_{\text{Maximum}}$ , measured in the far field at several wavelengths distant (at  $r \gg \lambda$ ). This is compared to the power density from an ideal isotropic antenna at the same distance with the same input power. Antenna gain is determined from

$$G_A = \frac{\text{Maximum radiated power per unit area}}{\text{Maximum radiated power per unit area for an isotropic antenna}} \tag{4.23}$$

$$= \frac{S_r|_{\text{Maximum}}}{S_r|_{\text{Isotropic}}} = 4\pi r^2 \frac{P_D|_{\text{Maximum}}}{P_{IN}}$$

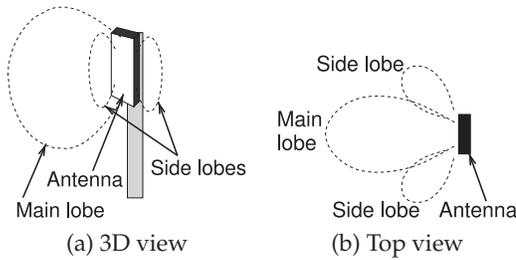
$$= 4\pi \frac{\text{Maximum radiated power per unit solid angle}}{\text{Total input power to the antenna}} \tag{4.24}$$

$$= 4\pi \frac{(dP_r/d\Omega)|_{\text{Maximum}}}{P_{IN}} = 4\pi r^2 \frac{(dP_r/dA)|_{\text{Maximum}}}{P_{IN}}$$

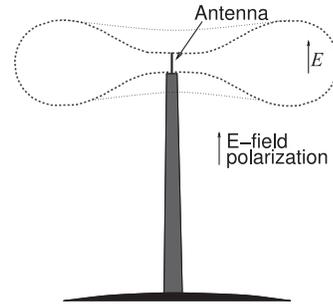
The antenna gains of common resonant and traveling-wave antennas are given in Table 4-1. In free space the antenna gain determined using Equation (4.22) is independent of distance. Antenna gain is measured on an antenna range using a calibrated receive antenna and care taken to avoid reflections from objects, especially from the ground.

The losses of an antenna are incorporated in the antenna gain which is defined in terms of the power input to the antenna, see Equation (4.24). Thus in calculations of radiated power using antenna gain, there is no need to separately account for resistive losses in the antenna.

In summary, antennas concentrate the radiated power in one direction so that the density of the power radiated in the direction of the peak field is higher than the power density from an isotropic antenna. Power radiated from a base station antenna, such as that shown in Figure 4-11, is concentrated in a region that looks like a toroid or, more closely, a balloon squashed at its north and south poles. Then the antenna does not radiate much power into space and will concentrate power in a region skimming the surface of



**Figure 4-10:** Field pattern produced by a microstrip antenna.



**Figure 4-11:** A base station transmitter pattern.

the earth. Antenna gain is a measure of the effectiveness of an antenna to concentrate power in one direction. Thus, in free space due to power spreading out the maximum power density (in SI units of  $W/m^2$ ) at a distance  $r$  is

$$P_D = \frac{G_A P_{IN}}{4\pi r^2}, \quad (4.25)$$

where  $4\pi r^2$  is the area of a sphere of radius  $r$  and  $P_{IN}$  is the input power.

Measurements of antenna gain are used to derive antenna efficiency. It is impossible to measure or simulate the resistive and dielectric losses of an antenna directly. Antenna efficiency is obtained using theoretical calculations of antenna gain assuming no losses in the antenna itself. This is compared to the measured antenna gain yielding the antenna efficiency.

#### EXAMPLE 4.2

#### Antenna Gain

A base station antenna has an antenna gain,  $G_A$ , of 11 dBi and a 40 W input. The transmitted power density falls off with distance  $d$  as  $1/d^2$ . What is the peak power density at 5 km?

#### Solution:

A sphere of radius 5 km has an area  $A = 4\pi r^2 = 3.142 \cdot 10^8 \text{ m}^2$ ;  $G_A = 11 \text{ dBi} = 12.6$ . In the direction of peak radiated power, the power density at 5 km is

$$P_D = \frac{P_{in} G_A}{A} = \frac{40 \cdot 12.6 \text{ W}}{3.142 \cdot 10^8 \text{ m}^2} = 1.603 \mu\text{W}/\text{m}^2.$$

#### EXAMPLE 4.3

#### Antenna Efficiency

A antenna has an antenna gain of 13 dBi and an antenna efficiency of 50% and all of the loss is due to resistive losses and resistance of metals is proportional to temperature. The RF signal input to the antenna has a power of 40 W.

(a) What is the input power in dBm?

$$P_{in} = 40 \text{ W} = 46.02 \text{ dBm}.$$

(b) What is the total power transmitted in dBm?

$$P_{\text{Radiated}} = 50\% \text{ of } P_{IN} = 20 \text{ W or } 43.01 \text{ dBm}.$$

$$\text{Alternatively, } P_{\text{Radiated}} = 46.02 \text{ dBm} - 3 \text{ dB} = 43.02 \text{ dBm}.$$

(c) If the antenna is cooled to near absolute zero so that it is lossless, what would the antenna gain be?

The antenna gain would increase by 3 dB and antenna gain incorporates both directivity and antenna losses. So the gain of the cooled antenna is 16 dBi.

### 4.5.3 Effective Isotropic Radiated Power

A transmit antenna does not radiate power equally in all directions and for a receiver in the main lobe of the transmit antenna it is as though there is an isotropic transmit antenna with a much higher input power. This concept is incorporated in the effective isotropic radiated power (EIRP):

$$\text{EIRP} = P_{\text{IN}}G_A. \quad (4.26)$$

This uses antenna gain to derive the total power that would be radiated by an isotropic antenna producing the same (peak) power density as the actual antenna. Regulatory limits on the power levels of transmitters are in terms of EIRP rather than the total radiated power. Sometimes **equivalent radiated power (ERP)** is used instead of EIRP with the same meaning.

### 4.5.4 Effective Aperture Size

Effective aperture size is defined so that the power density at a receive antenna when multiplied by its effective aperture size,  $A_R$ , yields the power output from the antenna at its connector. Since antennas are linear and reciprocal it is to be expected that there is a relationship between effective aperture size and antenna gain.

An antenna has an effective size that is more than its actual physical size. This is because of its influence on the EM fields around it. The power captured by an antenna is the effective aperture size (or area) multiplied by the transmitted power density. That is, the effective aperture size of an antenna is the area of the surface that captures all of the power passing through it and delivers this power to the output terminals of the antenna.

The effective aperture area of a receive antenna,  $A_R$ , is related to the receive antenna gain,  $G_R$ , as follows [2, 3] ((note that  $A_e$  is often used if it is not necessary to distinguish antennas):

$$A_R = \frac{G_R\lambda^2}{4\pi}, \quad (4.27)$$

where  $\lambda$  is the wavelength of the radio signal<sup>1</sup>. The effective aperture area of an antenna can have little to do with its physical size; e.g., a wire antenna has almost no physical size but has a significant effective aperture size.

If  $S_r$  is the transmitted power density at the receive antenna, the power received is

$$P_R = P_D A_R = P_D \frac{G_R\lambda^2}{4\pi}. \quad (4.28)$$

---

<sup>1</sup> The derivation of Equation (4.27) invokes a thought experiment with an antenna connected to a resistor in thermal equilibrium and each in separate thermally isolated chambers with black body walls [4, 5]. The available thermal (Johnson-Nyquist) noise power from the resistor is radiated by the antenna (incorporating antenna gain) and absorbed by the walls of the antenna's chamber. The same amount of power, as black body radiation from the walls of the antenna's chamber, must be received by the antenna (incorporating effective aperture size) and delivered to the resistor where it is dissipated as heat. At frequencies up to infrared and at room temperature (and hence at radio frequencies) the power density of black body radiation increases linearly with temperature and as the square of frequency (i.e. as  $1/\lambda^2$ ) whereas Johnson-Nyquist noise is independent of frequency but increases linearly with temperature. The derivation is exact for an isotropic antenna and assumed to apply to all antennas.

Here  $P_R$  is the power delivered at the output connector of the receive antenna, as loss in the receive antenna is incorporated in  $G_R$  and  $A_R$ . The total power radiated by the transmit antenna in the direction of maximum power density is given by multiplying the power input to the transmit antenna,  $P_T$ ,<sup>2</sup> by the antenna gain of the transmit antenna,  $G_T$ . This can be converted to the power density at a distance  $d$  (ignoring multipath effects),

$$S_r = \frac{P_T G_T}{4\pi d^2}, \quad (4.29)$$

and the power delivered by the receive antenna is

$$P_R = S_r A_R = \frac{P_T G_T}{4\pi d^2} \frac{G_R \lambda^2}{4\pi} = P_T G_T G_R \left( \frac{\lambda}{4\pi d} \right)^2. \quad (4.30)$$

That is,

$$P_R = P_T G_T G_R \left( \frac{\lambda}{4\pi d} \right)^2. \quad (4.31)$$

This equation is known as the **Friis transmission equation** or the **Friis transmission formula**.

Equation (4.31) provides a powerful tool for evaluating the power received in a communication system and modifications have been developed to account for impairments that are encountered. One of the assumptions in the development of Equation (4.31) is that the polarization of the receive antenna matches that of the transmitted signal. The transmit antenna will radiate a signal with an electric field with a particular polarization, i.e.  $E$  field orientation. Propagation through air has little effect on the polarization of a signal, although atmospheric conditions can rotate the polarization slightly. A polarization mismatch factor could be included in Equation (4.31). For ground-based communications, such as cellular communications, reflections and diffraction have a major impact on the power of the received signal and techniques for handling these effects are considered in the next section.

#### 4.5.5 Summary

This section introduced several metrics for characterizing antennas:

Metric	Equation	Description
$S_r$	(4.14)	Radiated power density, W/m <sup>2</sup>
$U$	(4.16)	Radiation intensity W/sr
$\eta_A$	(4.27)	Antenna efficiency
$D$	(4.21)	Antenna directivity
$G_A$	(4.23)	Antenna gain, used with a transmit antenna
$A_e$	(4.27)	Effective aperture area, used with a receive antenna
EIRP	(4.26)	Equivalent isotropic radiated power

<sup>2</sup> Note that this was  $P_{IN}$  previously and here the subscript  $T$  is used to indicate the power input to the transmit antenna.

**EXAMPLE 4.4** Point-to-Point Communication

In a point-to-point communication system, a parabolic receive antenna has an antenna gain of 60 dBi. If the signal is 60 GHz and the power density at the receive antenna is 1 pW/cm<sup>2</sup>, what is the power at the output of the receive antenna connected to the RF electronics?

**Solution:**

The first step is determining the effective aperture area,  $A_R$ , of the parabolic antenna. The frequency is 60 GHz, and so the wavelength ( $\lambda$ ) is 5 mm. Also note that  $G_R = 60$  dBi =  $10^6$ . From Equation (4.27),

$$A_R = \frac{G_R \lambda^2}{4\pi} = \frac{10^6 \cdot 0.005^2}{4\pi} = 1.989 \text{ m}^2. \quad (4.32)$$

Using Equation (4.28)  $P_D = 1$  pW/cm<sup>2</sup> = 10 nW/m<sup>2</sup>, the total power delivered to the RF receiver electronics (at the output of the receive antenna) is

$$P_R = P_D A_R = 10 \text{ nW} \cdot \text{m}^{-2} \cdot 1.989 \text{ m}^2 = 19.89 \text{ nW}. \quad (4.33)$$

## 4.6 The RF Link

The RF link is between a transmit antenna and a receive antenna. Sometimes the RF link includes the antenna and sometimes it does not, this will be clear from the context, but usually it includes the antennas. The principle source of link loss is the spreading out of the EM field as it propagates. In the absence of any other effects (such as atmospheric loss and reflections), the power density reduces as  $1/d^2$ , where  $d$  is distance, and this is called the **line-of-sight (LOS)** situation. In this section the propagation path is first described along with its impairments including propagation on multiple paths between a transmit antenna and a receive antenna. The resonant scattering is described as well as multiple fading effects due to reflections, diffraction, rain, and other atmospheric effects.

### 4.6.1 Propagation Path

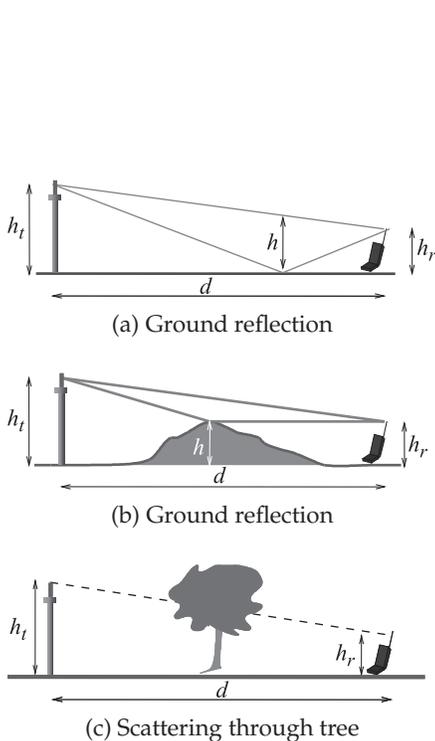
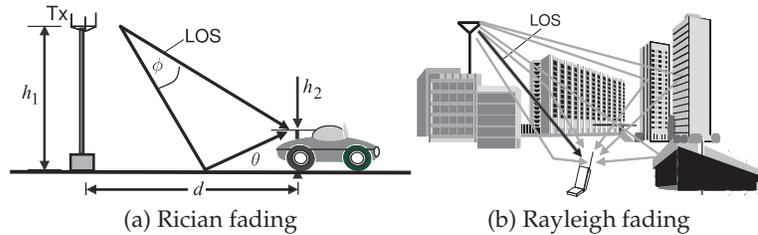
When the radiated signal reflects and diffracts there are multiple propagation paths that result in what is called fading, as the paths constructively and destructively combine at the receiver. The direct path and other paths, called multipaths, which are reflected or diffracted by ground, buildings, and other objects do not arrive at the receiver in phase leading to time-varying constructive and destructive combining, called multipath interference. Of these destructive combining is much worse as it can reduce a signal level below what it would be if propagation was in free space. In urban areas, 10 or 20 paths can have significant power in them and these combine at the receive antenna [6].

When the signal on one of the paths dominates, this is usually the LOS path, fading is called **Rician fading**. With LOS and a single ground reflection, the situation is the classic Rician fading as shown in Figure 4-12(a). This is a special situation, as the ground changes the phase of the signal upon reflection, generally by 180°. When the receiver is a long way from the base station, the lengths of the two paths are almost identical and the level of the signals in the two paths are almost the same. The net result is that these two signals almost cancel, and so instead of the power falling off by  $1/d^2$ , it

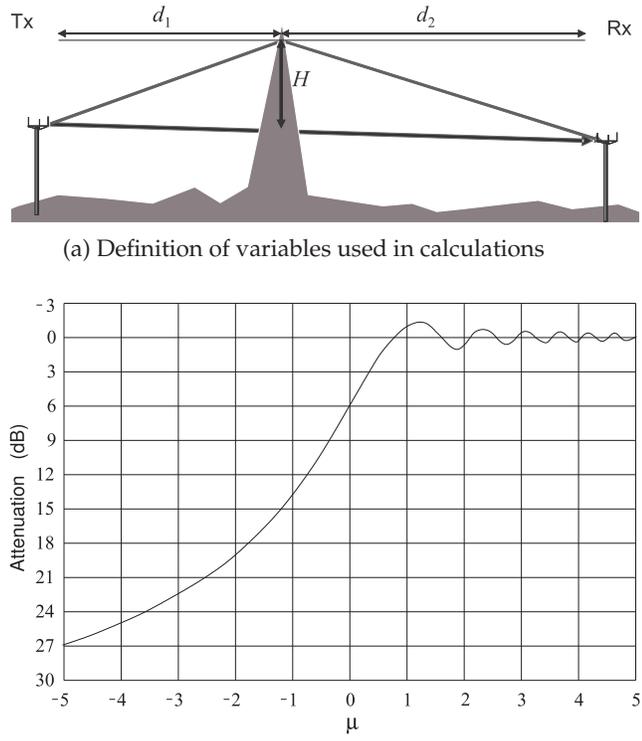
falls off by  $1/d^3$ . When there are many paths and all have similar amplitude signals, fading is called **Rayleigh fading**. In an urban area such as that shown in Figure 4-12(b), there are many significant multipaths and the power falls off by  $1/d^4$  and sometimes faster.

Common paths encountered in cellular radio are shown in Figure 4-13. As a rough guide, in the single-digit gigahertz range each diffraction and scattering event reduces the signal received by 20 dB. The knife-edge diffraction scenario is shown in more detail in Figure 4-14. This case is fairly easy to analyze and can be used to estimate the effects of individual obstructions. The diffraction model is derived from the theory of half-infinite screen diffraction [7]. First, calculate the parameter  $\nu$  from the geometry of

**Figure 4-12:** Multipath propagation: (a) line of sight (LOS) and ground reflection paths only; and (b) in an urban environment.

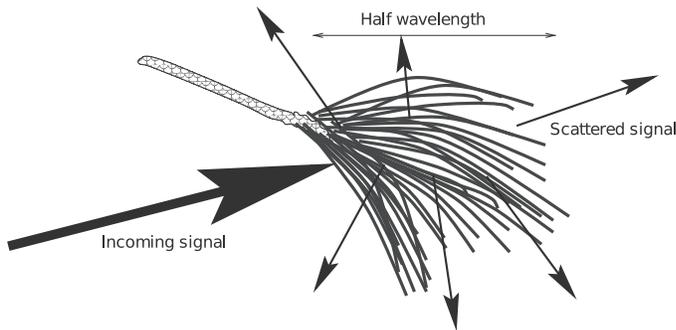


**Figure 4-13:** Common paths contributing to multipath propagation.



(b) Chart for determining attenuation using Equation (4.34)

**Figure 4-14:** Knife-edge diffraction.



**Figure 4-15:** Pine needles scattering an incoming EM signal.

the path using

$$\nu = -H \sqrt{\frac{2}{\lambda} \left( \frac{1}{d_1} + \frac{1}{d_2} \right)}. \quad (4.34)$$

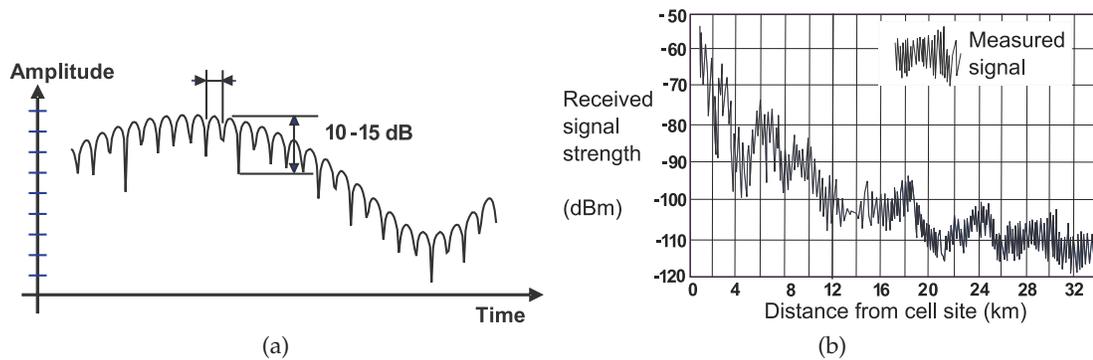
Next, consult the plot in Figure 4-14(b) to obtain the diffraction loss (or attenuation). This loss should be added (using decibels) to the otherwise determined path loss to obtain the total path loss. Other losses such as reflection cancellation still apply, but are computed independently for the path sections before and after the obstruction.

#### 4.6.2 Resonant Scattering

Propagation is rarely from point to point as the path is often obstructed and so non-LOS (NLOS). One type of event that reduces transmission is scattering. The level of the effect depends on the size of the objects causing scattering. Here the effect of the pine needles of Figure 4-15 will be considered. The pine needles (as most objects in the environment) conduct electricity, especially when wet. When an EM field is incident an individual needle acts as a wire antenna, with the current maximum when the pine needle is one-half wavelength long. At this length, the “needle” antenna supports a standing wave and will re-radiate the signal in all directions. This is scattering, and there is a considerable loss in the direction of propagation of the original fields. As well, there is loss due to a needle not being a very good conductor, so EM energy is lost as heat. The effect of scattering is frequency and size dependent. A typical pine needle is 15 cm long, which is exactly  $\lambda/2$  at 1 GHz, and so a stand of pine trees have an extraordinary impact on cellular communications at 1 GHz. As a rough guide, 20 dB of a signal is lost when passing through a small stand of pine trees. At 2 GHz, a similar impact comes from other leaves, and they do not need to look like wire antennas. Consider an oak leaf having a dimension of 7.5 cm. This is  $\lambda/2$  at 2 GHz, another dominant cellular frequency. So scattering results, but now the loss is season dependent, with the loss due to scattering being considerably smaller in winter (when trees such as oaks do not have leaves) than in summer.

#### 4.6.3 Fading

Fading refers to the variation of the received signal with time or when the position of transmit or receive antennas is changed. The variation in signal strength is much greater than would be expected from changes in path



**Figure 4-16:** Fast and slow fading: (a) in time as a radio and obstructions move; and (b) in distance.

length. Several forms of fading are identified and associated with different physical effects. The most important fading types are flat fading, multipath fading, and rain fading.

### Flat Fading

Temperature variations of the atmosphere between the transmit antenna and receive antenna give rise to what is called flat fading and sometimes called **thermal fading**. This fading is called flat because it is independent of frequency. One form of flat fading is due to refraction, which occurs when different layers of the atmosphere have different densities and thus dielectric permittivities increasing or decreasing away from the surface of the earth. The temperature profile can increase away from the earth surface or reduce depending on whether the temperature of the earth is higher than that of the air and is commonly associated with the beginning and end of the day. **Temperature inversions** can also occur where the temperature profile in air does not uniformly increase or decrease, thus causing a layer with a higher permittivity than that of the air above or below. RF energy gets trapped in this layer, reflecting from the top and bottom of the inversion layer. This is called **ducting**. In point-to-point communication systems the transmit and receive antennas are mounted high on towers and then reflection from ground objects is often small. In such cases flat fading is the most commonly observed phenomena and minor fluctuations of several decibels in receive signal level are common throughout the day. However, when temperature variations are extreme, ducting can severely impact communications reducing signal levels by up to 20 dB.

### Shadow Fading

Shadow fading occurs when the LOS path is blocked by an obstruction such as a building or hill. The full signal strength returns when the obstruction or receiver has moved to restore the direct path. Shadow fading is also called slow fading and the characteristic of the channel is regarded as being relatively constant over a short time. The amplitude response varies in time and distance and is shown in Figure 4-16. This figure shows both fast fades that are 10 to 15 dB deep and slow or shadow fades that are 20 to 30 dB deep.

### Multipath Fading

Multipath fading is the most common fading when either the transmit antenna or the receive antenna are close to the ground, near obstructing buildings or terrains, or inside a building. In such situations there are many reflections that combine destructively and constructively. Multipath fading is also called **fast fading**, as the characteristics of the channel can change significantly in a few milliseconds. Two types of multipath fading—Rician fading and Rayleigh fading—will be considered.

Multipath is principally a problem when the line of sight between transmit and receive antennas is obscured. However, where there is line of sight the signal reflecting from the ground immediately in front of a receive antenna can sometimes largely cancel the line-of-sight signal. In time, intervening objects can move and the propagation characteristics of the various paths can change due to thermal variations. All this adds to the randomness of fast fades. Multipath fading of 20 dB can occur for a small percentage of the time on time scales of many seconds when there are few propagation paths (e.g. in a rural area) to a large percentage of the time many times per second in a dense urban environment when there are many paths. Constructive combining does increase the signal level momentarily, but there is no advantage to this. Destructive combining can result in deep fades of 20 dB impacting communications and forcing the communication system to accommodate either by using higher average powers or using strategies such as multiple antennas or spreading the communication signal over a wide bandwidth since fades tend to be 500 kHz to 1 MHz wide at all frequencies.

### Rician Fading

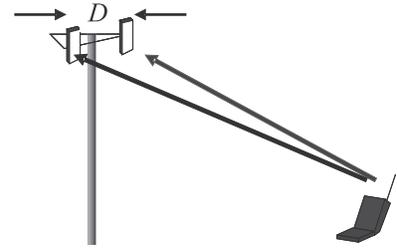
Rician fading occurs when there is one dominant RF path, usually the LOS path, and one or more other paths. The main situation is when there is an LOS path and a ground reflection as shown in Figure 4-12(a). The signal received is the sum of two signals:

$$v_r(t) = C [\cos(\omega_c t) + r \cos(\omega_c t + \phi)], \quad (4.35)$$

where  $C \cos(\omega_c t)$  is the LOS signal,  $r$  is the ratio of the amplitude of the signal reflected from the ground and the LOS signal, and  $\phi$  is the relative phase. When the distance between the transmitted and received signal is large, the ground reflection will be a glancing reflection and the LOS and reflected path will be almost exactly the same length. Ground reflection will usually introduce a  $180^\circ$  phase rotation in the reflected signal and  $r$  could be 0.8–0.9. Thus there will be nearly complete cancellation of the signal [8].

### Rayleigh Fading

Rayleigh fading, or **fast fading**, in a multipath environment results from destructive cancellation as individual paths of similar amplitude drift in and out of phase as the receiver and sources of multiple reflections, diffractions, and refractions move. With multiple paths between a transmitter and a receiver, different components of the received signal arrive at different times. When the line of sight is blocked, the received signal on as many as 20 of these paths can have appreciable signal power. The time between when the first significant received signal is received and the last is received is called the



**Figure 4-17:** Two receive antennas used to achieve **space diversity** and overcome the effects of Rayleigh fading.

delay spread. The delay spread is largest inside buildings. A typical office building has a delay spread of 30 ns, and an office building with highly reflective walls and large open spaces has a **delay spread** of up to 250 ns [9–11], similar delay spreads are obtained at all RF frequencies [12]. Rayleigh fading is the type of fading that occurs when the inverse of the delay spread is small compared to the bandwidth of the received signal. Then when the multipaths combine destructively the entire signal within a bandwidth of 500 kHz to 1 MHz is suppressed. This bandwidth is experimentally observed and is related to the delay spread of individual paths, each with random amplitude and phase. If a communication signal has a bandwidth that is entirely within the suppression bandwidth of a Rayleigh fade, the situation is called frequency flat fading. Some communication signals have bandwidths greater than the Rayleigh fade bandwidth, so it is possible to recover from Rayleigh fades if error correction is used.

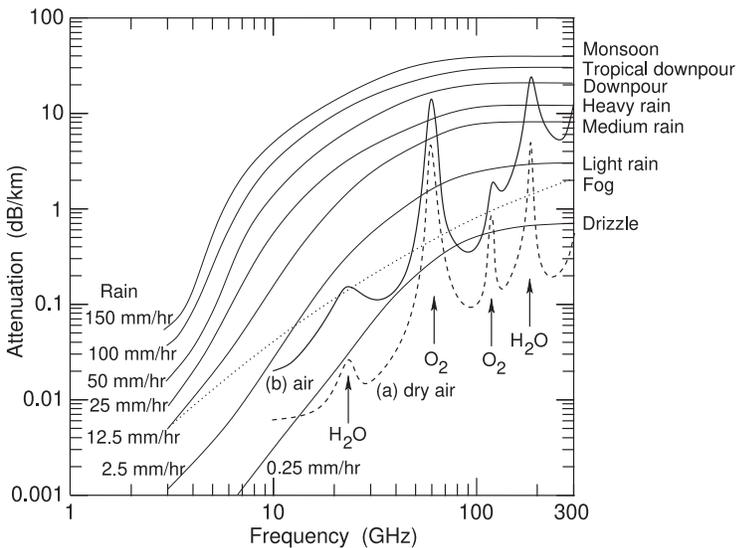
The result of Rayleigh fading is rapid amplitude variation with respect to frequency. In Figure 4-16(a) the amplitude varies in time depending on the speed of moving objects. This graph could just as well be a plot of amplitude versus distance or amplitude versus frequency. A measured amplitude response is shown in Figure 4-16(b). Although somewhat random, the fades occur over distances roughly  $\frac{1}{2}\lambda$  apart and 500 kHz wide. The 500 kHz width is almost independent of frequency from hundreds of megahertz to 100 GHz. The probability of receiving a signal  $x$  dB below the time-averaged received signal power is approximately  $10^{-x}$  [13]. Rayleigh fading is named after the statistical model that describes this.

Fortunately deep Rayleigh fades are very short, last a small percentage of the time, and slight changes to the propagation environment can circumvent their effects. One strategy for overcoming fades is to use two receive antennas as shown in Figure 4-17. Here the signals received by two antennas separated by several wavelengths will rarely fade concurrently. In practice, the required separation for good decorrelation is found to be 10 to  $20\lambda$ .

In an LOS wireless system (e.g. a point-to-point link without multipath), Rayleigh fading is due to rapidly changing atmospheric conditions, with the refractive index of small regions varying. These fades occur over a few seconds.

### Rain Fading

Rain fading is due to both the amount of rain and the size of individual rain drops and fading occurs over periods of minutes to hours. Propagation through the atmosphere is affected by absorption by molecules in the air, fog, and rain, and by scattering by rain drops. Figure 4-18 shows the attenuation in decibels per kilometer from 3 GHz to 300 GHz. Curve (a) shows the



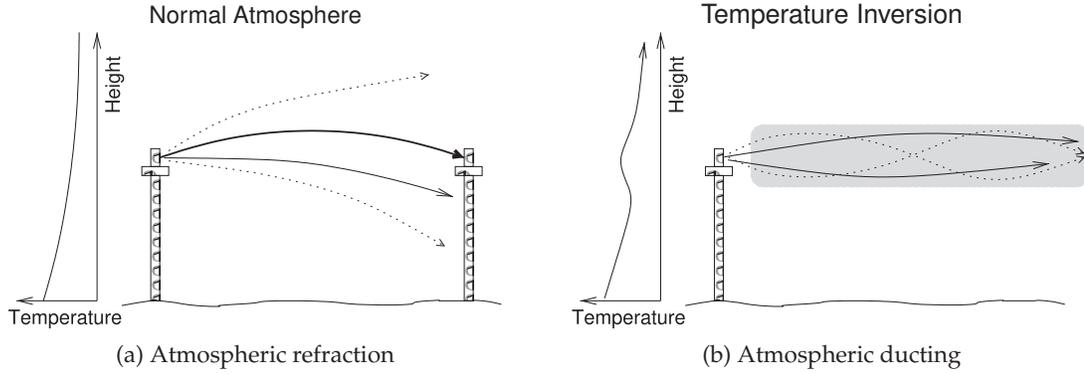
**Figure 4-18:** Excess attenuation due to atmospheric conditions showing the effect of rain on RF transmission at sea level. Curve (a) is atmospheric attenuation, due to excitation of molecular resonances, of very dry air at 0 °C, curve (b) is for typical air (i.e. less dry) at 20 °C. The attenuation shown for fog and rain is additional (in dB) to the atmospheric absorption shown as curve (b).

attenuation in dry air at 0°C with very low attenuation at a few gigahertz with several attenuation peaks as frequency increases. The first peak in attenuation is due to the resonance of water molecules in water vapor. The resonance peaks at 23 GHz but loss starts increasing before that. The next absorption peak is at 60 GHz due to the resonance of oxygen molecules. Two other absorption peaks are observed going up to 300 GHz. Curve (b) is for a less dry atmosphere at 20°C. The dotted curve shows the additional effect of fog on attenuation and then a family of curves shows the effect of rain on propagation. Below 10 GHz the effect of rain is very little and is safely ignored below 5 GHz. The attenuation due to rain increases with frequency and this derives largely from scattering and is related to the RF signal's wavelength relative to the circumference of rain drops.

### Fading Due to Ducting

Normally the temperature and density of air, and thus its refractive index, drops with increasing height above earth. As a result, radio waves will be refracted toward the earth as shown in Figure 4-19(a). However, during periods of stable weather characterized by a high-pressure system, the temperature can rise with increasing height before eventually falling, creating what is called temperature inversion. This can occur at heights of tens to hundreds of meters. Temperature inversion is most common in summer, but can also occur when the temperature is dropping quickly, such as at sunset. At sunset the inversion layer may occur at a meter or so above the ground as the earth cools rapidly. The denser colder air above the ground has a higher permittivity than the air above and this results in an inversion layer that has a higher refractive index than the air above and below. RF waves can become trapped in the inversion layer, as any RF energy leaving the temperature inversion layer is refracted back into the layer. This effect is called ducting. Ducting can also occur when a cold air mass is overrun by warm air.

Fading due to ducting occurs when the receiver wanders in and out of the ducting layer, as the ducting layer is not stable, increasing and decreasing in



**Figure 4-19:** Fading resulting from ducting: (a) normal atmospheric refraction (normally the temperature of air drops with increasing height and the lower refractive index at high heights results in a concave refraction); and (b) atmospheric ducting (resulting from temperature inversion inducing an air layer with higher dielectric constant than the surrounding air).

strength and in geometry.

### Summary of Fading

The fades of most concern in a mobile wireless system are the deep fades resulting from destructive interference of multiple reflections. These fades vary rapidly (over a few milliseconds) if a handset is moving at vehicular speeds but occur slowly when the transmitter and receiver are fixed. Fades can be viewed as deep amplitude modulation, and so so modulation is restricted to phase shift keying schemes when a transmitter and receive moving at vehicular speeds relative to each other.

#### 4.6.4 Link Loss and Path Loss

With transmit and receive antennas included in the RF link, the usual case, link loss is defined as the ratio of the power input to the transmit antenna ( $P_T$ ) to the power delivered by the receive antenna ( $P_R$ ). Rearranging Equation (4.31) and taking logarithms yields the total line-of-sight (LOS) link loss,  $L_{\text{LINK,LOS}}$ , between the input of the transmit antenna and the output of the receive antenna separated by distance  $d$  is (in decibels):

$$L_{\text{LINK,LOS}}|_{\text{dB}} = 10 \log \left( \frac{P_T}{P_R} \right) = 10 \log \left( \frac{P_T}{P_D A_R} \right) \quad (4.36)$$

$$= 10 \log \left[ P_T \left( \frac{4\pi d^2}{P_T G_T} \right) \left( \frac{4\pi}{\lambda^2 G_R} \right) \right] \quad (4.37)$$

$$= 10 \log \left[ \left( \frac{1}{G_T G_R} \right) \left( \frac{4\pi d}{\lambda} \right)^2 \right] \quad (4.38)$$

$$= -10 \log G_T - 10 \log G_R + 20 \log \left( \frac{4\pi d}{\lambda} \right). \quad (4.39)$$

The last term includes  $d$  and is called the LOS path loss (in decibels):

$$L_{\text{PATH,LOS}}|_{\text{dB}} = 20 \log \left( \frac{4\pi d}{\lambda} \right). \quad (4.40)$$

This is the preferred form of the expression for path loss, as it can be used directly in calculating link loss using the antenna gains of the transmit and receive antennas without the exercise of calculating the **effective aperture size** of the receive antenna.

An alternative definition of path loss comes directly from Equation (4.25) and is called the LOS path loss of the first kind:

$${}^1L_{\text{PATH,LOS}} = 4\pi d^2, \quad (4.41)$$

but this is not commonly used.

Multipath effects result in losses that are proportional to  $d^n$  [14, 15] so that the general path loss, including multipath effects, is (in decibels)

$$\begin{aligned} L_{\text{PATH}}|_{\text{dB}} &= L_{\text{PATH,LOS}}|_{\text{dB}} + \text{excess loss}|_{\text{dB}} \\ &= 20 \log \left( \frac{4\pi d}{\lambda} \right) + 10(n-2) \log \left( \frac{d}{1 \text{ m}} \right) \\ &= 20 \log \left[ \frac{4\pi(1 \text{ m})}{\lambda} \right] + 10(2) \log \left( \frac{d}{1 \text{ m}} \right) + 10(n-2) \log \left( \frac{d}{1 \text{ m}} \right) \\ &= 10n \log[d/(1 \text{ m})] + C, \end{aligned} \quad (4.42)$$

where the distance  $d$  and wavelength  $\lambda$  are in meters, and  $C$  is a constant that captures the effect of wavelength. Here,

$$C = 20 \log [4\pi(1 \text{ m})/\lambda]. \quad (4.43)$$

Combining this with Equation (4.39) yields the link loss between the input of the transmit antenna and the output of the receive antenna:

$$L_{\text{LINK}}|_{\text{dB}} = -G_T|_{\text{dB}} - G_R|_{\text{dB}} + 10n \log[d/(1 \text{ m})] + C. \quad (4.44)$$

As you can imagine, a few constants, here  $n$  and  $C$ , cannot capture the full complexity of the propagation environment. Many models have been developed to capture particular environments better and incorporate mast height, experimental correction factors, and statistical parameters.

The path loss between two antennas is exactly the same in both directions when the frequency of the signal in each direction is the same. This is radio link reciprocity. Many communication systems use different frequencies in the two directions, and then the links are not reciprocal.

**EXAMPLE 4.5** Power Density

A communication system operating in a dense urban environment has a power density roll-off of  $1/d^{3.5}$  between the base station transmit antenna and the mobile receive antenna. At 10 m from the transmit antenna, the power density is  $0.3167 \text{ W/m}^2$ . What is the power density at the receive antenna located at 1 km from the base station?

**Solution:**

$P_D(10 \text{ m}) = 0.3167 \text{ W/m}^2$  and let  $d_c = 10 \text{ m}$ , so at  $d = 1 \text{ km}$ , the power density  $P_D(1 \text{ km})$  is obtained from

$$\frac{P_D(1 \text{ km})}{P_D(10 \text{ m})} = \frac{d_c^{3.5}}{d^{3.5}} = \frac{10^{3.5}}{1000^{3.5}} = 10^{-7}, \quad (4.45)$$

so 
$$P_D(1 \text{ km}) = P_D(10 \text{ m}) \cdot 10^{-7} = 31.7 \text{ nW/m}^2. \quad (4.46)$$

**EXAMPLE 4.6** Link Loss

A 5.6 GHz communication system uses a transmit antenna with an antenna gain  $G_T$  of 35 dB and a receive antenna with an antenna gain  $G_R$  of 6 dB. If the distance between the antennas is 200 m, what is the link loss if the power density reduces as  $1/d^3$ ? The link loss here is between the input to the transmit antenna and the output from the receive antenna.

**Solution:**

The link loss is provided by Equation (4.44),

$$L_{\text{LINK}}|_{\text{dB}} = -G_T - G_R + 10n \log[d/(1 \text{ m})] + C,$$

and  $C$  comes from Equation (4.43), where  $\lambda = 5.36 \text{ cm}$ . So

$$C = 20 \log \left( \frac{4\pi}{\lambda} \right) = 20 \log \left( \frac{4\pi}{0.0536} \right) = 47.4 \text{ dB}.$$

With  $n = 3$  and  $d = 200 \text{ m}$ ,

$$L_{\text{LINK}}|_{\text{dB}} = -35 - 6 + 10 \cdot 3 \cdot \log(200) + 47.4 \text{ dB} = 75.4 \text{ dB}.$$

**EXAMPLE 4.7** Radiated Power Density

In free space, radiated power density drops off with distance  $d$  as  $1/d^2$ . However, in a terrestrial environment there are multiple paths between a transmitter and a receiver, with the dominant paths being the direct LOS path and the path involving reflection off the ground. Reflection from the ground partially cancels the signal in the direct path, and in a semi-urban environment results in an attenuation loss of 40 dB per decade of distance (instead of the 20 dB per decade of distance roll-off in free space). Consider a transmitter that has a power density of  $1 \text{ W/m}^2$  at a distance of 1 m from the transmitter.

- The power density falls off as  $1/d^n$ , where  $d$  is distance and  $n$  is an index. What is  $n$ ?
- At what distance from the transmit antenna will the power density reach  $1 \mu\text{W} \cdot \text{m}^{-2}$ ?

**Solution:**

- Power drops off by 40 dB per decade of distance. 40 dB corresponds to a factor of 10,000 ( $= 10^4$ ). So, at distance  $d$ , the power density  $P_D(d) = k/d^n$  ( $k$  is a constant).

At a decade of distance,  $10d$ ,  $P_D(10d) = k/(10d)^n = P_D(d)/10000$ , thus

$$\frac{k}{10^n d^n} = \frac{1}{10,000} \frac{k}{d^n}; \quad 10^n = 10,000 \Rightarrow n = 4.$$

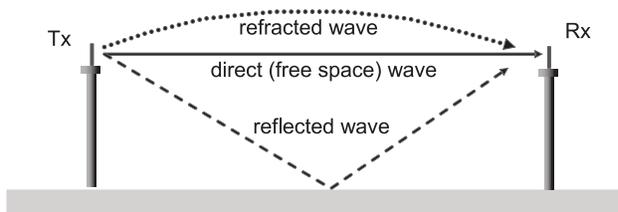
(b) At  $d = 1 \text{ m}$ ,  $P_D(1 \text{ m}) = 1 \text{ W/m}^2$ . At a distance  $x$ ,

$$P_D(x) = 1 \mu\text{W/m}^2 = \frac{k}{x^4} \text{m}^2 = \frac{k}{x^4} \rightarrow x^4 = \frac{1}{10^{-6}} \text{ and so } x = 31.6 \text{ m.}$$

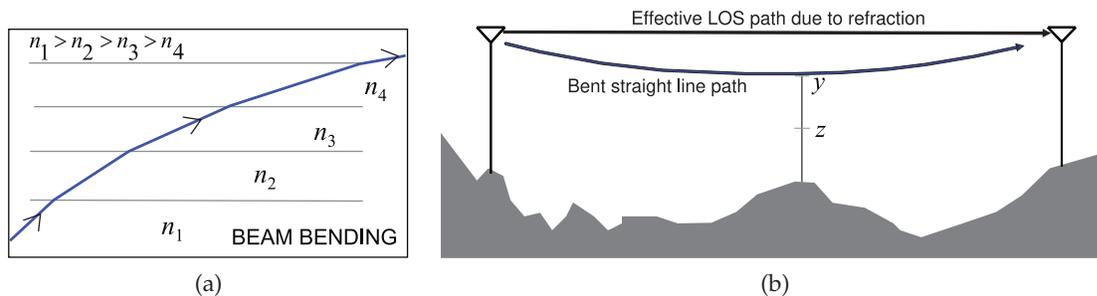
### 4.6.5 Fresnel Zones

In long-distance wireless communications from one fixed site to another, the intent is to use the LOS path and avoid reflections. Such systems are called **point-to-point links**. Thus avoiding reflections is an important consideration in design. The dominant propagation paths in a point-to-point system are shown in Figure 4-20. The refracted wave path arises because of density variations in the air producing a permittivity profile that varies with height, as shown in Figure 4-21(a). This effect is called **beam bending**. The other important propagation path to consider is the reflected wave from the ground, which can be important if the ground is too close to the propagation path. Both of these effects will be considered in this section.

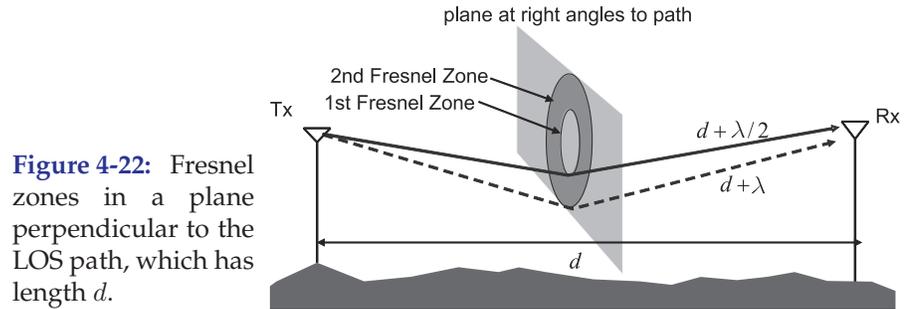
As radio waves propagate, they spread out in a plane perpendicular to the direction of propagation, the power density of the radio waves then reduces with distance from the centerline. One of the consequences of this is that an obstruction that is not in the LOS path can still interfere with signal propagation. The appropriate clearance is determined from the Fresnel zones, which are shown in Figure 4-22. The direct LOS path between the antennas has a length  $d$ . If there is a reflecting object near the LOS path, then there can be a second path between the transmit and receive antennas. The path from the first antenna to the circle defined by the first Fresnel zone and then to the second antenna has a path length  $d + \lambda/2$ , and so at the receive antenna this signal is  $180^\circ$  out of phase with the LOS signal and there will be



**Figure 4-20:** Three point-to-point characteristic propagation paths: line of sight, reflection, and refraction.



**Figure 4-21:** Beam bending by density variation in air: (a) refraction index profile with the air density reducing with height; and (b) incorporating beam bending in a curved-earth model.



**Figure 4-22:** Fresnel zones in a plane perpendicular to the LOS path, which has length  $d$ .

cancellation. The radius of the  $n$ th Fresnel zone at a point  $P$  is

$$F_n = \sqrt{\frac{n\lambda d_1 d_2}{d_1 + d_2}}, \quad (4.47)$$

where  $d_1$  is the distance from the first antenna to  $P$ ,  $d_2$  is the distance from the second antenna to  $P$  (so  $d = d_1 + d_2$ ), and  $\lambda$  is the wavelength of the propagating signal. Ninety percent of the energy in the wave is in the first Fresnel zone. A guideline is that an obstacle should be separated from the direct path by a distance more than the radius of the first Fresnel zone. Antenna heights are increased so that the beam achieves one or more Fresnel zone clearances.

The above analysis can be used even with beam bending. A convenient way of accommodating beam bending is to use a curved-earth model, as shown in Figure 4-21(b), so that subsequent calculations can use LOS considerations [16]. The amount of beam bending to account for comes from experimental surveys. In Figure 4-21(b) the original clearance from a hill to the first Fresnel zone is  $z$ . With beam bending the clearance increases to  $y$ .

Since a receive antenna is unlikely to be large enough to capture all of the energy contained in the first Fresnel zone, the effect of signal blocking by an obstruction that encroaches the first Fresnel zone is of little concern. What is of concern is the destructive combining of a reflected signal with the signal in the main LOS beam.

#### EXAMPLE 4.8

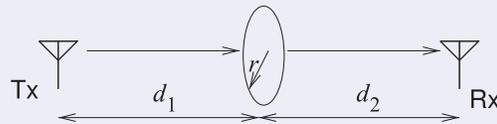
#### Fresnel Zone Clearance

A transmit antenna and a receive antenna are separated by 10 km and operate at 2 GHz.

- What is the radius of the first Fresnel zone?
- What is the radius of the second Fresnel zone?
- To ensure LOS propagation, what should the clearance be from the direct line between the antennas and obstructions such as hills and vegetation?

#### Solution:

The radius of the Fresnel zone is calculated at the midpoint so that  $d_1 = d_2 = d/2 = 5$  km. Also  $f = 2$  GHz,  $\lambda = 15$  cm.



- The radius of the first Fresnel zone is, from Equation (4.47),

$$r_1 = F_1 = \sqrt{\frac{\lambda d_1 d_2}{d_1 + d_2}} = \sqrt{\frac{\lambda d}{4}} = \sqrt{\frac{(0.15)(10^4)}{4}} = 19.36 \text{ m.}$$

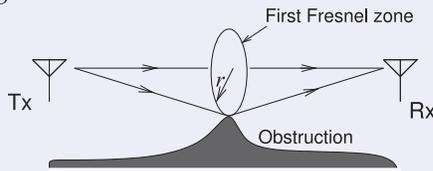
(b) At the midpoint the radius of the second Fresnel zone is

$$r_2 = F_2 = \sqrt{\frac{2\lambda d_1 d_2}{d_1 + d_2}} = \sqrt{\frac{\lambda d}{2}} = \sqrt{\frac{(0.15)(10^4)}{2}} = 27.39 \text{ m.}$$

(c) Ninety percent of the energy in the beam is contained within the first Fresnel zone of radius  $r_1$ . Obstructions within the first Fresnel zone will result in a significant fraction of the beam being blocked. In addition, reflections from the obstruction could destructively combine with the main beam and reduce the received signal level even beyond the reduction caused by part of the beam being blocked. Two criteria are commonly used to determine the clearance required to avoid signal obstruction.

One criterion used is that the minimum clearance between the direct beam and an obstruction is  $0.6r_1$ . So the minimum clearance is  $0.6r_1 = 11.6 \text{ m}$ .

A more conservative criterion is that the minimum clearance should be  $r_1 = 19.36 \text{ m}$ .



#### 4.6.6 Propagation Model in the Mobile Environment

RF propagation in the mobile environment cannot be accurately derived. Instead, a fit to measurements is often used. One of the models is the Okumura–Hata model [17], which calculates the path loss as

$$L_{\text{PATH}}|_{(\text{dB})} = 69.55 + 26.16 \log f - 13.82 \log H + (44.9 - 6.55 \log H) \cdot \log d + c, \quad (4.48)$$

where  $f$  is the frequency (in MHz),  $d$  is the distance between the base station and terminal (in km),  $H$  is the effective height of the base station antenna (in m), and  $c$  is an environment correction factor ( $c = 0 \text{ dB}$  in a dense urban area,  $c = -5 \text{ dB}$  in an urban area,  $c = -10 \text{ dB}$  in a suburban area, and  $c = -17 \text{ dB}$  in a rural area, for  $f = 1 \text{ GHz}$  and  $H = 1.5 \text{ m}$ ).

More sophisticated characterization of the propagation environment uses ray-tracing models to follow individual propagation paths. The ray-tracing models are based on deterministic methods using terrain data and calculate paths accounting for obstruction and reflection analyses. Each refraction and reflection event is characterized either experimentally or through detailed EM simulations. Appropriate algorithms are applied for best compliance with radio physics. Commonly required inputs to these models include frequency, distance from the transmitter to the receiver, effective base station height, obstacle height and geometry, radius of the first Fresnel zone, forest height/roof height, distance between buildings, arbitrary loss allowances based on land use (forest, water, etc.), and loss allowances for penetration of buildings and vehicles. Such a technique was used to calculate the radio coverage diagrams on the inside front cover of this book.

There are many propagation models for different frequency ranges and different environments. The considerable effort put into developing reliable models is because being able to predict signal coverage is essential to efficient design of basestation layout.

## 4.7 Multipath and Delay Spread

Multipath was discussed previously in Sections 4.6.1 and 4.6.3 in the context of reduction in signal amplitude due to the destructive interference of a signal traveling on several paths. Another effect of multipath is that instances of a signal traveling from a transmitter to a receiver on different paths will arrive at a receiver with various delays. These delays can differ by tens or thousands of nanoseconds which is much longer than the period of a microwave signal, e.g. a 1 GHz signal has a period of 1 ns.

### 4.7.1 Delay Spread

Each of the paths supports what is called a signal instance. For a single-tone signal the effect of various delays is for the signal instances to arrive at the receiver with very different phases. As well, the signal instances invariably have different amplitudes resulting in a composite received signal that is an unpredictable combination of constructive and destructive interference. Sometimes the constructive interference will result in a stronger composite received signal but this is not that important. What is much more significant is the destructive interference which will result in signal fades with the signal received being very small at times. Consider two signal instances  $y_1(t)$  and  $y_2(t)$ , traveling from a transmitter to a receiver on paths 1 and 2 respectively. If  $y_1(t)$  and  $y_2(t)$  have the same amplitude and phase then the combined signal will only be increased by 6 dB over a signal that travels on just one path. However if  $y_1(t)$  and  $y_2(t)$  have the same amplitude but phases that differ by  $180^\circ$  then there will be total cancellation and no composite signal will be received. So destructive interference is much more significant than constructive interference.

For cellular radio, in a rural area there can be two or three significant paths, usually taken as paths having signal instances that are within 20 dB of the largest signal instance. In an urban environment there can be tens of significant paths because there are many reflections from buildings and often there is not an LOS path which would otherwise have the largest signal instance. With multiple paths it is unlikely that there will be total cancellation of the received signal but it is very likely that the received signal will be smaller than if there was one path.

Early cellular radio had narrow bandwidths, e.g. 200 kHz for 2G's GSM cellular system, and low-order modulation, e.g. 2 bits per symbol and it was sufficient to use the phase-based concepts of constructive and destructive interference to understand the effects of multipath. As channel bandwidths increased and the order of modulation increased it was necessary to consider the actual physical effect of the paths having different delays.

If there is an LOS path, say with no building in the way, then the LOS signal instance will be stronger than a signal instance that travels on any other path. It will also have the smallest delay. This is partly because each non-LOS (NLOS) path will travel further and hence spread out more, and also because there will be signal loss at reflections. If there is scattering or diffraction on a path the amplitude of this signal instance on that path will be even smaller. Even if there is an LOS path, if there are multiple NLOS paths then the combined NLOS signal instances could indeed be larger than the LOS signal instance. If there is not an LOS path then one of the NLOS paths will become

Freq- uency (GHz)	Tx-Rx distance (km)	Environ- ment	max. excess delay $\tau_{d,max}$ (ns)	median excess delay $\tau_{d,med}$ (ns)	
0.43	3.8 <sup>†</sup>	urban	1300	900	[18]
0.90	2.2	rural	800	–	[19]
0.90	5	rural	1900	–	[20]
0.90	2	suburban	3700	300	[21]
0.90	2.2	suburban	900	300	[19]
0.90	0.6	urban	700	200	[22]
0.90	3.5	urban	5000	1100	[23]
0.90	7.0	urban	15000	1900	[24]
0.90	1.0	urban	2000	700	[21]
0.90	13	mountain <sup>§</sup>	3800	–	[25]
0.90	6.0	mountain <sup>§</sup>	1800	–	[26]
1.35	3.8 <sup>†</sup>	urban	1400	850	[18]
2.26	3.8 <sup>†</sup>	urban	1400	800	[18]
5.75	3.8 <sup>†</sup>	urban	1000	300	[18]
28	0.052	urban LOS	754	<200	[27]
28	0.097	urban NLOS	1388	200	[27]
38	0.2 <sup>‡</sup>	urban LOS	12	1.5	[28]
38	0.2 <sup>‡</sup>	urban NLOS	133	14	[28]

**Table 4-2:** Measured excess delays for various carrier frequencies, Tx-Rx distances and multipath environments. Data has been normalized where needed so that  $\tau_{d,max}$  is the maximum excess delay that is exceeded only 1% of the time.

<sup>†</sup>  $d$  ranges from 0.04 km to 3.8 km and here  $\tau_{d,max}$  is the maximum delay exceeded only 10% of the time.

<sup>‡</sup> high gain steerable antennas.

<sup>§</sup> urban area surrounded by mountains

the largest received signal instance and it may not have the shortest delay. Generally only signal instances that are within 20 dB of the largest signal instance must be considered and then the delay, also referred to as ‘delay (20 dB)’, that matters is the excess delay,  $\tau_d$ , which is the difference between the earliest and latest arriving signal instances within that 20 dB window. It is essential that a communications link be maintained a certain percentage of the time, e.g. 99%, for a specified level of service. This minimizes the number of times packets must be re-transmitted and assures acceptable operation of a system. Thus the maximum excess delay,  $\tau_{d,max}$  becomes an important metric and determines the interval for which the received signal is not valid, the required guard band to prevent signals from an earlier symbol interval overlapping with the current symbol interval, and also the maximum useful bandwidth of a signal.

It is found that the delay spread increases with distance between the transmitter and receiver, the Tx-RX distance, increases with frequency, and increases with a richer multipath environment. For example, a rural area has few paths but large cell sizes are used and the Tx-Rx distance is relatively large resulting in a large excess delay but a dense urban area has small cells but a rich multipath environment with many reflections from buildings. Measured excess delays for various situations are given in Table 4-2. There is a lot of variation in excess delays but it is clearly reduced with smaller cells and a less rich multipath environment.

It is also found that as a mobile unit moves the excess delay and the number of significant paths varies tens or hundreds of times per second depending on the environment but also the mobility, e.g. pedestrian versus highway speeds. However over a millisecond the propagation is found to be effectively fixed and consequently systems are designed with data sent in packets that are usually no longer than a millisecond.

Examining Table 4-2 further it is observed that the 2G and 3G cellular bands are below 2 GHz and the delay spreads can be several milliseconds

and even larger with larger Tx-Rx distances in an urban area. At the higher microwave frequencies up to 6 GHz covering the operating frequency range of 4G, the delay spread is smaller.

The 28 GHz and 38 GHz bands are where 5G operates with short Tx-Rx distances. Here the number of paths above 20 dB is around 4 to 8 for Tx-Rx separations of 35 to 193 m [27, 28].

The spreads are relatively low and this is largely a result of using high gain steerable antennas which are essential to 5G operation. The high gain antennas reduce the spread of multipath path lengths, and hence reduce the excess delay spread, concentrating the reflected paths close to the transmit and receive antennas.

The prior discussions on multipath effects and the RF link consider a single frequency tone propagating from a transmit antenna and a receive antenna. When a single-tone signal travels on multiple paths there will be destructive and/or constructive interference.

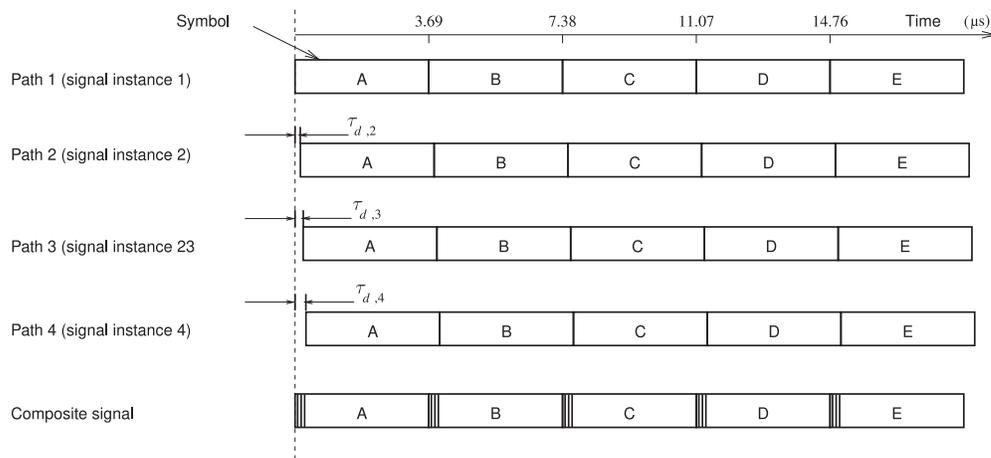
#### 4.7.2 Intersymbol Interference

Intersymbol interference (ISI) occurs when a symbol traveling on one path interferes with the symbol traveling on another path. This is the result of the paths having different delays. The effect of excess delay on the received composite signal is indicated in Figure 4-23. Consider the situation where there are four paths with the excess delay of path 1 being 0 by definition, i.e.  $\tau_{d,1} = 0$ , the excess delay of path 2,  $\tau_{d,2} = 100$  ns, and so on,  $\tau_{d,3} = 200$  ns., and  $\tau_{d,4} = 300$  ns. The 2G–5G cellular systems transmit a series of symbols, here A, B, C, D, and E, which have a symbol duration which differs by standard. In Figure 4-23 four paths are considered and the four signal instances each travel on a different path and are combined in the receiver's antenna to yield a composite received signal. The symbols of each of the signal instances overlap causing ISI.

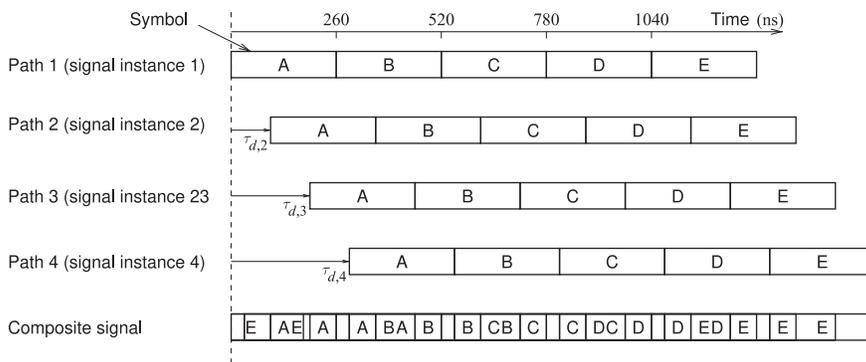
With the 2G GSM system symbol lengths are 3690 ns and the first 300 ns of a symbol instance traveling on the first path to arrive has interference from the previous symbol traveling on the other three paths, see Figure 4-23(a). In this case the first 300 ns of the 3690 ns of the symbol is garbled. In this example the delay spread is just 300 ns but this is conservative. Table 4-2 indicates that in some environments that the excess delay spread 1% of the time can be as long as the symbol interval for 2G/GSM (operating below 2 GHz). So even if the signal strength is high reception may not be possible because of ISI.

With 3G the packet length is 260 ns so with an excess delay spread of 300 ns for the four path example the symbol interference is severe, see Figure 4-23(b), and as noted 300 ns is a very conservative estimate of the excess delay spread. One of the features of 3G is a method for separating out the individual signal instances. Only when there are a very large number of significant paths does 3G have problems. In an urban area sometimes a communication link cannot be established with 2G and 3G even though the signal strength is high.

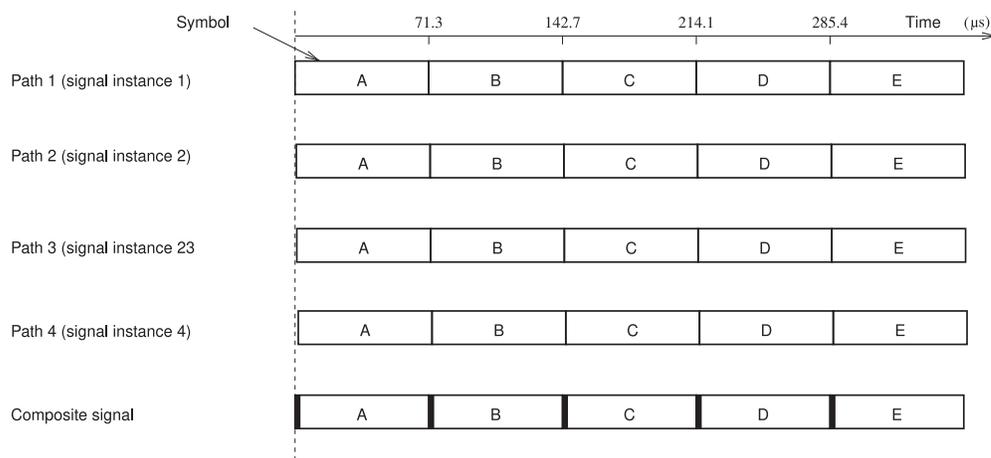
The 4G system minimizes the problem of ISI by having a very long symbol interval of 71.35  $\mu$ s and a 300 ns excess delay spread is a small fraction of the symbol length. The symbol interval is considerably longer than any of the maximum delay spreads given in Table 4-2. Also the 4G standard



(a) Symbols in 2G GSM



(b) Symbols in 3G WCDMA



(c) Symbols in 4G

**Figure 4-23:** The effect of excess delay on symbol interference when there are four paths. The excess delay of each path,  $\tau_{d,x}$ , increases by 100 ns. Five symbols are shown: A, B, C, D, and E. The 2G GSM standard has a symbol duration of 3.6928  $\mu\text{s}$ , the 3G WCDMA standard has a symbol duration of 260 ns, and the 4G OFDM standard has a symbol duration of 71.35  $\mu\text{s}$ . The situation with 5G is similar to that for 4G.

has a guard time band known as a cyclic prefix. (The cyclic prefix is a bit more than a guard band as it involves repeating the end of symbol but this extra feature will be discussed later.) The cyclic prefix can be either 4.7  $\mu\text{s}$ , 5.2  $\mu\text{s}$ , 16.7  $\mu\text{s}$ , or 33.3  $\mu\text{s}$  accommodating various excess delay spreads. This is more than enough to avoid the problem of ISI as any of the maximum excess delays listed in Table 4-2 are less than the maximum cyclic prefix length. The long symbol interval implies a very narrow bandwidth or sub-channel bandwidth. A 4G system communicates with a user using a very large number of sub-channels thus supporting high data rates.

### 4.7.3 Summary

Various excess delays on different paths causes intersymbol interference. High levels of intersymbol interference results in the failure to establish a communication link. Excess delay is a fundamental limitation with the 2G/GSM system and the only solution is to use very small cells in urban environments which have many significant signal paths. The 3G system employs a coding technique that enables the first few different paths to be resolved thus limiting the problem of ISI but not eliminating it. The 4G system, and 5G operates similarly, employs a long guard time, the cyclic prefix, to avoid ISI completely.

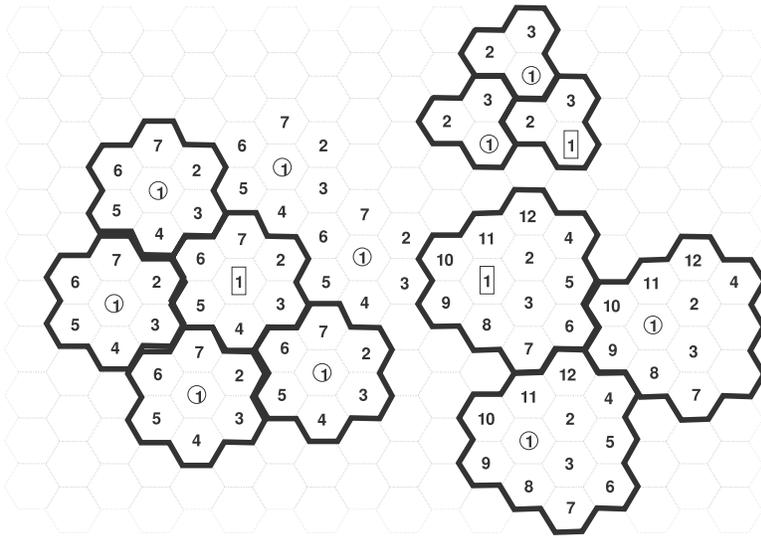
## 4.8 Radio Link Interference

The two major source of interference in a cellular system are self-interference resulting from excess delay spreads as discussed in the previous section, and interference, called radio link interference, from radios in other cells operating on the same frequency channel. This section discusses radio link interference and how it is controlled by using a frequency reuse plan.

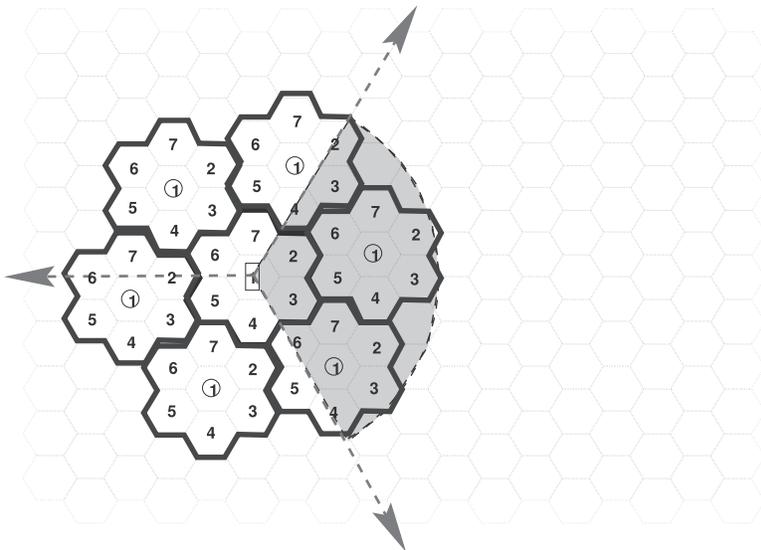
### 4.8.1 Frequency Reuse Plan

Radio systems using the same channel are geographically spaced to control interference. In a cellular system, the coverage areas are arranged in cells represented by the hexagons in Figure 4-24. The actual shape of the cells is influenced by obstructions such as hills and buildings. In 1G and 2G cellular systems the cells are arranged in clusters such as the 3-, 7-, and 12-cell clusters shown, and the total number of channels available is divided among the cells in a cluster with the full set of frequency channels repeated in each cluster. (More about 3G–5G arrangements later.) The size of the clusters is the major component of what is called the **frequency reuse** plan. So a three-cell cluster has a spacing of approximately one cell diameter to the next cell using the same frequency channels. The signal level transmitted from the original cell will interfere with the signal in its corresponding cell in adjacent clusters. The level of interference is reduced with 7-cell and then 12-cell clusters. There is also background noise coming from cosmic sources as well as artificial sources, but in a cellular system, interference from other radios operating in the same system nearly always dominates.

The use of directional antennas at the base station increases the SIR. The interference pattern obtained using a **triselector antenna** (each segment of the antenna providing  $120^\circ$  of coverage) is shown in Figure 4-25. The triselector antenna can be arranged so that the transmit and receive antennas are

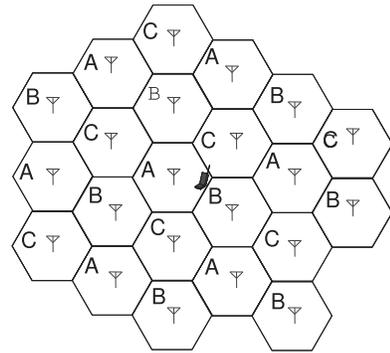
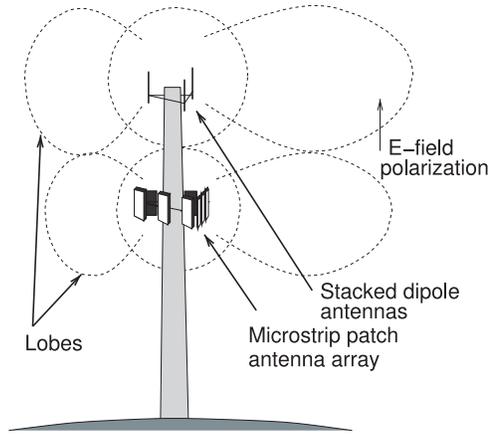


**Figure 4-24:** Cells arranged in clusters. Shown are 3-cell clusters (top right), 7-cell clusters (bottom left), and 12-cell clusters (bottom right).



**Figure 4-25:** Cells in a cellular radio system with a trisector antenna, such as in Figure 4-26, showing the area of possible interference as a shaded region.

separated (see Figure 4-26), and can also be arranged to tilt the coverage toward the ground (e.g., see the antenna in Figure 4-2(1)). It is also possible to use higher-order antenna sectoring and to have smaller cells (achieved possibly by relatively low-power transmissions) to provide higher levels of coverage at critical regions such as intersections of roads. These smaller cells are referred to as microcells or picocells.



**Figure 4-26:** Base station with trisector stacked-dipole transmit and microstrip patch receive antennas.

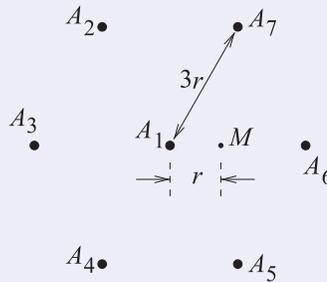
**Figure 4-27:** Three-cell cluster with mobile unit at the edge of cells A, B, and C.

**EXAMPLE 4.9 Cellular Interference**

In a cellular system, a signal is intentionally transmitted from a base station nominally located in the center of a cell to a mobile unit in the same cell. However, nearby transmitters using the same channel cause interference. In Figure 4-27, a mobile unit is located at the edge of a cell and uses frequency channel A. Many nearby transmitters also operate using channel A and the six nearest transmitters can be considered as causing significant interference. Consider that the mobile unit is a distance  $r$  from its cell's transmitter along the line connecting two channel A base stations, that the transmitters all operate at the same power level, and that the distance between base stations operating using channel A is  $3r$ . This three-cell cluster operates in a suburban area and the power density drops off with distance  $d$  as  $1/d^3$  due to multipath effects. What is the SIR at the mobile unit?

**Solution:**

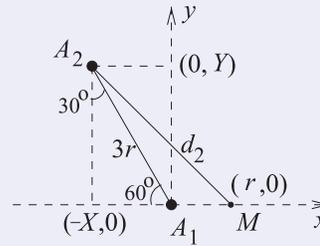
There are seven close towers transmitting signals on the same channel. Call these  $A_1$ – $A_7$  ( $M$  = mobile unit,  $A$  = transmitters).  $A_1$  is the desired signal and  $A_2$ – $A_7$  are interferers. The distances from the transmitters to the mobile unit are  $d_1$ – $d_7$ , and the powers from the transmitters received at the mobile unit are  $P_1$ – $P_7$ .



$$\text{Now SIR} = \frac{P_1}{P_2 + P_3 + P_4 + P_5 + P_6 + P_7}.$$

This problem requires  $d_2$ – $d_7$  to be determined.

Consider the diagram to the right with  $d_1 = r$ . Also, the distance from  $A_2$  to  $A_7$  is  $3r$ . Now  $d_3 = 3r + r = 4r$ ,  $d_6 = 3r - r = 2r$ . So  $d_2 = [(X + r)^2 + Y^2]^{\frac{1}{2}} = (2.5^2 + 2.6^2)r = 3.607r = d_4$ . ( $X = 3r \sin 30^\circ = 1.5r$  and  $Y = 3r \cos 30^\circ = 2.6r$ .) Similarly,  $d_7 = d_5 = [(1.5 - 1)^2 + 2.6^2]^{\frac{1}{2}} = 2.648r$ .



$$\frac{P_2}{P_1} = \frac{d_1}{d_2} = \frac{1}{3.607^3} = 0.0213 = \frac{P_4}{P_1} \quad \frac{P_3}{P_1} = \frac{d_1}{d_3} = \frac{1}{4^3} = 0.0156.$$

$$\frac{P_7}{P_1} = \frac{d_1}{d_7} = \frac{1}{2.648^3} = 0.0539 = \frac{P_5}{P_1} \quad \frac{P_6}{P_1} = \frac{d_1}{d_6} = \frac{1}{2^3} = 0.1250.$$

$$SIR = (0.0213 + 0.0156 + 0.0213 + 0.0539 + 0.1250 + 0.0539)^{-1} = 3.44 = 5.36 \text{ dB.}$$

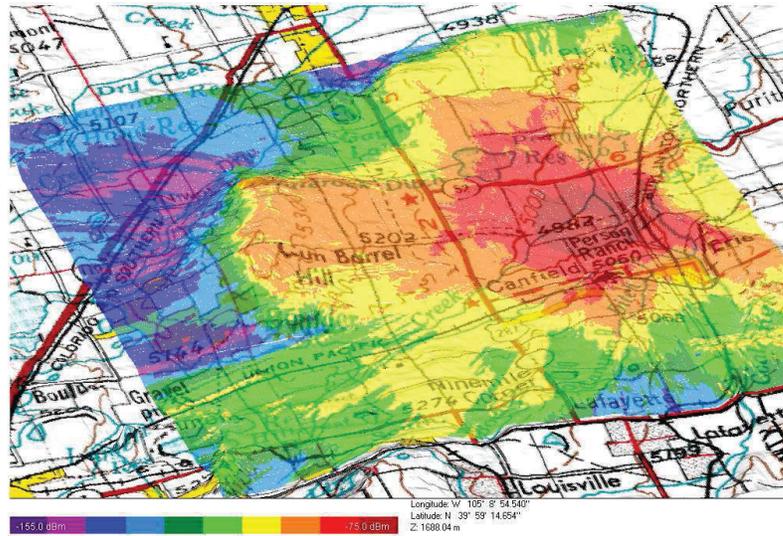
### 4.8.2 Summary

The clustering of cells as described above applies most directly to 1G and 2G cellular systems. The 3G system does not use clusters of cells but instead uses high frequency codes to separate users. There is a system awareness of managing interference but more complicated than the clustering arrangements. The 4G and 5G systems take frequency reuse algorithms to a whole new level performing global minimization of interference of wide areas. Still there is the basic concept that radios in adjacent cells operating on the same frequency channel can interfere with each other and careful attention must be given to frequency reuse.

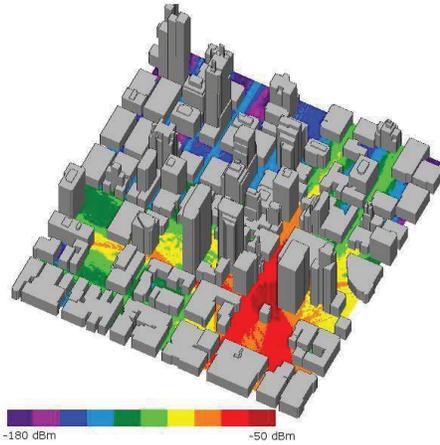
The region of signal coverage by a basestation is greatly influenced by terrain and is not a neat hexagon. See for example the coverage map for a rural area shown in Figure 4-28(a). To ensure adequate coverage basestations must be placed on a terrain-aware grid and not a regular grid. A wasted expense in operating a cellular system would be installing more basestations than are needed for optimum coverage. Thus cellular service providers simulate coverage before installing basestations. A significant characteristic of urban coverage as shown in Figures 4-28(b and c) is the “**urban canyon effect**” whereby coverage is concentrated along roads with signals bouncing off buildings directing signal propagation down the roadway. Even areas close to a transmitter but behind buildings are shadowed while distant areas near roads can have good coverage. This effect is also known as the urban **wave-guiding effect** or **wave-channeling effect**. The signal blocking by buildings is shown in Figure 4-28(e) where it is seen that there is some diffraction over buildings but shadowing is severe after two or more diffraction events. The indoor environment also affects coverage as shown in Figure 4-28(d) for WiFi.

### 4.9 Antenna array

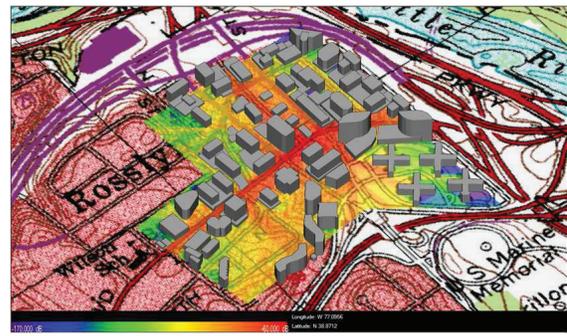
An antenna array comprises multiple radiating elements, i.e. individual antennas, and focuses a transmit beam in a desired direction. In Figure 4-29 the field pattern in the plane of the earth produced by an array of 30 antenna elements arranged horizontally is shown. The fields from each



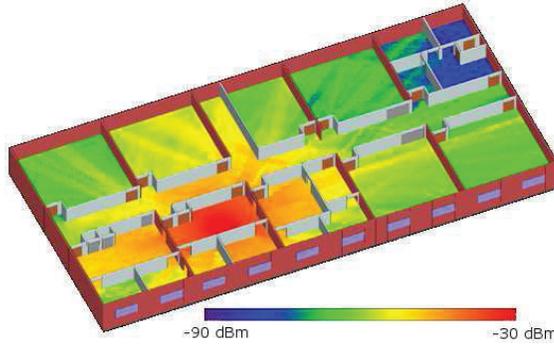
(a) Colorado plain east of Boulder and northwest of Denver (rural, gently rolling terrain).



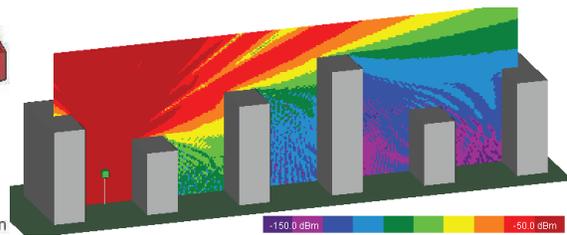
(c) Central Chicago (dense urban)



(c) Rosslyn Virginia (urban)

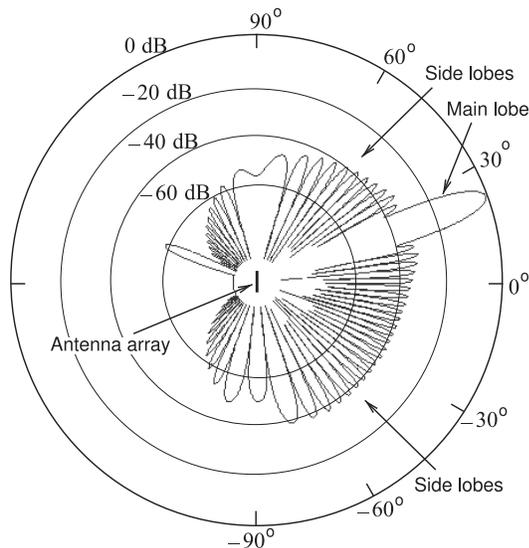


(d) Indoor coverage



(e) Vertical profile with low transmit antenna

**Figure 4-28:** Radio coverage profiles with 900 MHz transmitter except for (d) which is for 2.45 GHz. Calculated using Wireless Insite®. Copyright Remcom, Inc. Used with permission.



**Figure 4-29:** Electric field pattern from a 30 element array of antennas spaced  $0.65\lambda$  apart. The sidelobe levels are about 40 dB below the power level of the main lobe. The same signal is presented to the antenna elements except that phases of the signal at each antennas is adjusted to produce a main beam directed at 20 degrees. The signals to each antenna are thus correlated. After [29].

antenna element combined to narrow and strengthen the main beam. Side (or grating) lobes are produced and these will result in some interference but their level here is 40 dB below that of the main beam. This is another way of managing interference in a cellular system but the direction to the mobile unit must be known. Antenna arrays are used in 5G. The affect of the array is to increase the power density of the main beam and the density relative to that of an isotropic antenna is called the directional gain of the array,  $D_{\text{Array}}$ , and is the product of the antenna gain,  $G_A$ , of an individual antenna element and the array gain,  $G_{\text{Array}}$ :

$$D_{\text{Array}} = G_{\text{Array}} G_A. \quad (4.49)$$

The maximum value of  $G_{\text{Array}}$  is  $N$  for an  $N$  element array. So the maximum directional gain of the array in dBi is

$$D_{\text{Array}}|_{\text{dBi}} = G_A|_{\text{dBi}} + 10 \log N. \quad (4.50)$$

## 4.10 Summary

This chapter discussed the impacts on communication integrity of imperfections in the RF link from the output of the transmitter to the input of the receiver. Prior to cellular communications becoming so important, only the LOS communication path was considered and the impact of fading, multiple reflections, and delay spread was regarded as bothersome. Most commonly the solution to the problems introduced by these effects was either to shift to another frequency, to increase the operating power, or install more basestations. With the advent of cellular communications, and digital communications in general, techniques have been developed to overcome the negative impacts of fading without boosting signal levels to the level where they severely impact the operation of other radios. So while many aspects of RF propagation are random, concepts and statistical models have been developed that enable design choices to be made that permit digital communication systems to operate in what would have once been regarded as a hostile environment.

#### 4.11 References

- [1] S. Ramo, J. Whinnery, and T. Van Duzer, *Fields and Waves in Communication Electronics*. John Wiley & Sons, 1965.
- [2] C. A. Balanis, "Antenna theory: A review," *Proc. IEEE*, vol. 80, no. 1, pp. 7–23, 1992.
- [3] C. Balanis, *Antenna Theory: Analysis and Design*. John Wiley & Sons, 2005.
- [4] T. L. Wilson, K. Rohlf, and S. Hüttemeister, *Tools of radio astronomy*, 6th ed. Springer, 2013.
- [5] J. L. Pawsey and R. N. Bracewell, "Radio astronomy," *Oxford, Clarendon Press*, 1955.
- [6] J. Doble, *Introduction to Radio Propagation for Fixed and Mobile Communications*. Norwood, MA, USA: Artech House, Inc., 1996.
- [7] M. Born and E. Wolf, "Two-dimensional diffraction of a plane wave by a half-plane," in *Principles of Optics: Electromagnetic Theory of Propagation, Interference, and Diffraction of Light*, 7th ed. Cambridge University Press, 1999, section 11.5.
- [8] A. Abdi, C. Tepedelenlioglu, M. Kaveh, and G. Giannakis, "On the estimation of the  $k$  parameter for the rice fading distribution," *IEEE Communications Letters*, vol. 5, no. 3, pp. 92–94, Mar. 2001.
- [9] A. Saleh and R. Valenzuela, "A statistical model for indoor multipath propagation," *IEEE J. on Selected Areas in Communications*, vol. 5, no. 2, pp. 128–137, Feb. 1987.
- [10] D. Devasirvatham, "Time delay spread measurements of wideband radio signals within a building," *Electronics Letters*, vol. 20, no. 23, pp. 950–951, 8 1984.
- [11] A. Acampora and J. Winters, "System applications for wireless indoor communications," *IEEE Communications Magazine*, vol. 25, no. 8, pp. 11–20, Aug. 1987.
- [12] A. Siamarou and M. Al-Nuaimi, "Multipath delay spread and signal level measurements for indoor wireless radio channels at 62.4 ghz," in *2001 IEEE VTS 53rd Vehicular Technology Conf.*, 2001, pp. 454–458.
- [13] R. Gibson and R. Murray, "Systems organization for multichannel cordless telephones," in *Int. Zurich Seminar on Digital Communication*, 1982, pp. 59–61.
- [14] J. Andersen, T. Rappaport, and S. Yoshida, "Propagation measurements and models for wireless communications channels," *IEEE Communications Magazine*, vol. 33, no. 1, pp. 42–49, Jan. 1995.
- [15] T. Sarkar, Z. Ji, K. Kim, A. Medouri, and M. Salazar-Palma, "A survey of various propagation models for mobile communication," *IEEE Antennas and Propagation Magazine*, vol. 45, no. 3, pp. 51–82, Jun. 2003.
- [16] M. Murotani and K. Kohriyama, "Representation by approximate numerical formulas of radio-relay antenna beam bending due to atmospheric refraction," *IEICE Trans.*, vol. E72-E, no. 2, pp. 101–103, Feb. 1989.
- [17] J. Seybold, *Introduction to RF Propagation*. John Wiley & Sons, 2005.
- [18] P. Papazian, "Basic transmission loss and delay spread measurements for frequencies between 430 and 5750 mhz," *IEEE transactions on antennas and propagation*, vol. 53, no. 2, pp. 694–701, 2005.
- [19] J. Van Rees, "Measurements of the wideband radio channel characteristics for rural, residential, and suburban areas," *IEEE Transactions on Vehicular Technology*, vol. 36, no. 1, pp. 2–6, 1987.
- [20] J. A. Wepman, J. R. Hoffman, and L. H. Loew, "Characterization of macrocellular pcs propagation channels in the 1850-1990 MHz band," in *Proceedings of 1994 3rd IEEE International Conference on Universal Personal Communications*, 1994, pp. 165–170.
- [21] E. S. Sousa, V. M. Jovanovic, and C. Daigneault, "Delay spread measurements for the digital cellular channel in toronto," *IEEE Transactions on Vehicular Technology*, vol. 43, no. 4, pp. 837–847, 1994.
- [22] D. Devasirvatham, "Radio propagation studies in a small city for universal portable communications," in *38th IEEE Vehicular Technology Conference*, 1988, pp. 100–104.
- [23] D. Cox and R. Leck, "Distributions of multipath delay spread and average excess delay for 910-MHz urban mobile radio paths," *IEEE Transactions on Antennas and Propagation*, vol. 23, no. 2, pp. 206–213, 1975.
- [24] T. S. Rappaport, S. Y. Seidel, and R. Singh, "900-MHz multipath propagation measurements for us digital cellular radiotelephone," *IEEE Transactions on Vehicular Technology*, vol. 39, no. 2, pp. 132–139, 1990.
- [25] J.-P. De Weck, P. Merki, and R. W. Lorenz, "Power delay profiles measured in mountainous terrain (radiowave propagation)," in *38th IEEE Vehicular Technology Conference*, 1988, pp. 105–112.
- [26] T. Tanaka, S. Kozono, and A. Akeyama, "Urban multipath propagation delay characteristics in mobile communications," *Electronics and Communications in Japan (Part I: Communications)*, vol. 74, no. 8, pp. 80–88, 1991.
- [27] Y. Azar, G. N. Wong, K. Wang, R. Mayzus, J. K. Schulz, H. Zhao, F. Gutierrez Jr, D. Hwang, and T. S. Rappaport, "28 GHz

propagation measurements for outdoor cellular communications using steerable beam antennas in new york city." in *ICC*, 2013, pp. 5143–5147.

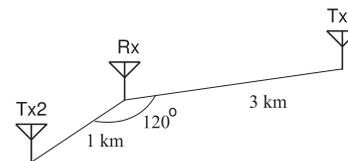
- [28] T. S. Rappaport, Y. Qiao, J. I. Tamir, J. N. Murdock, and E. Ben-Dor, "Cellular broadband millimeter wave propagation and angle of arrival for adaptive beam steering systems,"

in *2012 IEEE Radio and Wireless Symposium*. IEEE, 2012, pp. 151–154.

- [29] Phased array radiation pattern, *Phased\_array\_radiation\_pattern.gif*, By Maxter315 [CC BY-SA 4.0 (<https://creativecommons.org/licenses/by-sa/4.0/>)], from Wikimedia Commons.

## 4.12 Exercises

- An antenna only radiates 45% of the power input to it. The rest is lost as heat. What input power (in dBm) is required to radiate 30 dBm?
- The output stage of an RF front end consists of an amplifier followed by a filter and then an antenna. The amplifier has a gain of 27 dB, the filter has a loss of 1.9 dB, and of the power input to the antenna, 35% is lost as heat due to resistive losses. If the power input to the amplifier is 30 dBm, calculate the following:
  - What is the power input to the amplifier in watts?
  - Express the loss of the antenna in dB.
  - What is the total gain of the RF front end (amplifier + filter)?
  - What is the total power radiated by the antenna in dBm?
  - What is the total power radiated by the antenna in mW?
- The output stage of an RF front end consists of an amplifier followed by a filter and then an antenna. The amplifier has a gain of 27 dB, the filter has a loss of 1.9 dB, and of the power input to the antenna, 45% is lost as heat due to resistive losses. If the power input to the amplifier is 30 dBm, calculate the following:
  - What is the power input to the amplifier?
  - Express the loss of the antenna in decibels.
  - What is the total gain of the RF front end (amplifier + filter)?
  - What is the total power radiated by the antenna in dBm?
  - What is the total power radiated by the antenna in milliwatts?
- In the figure below there are two transmitters, Tx<sub>1</sub> and Tx<sub>2</sub>, operating at the same power level, and one receiver, Rx. Tx<sub>2</sub> is an intentional transmitter and its signal is intended to be received at Rx. Tx<sub>1</sub> uses the same frequency channel as Tx<sub>2</sub>, but it transmits an interfering signal. [Parallels Example 4.1]



- Assume that the antennas are omnidirectional and, being in a semiurban area, that the transmitted power density drops off as  $1/d^3$ , where  $d$  is the distance from the transmitter. Calculate the SIR at Rx. Express your answer in decibels.
  - Now consider a directional antenna at Rx while the transmit antennas remain omnidirectional. The antenna at Rx is directed toward the transmitter Tx<sub>2</sub> and the antenna gain is 6 dB. In the direction of Tx<sub>1</sub> the effective antenna gain is  $-3$  dB. Now recalculate the SIR. Express your answer in dB.
- Thirty five percent of the power input to an antenna is lost as heat, what is the loss of the antenna in dB.
  - Only 65% of the power input to an antenna is radiated with the rest lost to dissipation in the antenna, what is the gain of the antenna in dB? (This is not the antenna gain.)
  - The efficiency of an antenna is 66%. If the power input to the antenna is 10 W what is the power radiated by the antenna in dBm?
  - An antenna with an input of 1 W operates in free space and has an antenna gain of 12 dBi. What is the maximum power density at 100 m from the antenna?
  - A transmitter has an antenna with an antenna gain of 10 dBi, the resistive losses of the antenna are 50%, and the power input to the antenna is 1 W. What is the EIRP in watts?
  - A transmitter has an antenna with an antenna gain of 20 dBi, the resistive losses of the antenna are 50%, and the power input to the antenna is 100 mW. What is the EIRP in watts?
  - An antenna with an antenna gain of 8 dBi radiates 6.67 W. What is the EIRP in watts? Assume

- that the antenna is 100% efficient.
12. An antenna has an antenna gain of 10 dBi and a 40 W input signal. What is the EIRP in watts?
  13. An antenna with 5 W of input power has an antenna gain of 20 dBi and an antenna efficiency of 25% and all of the loss is due to resistive losses in the antenna. [Parallels Example 4.3]
    - (a) How much power in dBm is lost as heat in the antenna?
    - (b) How much power in dBm is radiated by the antenna?
    - (c) What is the EIRP in dBW?
  14. An antenna with an efficiency of 50% has an antenna gain of 12 dBi and radiates 100 W. What is the EIRP in watts?
  15. An antenna with an efficiency of 75% and an antenna gain of 10 dBi. If the power input to the antenna is 100 W,
    - (a) what is the total power in dBm radiated by the antenna?
    - (b) what is the EIRP in dBm?
  16. On a resonant antenna a large current is established by creating a standing wave. The current peaking that thus results establishes a strong electric field (and hence magnetic field) that radiates away from the antenna. A typical dipole loses 15% of the power input to it as resistive ( $I^2 R$ ) losses and has an antenna gain of 10 dBi measured at 50 m. Consider a base station dipole antenna that has 100 W input to it. Also consider that the transmitted power density falls off with distance  $d$  as  $1/d^3$ . Hint, calculate the power density at 50 m. [Parallels Example 4.2]
    - (a) What is the input power in dBm?
    - (b) What is the power transmitted in dBm?
    - (c) What is the power density at 1 km? Express your answer as  $W/m^2$ .
    - (d) What is the power captured by a receive antenna (at 1 km) that has an effective antenna aperture of  $6\text{ cm}^2$ ? Express your answer in first dBm and then watts.
    - (e) If the background noise level captured by the antenna is 1 pW, what is the SNR in decibels? Ignore interference that comes from other transmitters.
  17. A communication system operating at 2.5 GHz includes a transmit antenna with an antenna gain of 12 dBi and a receive antenna with an effective aperture area of  $20\text{ cm}^2$ . The distance between the two antennas is 100 m.
    - (a) What is the antenna gain of the receive antenna?
    - (b) If the input to the transmit antenna is 1 W, what is the power density at the receive antenna if the power falls off as  $1/d^2$ , where  $d$  is the distance from the transmit antenna?
    - (c) Thus what is the power delivered at the output of the receive antenna?
  18. Consider a point-to-point communication system. Parabolic antennas are mounted high on a mast so that ground effects do not exist, thus power falls off as  $1/d^2$ . The gain of the transmit antenna is 20 dBi and the gain of the receive antenna is 15 dBi. The distance between the antennas is 10 km. The effective area of the receive antenna is  $3\text{ cm}^2$ . If the power input to the transmit antenna is 600 mW, what is the power delivered at the output of the receive antenna?
  19. Consider a 28 GHz point-to-point communication system. Parabolic antennas are mounted high on a mast so that ground effects do not exist, thus power falls off as  $1/d^2$ . The gain of the transmit antenna is 20 dBi and the gain of the receive antenna is 15 dBi. The distance between the antennas is 10 km. If the power output from the receive antenna is 10 pW, what is the power input to the transmit antenna?
  20. An antenna has an effective aperture area of  $20\text{ cm}^2$ . What is the antenna gain of the antenna at 2.5 GHz?
  21. An antenna operating at 28 GHz has an antenna gain of 50 dBi. What is the effective aperture area of the antenna?
  22. A 15 GHz receive antenna has an antenna gain of 20 dBi. If the power density at the receive antenna is  $1\text{ nW/cm}^2$ , what is the power at the output of the antenna? [Parallels Example 4.7]
  23. A microstrip patch antenna operating at 2 GHz has an efficiency of 66% and an antenna gain of 8 dBi. The power input to the antenna is 10 W.
    - (a) What is the power, in dBm, radiated by the antenna?
    - (b) What is the equivalent isotropic radiated power (EIRP) in watts?
    - (c) What is the power density, in  $\mu\text{W/m}^2$ , at 1 km if ground effects are ignored?
    - (d) Because of multipath effects, the power density drops off as  $1/d^4$ , where  $d$  is distance. What is the power density, in  $\text{nW/m}^2$ , at 1 km if the power density is  $100\text{ mW/m}^2$  at 10 m from the transmit antenna?
  24. A communication system operating at 10 GHz uses a microstrip patch antenna as a transmit antenna and a dipole antenna as a receive antenna. The transmit antenna is directly connected to the

- transmitter and the output power of the transmitter is 30 W. The transmit antenna has an antenna gain of 9 dBi and an antenna efficiency of 60%. The receive antenna has an antenna gain of 2 dBi and a radiation efficiency of 80%. The receive antenna is connected to a receiver by a 10 m long cable with a loss of 0.1 dB/m. The link between the transmit and receive antenna is sufficiently elevated that ground effects and multipath effects are insignificant.
- What is the output power of the transmitter in dBm?
  - What is the EIRP in dBm?
  - The transmitted power will drop off as  $1/d^n$  ( $d$  is distance). What is  $n$ ?
  - What is the peak power density in  $\mu\text{W}/\text{m}^2$  at 1 km?
  - What is the effective aperture size of the receive antenna in  $\text{m}^2$ ?
  - If the radiated power density at the receive antenna is  $1 \mu\text{W}/\text{m}^2$ , what is the signal power at the output of the receive antenna in dBm?
  - What is the total cable loss in dB?
  - What is the power presented to the receiver in dBm?
25. A communication system operating at 10 GHz uses a microstrip patch antenna as a transmit antenna and a dipole antenna as a receive antenna. The transmit antenna is connected to the transmitter by a 20 m long cable with a loss of 0.2 dB/m and the output power of the transmitter is 30 W. The transmit antenna has an antenna gain of 9 dBi and an antenna efficiency of 60%. The link between the transmit and receive antenna is sufficiently elevated that ground effects and multipath effects are insignificant.
- What is the output power of the transmitter in dBm?
  - What is the cable loss between the transmitter and the antenna?
  - What is the total power radiated by the transmit antenna in dBm?
  - What is the power lost in the antenna as resistive losses and spurious radiation? Express your answer in dBm.
  - What is the EIRP of the transmitter in dBm?
  - The transmitted power will drop off as  $1/d^n$  ( $d$  is distance). What is  $n$ ?
  - What is the peak power density in  $\mu\text{W}/\text{m}^2$  at 1 km?
26. Stacked dipole antennas are often found at the top of cellphone masts, particularly for large cells and operating frequencies below 1 GHz. These antennas have an efficiency that is close to 90%. Consider an antenna that has 40 W of input power, an antenna gain of 10 dBi, and transmits a signal at 900 MHz.
- What is the EIRP in watts?
  - If the power density drops as  $1/d^3$ , where  $d$  is the distance from the transmit tower, what is the power density at 1 km if the power density is  $100 \text{ mW}/\text{m}^2$  at 10 m?
27. Consider an 18 GHz point-to-point communication system. Parabolic antennas are mounted on masts and the LOS between the antennas is just above the tree line. As a result, power falls off as  $1/d^3$ , where  $d$  is the distance between the antennas. The gain of the transmit antenna is 20 dBi and the gain of the receive antenna is 15 dBi. The antennas are aligned so that they are in each other's main beam. The distance between the antennas is 1 km. The transmit antenna is driven by a power amplifier with an output power of 100 W. The amplifier drives a coaxial cable that is connected between the amplifier and the transmit antenna. The cable loses 75% of its power due to resistive losses. On the receive side, the receive antenna is directly connected to a mast-head amplifier with a gain of 10 dB and then a short cable with a loss of 3 dB before entering the receive base station.
- Draw the signal path.
  - What is the loss and gain of the transmitter coaxial cable in decibels?
  - What percentage of the power input to the receive coaxial cable is lost in the receive cable?
  - Express the power of the transmit amplifier in dBW and dBm.
  - What is the propagation loss in decibels?
  - Determine the total power in watts delivered to the receive base station.
28. Consider a point-to-point communication system. Parabolic antennas are mounted high on a mast so that ground effects are minimal. Thus power density falls off as  $1/d^{2.3}$ , where  $d$  is the distance from the transmitter. The gain of the transmit antenna is 15 dBi and the gain of the receive antenna is 12 dBi. These antenna gains are normalized to a distance of 1 m. The distance between the antennas is 15 km. The output power of the receive antenna must be 1 pW. The RF frequency is 2 GHz; treat the antennas as lossless.
- What is the received power in dBm?
  - What is the path loss in decibels?
  - What is the link loss in decibels?
  - Using the link loss, calculate the input power,  $P_T$ , of the transmitter. Express the answer in dBm.

- (e) What is the aperture area of the receiver in square meters?
- (f) Determine the radiated power density at the receiver in terms of the transmitter input power. That is, if  $P_T$  is the power input to the transmit antenna, determine the power density,  $P_D$ , at the receive antenna where  $P_D = xP_T$ . What is  $x$  in units of  $\text{m}^{-2}$ ?
- (g) Using the power density calculation and the aperture area, calculate  $P_T$  in watts.
- (h) What is  $P_T$  in dBm? This should be the same as the answer you calculated in (d).
- (i) What is the total power radiated by the transmit antenna in dBm?
29. Two identical antennas are used in a point-to-point communication system, each having a gain of 50 dBi. The system has an operating frequency of 28 GHz and the antennas are at the top of masts 100 m tall. The RF link between the antennas consists only of the direct line-of-sight path.
- (a) What is the effective aperture area of each antenna?
- (b) How does the power density of the propagating signal rolloff with distance.
- (c) If the separation of the transmit and receive antennas is 10 km, what is the path loss in decibels?
- (d) If the separation of the transmit and receive antennas is 10 km, what is the link loss in decibels?
30. A transmitter and receiver operating at 2 GHz are at the same level, but the direct path between them is blocked by a building and the signal must diffract over the building for a communication link to be established. This is a classic knife-edge diffraction situation. The transmit and receive antennas are each separated from the building by 4 km and the building is 20 m higher than the antennas (which are at the same height). Consider that the building is very thin. It has been found that the path loss can be determined by considering loss due to free-space propagation and loss due to diffraction over the knife edge.
- (a) What is the additional attenuation (in decibels) due to diffraction?
- (b) If the operating frequency is 100 MHz, what is the attenuation (in decibels) due to diffraction?
- (c) If the operating frequency is 10 GHz, what is the attenuation (in decibels) due to diffraction?
31. A hill is 1 km from a transmit antenna and 2 km from a receive antenna. The receive and transmit antennas are at the same height and the hill is 20 m above the height of the antennas. What is the additional loss caused by diffraction over the top of the hill? Treat the hill as a knife-edge and the operating frequency is 1 GHz.
32. A 1 GHz point-to-point link has two major transmission paths. One is a LOS path and the other includes reflection from the ground so that the power density of the transmitted signal rollsoff as  $1/d^{2.5}$  where  $d$  is the distance from the transmit antenna. At 10 m the power density from the transmit antenna is  $100 \text{ mW/m}^2$ .
- (a) What is the power density at 1 km.
- (b) If the receive antenna has an antenna gain of 30 dBi, what is the effective aperture area of the receive antenna?
- (c) What is the power of the signal at the output of the receive antenna?
33. Two identical antennas are used in a point-to-point communication system, each having a gain of 30 dBi. The system has an operating frequency of 14 GHz and the antennas are at the top of masts 100 m tall. The RF link between the antennas consists only of the direct LOS path.
- (a) What is the effective aperture area of each antenna?
- (b) How does the power density of the propagating signal rolloff with distance?
- (c) If the separation of the transmit and receive antennas is 10 km, what is the path loss? Ignore atmospheric loss.
34. The three main cellular communication bands are centered around 450 MHz, 900 MHz, and 2 GHz. Compare these three bands in terms of multipath effects, diffraction around buildings, object (such as a wall) penetration, scattering from trees and parts of trees, and ability to follow the curvature of hills. Complete the table below with the relative attributes: high, medium, and low.
- | Characteristic      | 450 MHz | 900 MHz | 2 GHz |
|---------------------|---------|---------|-------|
| Multipath           |         |         |       |
| Scattering          |         |         |       |
| Penetration         |         |         |       |
| Following curvature |         |         |       |
| Range               |         |         |       |
| Antenna size        |         |         |       |
| Atmospheric loss    |         |         |       |
35. Describe the difference in multipath effects in a central city area compared to multipath effects

- in a desert. Your description should be approximately 4 lines long and not use a diagram
36. Wireless LAN systems can operate at 2.4 GHz, 5.6 GHz, 40 GHz and 60 GHz. Contrast with explanation the performance of these schemes inside a building in terms of range.
  37. At 60 GHz the atmosphere strongly attenuates a signal. Discuss the origin of this and indicate an advantage and a disadvantage.
  38. Short answer questions. Each part requires a short paragraph of about five lines and a figure, where appropriate, to illustrate your understanding.
    - (a) Cellular communications systems use two frequency bands to communicate between the basestation and the mobile unit. The bands are generally separated by 50 MHz or so. Which band (higher or lower) is used for the downlink from the basestation to the mobile unit and what are the reasons behind this choice?
    - (b) Describe at least two types of interference in a cellular system from the perspective of a mobile handset.
  39. The three main cellular communication bands are centered around 450 MHz, 900 MHz, and 2 GHz. Compare these three bands in terms of multipath effects, diffraction around buildings, object (such as a wall) penetration, scattering from trees and parts of trees, and the ability to follow the curvature of hills. Use a table and indicate the relative attributes: high, medium, and low. WAS 10(C)
  40. Describe Rayleigh fading in approximately 4 lines and without using a diagram.
  41. In several sentences and using a diagram describe Rayleigh fading and the impact it has on radio communications.
  42. A transmit antenna and a receive antenna are separated by 1 km and operate at 1 GHz. What is the radius of the first Fresnel zone at 0.5 km from each antenna? [Parallels Example 4.8]
  43. A transmit antenna and a receive antenna are separated by 40 km and operate at 10 GHz. [Parallels Example 4.8]
    - (a) What is the radius of the first Fresnel zone at the midpoint between the antennas?
    - (b) What is the radius of the second Fresnel zone?
    - (c) To ensure LOS propagation, what should the clearance be from the direct line between the antennas and obstructions such as hills?
  44. A transmitter and receiver operate at 100 MHz, are at the same level, and are separated by 4 km. The signal must diffract over a building half way between the antennas that is 20 m higher than the direct path between the antennas. What is the attenuation (in decibels) due to diffraction?
  45. A transmitter and receiver operate at 10 GHz, are at the same level, and are 4 km apart. The signal must diffract over a building that is half way between the antennas and is 20 m higher than the line between the antennas. What is the attenuation (in dB) due to diffraction?
  46. The path from a transmit antenna to a receive antenna is elevated so that ground and multipath effects are insignificant. The power radiated by the transmit antenna drop off as  $1/d^n$  where  $d$  is distance, what is  $n$ ?
  47. A communication system has a power density roll-off of  $1/d^{2.5}$  between a transmit antenna and a mobile receive antenna which are separated by 10 km. At 10 m from the transmit antenna, the power density is  $10 \text{ W/m}^2$ . What is the power density at the receive antenna? [Parallels Example 4.5]
  48. A 900 MHz communication system uses a transmit antenna with an antenna gain  $G_T$  of 3 dB and a receive antenna with an antenna gain  $G_R$  of 0 dB. If the distance between the antennas is 200 m, what is the link loss from the input to the transmit antenna and the output of the receive antenna if the power density reduces as  $1/d^{2.5}$ ? [Parallels Example 4.6]
  49. A transmitter has a power density of  $100 \text{ mW/m}^2$  at a distance of 1 m from the transmitter. The power density falls off as 33 dB per decade of distance. At what distance from the transmit antenna will the power density reach  $1 \mu\text{W}\cdot\text{m}^{-2}$ ? [Parallels Example 4.7]
  50. Describe at least two types of interference in a cellular system from the perspective of a mobile handset.
  51. In a cellular system, a signal is intentionally transmitted from a base station nominally located in the center of a cell to a mobile unit in the same cell. However, nearby transmitters using the same channel cause interference. In Figure 4-27, a mobile unit is located at the edge of a cell and uses frequency channel A. Many nearby transmitters also operate using channel A and the six nearest transmitters can be considered as causing significant interference. Consider that the mobile unit is a distance  $r$  from its cell's transmitter along the line connecting two

- channel A base stations, that the transmitters all operate at the same power level, and that the distance between base stations operating using channel A is  $3r$ . This three-cell cluster operates in a suburban area and the power density drops off with distance  $d$  as  $1/d^{3.5}$  due to multipath effects. What is the SIR at the mobile unit? Express your answer in decibels. [Parallels Example 4.9]
52. A cellular system uses a three-cell cluster.
- Treating cells as equal sized hexagons with towers in the center of each cell, draw the cell map including all cells within 3 cell diameters of the main channel. Label the main cell and other cells using the same frequencies as A.
  - Consider a mobile unit at the edge of a cell transmitting the intended signal from the tower at the center of that cell. Identify the interfering towers.
- If ground effects and multipath effects are negligible, what is the power roll-off factor if the distance between a tower and the mobile unit is  $d$ ? That is, what is  $n$  if power falls off as  $1/d^n$ ?
  - If trisector antennas are used, identify the interfering cells and approximately determine the improvement in SIR compared to using a nonsectorized antenna? You do not need to do detailed calculations.
53. Describe trisector antennas in 4 lines and without using a diagram.

### 4.12.1 Exercises By Section

†challenging, ‡very challenging

- |       |   |        |   |
|-------|---|--------|---|
| §14.3 | 1 <sup>†</sup>  | 19, 22 | 35 <sup>†</sup> , 36, 37, 38 <sup>†</sup> , 39, 40, 41,   |
| §14.5 | 2 <sup>†</sup> , 3 <sup>†</sup> , 4 <sup>‡</sup> , 5, 6, 7 <sup>†</sup> , 8 <sup>†</sup> , 9 <sup>†</sup> , 10, | §14.6  | 23 <sup>‡</sup> , 24 <sup>‡</sup> , 25 <sup>‡</sup> , 26 <sup>‡</sup> , 27 <sup>‡</sup> , 28 <sup>‡</sup> , |
|       | 11, 12, 13, 14, 15, 16 <sup>‡</sup> , 17, 18,   |        | 29 <sup>†</sup> , 30 <sup>†</sup> , 31 <sup>†</sup> , 32 <sup>‡</sup> , 33 <sup>†</sup> , 34 <sup>‡</sup> , |
|       |   | §14.8  | 50 <sup>†</sup> , 51 <sup>‡</sup> , 52 <sup>‡</sup> , 53 <sup>†</sup>                                       |

### 4.12.2 Answers to Selected Exercises

- |       |             |       |                        |    |         |
|-------|-------------|-------|------------------------|----|---------|
| 2(e)  | 53.23 dBm   | 25(e) | 49.77 dBm              | 42 | 8.66 m  |
| 23(b) | 63.1 W      | 27(f) | 708 pW                 | 51 | 7.33 dB |
|       | 17 0.251 μW |       | 28 107.65 dB           |    |         |
|       | 18 14.3 pW  |       | 32 7.16 m <sup>2</sup> |    |         |

# RF Systems

5.1	Introduction .....	173
5.2	Simplex and Duplex Operation .....	174
5.3	Cellular Communications .....	178
5.4	Multiple Access Schemes .....	182
5.5	Spectrum Efficiency .....	185
5.6	Processing Gain .....	187
5.7	Cellular Phone Systems .....	196
5.8	Early Generations of Radio .....	197
5.9	3G, Third Generation: Code Division Multiple Access (CDMA) .	201
5.10	4G, Fourth Generation Radio: Long-Term Evolution .....	208
5.11	5G, Fifth Generation Radio .....	219
5.12	6G, Sixth Generation Radio .....	223
5.13	Radar Systems .....	224
5.14	Summary .....	228
5.15	References .....	229
5.16	Exercises .....	231

## 5.1 Introduction

In the history of wireless communication, radar, and sensor systems there have been many standards and different types of systems. Following the development of the main cellular radio systems is a good proxy for the evolution of nearly all RF systems. Mobile radio up to and including 1G cellular radio was essentially analog with only simple signaling using tones or FSK modulation. There were many incompatible mobile radio and 1G systems as little thought was given to worldwide interoperability and radio companies were content with proprietary standards.

The situation continued with 2G cellular radio as there were many incompatible 2G systems with each only able to support one modulation method. Most of the radio functionality was in the analog hardware but VLSI was starting to become mature and it was possible to do limited error correction coding. At the basestation there was enough computing power to dynamically manage interference and implement simple system-level optimization. In the evolution to 3G there were many proposals all exploiting the increasing capability of VLSI. With the introduction of 3G,

communication providers enforced the adoption of a single world-wide standard. Perhaps 'single' is a stretch as there were still several variations of the 3G implementation but a core set of standards supporting worldwide connectivity was adopted.

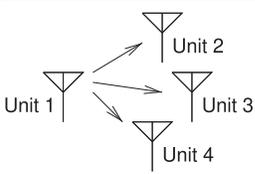
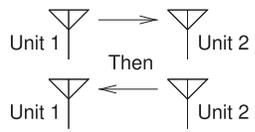
The path forward was becoming unmanageable so following the introduction of 3G an international consortium focused on developing a strategy to support cellular communications that was more capable and upwards compatible.

For the first time, with 4G a large number of modulation methods were supported with most of the modulation and demodulation functionality implemented in a DSP unit called the baseband processor. Changing modulation and demodulation formats was a simple procedure of running different code. This is software-defined radio with analog circuitry implementing just the translation from an analog baseband signal, really a modulated signal on a low-frequency intermediate frequency carrier, to the high-frequency RF signal. Modulation by the baseband data produced the baseband analog signal with the baseband processor implementing a digital version of an analog modulation format. With 5G additional capability is provided with the 4G standard becoming a subset of the 5G standard.

An understanding of systems is required to understand the specifications for RF hardware. This is particularly important as the actual performance required of hardware is often not directly related to the specifications developed by system designers. For example, one of the most important characteristics of digital radio systems is the bit error rate (BER). The BER is a quantity that cannot be determined until the components of a system are integrated. Thus in the design of subsystems, indirect measures such as intermodulation distortion (IMD) (referring to the generation of spurious signals when discrete tones are applied to a subsystem) are specified. The relationship between IMD and BER is weak. Clearly higher IMD tends to indicate a higher BER for the same technology, but the relative performance of different technologies cannot be evaluated this way. Thus an essential system design problem is developing sufficient and tractable criteria that enable subsystems to be locally designed and optimized, leading to an overall optimized system.

## **5.2 Broadcast, Simplex, Duplex, Diplex, and Multiplex Operations**

One would think that defining these terms would be a simple matter. However, there have been different conventions in various segments of the telecommunications industry. Now that telecommunications are converging, universal definitions that accommodate all earlier uses are not possible. There is now standardization by the International Telecommunications Union (ITU) of the terms broadcast, simplex, duplex, diplex, and multiplex, as they relate to wireless communications [1, 2] (see Table 5-1). National Standards, however, do not need to conform to the ITU definitions. Examples of slightly different definitions come from the Telecom Glossary established as part of an American National Standard (ANS) [3]. There is greater specificity to radio in the ITU usage of the terms. In general the ITU definition should be used if the context is radio communications, but be careful if the interpretation of the terms is critical.

Schematic	ANS definition (2005) [3]	ITU definition (2004) [1]	ITU: obsolete usage [1]
	Broadcast	Broadcast	Simplex
	Simplex	Simplex	Halfduplex
	Duplex	Duplex	Fullduplex

**Table 5-1:** Definitions of operating modes of wireless communication systems. American National Standard (ANS) definition, International Telecommunication Union (ITU) definition, and deprecated (i.e., obsolete) usage

### 5.2.1 International Telecommunications Union Definitions

The ITU is the agency of the United Nations that coordinates the shared global use of the radio spectrum and establishes worldwide standards.

**Broadcast operation** refers to one-way communication in which there is only one transmitter and at least one, and perhaps more, receivers. The ITU defines broadcasting as [1, 2]:

a form of unidirectional telecommunication intended for a large number of users having appropriate receiving facilities, and carried out by means of radio or cable networks.

**Simplex operation** is defined by the ITU as [1]

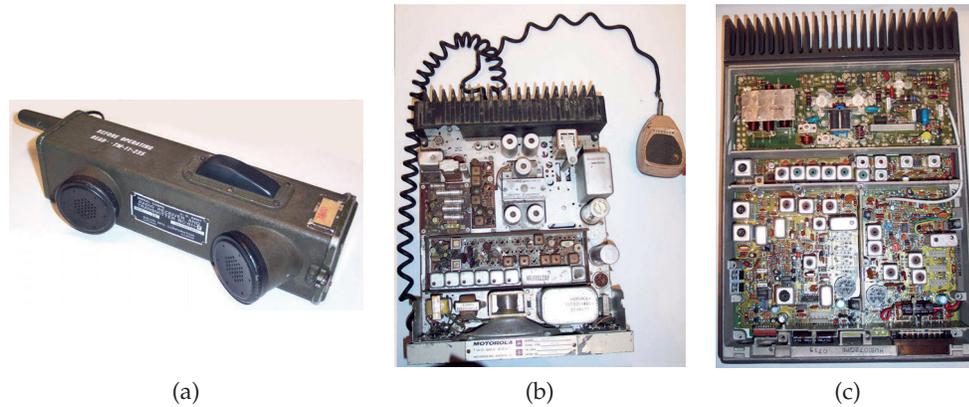
designating or pertaining to a method of operation in which information can be transmitted in either direction, but not simultaneously, between two points.

So simplex operation is when there is communication from a first unit to a second unit. This is followed by a hand-off procedure and then communication from the second unit to the first. So the push-to-talk Family Radio Service (FRS) is an example of simplex communication as two people on opposite ends of the link cannot talk simultaneously. Simplex operation may use either one or two frequencies. That is, the forward link (from user 1 to user 2) may use the same frequency channel as the reverse link (from user 2 to user 1). However, simplex operation may use different frequency channels for the forward and the reverse link. The ITU defines “half duplex” as being synonymous with simplex, but the use of the term “half duplex” is deprecated (not recommended) [1, 2]. An early radio using simplex operation is shown in Figure 5-1(a).

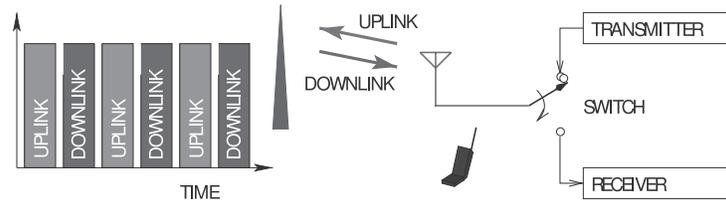
**Duplex Operation** is defined by the ITU as [2]

the operating method in which transmission is possible simultaneously in both directions of a telecommunication channel.

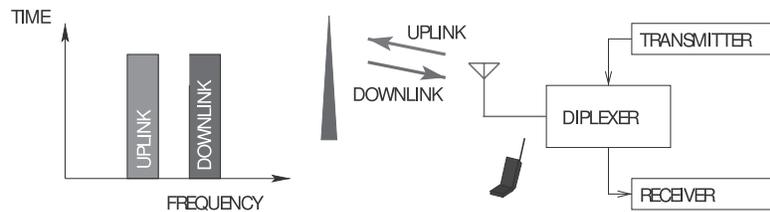
Simultaneous two-way communication uses duplexing. In analog radio the links must be simultaneous and so operation requires two frequencies. But in digital radio with compression and expansion of digitized speech and



**Figure 5-1:** Early radios: (a) walkie-talkie with a push-to-talk (PTT) switch on top and using simplex communication; (b) Motorola business dispatcher two-way radio operating at 33.220 MHz designed in the 1960s as a dash mount unit; and (c) the Mitrek two-way radio designed by Motorola in 1977. The radio has two PC boards and was crystal controlled with a channel scanning control head. The radio was trunk mounted, with the control head, microphone, and speaker mounted under the dash board. The units in (b) and (c) are approximately 20 cm long.



(a) Time-division duplex (TDD)

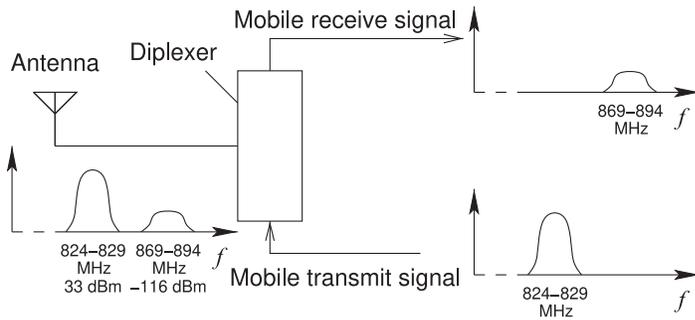


(b) Frequency-division duplex (FDD)

**Figure 5-2:**  
Duplex schemes.

transmission of data, duplex operation can involve transmission of data in time slots that may not be simultaneous in each direction. If the time slots are short enough, say less than ten milliseconds, the effect will be that two-way voice communication is simultaneous. This is still called duplex operation. Early examples of radios using duplex operation are shown in Figure 5-1(b and c). A modern example is the cell phone.

The two predominant duplex schemes are **frequency-division duplex (FDD)** and **time-division duplex (TDD)**, as illustrated in Figure 5-2. In these diagrams the mobile unit is shown communicating with a basestation, in which case transmission from the mobile unit to the basestation is called the **uplink (UL)**, or less commonly the **reverse link (RL)**, or **reverse path**.



**Figure 5-3:** A diplexer which separates receive and transmit signals. A diplexer is a type of duplexer. However not all duplexers are diplexers.

Communication from the basestation and received by the mobile unit is called the **downlink**, (**DL**), **forward link** (**FL**), or **forward path**. In FDD it is necessary to use a filter to separate the uplink and downlink signals, as shown in Figure 5-3, as the two links are in use simultaneously.

In TDD, see Figure 5-2(a), the uplink and downlink are separated in time and a switch connects the antenna to the transmitter to send the uplink and then to the receiver to receive the downlink. The sequence is then repeated.

In FDD the uplink and downlink are at different frequencies and the transmitter and receiver are connected permanently to the antenna, see Figures Figure 5-2(b) and 5-3. FDD requires a diplex filter, which is a special filter with three ports that looks like a lowpass filter (usually for a handset) for the uplink when the transmitter is connected to the antenna (the uplink is generally at a lower frequency than the downlink) and a highpass filter for the downlink when the receiver is connected to the antenna.

### 5.2.2 Duplex Versus Diplex

The term duplex refers to a way of handling two communication channels. It derives from the term **multiplexing** (or **MUXing**) defined as [3]

the combining of two or more information channels onto a common transmission medium. Note: In electrical communications, the two basic forms of multiplexing are time-division multiplexing (TDM) and frequency-division multiplexing (FDM). Synonym: multiplex.

**Duplexing**, or **duplex operation**, is used when there are two channels and in radio communications it nearly always refers to one transmit channel and one receive channel.

**Diplex operation** is defined as [3]:

the sharing of one common element, such as a single antenna or channel, for transmission or reception of two simultaneous, independent signals on two different frequencies. Note: An example of diplex operation is the use of one antenna for two radio transmitters on different frequencies.

A **diplexer** is defined as [3]

a three-port frequency-dependent device that may be used as a separator or a combiner of signals.

So in the context of cellular communications, a diplexer is a filter (see

Figure 5-3). It is a three-port filter that separates the transmitted and received signals, which are at different frequencies.

### 5.3 Cellular Communications

Cellular communications, as the name implies, are based on the concept of cells in which a terminal unit communicates with a basestation at the center of a cell. Each cell can be relatively small and a terminal unit travels smoothly from cell to cell with a connection transferring using what is called a hand-off process. For communication in closely spaced cells to work, interference from other radios must be managed. This is facilitated using the ability to recover from errors available with error correction schemes and the use of antenna beam-forming technology.

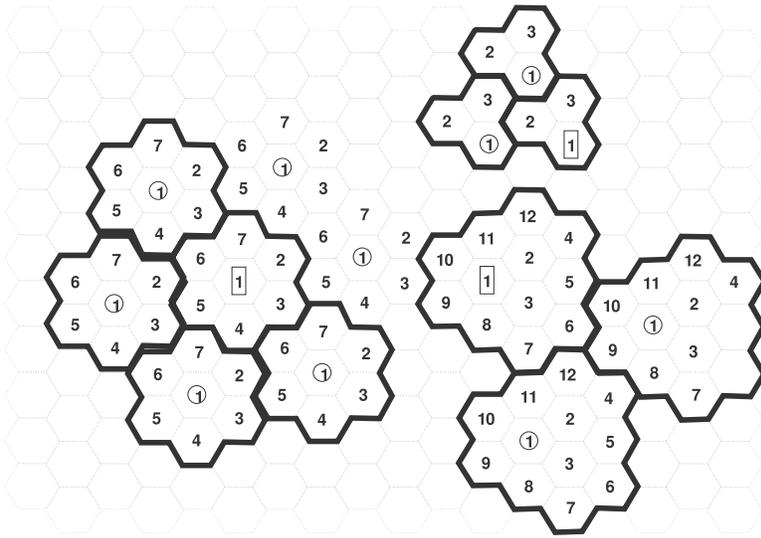
#### 5.3.1 Cellular Concept

The cellular concept was outlined in a 1947 Bell Laboratories technical memorandum [4]. It described a system of frequency reuse with small geographical cells, and this remains the key concept of cellular radio. This was elaborated on in two articles published in 1957 and 1960 [5, 6]. The first widespread cellular radio system was the **Advanced Mobile Phone System (AMPS)**, one of many 1G cellular radio systems, and was fully described by Bell Laboratories in a submission to the U.S. **Federal Communications Commission (FCC)** and in a patent filed on December 21, 1970 [7]. Bell Labs petitioned the FCC in 1958 for a frequency band around 800 MHz for a cellular system. The FCC, believing that it was better to allocate spectrum for the public good (such as radio, television, and emergency services) was reluctant to act on the petition. In 1968, pressure on the FCC became too great and an agreement was reached in principle to make frequencies available. Thus began the research and development of cellular systems in the United States. In 1961, Ericsson reorganized to address mobile radio, including cellular radio systems. Nokia did not begin developing 1G cellular systems until the 1970s. NTT was working away as well and began developing cellular radio systems in 1967 [8]. Meanwhile, in January 1969, the Bell System launched an experimental cellular radio system employing frequency reuse for the first time to achieve optimum use of a limited number of RF channels. The first commercial cellular system was launched by the Bahrain Telephone Company in May 1978 using Matsushita equipment. This was followed by the launch of AMPS by Illinois Bell and AT&T in the United States in July 1978.

In 1979 the **World Administrative Radio Conference (WARC)** allocated the 862–960 MHz band for mobile radio, leading to the FCC releasing, in 1981, 40 MHz in the 800–900 MHz band for “cellular land-mobile phone service.” The service, as defined in the original documents, is (and this is still the best definition of cellular radio)

a high capacity land mobile system in which assigned spectrum is divided into discrete channels which are assigned in groups to geographic cells covering a cellular geographic area. The discrete channels are capable of being reused within the service area.

The key attributes here are

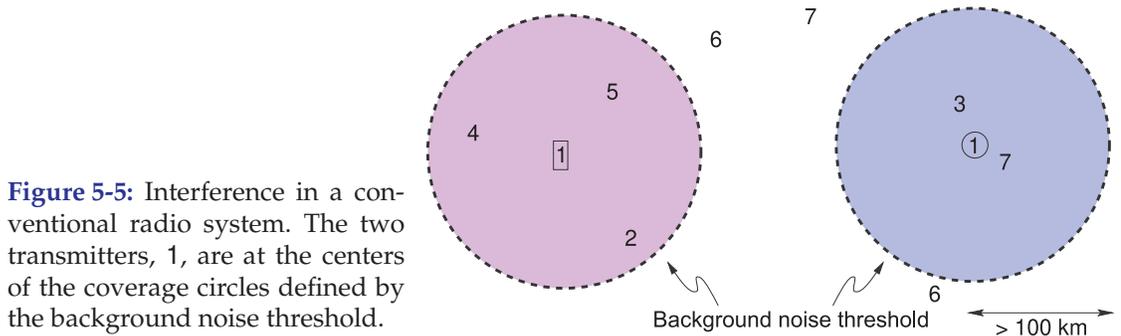


**Figure 5-4:** Cells arranged in clusters.

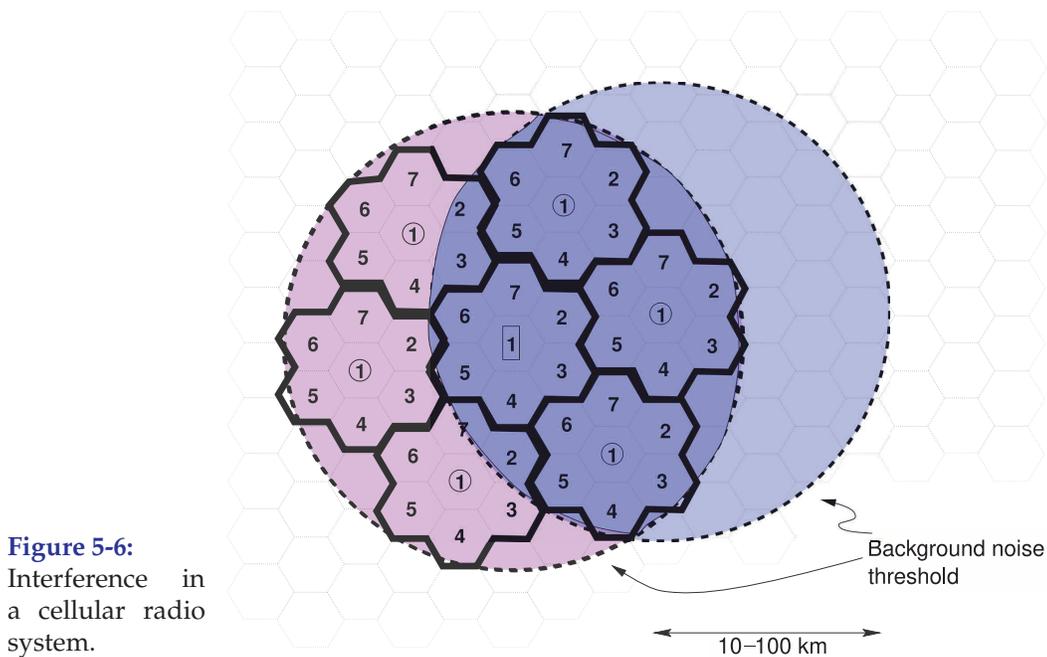
- High capacity. Prior to the cellular system, mobile radio users were not always able to gain access to the radio network and frequently access required multiple attempts.
- The concept of cells. The idea is to divide a large geographic area into cells, shown as the hexagons in Figure 5-4. The actual shape of the cells is influenced by obstructions such as hills and buildings, but the hexagonal shape is used to convey the concept of cells. The cells are arranged in clusters and the total number of channels available is divided among the cells in a cluster and the full set is repeated in each cluster. In Figure 5-4, 3-, 7-, and 12-cell clusters are shown. As will be explained later, the number of cells in a cluster affects both capacity (the fewer cells the better) and interference (the more cells per cluster, the further apart cells operating at the same frequency are, and so interference is less).
- Frequency reuse. Frequencies used in one cell are reused in the corresponding cell in another cluster. As the cells are relatively close, it is important to dynamically control the power radiated by each radio, as radios in one cell will produce interference in other clusters.

The shape of a cell depends on many factors. In a flat desert the coverage area of each basestation would be circular, so that with a cluster of cells there would be overlapping circles of coverage. (Power levels are adjusted to minimize the overlap of these circles.) Buildings, hills, lakes, etc., affect cell size. In a city, what is called the **urban canyon effect**, or **urban waveguide effect** greatly distorts cells and creates havoc in managing cellular systems [9–11]. In the urban canyon effect, good coverage extends for large distances down a street (see the inside front cover).

Achieving maximum frequency reuse is essential in achieving high capacity. In a conventional wireless system, be it broadcast or the mobile telephone service, basestations are separated by sufficient distance such that the signal levels fall below a noise threshold before the same frequencies are reused, as shown in Figure 5-5. There is clearly poor geographical use of the



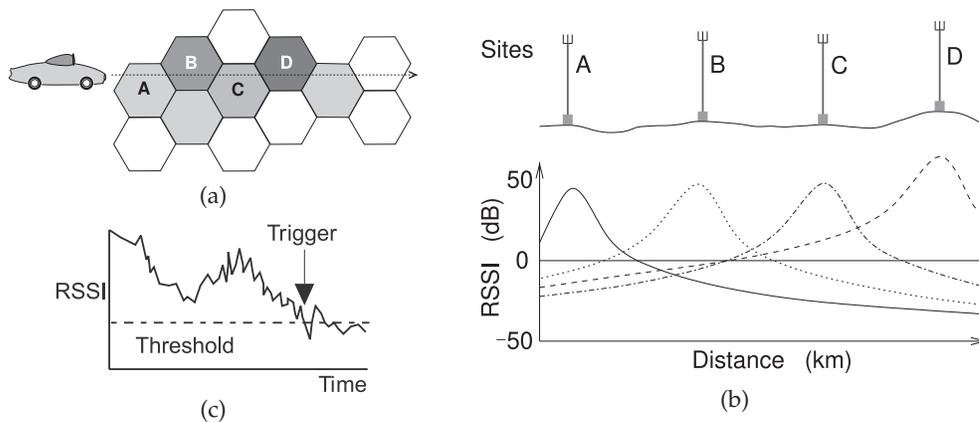
**Figure 5-5:** Interference in a conventional radio system. The two transmitters, 1, are at the centers of the coverage circles defined by the background noise threshold.



**Figure 5-6:** Interference in a cellular radio system.

spectrum here. The geographic areas could be pushed closer to each other at the expense of introducing what is called cochannel interference—a receiver could pick up transmissions from two or more basestations operating at the same frequency. This is strictly avoided so that interference in conventional radio systems results solely from background noise. In a cellular system, there is a radical departure in concept from this. Consider the interference in a cellular system as shown in Figure 5-6. The signals in corresponding cells in different clusters interfere with each other and the interference is much larger than that of the background noise. Generally only in rural areas and when mobile units are near the boundaries of cells will the background noise level be significant. Interference can also be controlled by dynamically adjusting the basestation and mobile transmit powers to the minimum acceptable level. Tolerating interference from neighboring clusters is a key concept in cellular radio.

The analog 1G AMPS system, which uses frequency modulation, has a qualitative minimum SIR of 17 dB (about a factor of 50) that was determined via subjective tests with a criterion that 75% of listeners ranked the voice quality as good or excellent. The seven-cell clustering shown in Figure 5-



**Figure 5-7:** Process of handoff: (a) movement of a mobile unit through cells; (b) received signal strength indicator (RSSI) during movement of a mobile unit through cells; and (c) RSSI of a mobile unit showing the handoff triggering event.

6 does not yield this required minimum SIR. So either a 12-cell cluster is required or directed antennas are used, as these provide enough SIR.

The 2G and 3G digital cellular systems use error correction coding and can tolerate high levels of interference and can reuse frequency channels more efficiently. Indeed, in the 3G CDMA system the tolerance to interference is so high that the concept of clustering is not required and every frequency channel is available in each cell.

In 4G and 5G cellular radio interference must be low to enable high-order modulation to be used. The high modulation efficiencies more than makes up for the reduced frequency reuse.

### 5.3.2 Personal Communication Services

The personal communication services (PCS) concept was implemented in the early 1990s and was a development in the thinking of cellular communications. In PCS, the concept is that communication is from person to person, whereas in cellular radio communication as originally conceived, it was from terminal-to-terminal. The idea is that when a call is placed, an individual is being contacted rather than a piece of hardware. One way this is achieved is by using a card, a **subscriber identification module (SIM card)**, to identify the user. A user can insert his or her SIM card into any (appropriate) handset and the handset becomes personalized. The term PCS is not commonly used now, as the concept has been incorporated in all evolved cellular systems.

### 5.3.3 Call Flow and Handoff

Mobile users can be expected to move frequently between cells and thus handoff procedures for transferring connections from one cell to the next are necessary. The triggering events that initiate handoff are shown in Figure 5-7. The main aspect is monitoring of the signal strength, the **received signal strength indicator (RSSI)**, both in the handset for the signal received from a

basestation and in the basestation for the signal received from a handset. If either of these falls below a threshold, computers in the basestation initiate a handoff procedure by polling nearby basestations for the RSSI they have for the user. If a suitable RSSI is found handoff proceeds and the other basestation takes control of the RF link to the mobile terminal.

#### 5.3.4 *Cochannel Interference*

The minimum signal detectable in conventional wireless systems is determined by the received SIR. In cellular wireless systems the dominant interference is due to other transmitters in the cell and adjacent cells. The noise that is produced in the signal band from other transmitters operating at the same frequency is called cochannel interference. The level of cochannel interference is dependent on cell placement and the frequency reuse pattern. The degree to which cochannel interference can be controlled has a large effect on system capacity. Control of cochannel interference is largely achieved by controlling the power levels at the basestation and at the mobile units.

### 5.4 Multiple Access Schemes

Many schemes have been developed to enable multiple users to share a frequency band. The simplest scheme requiring the least sophistication in channel management is the **frequency division multiple access (FDMA)** scheme, shown in Figure 5-8(a), where the numerals indicate a particular user. In FDMA the available spectrum is divided by frequency, with each user assigned a narrow frequency band that is kept for the duration of a call. This can be conveniently implemented in analog radio. In duplex operation a frequency channel is assigned for the uplink as well as the downlink. An example of where FDMA is used is AMPS, where 30 kHz is assigned to each channel. Clearly this is wasteful of spectrum, as not all of the band is used continuously.

In the first digital access technique, **time division multiple access (TDMA)**, shown in Figure 5-8(b), a bitstream is divided among a few users using the same physical channel. In the **GSM** mobile phone system,<sup>1</sup> a physical channel is divided into eight time slots and a user is allocated to each so that eight users can be supported in each physical channel. Thus the logical channels are divided in both frequency and time. In a TDMA system, the basestation transmits a continuous stream of data containing frames of time slots for multiple users. The mobile unit listens to this continuous stream, extracting and processing only the time slots assigned to it. On the reverse transmission, the mobile unit transmits to the basestation in bursts only in its assigned time slots. This is yet another complication to the RF design of TDMA systems. The RF circuitry operates in a burst mode with constraints on settling time. One of the advantages of GSM is that multiple slots can be assigned to the same user to support high data rates and time

---

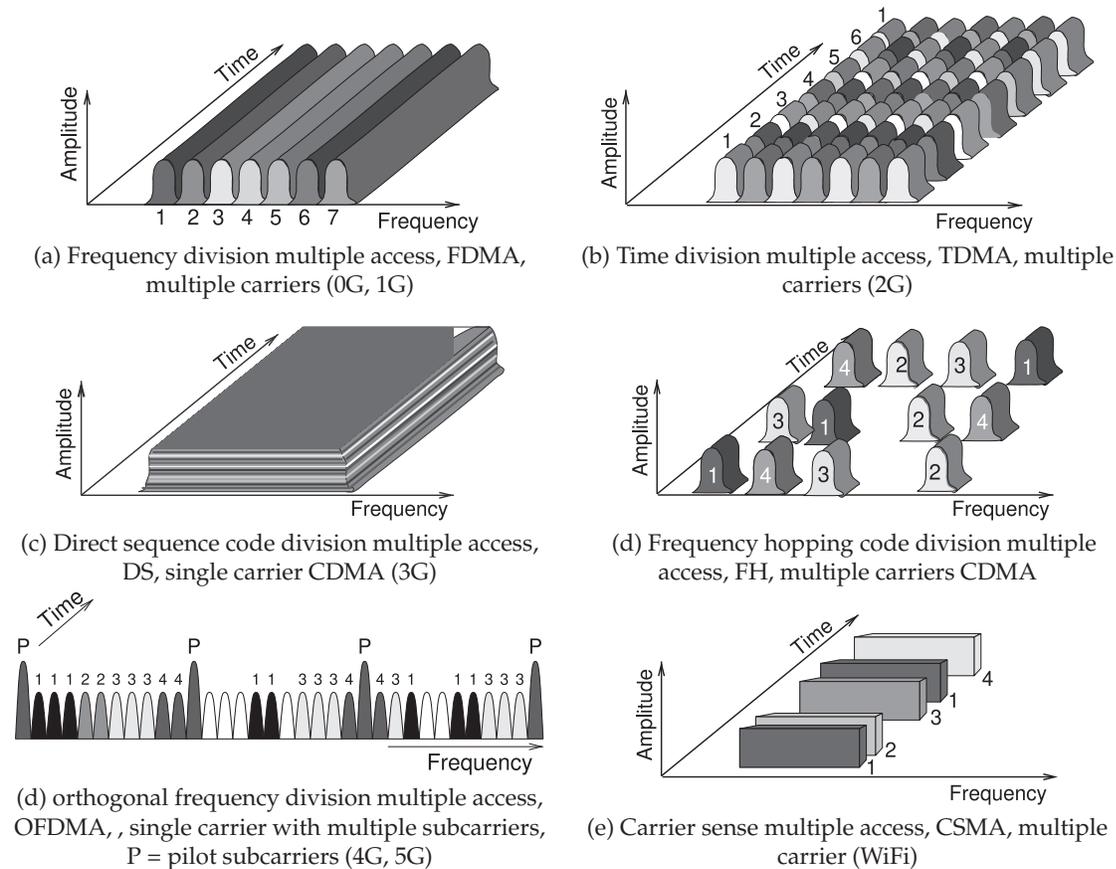
<sup>1</sup> GSM stands for the 2G Global System for Mobile Communications and was formerly known as Groupe Spécial Mobile. The system was deployed worldwide in 1991. An interesting footnote is that the GSM group began discussions leading to the system in 1982. By the mid-1980s there were many different versions of the GSM system in Europe. The European Union (EU) intervened and all member countries adopted a single standard.

slots can be skipped.

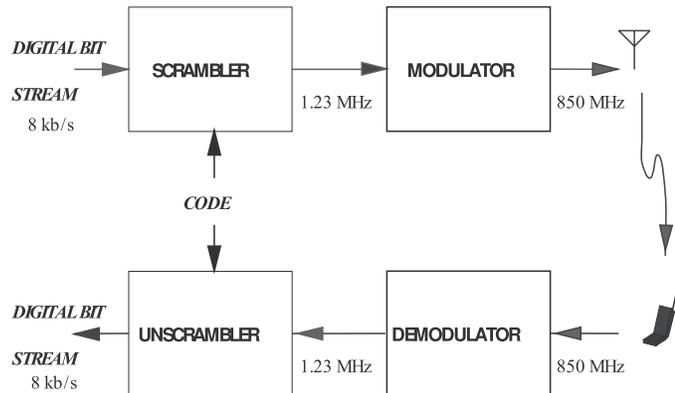
Another technique that can be used in digital radio is **spread spectrum** (SS [12]), the invention of which is attributed to Hedwig Kiesler Markey<sup>2</sup>. Most of the development of spread spectrum was secret up until the mid-1970s. There are two main types of spread-spectrum techniques used in multiuser communication: direct sequence, and frequency hopping. It is both an access technique and a way to secure communications. **Direct sequence code division multiple access (DS-CDMA)**<sup>3</sup> is shown in Figure 5-8(c). DS-CDMA mixes the baseband signal with a broadband spreading code signal to produce a broadband signal that is then used to modulate an RF signal [13–15]. This process is illustrated in Figure 5-9. The rate of the spreading code is referred to as the **chip rate** rather than the bit rate, which is reserved for the rate of the information-bearing signal. A unique code is used to

<sup>2</sup> The 1940 invention is described in [12]. The patent described a frequency-hopping scheme to render radio-guided torpedoes immune from jamming using a piano roll to hop among 88 carrier frequencies.

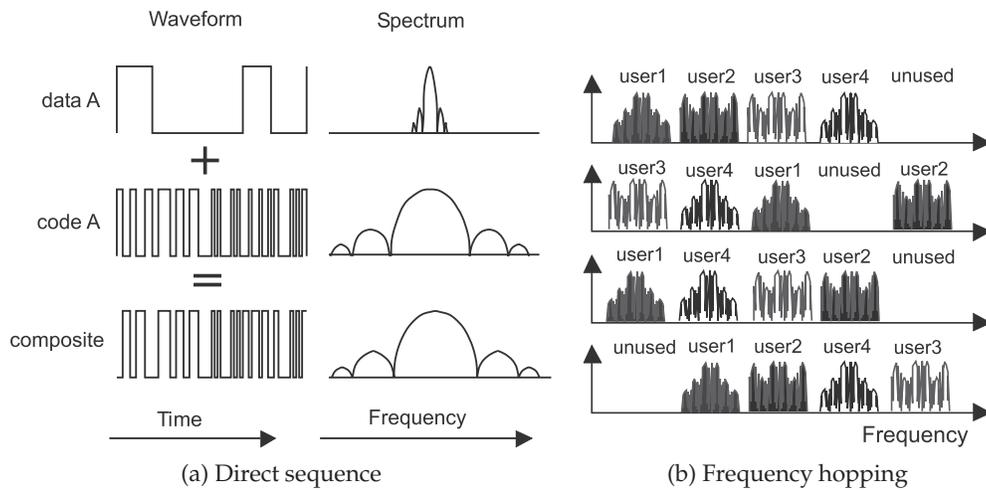
<sup>3</sup> DS-CDMA is usually referred to as just CDMA.



**Figure 5-8:** Multiple access schemes for supporting multiple users with each user indicated by a numeral except for (b), TDMA, where the numerals indicate time slots.



**Figure 5-9:** Block schematic of the operation of a CDMAOne spread-spectrum radio.



**Figure 5-10:** Basis of multiuser spread-spectrum communication.

“scramble” the original baseband signal, and the original signal can only be recovered using that particular code. The system can support many users, with each user assigned a unique code, hence the term code division multiple access (CDMA). Another interpretation of the DS-SS process is shown in Figure 5-10(a). One of the unique characteristics is that different users are using the same frequency band at the same time. CDMA is used in the 3G **wideband CDMA (WCDMA)** system.

Another form of CDMA access is the **frequency-hopping CDMA** scheme (**FH-CDMA**) shown in Figures 5-8(d) and 5-10(b). This scheme uses **frequency-hopping spread spectrum (FHSS)** to separate users. For one particular user, a code is used to determine the hopping sequence and hops can occur faster than the data rate. The signal can only be reassembled at the receiver if the hopping code is known. FH-CDMA is useful in unregulated environments, such as an ad hoc radio network, where it is not possible to coordinate multiple users. In this case, two users can transmit information in the same frequency band at the same time with error correction for occasional errors. FH-CDMA is commonly used in radios for emergency workers and by the military.

Another multiple access scheme is **carrier-sense multiple access (CSMA)**, which is used in the WiFi system (IEEE 802.11). The access scheme is illustrated in Figure 5-8(e). This scheme is also used in hostile environments where various radios cannot be well coordinated. In CSMA, users transmit at different times, but without coordination. A terminal unit listens to the channel and transmits a data packet when it is free. Collisions are unavoidable and data are lost, requiring re-transmission of data, but the terminals use a random delay before re-transmitting data. In the case of WiFi there are several channels and if collisions are excessive another channel can be selected either as the preferred start-up channel or in evolved systems under central unit control.

**Orthogonal frequency division multiple access (OFDMA)** is used in **worldwide interoperability for microwave access (WiMAX)** and in 4G and 5G cellular radio. It overcomes many of the performance limitations encountered with CSMA as assignment of effective channels is used. OFDMA builds on orthogonal frequency division multiplexing (**OFDM**) introduced in Section 5.10.2.

## 5.5 Spectrum Efficiency

The concept of spectral efficiency is important in contrasting digital radio systems. Spectral efficiency has its origins in **Shannon's theorem**, which expresses the information-carrying capacity of a channel as [16–18]

$$\hat{C} = B_c \log_2(1 + \text{SNR}), \quad (5.1)$$

where  $\hat{C}$  is the capacity in units of bits per second (bit/s),  $B_c$  is the channel bandwidth in hertz, and SNR is the signal-to-noise ratio.  $N$  is assumed to be Gaussian noise, so interference that can be approximated as Gaussian can be incorporated by adding the noise and interference powers, and then it is more appropriate to use the SIR. Thus Equation (5.1) becomes

$$\hat{C} = B_c \log_2[1 + S/(N + I)] = B_c \log_2(1 + \text{SIR}). \quad (5.2)$$

Shannon's theorem is widely accepted as the upper limit on the information-carrying capacity of a channel. So the stronger the signal, or the lower the interfering signal, the greater a channel's information-carrying capacity. Indeed, if there is no noise and no interference, the information-carrying capacity is infinite. Shannon's capacity formula indicates that increasing the interference level (lower SIR) has a more weakened effect on the decrease in capacity than may initially be expected; that is, doubling the interference level does not halve  $\hat{C}$ . This is the conceptual insight that supports the use of closely packed cells and frequency reuse, as the resulting increase in interference, and its moderated effect on capacity, is offset by having more cells and supporting more users.

Shannon's carrying capacity limit has not been reached, but today's radio systems are very close. Current systems operate with SNRs only a few decibels away from the limit [16]. Different modulation and radio schemes come closer to the limit, and two quantities will be introduced here to describe the performance of different schemes. From the capacity formula, a useful metric for the performance of modulation schemes can be defined. This is the modulation efficiency (also referred to as the **channel efficiency**,

channel spectrum efficiency, and channel spectral efficiency),

$$\eta_c = R_c/B_c, \quad (5.3)$$

where  $R_c$  (in bit/s) is the bit rate transmitted on the channel, so  $\eta_c$  has the units of bit/s/Hz. The unit is dimensionless, as hertz has the units of  $s^{-1}$ .

In a cellular system, the number of cells in a cluster must be incorporated to obtain a system metric [19]. The available channels are divided among the cells in a cluster, and a channel in one cell appears as interference to a corresponding cell in another cluster. Thus the SIR is increased and the capacity of the channel drops. System throughput increases, however, because of closely packed cells. So the system throughput is a function of the frequency reuse pattern. The appropriate system-level metric is the radio spectrum efficiency,  $\eta_r$ , which incorporates the number of cells,  $K$ , in a cluster:

$$\eta_r = \frac{R_b}{B_c K} = \frac{\eta_c R_b}{K R_c}, \quad (5.4)$$

where  $R_b$  is the bit rate of useful information  $R_c$ , ( $R_c$  is higher than  $R_b$  because of coding). Coding is used to enable error correction, assist in identifying the start and end of a packet, and also to provide orthogonality of users in some systems that overlap users as with CDMA. The units of  $\eta_r$  are bit/s/Hz/cell. The decrease in channel capacity resulting from the increased SIR associated with fewer cells in a cluster (i.e., lower  $K$ ) is more than offset by the increased system throughput.

So there are two definitions of spectral efficiency: the channel spectrum efficiency (also known as the modulation efficiency),  $\eta_c$ , which characterizes the efficiency of a modulation scheme, and the **radio spectrum efficiency**,  $\eta_r$ , which incorporates the added interference that comes from frequency reuse. Commonly both measures of efficiency are referred to as spectral efficiency, and then only the units identify which is being referred to.

One may well ask why efficiency is not expressed as a ratio of actual bit rate to Shannon's limit for a given set of conditions. There are several reasons:

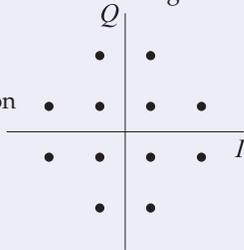
1. The historical use of bits per hertz to characterize a modulation scheme was used long before cellular systems came about.
2. Shannon's capacity limit is so high that back in the 1950s people would have been talking about extremely low efficiencies if performance was referred to the capacity limit.
3. Only additive white noise is considered, but this does not capture all types of interference, which can be multiplicative or partially correlated with the signal. Shannon's limit is not really a theoretical limit, there is no proof. In the digital communications world, much has been published about how close Shannon's capacity limit can be approached. In a direct LOS system, such as a point-to-point microwave link, the limit is now approached within a few percent. In MIMO systems, (see Section 5.10.5) the limit has been exceeded, prompting a redefinition of the limit when multiple transmit and receive antennas are used.

**EXAMPLE 5.1** Modulation Efficiency

A radio uses a modulation scheme based on 16-QAM but the four constellation points corresponding to the largest signal are not used. Consequently the distortion that would occur in RF amplifiers is reduced. Even though there are 4 bits of information per symbol for the symbols that are actually used, not every possible combination of the bits is used. Ignoring error correction coding all of the bits modulated are information bits.

**Solution:**

(a) Draw the constellation diagram.



(b) How many symbols are there?  
12.

(c) On average, how many bits of information are transmitted per symbol?

It takes 8 symbols to transmit 3 bits of information and 16 symbols to transmit 4 bits. With 12QAM 8 symbols are sent in the first clock tick interval and 4 symbols are borrowed from the second clock tick interval to provide 16 symbols combined and hence 4 bits of information. There are 8 symbols left over in the second clock tick interval and these can be used to send 3 bits. Thus over two clock tick intervals 7 bits of information are sent. Thus in one clock tick interval 3.5 bits are transmitted. So the number of bits of information sent by each symbol is 3.5 bits.

(d) What is the maximum possible modulation efficiency,  $\eta_c$ , in bit/s/Hz?

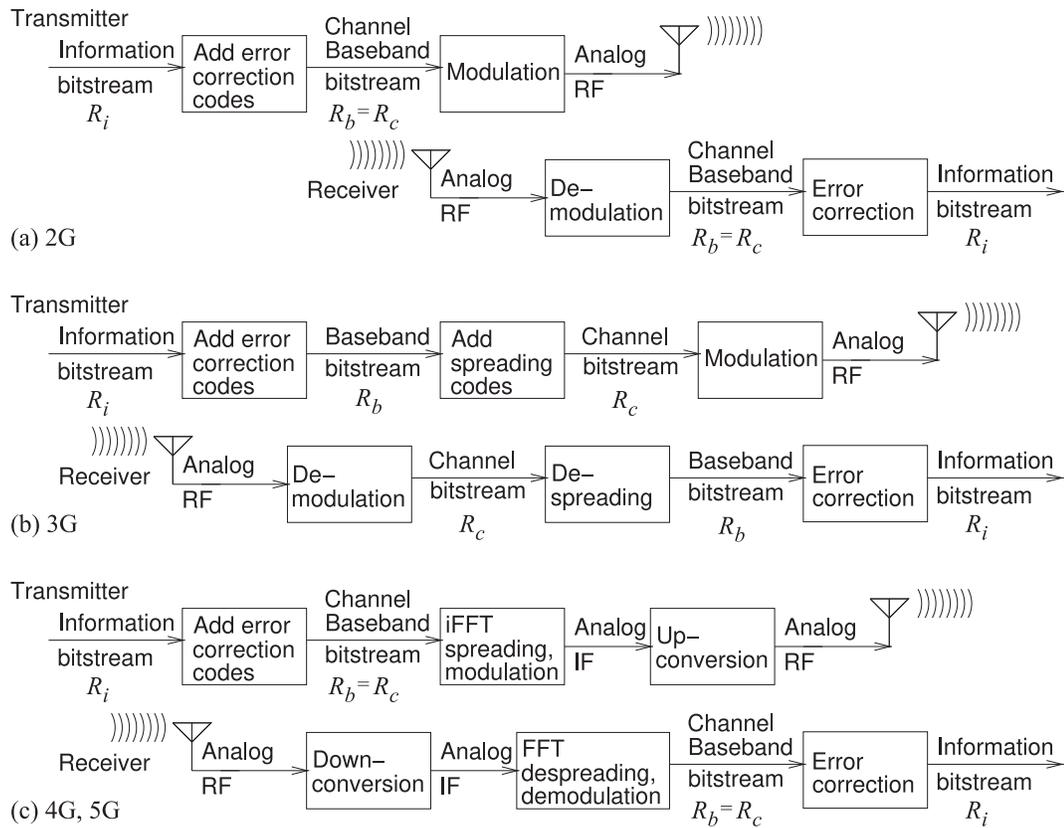
Ideally  $\eta_c$  is equal to the number of bits per symbol. However not all of the symbol transitions are of equal length and the bandwidth must be high enough to allow the longer transitions to take place in the same amount of time as the short transitions. However the maximum possible modulation efficiency is 3.5 bit/s/Hz.

## 5.6 Processing Gain

In digital radio the demodulated signal is a binary signal that includes information bits, coding bits, as well as bits that are in error which contributes to a raw bit error rate (BER). In the receiver decoding creates a smaller bitstream with just the information bits and which has a relatively low BER. Another way of looking at what happens is that the signal-to-interference ratio (SIR) of the decoded signal divided by the SIR of the raw signal represents a gain and this gain is called the processing gain. In 2G cellular radio this processing gain is also called the coding gain. With 3G spreading codes are used in addition to error correction codes with both reducing the bit error rate but now the gain from despreading is called spreading gain and the gain from error correction coding is called coding gain. The processing gain is now the product of spreading gain and coding gain. With 4G and 5G the functions of the spreading codes and error correction codes merge so just the term processing gain is used.

### 5.6.1 Energy of a Bit

Processing gain ( $G_P$ ) occurs when demodulating and recovering an information bitstream and is the ratio of the SIR of the processed signal to the SIR of the unprocessed signal so  $G_P$  captures the amount that the SIR increases.  $G_P$  is a measure of the additional noise immunity obtained



**Figure 5-11:** Flow of information in digital cellular radio with  $R_i$ ,  $R_b$ , and  $R_c$  being the information, baseband, and channel bit rates, respectively.

through decoding and demodulation. There are various formulas for  $G_P$  depending on whether the processed signal is the baseband bitstream with or without error correction coding, and whether the unprocessed signal is the analog RF signal, the channel bitstream, or the channel symbol (i.e. chip) stream. Processing gain is an essential concept in cellular radio and really is what makes it work. In the following, reference will be made to three bitstreams in the transmitter which will be recovered in the receiver. The first is the **information bitstream** with a bit rate  $R_i$  (in units of bit/s). To this error correction coding can be used to produce a **baseband bitstream** with baseband bit rate  $R_b > R_i$ . In some generations of radio the baseband bitstream is modulated on a carrier and the baseband bitstream is the same as the channel bitstream. The bitstreams are shown in the depiction in Figure 5-11(a) of the information flow in 2G cellular radio.

In 3G cellular radio relatively fast spreading codes merged with the baseband bitstream to produce a much faster **channel bitstream** with bit rate  $R_c \gg R_b$ . The introduction of this separate third bitstream is the depiction in Figure 5-11(b) of the information flow in 3G cellular radio. The third bitstream is the **channel bitstream** with bit rate  $R_c$  and this is the bitstream that is modulated to produce a modulated signal with a

constellation diagram where each symbol representing  $b$  bits.

In 4G and 5G cellular radio error correction codes are merged with the information bitstream to produce a baseband bitstream and this is spread during a first stage of the modulation process which produces an analog intermediate frequency modulated signal. Details of this process will be given later. There is not a separate spreading code so that the baseband and channel bitstreams are the same and  $R_c = R_b$ .

Derivation of the coding, spreading, and processing gains is based on the energy of a bit and the energy of the noise in the interval corresponding to a bit. The energy of a bit in the  $x$ th bitstream ( $x = i, b, c$  for information, baseband, and channel respectively) is denoted  $E_{b,x}$  and the energy of the noise corresponding to the duration of the bit is denoted  $N_{o,x}$ . The digital equivalent of the analog SIR is  $E_b/N_o$  (pronounced E-B-N-O for **EBNO**) and the effective SIR of the  $x$ th bitstream is

$$\text{SIR}_{\text{eff},x} = \frac{E_{b,x}}{N_{o,x}}. \quad (5.5)$$

A bitstream of course is a binary signal and the noise in the bitstream is manifested as binary errors in the bitstream. Consider a sequence of 7 bits, ideally 1001110, with one of the bits being in error so that the bitstream recovered is 1001010 where 1 bit in 7 is in error so  $\text{SIR}_{\text{eff}} = 7$ .

### 5.6.2 Coding Gain

There are many types of error correction coding schemes with some schemes better for certain types of errors<sup>4</sup>. Generalizing, with each extra error correcting bit added to an information bitstream to produce the baseband bitstream, one error in the recovered bitstream can be corrected. Thus the processing gain due to coding (and often called just coding gain) in going from the baseband to the information bitstream is

$$G_{PC} = \frac{\text{SIR}_{\text{eff},i}}{\text{SIR}_{\text{eff},b}} = \frac{E_{b,i}/N_{o,i}}{E_{b,b}/N_{o,b}} = \frac{R_b}{R_i}. \quad (5.6)$$

This is a bit-wise processing gain as it applies to bitstreams. The coding gain here is not restricted to error correction coding as there are other types of codes called spreading codes that have a similar property of providing redundancy that can be used to remove errors. Not codes are 100% effective. However Equation (5.6) is the best simple measure of processing gain when only bitstreams are considered. It provides the coding gain of one bitstream derived from a second bitstream to which coding has been used to randomize and increase the rate of a bitstream and thus introduce redundancy.

---

<sup>4</sup> The error codes used in cellular radio are **forward error correction (FEC) codes** also called **channel codes**. There are many FEC codes broadly categorized as either block codes (because they work on blocks of data) or convolution codes (because they work on arbitrary-length bitstreams). The various FEC codes use different models (i.e. assumptions) about the types of errors encountered. The selection of the FEC code to use is a design choice based on the available computing power and the nature of the errors, e.g. random or long strings of errors.

### 5.6.3 Spreading Gain

Equation (5.6) can be rearranged so that the information EBNO is

$$\frac{E_{b,i}}{N_{o,i}} = G_{PC} \frac{E_{b,b}}{N_{o,b}}. \quad (5.7)$$

The processing gain determined in Equation (5.6) applies to bitstreams and does not include the effect of modulation. A second form of the processing gain relates the SIR of the analog RF signal, i.e.  $SIR_{RF}$ , to the EBNO of the baseband bitstream. To include the effect of modulation it is recognized that what is modulated and transmitted are symbols and the transitions from one symbol to another. The energy of a symbol is denoted  $E_s$  and this energy is shared by the  $b$  bits associated with the symbol. Thus a channel bit will have the energy  $E_{b,c} = E_s/b$ . In digital radio power levels of transmitted signals are adjusted so that at most one received channel bit can be in error per received symbol. Thus if received noise results in a symbol error it will affect just one bit. That is, the symbol noise will be the same as the bit noise,  $N_{o,c} = N_o$  where  $N_o$  is the noise energy in the duration of one symbol. Thus the effective SIR of the channel bitstream is

$$SIR_{\text{eff},c} = \frac{E_{b,c}}{N_{o,c}} = \frac{E_{b,c}}{N_o} = \frac{E_s/b}{N_o}. \quad (5.8)$$

Development is completed by relating the energy of a symbol to the RF signal energy. For PSK modulation all of the symbols have the same energy. (The energy of each symbol is not the same for modulation formats such as QAM where symbols have different energy levels and a more sophisticated derivation is required than that provided here.) Also the noise energy in the duration of a symbol is the same for all symbols for all modulation formats. Thus the effective SIR of a symbol is

$$SIR_{RF} = \frac{E_s}{N_o}. \quad (5.9)$$

Thus the effective SIR of the channel bitstream is (combining Equations (5.8) and (5.9)),

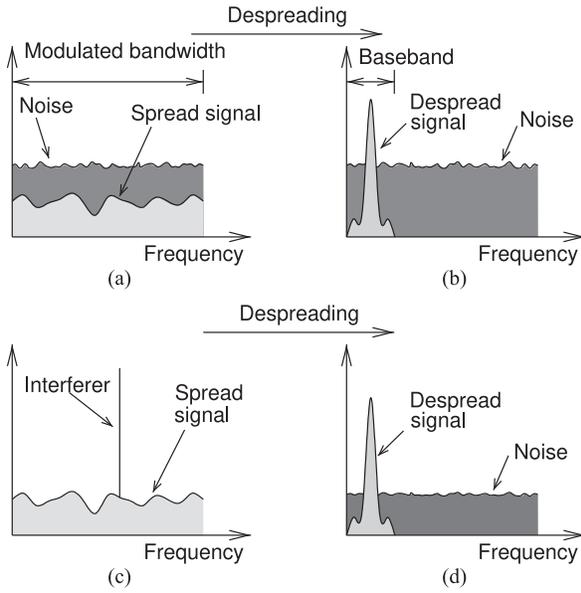
$$SIR_{\text{eff},c} = \frac{E_{b,c}}{N_{o,c}} = \frac{1}{b} SIR_{RF}. \quad (5.10)$$

The processing gain relating the EBNO of the recovered baseband bitstream to the SIR of the RF signal can be defined as the processing gain due to spreading and modulation and this is sometimes called just the **spreading gain**  $G_{PS}$ . It is defined here as the ratio of the effective SIR of the baseband bitstream and the RF SIR (and using Equation (5.10)):

$$G_{PS} = \frac{E_{b,b}/N_{o,b}}{SIR_{RF}} = \frac{1}{b} \frac{E_{b,b}/N_{o,b}}{SIR_{\text{eff},c}} = \frac{1}{b} \frac{SIR_{\text{eff},b}}{SIR_{\text{eff},c}}. \quad (5.11)$$

This can be rearranged so that the EBNO of the baseband bitstream is

$$\frac{E_{b,b}}{N_{o,b}} = G_{PS} SIR_{RF}. \quad (5.12)$$



**Figure 5-12:** Increase in SNR obtained by despreading a spread-spectrum signal: (a) the spectrum of noise and a spread-spectrum signal that has a power below the level of the noise; (b) the despread signal, where the power distribution of the noise is not changed but the despread signal is confined to a narrower bandwidth called the baseband bandwidth; (c) a signal with a single-tone interferer; and (d) after despreading where the power of the interferer spread across the the original bandwidth.

Equation (5.11) includes the ratio of the effective SIRs of two bitstreams and this was related to the bit rates of the bitstreams in Equation (5.6) where the faster bitstream is coded (or spread) to enable recovery from errors. Using the result in Equation (5.6), Equation (5.11) can be written as

$$G_{PS} = \frac{1}{b} \frac{SIR_{eff,b}}{SIR_{eff,c}} = \frac{1}{b} \frac{R_c}{R_b}. \quad (5.13)$$

Using this the EBNO of the baseband bitstream is

$$\frac{E_{b,b}}{N_{o,b}} = SIR_{eff,b} = \frac{1}{b} \frac{R_c}{R_b} SIR_{RF} = \frac{1}{b} G_{PC} SIR_{RF}. \quad (5.14)$$

Then the EBNO of the information bitstream is

$$\frac{E_{b,i}}{N_{o,i}} = G_{PC} \frac{E_{b,b}}{N_{o,b}} = G_{PC} G_{PS} SIR_{RF} = G_P SIR_{RF} = \frac{1}{b} \frac{R_c}{R_i} SIR_{RF} \quad (5.15)$$

where the processing gain  $G_P = G_{PC} G_{PS}$ . (5.16)

### 5.6.4 Spreading Gain in Terms of Bandwidth

The spreading gain that results from spreading a baseband bitstream in a transmitter and then despreading in a receiver is illustrated graphically in Figure 5-12. The modulated signal shown in Figure 5-12(a) is the “spread signal” which here has a power density which is below that of the noise. Following despreading all of the energy in the modulated RF signal is correlated with the spreading code (it was spread using the spreading code) and is collapsed to the despread signal shown in Figure 5-12(b). Only the signal correlated to the despreading code is collapsed to the smaller baseband bandwidth. Noise is rearranged with the spectral density of the noise unchanged so that the total noise energy in the baseband for the duration of a bit of data, i.e. the noise in the baseband bandwidth  $B_b$ , is

greatly reduced by the ratio of the modulated bandwidth to the baseband bandwidth.

In Equation (5.13)  $R_c$  is related to the bandwidth  $B_m$  of the modulated carrier by the modulation efficiency as  $\eta_c = R_c/B_m$  (see Equation (5.3)). Also the minimum bandwidth of a baseband bitstream with a bit rate  $R_b$  is the baseband bandwidth  $B_b = R_c$ . So Equation (5.13) becomes

$$G_{PS} = \frac{\eta_c B_m}{b B_b}. \quad (5.17)$$

Ideally for a modulation scheme  $\eta_c = b$  so that

$$G_{PS,ideal} = \frac{B_m}{B_b}. \quad (5.18)$$

For all modulation schemes other than BPSK where  $\eta_c$  can be very close to  $b$ , the spreading gain will be less than  $G_{PS,ideal}$  as in reality  $\eta_c < b$ , see Table 2-2. If both spreading and error correction coding are used the two processing gains in Equations (5.6) and (5.17) can be multiplied.

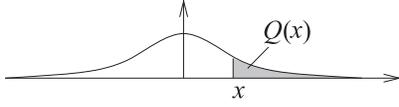
### 5.6.5 Symbol Error Rate and Bit Error Rate

With digital modulation, symbols (i.e., groups of bits) are transmitted rather than individual bits. As the modulation order increases, the number of symbols increases and there is a smaller margin between the symbols. That is, on a constellation diagram there are more symbols and the symbols are closer together. For the same signal and noise levels, as the modulation order increases the probability of a symbol error will increase and thus the symbol error rate (SER) increases. However the bit error rate (BER) increases more slowly than the SER. This is because if the SIR is high enough, the only possible errors will be nearest-neighbor errors (on a constellation diagram). If there are at least  $b$  bits per symbol the BER will be less than the **symbol error rate** by a factor of  $b$  as errors are at most one bit per symbol since the modulation order or the signal power level are adjusted to ensure this. If it is not possible to achieve a maximum of one bit error per symbol error then communication is lost. The rest of this section treats this analysis more mathematically.

The discussion begins with the SIR of the incoming RF signal. Through sampling of the RF signal at appropriate times (determined by the recovered carrier) discrete symbols are obtained. The digital form of SIR relates the energy of a symbol,  $E_s$ , to the noise and interference energy corresponding to the symbol (the noise and interference in the duration of the symbol). The height of the double-sided noise spectral density is conventionally taken as  $N_o/2$ , so the noise power corresponding to a symbol is  $N_o$ . Thus [16–18]

$$\frac{E_s}{N_o} = \text{SIR}. \quad (5.19)$$

Error in a digitally modulated communications system is first manifested as a symbol error that occurs when a symbol selected in a receiver is not the symbol transmitted. The probability of a symbol error is a function of  $E_s/N_o$ . However it is not always possible to develop a closed-form expression relating the two. For a BPSK system it can be derived. The probability of



**Figure 5-13:** Gaussian distribution function showing the  $Q$  function as the area of the shaded region.

a symbol error, the SER, is [16, p. 187]

$$\text{SER}_{\text{BPSK}} = \Pr[\text{symbol error}] = \Pr_s^{\text{BPSK}} = Q\left(\sqrt{\frac{2E_s}{N_o}}\right) = Q\left(\sqrt{2 \cdot \text{SIR}}\right), \quad (5.20)$$

where  $Q(x)$  is known as the  $Q$  function [16] and is the integral of the tail of the Gaussian density function (see Figure 5-13). It can be expressed in terms of the error function  $\text{erf}(x)$  and complementary function  $\text{erfc}(x)$  [16, page 63]:

$$Q(x) = \frac{1}{2}\text{erfc}\left(\frac{x}{\sqrt{2}}\right) = \frac{1}{2}\left[1 - \text{erf}\left(\frac{x}{\sqrt{2}}\right)\right]. \quad (5.21)$$

For  $M$ -PSK [16, page 191] the SER is

$$\text{SER}_{M\text{-PSK}} = \Pr_s^{M\text{-PSK}} \approx 2Q\left[\sqrt{\frac{2E_s}{N_o}} \sin\left(\frac{\pi}{M}\right)\right], \quad (5.22)$$

where  $M$  is the number of symbols (e.g. for 8-QPSK  $M = 8$ ). To see the effect of higher-order modulation (i.e., higher  $M$ ) consider Equation (5.22). As the order of modulation increases,  $M$  increases, and the argument of  $Q$  reduces, thus increasing SER.

The per-bit SIR is obtained after noting that each symbol can represent several bits. With a uniform constellation and the same number of bits per symbol,  $b$ , the signal energy received per bit is

$$E_b = E_s/b, \quad b = \log_2 M, \quad (5.23)$$

$$\text{and } M = 2^b. \quad (5.24)$$

For high SIR, a symbol error is the erroneous selection of a nearest neighbor symbol by the receiver. So with **Gray code mapping** (also called **Gray mapping**), such a symbol error results in only a single bit being in error. Thus the probability of a bit error, the BER, is

$$\Pr[\text{bit error}] = \Pr_b = \frac{1}{b}\Pr_s. \quad (5.25)$$

The final results are the bit error probabilities for BPSK and M-PSK [16, page 193]:

$$\text{BER}_{|\text{BPSK}} = \Pr_{b,\text{BPSK}} = Q\left(\sqrt{2E_s/N_o}\right) = Q\left(\sqrt{2 \cdot \text{SIR}}\right) \quad (5.26)$$

$$\begin{aligned} \text{BER}_{|M\text{-PSK}} &= \Pr_{b,M\text{-PSK}} = \frac{2}{b}Q\left[\sqrt{2E_s/N_o} \sin\left(\frac{\pi}{M}\right)\right] \\ &= \frac{2}{b}Q\left[\sqrt{2 \cdot \text{SIR}} \sin\left(\frac{\pi}{M}\right)\right]. \end{aligned} \quad (5.27)$$

Thus for the same SIR the probability of bit errors, the BER, grows with higher-order modulation (i.e., larger  $M$ ), but not as fast as SER. At the same time the number of bits transmitted increases. These extra bits are use to

embed error correction codes. The net gain in throughput can be tremendous as long as the SIR is high enough.

### EXAMPLE 5.2 Symbol Error Rate

Calculate the SER and BER for QPSK and 8-PSK if the SIR is 10 dB.

#### Solution:

SER is found by evaluating Equation (5.22). For SIR = 10 dB,  $E_s/N_o = 10^{(\text{SIR}_{\text{dB}}/10)} = 10$ . For QPSK,  $M = 4$  and  $b = 2$ , and the SER is

$$\begin{aligned} \text{SER}_{\text{QPSK}} = \text{Pr}_s^{M-\text{PSK}} &\approx 2Q \left( \sqrt{\frac{2E_s}{N_o}} \sin \left( \frac{\pi}{M} \right) \right) = 2Q \left[ \sqrt{20} \sin \left( \frac{\pi}{4} \right) \right] \\ &= 2Q(3.162) = \left[ 1 - \text{erf}(3.162/\sqrt{2}) \right] = 0.001565. \end{aligned} \quad (5.28)$$

For 8-PSK,  $M = 8$  and  $b = 3$ , and the SER is

$$\begin{aligned} \text{SER}_{8-\text{PSK}} = \text{Pr}_s^{M-\text{PSK}} &\approx 2Q \left[ \sqrt{\frac{2E_s}{N_o}} \sin \left( \frac{\pi}{M} \right) \right] = 2Q \left( \sqrt{20} \sin \left( \frac{\pi}{8} \right) \right) \\ &= 2Q(1.711) = 0.08701. \end{aligned} \quad (5.29)$$

The corresponding BERs are

$$\text{BER}_{\text{QPSK}} = \frac{1}{2} \text{SER}_{\text{QPSK}} = 0.000783 \quad \text{and} \quad \text{BER}_{8-\text{PSK}} = \frac{1}{3} \text{SER}_{8-\text{PSK}} = 0.0290. \quad (5.30)$$

### 5.6.6 Summary

The formulas for processing gain developed above are only approximate as many simplifications were made. One assumption is that the error correction coding is fully randomized bits but in practice there are many error correction codes that address particular types of errors. The aim of error correction coding is to recover from errors not necessarily to achieve processing gain.

Another caveat to the processing gain formulas developed here is that it considered the energy of a symbol as being the same for all symbols. While this is true for PSK modulation it is not true for all modulation formats. Considerable effort is put into developing better estimates that capture the essence of processing gain but in a more rigorous manner. Still the simple formulas provide the required insight into the effect of using error correction and spreading codes.

In 2G cellular radio processing gain is achieved using error correction codes alone and a separate spreading code is not used. Thus in 2G the baseband bits are also the channel bits used directly in modulating the RF carrier. The extra bits provide redundancy and the ability to recover from some errors resulting from noise. In 3G spreading codes as well as error correcting codes are used. A spreading code greatly increases the bit rate of the channel bitstream above that of the baseband bitstream and the bandwidth of the modulated carrier is much greater than that needed to transmit the baseband bitstream. Generally error correcting codes do not increase the bit rate by more than a factor of 2. Thus in 3G the processing gain achieved with the error correction codes is almost insignificant compared to

**Table 5-2:** Processing gain definitions.  $R_i$  is the information bit rate,  $R_b > R_i$  is the baseband bit rate after error correction coding of the information bitstream, and  $R_c > R_b$  is the channel bit rate after spreading the baseband bitstream. Reference is made to the  $x$ th bitstream with  $x = i, b,$  or  $c$  indicating the information, baseband, or channel bitstreams respectively.  $E_{b,x}$  is the energy of a bit in the  $x$ th bitstream, and  $N_{o,x}$  is equivalent noise energy corresponding to a bit in the  $x$ th bitstream.  $E_s$  is the energy of a symbol with  $b$  bits per symbol.  $B_m$  is the bandwidth of the modulated carrier and  $B_b$  is the bandwidth of the baseband signal.

Description	Definition	Formula	Equation
Coding gain, processing gain from error correction coding. Calculated for bitstreams.	$G_{PC} = \frac{E_{b,i}/N_{o,i}}{E_{b,b}/N_{o,b}}$	$G_{PC} = \frac{R_b}{R_i}$	(5.6)
Processing gain from spreading calculated on a bitstream basis	$G_{PS} = \frac{E_{b,c}/N_{o,c}}{E_{b,b}/N_{o,b}}$	$G_{PS} = \frac{1}{b} \frac{R_c}{R_b}$	(5.13)
Processing gain from spreading in terms of bandwidth.	$G_{PS} = \frac{E_{b,b}/N_{o,b}}{\text{SIR}_{\text{RF}}}$	$G_{PS} = \frac{\eta_c}{b} \frac{B_m}{B_b}$	(5.17)
Processing gain	$G_P = \frac{E_{b,c}/N_{o,c}}{E_{b,i}/N_{o,i}}$	$G_P = G_{PC}G_{PS}$	(5.16)
EBNO	$E_{b,i}/N_{b,i}$	$E_{b,i}/N_{b,i} = G_P \text{SIR}_{\text{RF}}$	(5.12)

that obtained from spreading the signal.

In 4G and 5G error correction coding also provides spreading and there is not a separate spreading code. The 4G and 5G systems have great complexity.

A summary of the key results for processing gain is given in Table 5-2.

### EXAMPLE 5.3 Processing Gain

A new communication system is being investigated for sending data to a printer. The system will use GMSK modulation and a channel with 10 MHz bandwidth and the baseband bit rate will be 1 Mbit/s. The modulation format will result in a spectrum that distributes (i.e. spreads) power almost uniformly over the 10 MHz bandwidth.

(a) What is the processing gain?

The most efficient way to spread the signal across the 10 MHz of available bandwidth is to use a combination of error correction and spreading. Referring to Table 2-2, GMSK has a modulation efficiency  $\eta_c$  of 1.354 bit/s/Hz, so that the channel bit rate  $R_c = 13.54$  Mbit/s. Each symbol in GMSK represents two bits so  $b = 2$ . The information bit rate  $R_i = 1$  Mbit/s and so coding at the rate of 12.54 Mbit/s will be used and  $R_b = (R_i + 12.54 \text{ Mbit/s}) = 13.54 \text{ Mbit/s}$ .

There are two processing gains. The coding gain is determined using Equation (5.6):

$$G_{PC} = \frac{R_b}{R_i} = \frac{13.54 \text{ Mbit/s}}{1 \text{ Mbit/s}} = 13.54 = 11.31 \text{ dB.}$$

The spreading gain, using Equation (5.13), is

$$G_{PS} = \frac{1}{b} \frac{R_c}{R_b} = \frac{1}{2} \frac{13.54 \text{ Mbit/s}}{13.54 \text{ Mbit/s}} = 0.5 = -3 \text{ dB.}$$

Thus the processing gain

$$G_P = G_{PC}G_{PS} = 6.77 = 8.31 \text{ dB.}$$

Another way of calculating the spreading gain is to consider the modulated bandwidth,  $B_m = 10$  MHz and the baseband bandwidth  $B_b = 1$  MHz (taken to be numerically equal to  $R_b$ ). The processing gain due to spreading, from Equation (5.17):

$$G_{PS} = \frac{\eta_c B_m}{b B_b} = \frac{1.354}{2} \frac{10 \text{ MHz}}{1 \text{ MHz}} = 6.770 = 8.31 \text{ dB.}$$

- (b) If the received RF SIR is  $\text{SIR}_{\text{RF}} = 6$  dB, what is the effective system SIR (or  $E_b/N_o$ ) after the digital signal processor?

$$\text{Effective SIR} = \frac{E_{b,i}}{N_{o,i}} = G_P \text{SIR}_{\text{RF}} = 8.31 \text{ dB} + 6 \text{ dB} = 14.32 \text{ dB.}$$

#### EXAMPLE 5.4

#### Signal-to-Interference Ratio

At the output of a receiver antenna, the level of interfering signals is 1 pW, the level of background noise is 500 fW, and the level of the desired signal is 4 pW.

- (a) What is the SIR? Note that SIR includes the effect of the signal, interference, and noise.  
 (b) If the processing gain is 20 dB and 16-QAM is the modulation scheme used, what is the effective system SIR, that is, what is the signal energy in a bit versus the noise energy in the duration of the bit (i.e.,  $E_b/N_o$ )?

#### Solution:

- (a) Interfering signal  $P_I = 1$  pW, noise signal  $P_N = 500$  fW, and signal  $P_S = 4$  pW.

$$\begin{aligned} \text{SIR}_{\text{RF}} &= P_S / (P_I + P_N) \\ &= 4 \text{ pW} / (1 \text{ pW} + 0.5 \text{ pW}) = 4 / 1.5 = 2.667 = 4.26 \text{ dB.} \end{aligned}$$

- (b) Processing gain  $G_P = 20$  dB so, from Equation (5.15),

$$\text{Effective SIR} = \frac{E_{b,i}}{N_{o,i}} = G_P \cdot \text{SIR}_{\text{RF}} = 20 \text{ dB} + 4.26 \text{ dB} = 24.26 \text{ dB.}$$

## 5.7 Early Generations of Cellular Phone Systems

The evolution of cellular communications is described by generations of radio. The major mobile communication systems are outlined in Table 5-3. No 0G and 1G systems remain and 2G systems are being phased out with only the 2G GSM system still in use. Third generation (3G) offers a significant increase in capacity and is optimized for broadband data access. The 0G–2G systems were plagued by many incompatible systems. Even with 3G there were several incompatible systems but the system known as WCDMA became overwhelmingly dominant. With 4G the evolution coalesced and while there were many modulation methods a standardized interface was developed with the physical interface defined as common. The 4G radio, and also 5G, support many modulation methods and frequency bands and cellular radios support a large number of these combinations. In part this is because the bulk of the modulation is now performed in a DSP unit making the switching of modulation methods a simple process of using a different algorithm. The development of 4G, 5G (and beyond) has been guided by the Third Generation Partnership, 3GPP, which developed a plan called long term evolution (LTE) that ensured upwards compatibility as the capability of the systems were enhanced every year or two eventually becoming 5G.

## 5.8 Early Generations of Radio

### 5.8.1 1G, First Generation: Analog Radio

The initial cellular radio system was analog, with the dominant system being AMPS, the attributes of which are given in Table 5-4. This is a relatively simple system, but appropriate for the low levels of integration of the 1980s, as most of the functionality could be realized using analog circuits. The first-

**Table 5-3:** Major mobile communication systems with the year of first widespread use.

System	Year	Description
0G MTS	1946	Broadcast, no cells, few users, analog modulation Mobile telephone service, halfduplex, operator assist to establish call, push to talk
AMTS	1965	Advanced mobile telephone system, Japan, duplex, 900 MHz
IMTS	1969	Improved mobile telephone service, duplex, up to 13 channels, 60–100 km (40–60 mile) radius, direct dial using dual-tone multifrequency (DTMF) keypad
0.5G PALM	1971	FDMA, analog modulation (also Autotel) Public automated land mobile radiotelephone service, used digital signaling for supervisory messages, technology link between IMTS and AMPS
ARP	1971	Autoradiopuhelin (car radio phone), obsolete in 2000, used cells (30 km radius) but not hand-off, 80 channels at 150 MHz, simplex and later duplex
1G NMT	1981	Analog modulation, FSK for signaling, cellular, FDMA Nordic mobile telephone, 12.5 kHz channel, 450 MHz, 900 MHz
AMPS	1983	Advanced mobile phone system, 30 kHz channel
TACS	1985	Total access communication systems, 25 kHz channel, widely used in Europe until 1990s, similar to AMPS
Hicap	1988	NTT's mobile radiotelephone service in Japan
Mobitex	1990	National public access wireless data network, first public access wireless data communication services including two-way paging network services, 12.5 kHz channel, GMSK
DataTac	1990	Point-to-point wireless data communications standard (like Mobitex), wireless wide area network, 25 kHz channels, maximum bandwidth 19.2 kbit/s (used by the original Blackberry device)
2G PHS	1990	Digital modulation Personal handyphone system, originally a cordless phone, now functions as both a cordless phone and as a mobile phone elsewhere
GSM	1991	Global system for mobile communications (formerly Groupe Spécial Mobile), TDMA, GMSK, constant envelope, 200 kHz channel, maximum 13.4 kbits per time slot (at 1900 MHz), 2 billion customers in 210 countries
DAMPS	1991	Digital AMPS (formerly NADC [North American digital cellular] and prior to that as U.S. Digital Cellular [USDC]), narrowband, $\pi/4$ DQPSK, 30 kHz channel
PDC	1992	Personal Digital Cellular, Japan, 25 kHz channel
CDMAOne	1995	Brand name of first CDMA system known as IS-95, spread spectrum, CDMA, 1.25 MHz channel, QPSK
CSD	1997	Circuit switched data, original data transmission format developed for GSM, maximum bandwidth 9.6 kbit/s, used a single time slot

Table 5-3 continued.

System	Year	Description
2.5G WiDEN	1996	Higher data rates Wideband integrated dispatch enhanced network, combines four 25 kHz channels, maximum bandwidth 100 kbit/s
GPRS	2000	General packet radio system, compatible with GSM network, used GSM time slot and higher-order modulation to send 60 kbits per time slot, 200 kHz channel, maximum bandwidth 171.2 kbit/s
HSCSD	2000	High-speed circuit-switched data, compatible with GSM network, maximum bandwidth 57.6 kbit/s, higher quality of service than GPRS
2.75G CDMA2000 EDGE	2000 2003	Medium bandwidth data—1 Mbit/s CDMA, upgraded CDMAOne, double data rate, 1.25 MHz channel Enhanced data rate for GSM Evolution, compatible with GSM network, 8-PSK, TDMA, maximum bandwidth 384 kbit/s, 200 kHz channel
3G FOMA	2001	Spread spectrum Freedom of mobile multimedia access, first 3G service, NTT's implementation of WCDMA
UMTS WCDMA	2004	Universal mobile telephone service, 5 MHz channel, data up to 2 Mbit/s Main 3G outside China
OFDMA 1xEV-DO	2007	Evolution to 4G (downlink high bandwidth data) (IS-856) Evolution of CDMA2000, maximum downlink bandwidth 307 kbit/s, maximum uplink bandwidth 153 kbit/s
TD-SCDMA	2006	Time division synchronous CDMA, China, uses the same band for transmit and receive, basestations and mobiles use different time slots to communicate, 1.6 MHz channel
GAN/UMA	2006	Generic access network, formerly known as Unlicensed Mobile Access, provides GSM and GPRS mobile services over unlicensed spectrum technologies (e.g., Bluetooth and WiFi)
3.5G UMTS/HSDPA	2006	Upgraded WCDMA, High-speed downlink packet access, download of 7.2 Mbit/s
EV-DO Rev A	2006	CDMA2000 EV-DO revision A, downlink to 3.1 Mbit/s, uplink to 1.8 Mbit/s
3.75G UMTS/HSUPA EV-DO Rev B	2007 2008	High-speed uplink packet access, upload speeds to 5.76 Mbit/s CDMA2000 EV-DO revision B, downlink to 73 Mbit/s, uplink to 27 Mbit/s
UMTS/HSPA	2009	Upgraded WCDMA, High-speed packet access, downlink to 40 Mbit/s, upload to 10 Mbit/s. Eventually added CA and MIMO
3.9G	2009	WiMAX 1 (IEEE 802.16), 10 MHz bandwidth; IP-based; branded as 4G by service providers; MIMO + OFDMA, downlink of 37 Mbit/s, uplink 17 Mbit/s (for 2×2 MIMO); WiMAX 2, IEEE 802.16m, 20 MHz bandwidth, downlink of 110 Mbit/s, uplink 70 Mbit/s; not allowed in many countries
3.9G	2011	Long-term evolution (LTE); up to 20 MHz channel bandwidth, IP-based; branded as 4G by service providers Low latency (for VoIP) + MIMO + OFDMA, downlink of 100 Mbit/s
4G	2013	LTE-advanced, downlink of 1 Gbit/s fixed, 100 Mbit/s mobile, variable bandwidths of 5–40 MHz
5G	2019	Millimeter waves with beam steering and massive MIMO; Mesh networks and cognitive radio

**Table 5-4:** Attributes of AMPS.

Property	Attribute
Number of physical channels	832; 2 groups of 416 channels, each group has 21 signaling channels and 395 traffic or voice channels
Bandwidth per channel	30 kHz
Cell radius	2–20 km
Base-to-mobile frequency	869–894 MHz (downlink)
Mobile-to-base frequency	824–849 MHz (uplink)
Channel spacing	45 MHz between transmit and receive channels
Modulation	30 kHz FM with peak frequency deviation of $\pm 12$ kHz Signaling channel uses FSK Can send data at 10 kbit/s
Access method	FDMA
basestation ERP	100 W per channel (maximum)
Channel coding	None
RF specifications of mobile unit	
Transmit RF power	3 W maximum (33 dBm) (600 mW for hand-held)
Transmit power control	10 steps of 4 dB attenuation each, minimum power is $-4$ dBm
Receive sensitivity	$-116$ dBm
Receive noise figure	6 dB measured at antenna terminals
Receive spurious response	$-60$ dB from center of the passband

generation systems handled analog 3 kHz voice transmissions with very limited ability to transmit digital information limited to signaling.

### 5.8.2 2G, Second Generation: Digital Radio

The second generation (2G) of cellular radio is characterized by digitization. Many different types of 2G digital systems were installed around the world. The 2G systems can transmit data and voice at rates of 8–14.4 kbit/s. This can be contrasted to the wireline phone system where, once signals reach the exchange, the 3 kHz analog signals are sampled at 64 kbit/s to achieve an undistorted signal representation. These cellular systems sacrifice some voice quality but use reasonably sophisticated algorithms that use the characteristics of speech to achieve greater than a factor of four compression.

In North America the first digital system introduced was the **digital advanced mobile phone system (DAMPS)** (originally known as **North American Digital Cellular [NADC]** and as the EIA/TIA interim standard **IS-54**). The system was designed to provide a transition from the then current 1G analog system to a fully digital system by reusing existing spectrum. The idea was that system providers could allocate a few of their channels for digital radio out of the total available. As analog radio was phased out, more of the channels could be committed to digital radio. The main motivation behind this system is that it provided three to five times the capacity of the analog system for the same bandwidth. The 2G GSM system provides a similar increase in capacity, and is compatible with the **Integrated Services Digital Network (ISDN)** which was the protocol used with the wired telephone system. The GSM system was initially (early 1990s) dominant in Europe and had the advantage that it did not need to coexist

**Table 5-5:** Attributes of the GSM system. Uplink and downlink frequencies are for the GSM-900 implementation, see Table 5-6 for other GSM implementations. Slow frequency hopping improves robustness.

Property	Attribute
Number of channels	125 (for GSM-900)
Bandwidth per channel	200 kHz
Channel spacing	200 kHz
Cell radius	2–20 km
Base-to-mobile frequency	935–960 MHz (GSM-900)
Mobile-to-base frequency	890–915 MHz (GSM-900)
Modulation	GMSK, Slow frequency hopping (217 hops/s)
Access method	TDMA, 8 slots per frame, user has one slot, each frame is 4.615 ms and each slot is 577 $\mu$ s. There are 148 data bits and 8.25 guard bits in a slot.
Symbol duration	3.6828 $\mu$ s
Transmit rate	270.833 kbit/s

**Table 5-6:** GSM frequency bands. GSM channels have a bandwidth of 200 kHz. The base-to-mobile transmission is the downlink and the mobile-to-basestation transmission is the uplink. GSM-900 and GSM-1800 are used in most of the world.

Band	Uplink (MHz)	Downlink (MHz)	Duplex spacing (MHz)
GSM-900	890–915	935–960	45
GSM-1800	1710–1785	1805–1880	95
GSM-900 extended	876–915	921–960	45
PCS-1900	1850–1910	1930–1990	80
GSM-850 (Americas)	824–849	869–894	45
GSM-450	450.4–457.6	460.4–467.6	10
GSM-480 (Nordic, Eastern Europe, Russia)	478.8–486	488.8–496	10

with uncoordinated 1G analog phone systems. The attributes of the GSM system are shown in Tables 5-5 and 5-6.

From an RF design perspective, the main differences between analog and digital standards are

1. The RF envelope. In AMPS, FM was used, which produces a constant envelope RF signal. Consequently, high-efficiency saturation mode amplifiers (such as Class C) can be used. In most digital modulation schemes the modulation results in a nonconstant envelope. This is true for PSK modulation, as the transition from one symbol to another does not follow a circle on the constellation diagram. The information contained in the amplitude of the RF signal is just as important as the information contained in the phase or frequency of the signal. Consequently, with digital radio saturation mode amplifiers that severely distort the amplitude characteristic must be avoided.
2. Bursty RF. In an analog system, RF power is continually being transmitted. In a digital system, transmission is intermittent and the RF signal is bursty. Therefore an RF designer must be concerned about turn-on transients and thermal stability of the power amplifier.

## 5.9 3G, Third Generation: Code Division Multiple Access (CDMA)

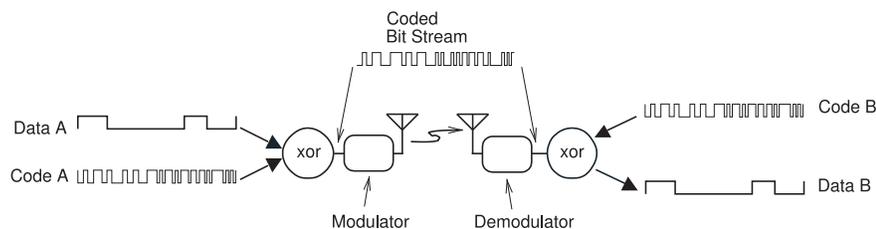
Originally there were multiple 3G cellular radio systems but the one that became dominant uses code division multiple access. This section begins with a precursor to 3G and which is now called 2.5G cellular radio and this is relatively narrowband, 1.23 MHz-wide, CDMA. The 3G system uses 5 MHz-wide wideband CDMA (WCDMA).

### 5.9.1 Generation 2.5: Direct Sequence Code Division Multiple Access

CDMA, or more specifically CDMAOne, was initially promoted as being third generation, but the definition now is that data rates of at least 2 Mbit/s must be supported in 3G. Thus CDMAOne is now referred to as 2.5G. A depiction of spread spectrum is shown in Figure 5-14, in which a very fast code is superimposed on a slower data sequence and the combined code is used to modulate a carrier. The same fast code is used to extract the baseband signal from the received bitstream. The effect of the fast code is to greatly spread out the baseband signal, transforming perhaps a 12 kbit/s baseband bitstream into an RF signal with a bandwidth of 1.23 MHz.

The key feature of the DS-CDMA system is the use of lengthy codes to spread the spectrum of the signal that is to be transmitted. In the case of voice, an 8 kbit/s bitstream, for example, with error correction coding becomes a 12.5 kbit/s baseband bitstream that is mixed with a much faster code that is unique to a particular user. Thus the 8 kbit/s bitstream becomes a 1.23 MHz-wide analog baseband signal. This signal is then modulated up to RF and transmitted. On the receiver side, the demodulated RF signal can only be decoded using the original fast code. Use of the original code to decode the signal rejects virtually all interference and noise in the received signal. Despreading distributes the signal and noise components differently. Upon despreading, the noise is still distributed uniformly in frequency while the information-bearing signal is concentrated in a narrow bandwidth, the bandwidth of the baseband signal. Tremendous processing gain is available using this spreading and despreading approach.

The mechanism that increases SNR in DS-CDMA is shown in Figure 5-12. The SNR is enhanced by grouping the signal energy in a narrower bandwidth and the noise is reduced, as only the noise in a narrower



**Figure 5-14:** Depiction of direct sequence CDMA transmission. If code B is code A, then data B is the same as data A with some corruption by noise.

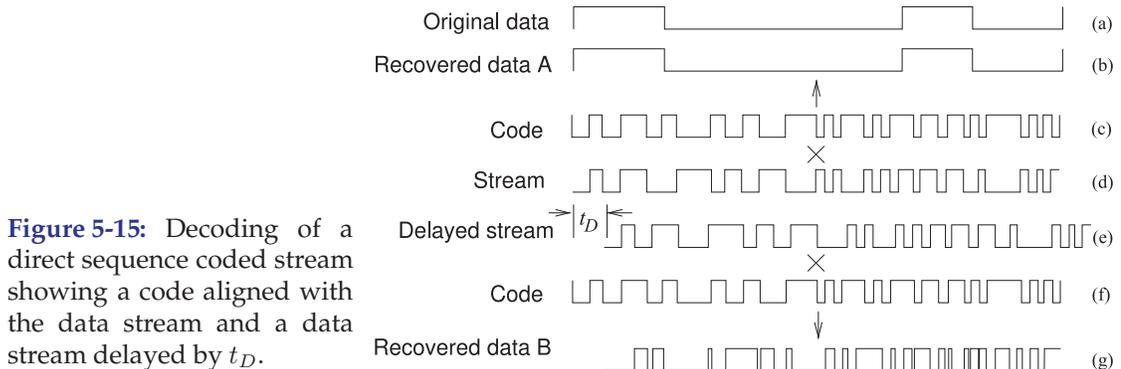
**Table 5-7:**  
Attributes of  
the CDMAOne  
system.

Property	Attribute
Bandwidth per channel	1.23 MHz
Channel spacing	1.25 MHz (20 kHz guard band)
Cell radius	2–20 km
Base-to-mobile frequency	869–894 MHz
Mobile-to-base frequency	824–849 MHz
	45 MHz between transmit and receive channels.
Modulation	QPSK
Access method	CDMA
	64 radio channels per physical channel
	Forward:
	55 traffic channels, 7 paging channels
	1 pilot channel 1 sync channel
	Reverse:
	55 traffic channels, 9 access channels
RF specifications of mobile unit	
Transmit power control	1 dB power control

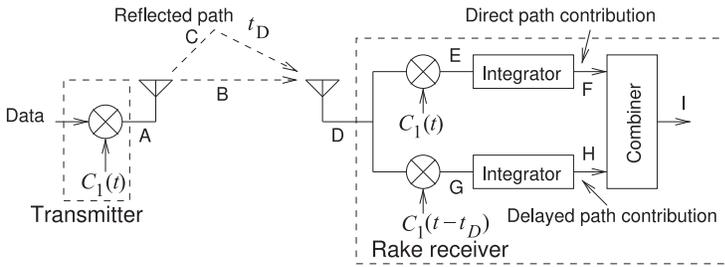
bandwidth is important. In despreading a signal with a single-tone interferer, as seen in Figure 5-12(c), the interferer is spread as noise while the signal energy is concentrated in the bandwidth of the baseband signal. Since orthogonal codes are used, many radio channels can be supported on the same radio link. CDMA can support approximately 120 radio channels on the same physical channel. Another important feature is that the same 120 channels can be reused in adjacent cells, as the information bitstream of each user can still be extracted. Thus there is no need for clustering as in the 2G systems. The attributes of the cellular 2.5G CDMAOne system, an immediate precursor to the 3G WCDMA system, are given in Table 5-7.

### 5.9.2 Multipath and Rake Receivers

In a line-of-sight CDMA system the code must be aligned with the data stream received so that the stream can be decoded correctly to reveal the original data. If there are multiple paths then the paths will, in general, have different delays and often the delay differences are more than the chip duration (the time allocated to transition from one symbol to another). The CDMA signal is then said to be transmitted over a **dispersive multipath channel**. This introduces complexity in aligning codes. The problem that arises is shown in Figure 5-15. The original data are shown in (a) and this



**Figure 5-15:** Decoding of a direct sequence coded stream showing a code aligned with the data stream and a data stream delayed by  $t_D$ .



**Figure 5-16:** Rake receiver used to decode two paths, a direct path and a path with excess delay  $t_D$ . The integrators eliminate incorrectly recovered data.

is coded using the code in (c) yielding the stream in (d). The stream is just the exclusive or (XOR) of the data and code. On receiving, the coded data stream, (d), is despread using the original code (c). If the code and the data stream are aligned, then an XOR operation results in the data being recovered correctly by XORing (c) and (d) to obtain (b), the original data A. However, if the stream is delayed, as shown in (e), then XORing the delayed stream and the code (which is not delayed) will result in recovered data B, (g), which has no relationship to the original data in (a). The solution to this problem is to appropriately delay the despreading code for each propagation-delay path.

When there are reflections of the transmitted signal, versions of the signal following different paths will arrive at the receiver with different delays. The **rake receiver**, introduced in the 1950s, recovers the data from each propagation delay path and combines the signal from each path to produce a combined signal with a higher SNR than that obtained using just the line-of-sight component of the signal [20–23].

A two-path version of the rake receiver is shown in Figure 5-16. The transmitter spreads the data using a code  $C_1(t)$  to produce a stream at A that is transmitted over two paths C and B. Path C is delayed with respect to path B by a time  $t_D$  and the signal at the receive antenna, D, is a combination of the signal following the two paths. On the top path of the receiver, called a *finger*, the direct signal (that followed path B) is despread by the original code  $C_1(t)$  and, if it is appropriately aligned with the signal that followed path B, it produces a signal at E that will contain the low-frequency user data plus a high-frequency component resulting from the wrong code being used to despread the delayed data stream that followed path C. The signal at E is lowpass filtered to produce the original data at F as the high frequency component is eliminated. Since the operation is implemented in DSP, an integrator is used over the duration of a data bit. The combined received signal, D, is despread using a delayed code,  $C_1(t - t_D)$ , to produce a signal at G that is integrated so that the data at H is the original D due to the component that followed path C. The signals at F and H, sampled after the duration of a bit, will both contain the original data plus some noise and the noise at the two points will be uncorrelated. When F and H are combined, the signal will combine coherently while the noise will combine incoherently, thus improving the SNR.

The rake receiver can be generalized to have many fingers. After despreading in each finger of the rake receiver, each delayed component is demodulated and the results are combined. The rake receiver is so named because each finger sweeps up information, resembling the tines on a garden rake collecting leaves. Since each component contains the original information, if the magnitude and time of arrival (phase) of each component is computed at the receiver (through a process called channel estimation),

**Table 5-8:** Attributes of the 3G WCDMA system. 3G operates in many frequency bands and just one representative band, Band 3, is considered here. The uplink modulation in WCDMA is unusual applying data to the I channel of an I/Q modulator and control data to the Q channel. Here the I channel (and Q channel) has two or four levels classified as **pulsed amplitude modulation (PAM)**. † per user.

\* shared, non-MIMO. ‡ shared, 2-cell carrier aggregation, MIMO. The 3G standards do not specify minimum information transmit rates per user.

Property	Attribute
Number of channels	
Bandwidth per channel	5 MHz
Channel spacing	5 MHz
Cell radius	2–20 km
Downlink frequency	1805–1880
Uplink frequency	1710–1785
Modulation (downlink)	QPSK, 16QAM, 64QAM
Modulation (uplink)	PAM on I and on Q
Access method	CDMA
Symbol rate	3.84 Msymbols/s
Modulation bandwidth	3.84 MHz
Symbol duration	260 ns
Information rate (original)†	200 kbit/s up- & down-link
Information rate (HSPA+)*	40 Mbit/s downlink
Information rate (HSPA+)‡	168 Mbit/s downlink
Information rate (HSPA+)†	10 Mbit/s uplink

then all the components can be added coherently to improve the information reliability. This method of combining is called **maximal-ratio combining (MRC)**, a method of **receiver diversity** combining in which

1. the signals from each channel are added together,
2. the gain of each channel is made proportional to the rms signal level and inversely proportional to the mean square noise level in that channel, and
3. different proportionality constants are used for each channel.

It is also known as **ratio-squared combining** and **predetection combining**. MRC is optimum combining for independent AWGN channels.

CDMA, WCDMA, and WLAN (WiFi) units use multipath signals and combine them to increase the SNR at the receivers. In contrast, the narrowband 2G systems cannot discriminate between the multipath arrivals, and multipath has negative impact. The rake receiver in Figure 5-16 has two fingers, but the practical limit on the number is based on the expected delay spread of the multipaths divided by the chip duration. Only delays that are an integer multiple of a chip duration are needed in a rake receiver [21, 22]. A trade-off must be made in terms of the amount of circuitry and the DC power available. In the original CDMA cell phone system (circa 2000), three fingers were used in handsets and four to five in basestations. Many more are used today.

### 5.9.3 3G, Wideband CDMA

**Third-generation radio, (3G)**, is coordinated by the **Third Generation Partnership Project (3GPP)**. This is a collaborative agreement of standards development organizations and other related bodies for the production of a complete set of globally applicable technical specifications for mobile communication systems (see <http://www.3gpp.org>). Initially the efforts of 3GPP were directed at establishing the 3G standards but the scope has expanded and it develops evolving standards for the transition to 4G and now 5G systems. Some of the attributes of 3G are given in Table 5-8.

The 3G mobile systems support variable data rates depending on demand and the level of mobility. Switched-packet radio techniques are required to

support this bandwidth-on-demand environment. Here the physical channel is shared (i.e., **packet switched**) rather than the user being assigned a physical channel for exclusive use (referred to as **circuit switched**).

The drive for 3G systems was partly fueled by the saturation of 2G systems in many places and a desire to increase revenues by supporting high-speed data. Prior to the rollout of 3G systems, the increased demand primarily resulted from an increased consumer base rather than the emergence of significant data traffic. The increased subscriber base was addressed by 2.5G systems, which have some of the 3G concepts but only partially implemented. The driving concept of 3G was the development of a standard that supports high-speed data, global roaming, and supports advanced features including two-way motion video and internet browsing.

Third-generation cellular radio is defined by the **International Telecommunications Union (ITU)** [24] and is formally called International Mobile Telecommunications 2000 (**IMT-2000**). The basic requirements are for a system that supports data rates up to 2 Mbit/s in fixed environments ranging down to 144 kbit/s in wide area mobile environments. In 1999 the ITU adopted five radio interfaces for IMT-2000:

1. IMT-DS direct-sequence CDMA, more commonly known as WCDMA;
2. IMT-MC multi-carrier CDMA, more commonly known as CDMA2000, the successor to CDMAOne (specifically international standard IS-95);
3. IMT-SC time-division CDMA, which includes time division CDMA (TD-CDMA) and time division synchronous CDMA (TD-SCDMA);
4. IMT-SC single carrier, more commonly known as EDGE; and
5. IMT-FT frequency time, more commonly known as DECT.

The dominant choice for 3G is WCDMA. In October 2007 the ITU Radio-communication Assembly included WiMAX-derived technology, specifically **orthogonal frequency division multiple access (OFDMA)**, see Figure 5-8(d) and MIMO, in the set of IMT-2000 standards as the sixth radio interface. 3GPP [25] provides a migration strategy for cellular communications through a process called **long-term evolution LTE** and through a number of releases each building on prior infrastructure and adding capabilities.

EDGE has intermediate data speeds between those of GSM and WCDMA. The following terms are also used to describe networks using the 3G WCDMA specification: Universal Mobile Telecommunication System (**UMTS**) (UMTS, in Europe), UMTS Terrestrial Radio Access Network (UTRAN), and **Freedom of Mobile Multimedia Access (FOMA)**, in Japan). UMTS is the 3G successor of the GSM standard, with the air interface now using WCDMA. The terminology used in UMTS, listed in part in Table 5-9, is based on the terminology used in GSM, with subtle differences. UMTS was first deployed in Japan in 2001. The term WCDMA describes the physical interface and protocols that support it, while UMTS refers to the whole network. A large number of frequency bands are designated for 3G, see Tables 5-10 and 5-11.

The 3GPP timeline is summarized in Figure 5-17. The CDMA2000 and WCDMA paths become the single LTE path beyond 3G. The CDMA2000 (the IS-2000 standard) path builds on the original CDMA system defined by the IS-95 standard and commonly known as CDMAOne (the **IS-95** standard). CDMA2000 1xEV-DO is the first evolution of CDMA2000 that meets the ITU basic specification for 3G. **Evolution-data optimized (EV-DO)**, combines CDMA and TDMA for higher data throughput.

**Table 5-9:**  
UMTS terminology.

Term	Description
AuC	Authentication center
GGSN	Gateway GPRS support node
GMSC	Gateway MSC
HLR	Home location register
ISDN	Integrated services digital network
MSC	Mobile switching center
Node B	Basestation
PSTN	Public switched telephone network
RNC	Radio network controller
SGSN	Serving GPRS support node
UE	User equipment
USIM	Universal subscriber identity module
VLR	Visitor location register

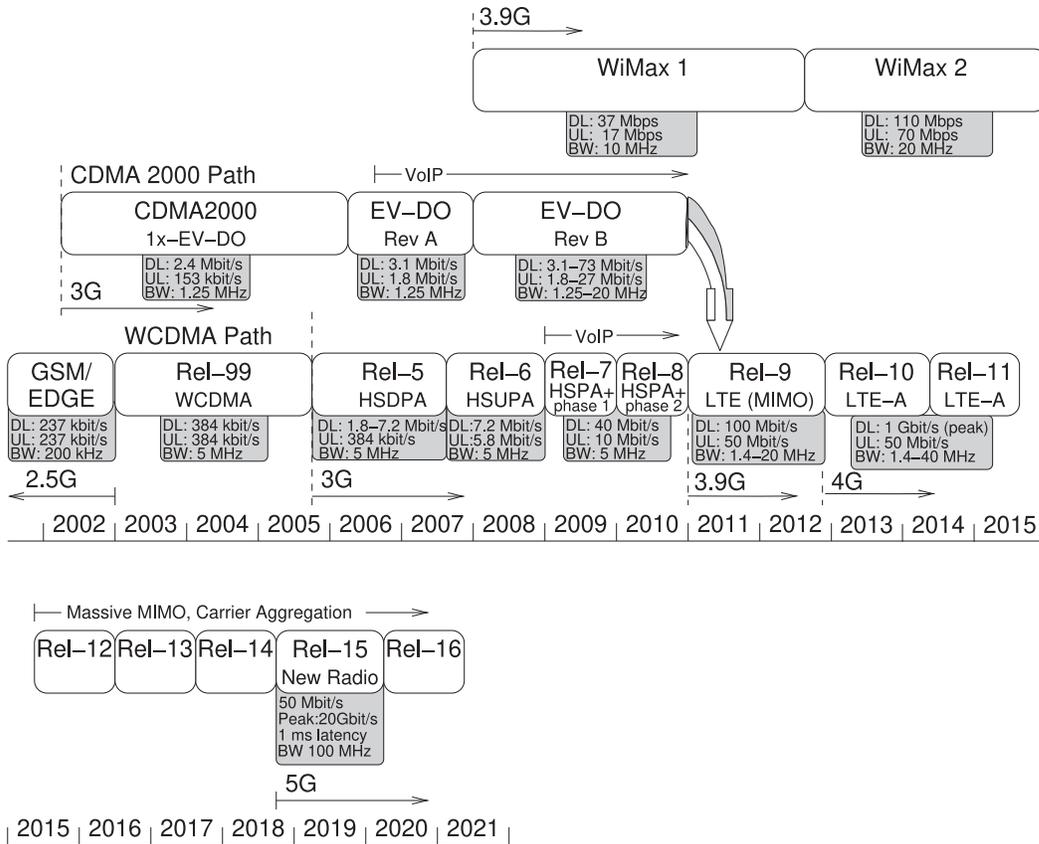
**Table 5-10:** Spectrum assignments for 3G, [25, Release 99].

Band	Uplink (MHz)	Downlink (MHz)	Available spectrum					
			NA	LA	EMEA	ASIA	Oceania	Japan
1	1920–1980	2110–2170						
2	1850–1910	1930–1990						
3	1710–1785	1805–1880						
4	1710–1755	2110–2155						
5	824–849	869–894						
6	830–840	875–885						
7	2500–2570	2620–2690						
8	880–915	925–960						
9	1749.9–1784.9	1844.9–1879.9						
10	1710–1770	2110–2170						
11	1427.9–1447.9	1475.9–1495.9						
12	698–716	728–746						
13	777–787	746–756						
14	788–798	758–768						
15–18	reserved							
19	832.4–842.6	877.4–887.6						
20	832–862	791–821						

The WCDMA and LTE evolution is defined by releases beginning with an initial release in 2000 known as Release 99 (Rel-99) [26]. This saw the beginning of the Third Generation Partnership Project (3GPP) specifying and controlling the evolution of cellular communications through 3G, 4G, and now 5G. The releases are designed to protect the installed investment in cellular systems while providing a migration path.

#### 5.9.4 Summary

WCDMA 3G and 4G are widely deployed. While 3G/WCDMA it will be gradually replaced by 4G and 5G it has particular aspects that could see the system remain for many years.



**Figure 5-17:** Timeline for implementation of 3G, 4G and 5G. DL indicates the downlink data rate; UL indicates the uplink data rate; BW indicates the channel bandwidth. Development supports Internet protocol (IP) and voice over IP (VoIP). The two 3G paths become a single long term evolution (LTE) path. LTE is the concept of smoothly evolving through 4G and into 5G utilizing existing infrastructure and adding capability [25]). Beginning with Release 99 (Rel-99) the timeline is controlled by 3GPP . The dates refer to the first significant commercial availability of the standard and the official release dates can be found at <http://www.3gpp.org>.

Band	Uplink (MHz)	Downlink (MHz)
33	1900–1920	1900–1920
34	2010–2025	2010–2025
35	1850–1910	1850–1920
36	1930–1990	1930–1990
37	1910–1930	1910–1930
38	2570–2620	2570–2620
39	1880–1920	1880–1920
40	2300–2400	2300–2400

**Table 5-11:** Spectrum assignments for TDD 3GPP [25, Release 8].

## 5.10 4G, Fourth Generation Radio

The 4G cellular system provides downlink data rates of 100 Mbit/s while mobile and 1 Gbit/s while stationary. The uplink data rates are much lower at 50 Mbit/s. The maximum rates are assigned to a cellular service area and are shared but rely on a user's average data rate being much lower. Providing these data rates relies on a number of significant advances over 3G:

- Orthogonal frequency division multiplexing (OFDM). Using a very large number of narrowband channels, e.g. 1200, with each narrowband channel having the highest-order modulation possible. Each narrowband channel has its own subcarrier which is individually modulated. Gets around the multipath problem where poor channel characteristics affecting a narrow range of frequencies does not limit overall bit rates.
- High-order modulation. A terminal unit supports a very large number of modulation formats, several PSK modulation methods and QAM methods ranging from 16-QAM to 1024-QAM. Modulation and demodulation are implemented in a DSP unit.
- Cyclic prefix (CP). Enables the use of efficient discrete Fourier transforms in the transmitter. Simplifies hardware by yielding lower PMEPR than would be obtained if a zero-level guard band was provided to overcome interference caused by multipath.
- Multiple input, multiple output (MIMO). MIMO uses multiple antennas to transmit and receive signals exploiting independent transmit-receive paths to increase the total capacity between a basestation and terminal unit. Multiplies data rates by the number of transmit or receive antennas (which ever is less).
- Pervasive access. Terminal units support many access technologies (e.g., 2G–4G, WiFi, Bluetooth etc.) and seamlessly switching among them. Since 4G is an IP-based system data can be sent on any network, e.g. WiFi utilization reduces the demands on the cellular network.

These advances require tremendous computing power made available by advanced low-power VLSI and developments in low-power RF modems. Another critical aspect of 4G (and now 5G) is deployment of a well thought out road map called Long Term Evolution (LTE) that enables incremental (although significant) advances building on previously implemented enhancements.

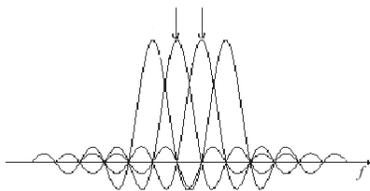
The LTE timeline is shown in Figure 5-17, 4G effectively began with Release-9 implemented first at the end of 2010 with this release introducing MIMO and OFDM (to be defined latter). There are fundamental changes to the physical interface, and 4G uses Internet Protocol (IP) exclusively including for voice (voice over IP—VoIP). Subsequent releases introduced more concepts and increased data coverage. 4G and now 5G can not coexist with 3G systems and separate bands must be allocated. Eventually 3G will go away. It is worth noting that 5G is upwards compatible with 4G.

The ITU World Radiocommunication Conference has allocated many bands to 4G ranging from 450 MHz to 5850 MHz. There is not a single band that can be used world-wide but cellular handsets are designed to support

**Table 5-12:** Selected LTE bands used in 4G.

Band	Mode	Uplink (MHz)	Downlink (MHz)	Bandwidths (MHz)	Regions
3	FDD	1710–1755	1805–1880	1.4, 3, 5, 10, 15, 20	Asia, Parts of Africa, Europe, Parts of Latin America, Oceania
12	FDD	699–716	729–746	1.4, 3, 5, 10	North America, Asia, Parts of Africa, Europe, Parts of Latin America
40	TDD	2300–2400		5, 10, 15, 20	Asia, Parts of Europe, Oceania
41	TDD	2496–2690		5, 10, 15, 20	Parts of Africa, Parts of Asia, USA

Bandwidth	Resource blocks	Subcarriers (downlink)	Subcarriers (uplink)
1.4 MHz	6	73	72
3 MHz	15	181	180
5 MHz	25	301	300
10 MHz	50	601	600
15 MHz	75	901	900
20 MHz	100	1201	1200

**Table 5-13:** Resources for different LTE bandwidths.**Figure 5-18:** OFDM spectrum with four subcarriers showing orthogonality.

multiple bands. The features of representative bands are shown in Table 5-12. Channel bandwidth, with a channel having a single carrier, range from 1.4 MHz to 20 MHz, see Table 5-13. Each carrier has correlated subcarriers each modulated to produce a subchannel with a 15 kHz-wide bandwidth.

### 5.10.1 Orthogonal Frequency Division Multiplexing

In OFDM, data is divided into many bitstreams with each bitstream modulating a **subcarrier** producing a relatively narrowband subchannel and a user being assigned multiple subchannels. In 4G these subchannels are 15 kHz-wide in normal operation and 7.5 kHz wide in a rich multipath environment with a large excess delay spread as described in Section 4.7. The concept behind this is that having a narrow bandwidth, the duration of a symbol is long and will exceed the excess delay spread. This minimizes the impact of intersymbol interference (ISI) when a symbol traveling on a fast path overlaps with the component of a previous symbol traveling on a longer path. Also each subchannel can be modulated with the maximum order modulation enabled by that particular channel's characteristics. For this to work the subcarriers must be orthogonal and precisely spaced in frequency and time. The orthogonality is shown in the spectrum of Figure 5-18, where the arrows at the top indicate sampling points for two subcarriers. The orthogonality of the subcarriers is seen by noticing that the peak of one subcarrier is at the zeros of the other subcarriers. When one subcarrier is sampled, the contribution from all other carriers is zero; they are orthogonal.

The spectra of the subcarriers overlap, but this does not matter.

In 4G groups of 12 subcarriers lasting 0.5 ms are grouped in a resource block which is the minimal allocation to a user. In communications with a handset, a basestation (called a B node in 4G and 5G) allocates a specific number of subcarriers. These may not be contiguous on the downlink, as shown in Figure 5-8(d), but are contiguous on the uplink. The number of subcarriers allocated to a particular user varies according to the data rate requirements and the order of modulation used with each subcarrier being adjusted to accommodate the subchannel characteristics including fading and interference.

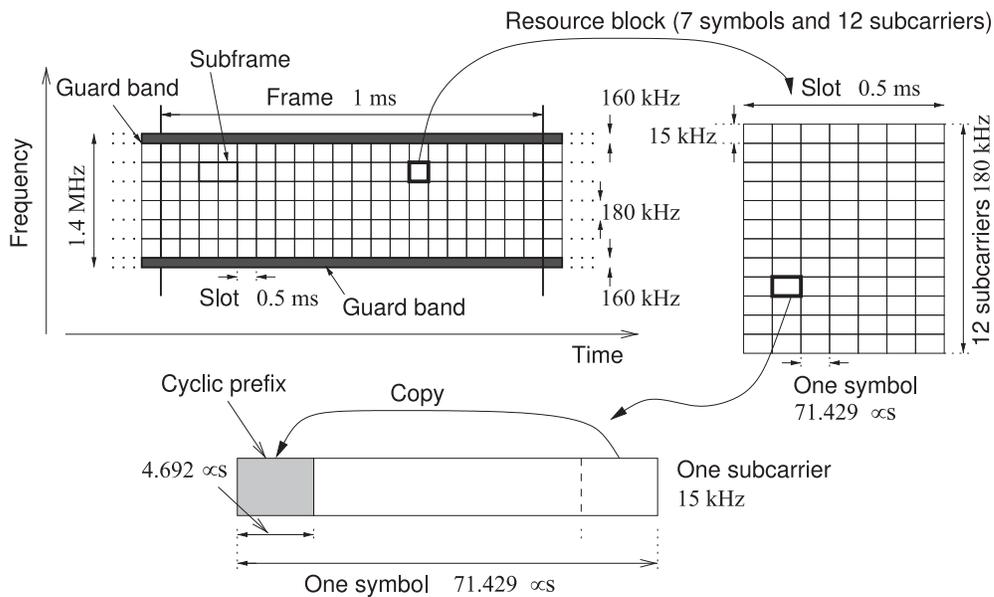
In communicating to a single user, the total bitstream is error encoded and then the bitstream is divided up so that different parts of the bitstream are sent over multiple subchannels thus spreading the coded data. As a result 4G implements spectrum spreading but at a much lower coding rate than the spreading in 3G. This further mitigates the impact of multipath which can affect the integrity of individual subchannels. Signal strength and interference, and hence SIR, can differ for each channel and this is compensated for by having different bit rates in each subchannel and spreading the data plus error correction coding.

A second effect of multipath is referred to as delay spread. Transmitted signal components following different paths arrive at a receiver at different times. If the time duration of a symbol is short, which is the case if wideband modulation is used, then the components of a signal, corresponding to one symbol, following longer paths could arrive at the same time as the components of the next symbol following shorter paths. This causes intersymbol interference. With the longer symbol duration with OFDM and a guard band interval (which is now relatively short) between symbols, intersymbol interference can be greatly reduced.

### Implementation

OFDM could be implemented by using separate modulators and demodulators for each subcarrier. Instead the separate modulators and demodulators are replaced by a **fast Fourier transform (FFT)** and an **inverse FFT (iFFT)**, respectively, implemented in a DSP unit called the baseband processor. In the transmitter the iFFT produces a modulated signal at a low intermediate carrier frequency. For example, a 1.4 MHz bandwidth OFDM signal could be a DSB-SC signal with a center frequency, i.e. IF carrier frequency, of 700 kHz but otherwise look identical to the final transmitted OFDM signal. Then a SSB-SC up-converter translates the OFDM signal to a radio frequency signal. Each of the time-varying frequency inputs of the iFFT is centered at a subcarrier frequency. Each of the frequency inputs is a sequence of time-samples of a slowly varying modulated signal implementing the selected modulation method for that subchannel. On receive the input of the FFT is the sequence of time samples of the modulated signal centered on the IF carrier which has been down-converted from RF. Each of the outputs of the FFT correspond to a time sample of an individual modulated subcarrier.

Since the OFDM signal combines many individually modulated subcarriers the PMEPR of OFDM is large but grows relatively slowly with the number of subcarriers. The higher PMEPR, compared to a single modulated carrier, introduces microwave design challenges especially for mixers and am-



**Figure 5-19:** Components of a frame in 4G for a 1.4 MHz-bandwidth channel. The 10 ms-long frame, top left, has an array of resource blocks each with a slot duration of 0.5 ms and a bandwidth of 180 kHz. This frame has a 1.08 MHz data bandwidth with 160 kHz guard bands for a total bandwidth of 1.4 MHz. (The guard-band differs for channels having other bandwidths. Each resource block has on seven OFDM symbols with each comprised of seven symbols during a symbol interval. There are 84 symbols (Seven OFDM symbols) per resource block (with 1024-QAM having 10 bits per subcarrier 840 bits are transmitted per resource block).

plifiers. However, since most of the heavy lifting is done in DSP which being numerical is distortion-free, the up-converter is simplified .

**Resource Blocks**

The way that OFDM is implemented in 4G is illustrated in Figure 5-19 for a 1.4 MHz-bandwidth modulated signal. The subchannels in 4G are organized in a frame that is 10 ms long. Within the frame subchannels are grouped into **resource blocks** each of which has a duration of 0.5 ms and a bandwidth of 180 kHz. The resource block is the smallest unit allocated to a user and a user can have many simultaneous resource blocks and many sequential resource blocks. A resource black has many subcarriers and symbol intervals and the exact length and bandwidth of a single symbol-subcarrier segment has a normal mode used most of the time, and two extended modes that handle severe multipath environments. In the normal mode, each resource block has seven 71.4 μs-long symbols per subcarrier and there are 12 subcarriers. (The extended modes have longer symbol duration and 7.5 kHz bandwidth.)

In OFDM there are subchannels that are reserved for pilot tones that send known data streams and these can be used to characterize the channel. The pilot subchannels are denoted as P in Figure 5-8(d). The dedicated pilot subchannels enable carrier recovery. Therefore it is not necessary to restrict modulation by avoiding transitions through the origin of the constellation

diagram (or more precisely the phasor of the modulated subcarrier). Thus near-ideal modulation efficiency is possible for the data subcarriers.

### Summary

OFDM has been tremendously successful and only had to wait for the capabilities of low-power VLSI to advance. It is now widely used in communications including WiFi, wired communications (e.g. digital subscriber line (DSL)) as well as in 4G and 5G.

### 5.10.2 Orthogonal Frequency Division Multiple Access

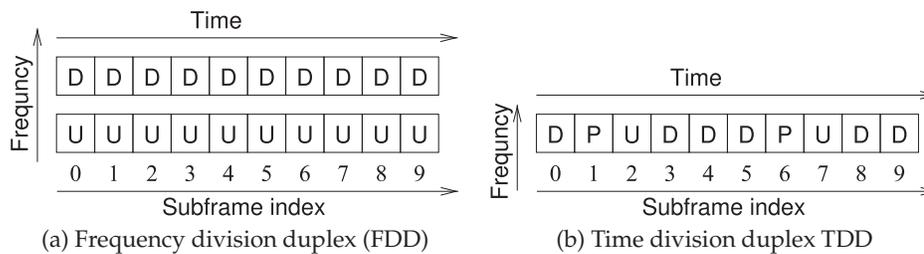
Orthogonal Frequency division multiple access (**OFDMA**), also known as **multiuser OFDM**, is the multiuser version of OFDM that supports simultaneous communication by multiple users and so is an access technology. The access scheme is illustrated in Figure 5-8(d) where one channel is illustrated and could have a bandwidth ranging from 1.4 MHz to 20 MHz. There could be 1201 subchannels so not all of these are shown. The subchannels are grouped in contiguous groups of 12 subchannels to form a 1 ms-long resource block (24 7.5 kHz-wide subchannels in extended mode used with a very rich multipath environment). So what are numbered in Figure 5-8(d) are resource blocks. In OFDMA, generally one or more resource blocks are assigned to a particular user and these do not need to be contiguous in the downlink but are contiguous in the uplink as then the PMEPR is lower

This lessens the demand on the power amplifier in the terminal unit. Effectively there is a single carrier on the uplink and hence the uplink access scheme is also called **single-carrier frequency division multiple access (SC-FDMA)** or **single-carrier orthogonal frequency division multiple access (SC-OFDMA)**. On the downlink the basestation is using all of the subchannels simultaneously communicating with multiple terminal units but in the uplink each terminal unit only uses a few. Users share pilot subchannels, designated as P in Figure 5-8(d).

### 5.10.3 Cyclic Prefix

The cyclic prefix (**CP**) refers to the scheme used to accommodate excess delay spread resulting from multipath. Since symbols sent on different paths arrive at the receive antenna with various delays, there can be overlap of a symbol arriving on a fast path with a previous symbol arriving on a longer path. This could result in intersymbol interference (ISI). In 4G there could have been a zeroed (i.e. no signal) guard interval to eliminate ISI. Instead, in 4G, a cyclic prefix (**CP**) is used whereby a symbol is prefixed repeating the end part of a symbol. One symbol of a resource block is shown at the bottom of Figure 5-19. Each symbol of the resource block is 66.7  $\mu\text{s}$  long with a cyclic prefix which normally is 4.7  $\mu\text{s}$  long but longer for the first symbol of the resource block and when in extended mode which is used in severe multipath environments.<sup>5</sup> The 4.7  $\mu\text{s}$  long CP corresponds to a

<sup>5</sup> In the normal mode the first CP is 5.2  $\mu\text{s}$  long and there are two extended modes. One is 16.7  $\mu\text{s}$  long and another is 33.3  $\mu\text{s}$  long and then the bandwidth of each subchannel is reduced to 7.5 kHz. In normal mode the overhead is 7% and in extended mode the CP overhead is 25%. If the first extended CP is used then there are 6 OFDM symbols per resource block.



**Figure 5-20:** Duplex schemes used in 4G with ten subframes in a 10 ms frame.

difference in path lengths of 1.4 km. This cyclic prefix provides the required guard interval. There are several advantages to copying the end of a symbol and repeating it before the symbol. This reduces the PMEPR that would result if the guard interval had no signal, and it enables a discrete Fourier transform to be used in processing the signal. This is essential to efficient VLSI implementation of OFDM.

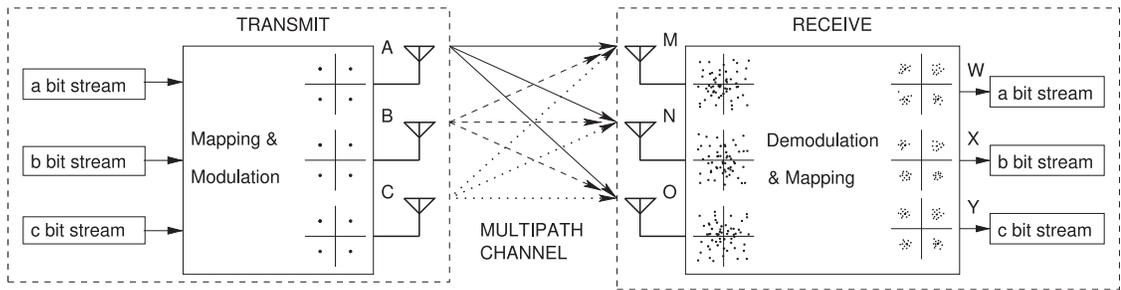
The repetition of the cyclic prefix for each symbol in a resource block is equivalent to repeating the end of each OFDM symbol at the start of an OFDM symbol (the aggregate of symbols across all 12 subcarriers constitutes an OFDM symbol). The repetition enables the start of each OFDM symbol to be determined as the CP is correlated to the end of the symbol. The cyclic prefix also helps in characterizing the channel and removing distortion within the OFDM symbol

#### 5.10.4 FDD versus TDD

Frequency division duplex (FDD) and time division duplex (TDD) are two duplex schemes supported in 4G although only one mode or the other is supported in a particular band, e.g. see Table 5-12. FDD and TDD are arranged in frames of 10 ms comprising ten subframes of 1 ms duration, see Figure 5-20. In FDD, Figure 5-20(a), there are separate uplink, denoted U, and downlink, denoted D, frequency bands and uplink and downlink transmissions are simultaneous. The use of paired spectrum requires a good diplexer to isolate the receiver and transmitter. In TDD, Figure 5-20(b), uplink and downlink use the same band and there are separate uplink, denoted U, and downlink, denoted D, transmissions in different frequency bands. In TDD, uplink and downlink transmissions are simultaneous. Channel propagation is the same in both directions at least over 10 ms as long as mobility is not too high. TDD allows dynamic allocation of uplink and downlink capacity.

#### 5.10.5 Multiple Input, Multiple Output

Multiple input, multiple output (MIMO, pronounced my-moe) technology uses multiple antennas to transmit and receive signals. The MIMO concept was developed in the 1990s [27, 28] and implemented in 4G and 5G, and several WLAN systems. There are several aspects to MIMO. First, each transmit antenna sends different data streams simultaneously on the same frequency channel as other transmit antennas. The most interesting feature



**Figure 5-21:** A MIMO system showing multiple paths between each transmit antenna and each receive antenna.

is that MIMO relies on signals traveling on multiple paths between an array of transmit antennas and an array of receive antennas. In a conventional communications system the various paths result in interference and fading, but in MIMO these paths are used to carry more information. In a MIMO system, each path propagates an image of one transmitted signal (from one antenna) that differs in both amplitude and phase from the images following other paths. Effectively there are multiple connections between each transmit antenna and each receive antenna, see Figure 5-21. For simplicity, three transmit antennas and three receive antennas are shown. However, MIMO can work with as few as two transmit antennas and two receive antennas. In MIMO a high-speed data-stream is split into several slower data streams, shown in Figure 5-21 as the a, b, and c bitstreams. The distinct bitstreams are separately modulated and sent from their own transmit antenna, with the constellation diagrams of the transmitted modulated signals labeled A, B, and C. The signals from each of the transmit antenna reaches all of the receive antennas by following different uncorrelated paths.

The output of each receive antenna is a linear combination of the multiple transmitted data streams, with the sampled RF phasor diagrams labeled M, N, and O. (It is not really appropriate to call these constellation diagrams.) That is, each receive antenna has a different linear combination of the multiple images. In effect, the output from each receive antenna can be thought of as the solution of linear equations, with each transmit antenna-receive antenna link corresponding to an equation. Continuing the analogy, the signal from each transmit antenna represents a variable. So a set of simultaneous equations can be solved to obtain the original bitstreams. This is accomplished by demodulation and mapping using knowledge of the channel characteristics to yield the original transmitted signals modified by interference. The result is that the constellation diagrams W, X, and Y are obtained. The composite channel can be characterized using known test signals. Special coding called **space-time** (or **spatio-temporal**) coding embedded in the transmitted data-stream also enables estimation of the communication matrix. Space-time coding encodes each transmitted data-stream with information that can be used to update the channel characterization.

The capacity of a MIMO system with high SIR scales approximately linearly with the minimum of  $M$  and  $N$ ,  $\min(M, N)$ , where  $M$  is the number

Modulation scheme	Capacity bit/s/Hz (bits per second per hertz)				
	Non-MIMO $M = 1, N = 1$ maximum	MIMO with $M = 2, N = 2$			
		SIR 0 dB	SIR 10 dB	SIR 20 dB	SIR 30 dB
BPSK	1	1.2	2	2	2
QPSK	2	1.6	3.7	4	4
8-PSK	3	1.6	4.8	6	6
16-PSK	4	1.6	4.9	7.5	8

**Table 5-14:** Capacity of MIMO schemes with PSK modulation for different received SIRs compared to the maximum capacity of a conventional non-MIMO scheme.  $M$  is the number of transmit antennas,  $N$  is the number of receive antennas. Data from [31].

of transmit antennas and  $N$  is the number of receive antennas (provided that there is a rich set of paths) [29, 30]. So a system with  $M = N = 4$  will have four times the capacity of a system with just one transmit antenna and one receive antenna. Table 5-14 presents the capacity of a MIMO system with ideal PSK modulation and two transmit and two receive antennas. This is compared to the capacity of a conventional (non-MIMO) system. The capacity is presented in bits per second per hertz and it is seen that significant increases in throughput are obtained when SIR is high. MIMO is a successful way to increase capacity and is included in modern WiFi, radars, and other communication systems.

In summary, MIMO systems achieve throughput and range improvements through four gains achieved simultaneously:

1. Array gain resulting from increased average received SIR obtained by coherently combining signals. To exploit this the channel must be characterized. This increases coverage and quality of service (QoS).
2. Diversity gain obtained by presenting the receiver with multiple identical copies of a given signal. This combats fading. This also increases coverage and QoS.
3. Multiplexing gain by transmitting independent data signals from different antennas to increase throughput. This increases spectral efficiency.
4. Cochannel interference reduction. This increases cellular capacity.

### 5.10.6 Carrier Aggregation

Carrier aggregation (CA) is one of the main features that 4G introduced. In carrier aggregation bitstreams on different carriers are combined to yield a higher overall bit rate than one carrier can support. This enables short-term high bit rates to be transmitted to a terminal. The full concept supports combining of bit-streams with carriers that are in different frequency bands (called inter-band CA), the same frequency band (called intra-band CA), different cells, and even in a combination of licensed and unlicensed (think WiFi) bands. It is 4G LTE-A combining carrier aggregation and MIMO that achieves, and can even surpass, the peak 1Gbit/s goal of 4G. A typical goal of 4G LTE A is to combine five carriers. The 4G standards, 3GPP release 13, supports 32-carrier aggregation and up to 3 Gbit/s.

### 5.10.7 IEEE 802.11n

The IEEE 802.11n WiFi system is not 4G but this is discussed here to indicate that many of the concepts that are incorporated in wireless systems

**Table 5-15:** Data rates used in the 802.11n standard. MCS is the modulation and coding scheme index. GI is the guard interval. With a 20 MHz-wide channel there are 52 data subcarriers, and an additional 4 subcarriers, called pilots, enable carrier recovery and monitor the channel characteristics. With a 40 MHz-wide channel there are 108 data subcarriers and 6 pilot subcarriers.  $R_i$  = information bit rate,  $R_c$  = coded bit rate.  $R_c$  is the bit rate of information bits plus bits added for error correction. Thus a coding rate of 5/6 means that for every 5 information bits there is an additional (redundant) code bit.

MCS index	Spatial streams	Modulation type	Coding rate $R_i/R_c$	Data rate (Mbit/s)			
				20 MHz channel		40 MHz channel	
				800 ns GI	400 ns GI	800 ns GI	400 ns GI
0	1	BPSK	1/2	6.50	7.20	13.50	15.00
1	1	QPSK	1/2	13.00	14.40	27.00	30.00
2	1	QPSK	3/4	19.50	21.70	40.50	45.00
3	1	16-QAM	1/2	26.00	28.90	54.00	60.00
4	1	16-QAM	3/4	39.00	43.30	81.00	90.00
5	1	64-QAM	2/3	52.00	57.80	108.00	120.00
6	1	64-QAM	3/4	58.50	65.00	121.50	135.00
7	1	64-QAM	5/6	65.00	72.20	135.00	150.00
8	2	BPSK	1/2	13.00	14.40	27.00	30.00
9	2	QPSK	1/2	26.00	28.90	54.00	60.00
10	2	QPSK	3/4	39.00	43.30	81.00	90.00
11	2	16-QAM	1/2	52.00	57.80	108.00	120.00
12	2	16-QAM	3/4	78.00	86.70	162.00	180.00
13	2	64-QAM	2/3	104.00	115.60	216.00	240.00
14	2	64-QAM	3/4	117.00	130.00	243.00	270.00
15	2	64-QAM	5/6	130.00	144.40	270.00	300.00
16	3	BPSK	1/2	19.50	21.70	40.50	45.00
17	3	QPSK	1/2	39.00	43.30	81.00	90.00
18	3	QPSK	3/4	58.50	65.00	121.50	135.00
19	3	16-QAM	1/2	78.00	86.70	162.00	180.00
20	3	16-QAM	3/4	117.00	130.70	243.00	270.00
21	3	64-QAM	2/3	156.00	173.30	324.00	360.00
22	3	64-QAM	3/4	175.50	195.00	364.50	405.00
23	3	64-QAM	5/6	195.00	216.70	405.00	450.00
...	4	...	...	...	...	...	...
31	4	64-QAM	5/6	260.00	288.90	540.00	600.00

first appear in WiFi. IEEE 802.11n is an example of a data communication system that combines many advanced concepts. The 802.11n system uses MIMO, OFDM, and many dedicated subcarriers for continuous channel characterization and carrier recovery. The 802.11n standard uses several modulation types, as shown in Table 5-15, and it supports the use of either a 20 MHz- or 40 MHz-wide channel. There can be one, two, three or four spatial streams and these are equal to the minimum of the number of transmit and receive antennas. If the number of transmit or receive antennas is more than the number of spatial streams, the extra antennas are used for receiver or transmit diversity. Under the best conditions with high SNR, negligible correlation of the transmit-to-receive antenna paths, and four transmit and four receive antennas, the **modulation and coding scheme (MCS)** index 31 is used. With this index, modulation uses 64-QAM and 5/6

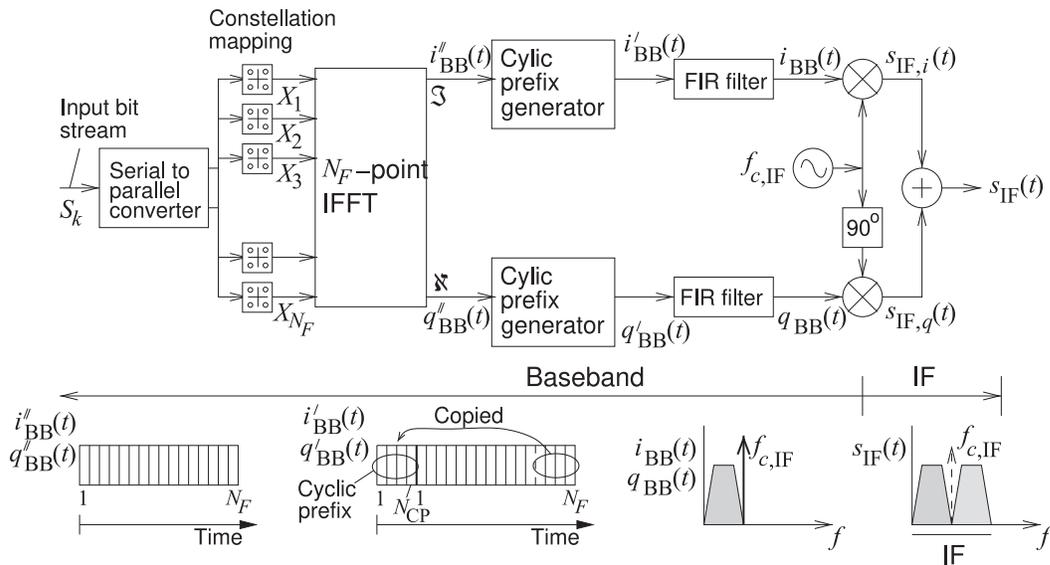


Figure 5-22: Block diagram of an OFDM modulator.

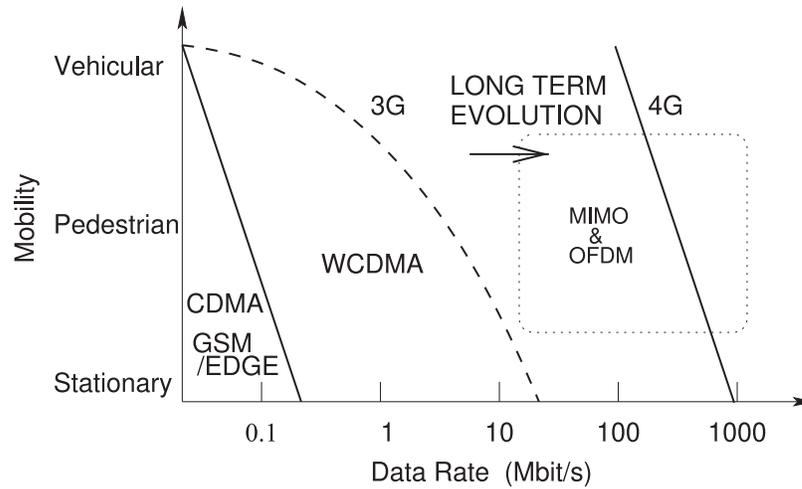
coding. This mode supports a data rate of 600 Mbit/s if a 400 ns guard interval is used. This guard interval, rather than the 800 ns guard interval, is used when there is low delay spread. If the SNR is lower, a lower-order modulation scheme is used.

### 5.10.8 OFDM Modulator

The basic structure of an OFDM modulator is shown in Figure 5-22. From the perspective of a single user, a bitstream  $S_k$  is divided up by a serial-to-parallel converter to produce multiple slower bitstreams each of which is mapped onto complex signals  $x_i$  which become the frequency inputs of an iFFT. The outputs of the iFFT are real and imaginary time-domain baseband signals  $i'_{BB}(t)$  and  $q'_{BB}(t)$ . At each symbol interval  $(i'_{BB}(t), q'_{BB}(t))$  indicates a symbol. The  $i'_{BB}(t)$  and  $q'_{BB}(t)$  pass through a cyclic prefix generator that copies the end of the symbol and prefixes it to the beginning of the symbol. Then the output of the cyclic prefix generators are shaped by an FIR filter (generally a raised cosine filter) to produce the shaped I/Q signals  $i_{BB}(t)$  and  $q_{BB}(t)$  which are input to a quadrature modulator with an intermediate carrier frequency to produce a DSB-SC intermediate frequency signal modulated signal  $s_{IF}(t)$ . For example, if the frequency range of  $i_{BB}(t)$  and  $q_{BB}(t)$  range from just above DC to just below 700 kHz the choice of  $f_{c,IF} = 700$  kHz results in a 1.4 MHz wide modulated signal  $s_{IF}(t)$  which is a DSB-SC OFDM modulated signal. The entirety of the OFDM modulator in Figure 5-22 is implemented digitally in a DSP unit.

### 5.10.9 Summary of 4G

At the physical layer the 4G standard introduces and combines OFDM, carrier aggregation (CA), and MIMO to achieve tremendous data rates that,



**Figure 5-23:** Data rate capacity of 4G as the long term evolution of 3G.

in a fully implemented 4G system achieves low-latency data download at mobile rates up to 100 Mbit/s and rates while stationary of 1 Gbit/s. OFDM is the technology that sends many relatively slow bitstreams over narrow bandwidth sub-channels to efficiently support multiple users making best use of the channels available and circumventing many multipath effects. It replaces the use of spreading codes in CDMA to support multiple access and circumvent some multipath problems. CA combines multiple bitstreams sent on different carriers and even from different basestations to increase the short term overall bit rate. MIMO uses multiple, i.e.  $M$ , transmit antennas, at the basestation, and multiple, i.e.  $N$ , receive antennas, at the handset, and exploits uncorrelated communication paths between pairs of transmit and receive antennas. The (ideally) zero correlation is made possible by multipath and the correlation is lowest if there is no line-of-sight path. If the de-correlation is complete the data capacity is increased by a function of the minimum of the number of transmit antennas, i.e.  $\text{MIN}(M, N)$ . If there is some correlation, perhaps due to coupling between the receive antennas, then a MIMO capacity factor  $H \leq \text{MIN}(M, N)$  is defined. With MIMO it is necessary to modify Shannon's capacity limit as MIMO systems can exceed the limit defined for a single channel. Shannon's capacity limit for a MIMO system becomes [29]

$$\hat{C} = B_c \log_2(1 + \text{SIR} \cdot H). \quad (5.31)$$

This indicates that the channel carrying capacity is greatly increased, especially when the SIR is high. The data rate capacity of 2G systems (GSM and CDMA) is contrasted with the capacity of 3G (WCDMA) and 4G systems in Figure 5-23.

## 5.11 5G, Fifth Generation Radio

As with 4G there will be long term evolution of the 5G standard and this will be upwards compatible with 4G. Even though 4G was touted as the long term evolution of 3G, at the physical layer it was fundamentally incompatible. 5G is fundamentally compatible with 4G at the physical layer and it is expected that cellular phone services will continue to function much the same as if 4G was being used. The first iteration (i.e. evolution) of 5G is called Next Radio (NR) and begins with 3GPP Release 15, see Figure 5-17, and operation began in the second half of 2018.

The full 5G standard supports three use cases and a significantly improved level of service over 4G. One of the uses cases is **enhanced mobile broadband (eMBB)** with peak data rates of 20 Gbit/s and sustained data rates of 100 Mbit/s. High mobility of up to 500 km/s is supported. Another use case is **massive machine-to-machine communication (mMTC)** supporting the **internet of things (IoT)** with high density of devices (up to  $10^6/\text{km}^2$ , long range, low data rate (1 kbit/s–100 kbit/s), and enabling 10 year battery life. The third use case is ultra reliability and low latency (URLLC) with less than 1 ms air interface latency and less than 5 ms overall end-to-end latency, 99.9999% reliability, data rates up to 10 Mbit/s, and supporting high mobility.

Many scenarios require ultra-reliable and very low end-to-end latency, perhaps as low as 1 ms. The early deployments of 4G had latencies of 80 ms or more and varied considerably across service providers. Over the years since 4G was launched latencies have reduced considerably but in 4G it is not possible to have air-interface latencies less than the 10 ms frame rate of 4G. A fundamental limit in 4G is due to the 0.5 ms time duration of a slot, see Figure 5-19. One concept in 5G is to have more slots per 1 ms subframe maintaining the subframe for upwards compatibility. In 4G there are two slots per subframe. Possible developments in 5G are 1) to have more slots per subframe and 2) (what seems to be similar) to have more symbols per slot. With 16 slots per subframe the slot duration is 0.0625 ms. The long-term evolution of 5G will continue and achieve overall latencies of 1 ms or less.

An important concept of 5G is to be fully compatible with the 4G infrastructure and the only new infrastructure required is that to support new services. There were two central ideas behind 4G: OFDM and MIMO and these are used in 5G OFDM. OFDM in 4G used two fixed sub-carrier spacing, but 5G has a number of different sub-carrier bandwidths including a subchannel bandwidth of 480 kHz, and a low bandwidth subchannel bandwidth of 7.5 kHz to accommodate high speed mobility. As with 4G, carrier aggregation is used but now the bitstreams on up to 32 carriers are combined. Also channel bandwidths are up to 400 MHz compared to the 20 MHz RF bandwidth maximum of 4G.

### 5.11.1 Mesh Radio

In 5G a wireless mesh removes the need for a fixed basestation to communicate directly to an end unit [32]. If there is a node between the basestation and a mobile user, an intervening node, possibly another mobile user, can be used effectively as a relay. This reduces overall power requirements and levels of interference, which in turn leads to greater data

carrying capacity. This concept can be extended to use multiple intervening nodes forming a dense mesh with considerable tolerance to multipath and interference effects. One significant benefit of this can be seen by noting that in an urban environment power can fall off by the fourth power of distance. Another benefit is that the impact of fading can be greatly reduced, as the paths in the mesh will fade independently. This will be augmented by many small cells called femtocells and picocells handling traffic in small geographic areas such as buildings, airports, and sporting facilities.

### 5.11.2 Cognitive Radio

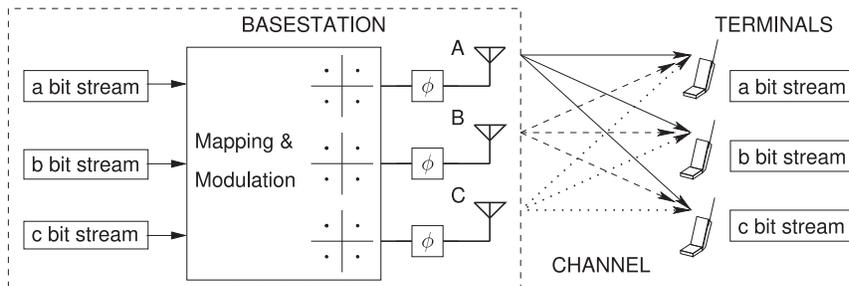
Cognitive radio exploits the fact that much of the EM spectrum remains unused even while the bands reserved for cellular communications have reached capacity. This situation is a consequence of the allocation of bands for dedicated services that may not be active at a particular time and place. Examples are unused bands reserved for broadcast television channels and bands reserved for communication in times of emergencies. A cognitive radio senses its environment and adapts in real time to a user's communication needs by temporarily borrowing unused spectrum [33, 34]. As a result spectrum is used more efficiently. A cognitive radio avoids causing interference with the communications of other users by sensing spectrum use, changing frequency, adjusting power level, and altering transmission protocols. If a licensed band is borrowed, then a cognitive radio discontinues use of this part of the spectrum if the licensee of the band becomes active. This behavior is also called **dynamic spectrum management**.

### 5.11.3 Massive MIMO

MIMO in 4G uses multiple antennas at the basestation and at terminals to increase overall data rates provided that there are multiple paths between the transmitter and receiver. For example, with two basestation antennas transmitting separate bitstreams but in the same frequency band and two uncorrelated receive antennas, known as 2x2 MIMO, the theoretical maximum data rate is twice what could be supported by line-of-sight communication. In 5G there can be a very large number of transmit antennas even though there are few receiver antennas per terminal unit as multiple terminal units mean that there can effectively be a very large number of receive antennas.

5G operates at frequencies below 6 GHz and at millimeterwave frequencies. Below 6 GHz 8x8 MIMO is supported in the standard but there can be a hundred or more basestation antennas. For example, with 128 basestation antennas and four receiver antennas per terminal unit, 32 terminal units can be supported in the same frequency band. Overall the system has much higher throughput, theoretically a maximum of 128 times more than a single line-of-sight link could support. This concept, i.e. not all of the receive antennas need be on the same platform, is called massive MIMO.

Massive MIMO uses characterization of the channel as in conventional MIMO but uses this information to form multiple beams each directed at each terminal unit. This situation is shown in Figure 5-24 where there are multiple transmit/receive antennas at each basestation but relatively



**Figure 5-24:** A massive MIMO system with beams from the basestation antennas to each terminal unit.

few transmit/receive antennas at each terminal unit. At the start of each correlation interval, i.e. the time over which a channel does not change significantly, each terminal unit transmits a code to the basestation. Each terminal unit uses an orthogonal code and so the basestation is able to characterize the channels from each of the terminal units to the basestation. The basestation is then able to use its antennas in a phased array to form a beam to transmit data just to the intended terminal. The real situation is a little different as a single beam is not formed but instead a weighted signal is broadcast from each of the basestation antennas with the effect being that all of the power following multiple paths comes together at the intended terminal unit. Still this process is called beamforming. Different beams are used for each basestation. In the second half of the correlation interval the terminal unit transmits back to the basestation and the basestation uses its beamforming in reverse to effectively create multiple receive beams each directed at a terminal units. Of course these are not strictly beams but use the channel model to combine the signals from each of the basestation receive antennas to effectively receive signals with the highest signal-to-interference ratio from each terminal.

Massive MIMO has a performance advantage even with line-of-sight as long as the terminal units are not on top of each other as the paths from each of the basestation antennas to each of the terminals will be uncorrelated. Beam steering can be also used to concentrate transmitted energy in a particular direction.

#### 5.11.4 Active Antenna Systems

**Adaptive antennas**, or **active antenna systems (AAS)**, are also known as **smart antennas** and adapt an array of antennas to either direct an antenna beam between a transmitter and a receiver, eliminate interference, or use the **spatial signatures** provided by diverse propagation paths to transfer more information than can be sent over a single link [35]. All this is done using signal processing rather than hardware reconfiguration.

In one approach a smart antenna is used to switch between a number of fixed, usually narrow, antenna beams in what is analogous to a highly sectored antenna array. Switching requires knowledge of the angle of arrival and the beams track a mobile unit. A simple version of this concept is used in 3G and 4G radio. As well as focusing the available power on the intended communication nodes, beam forming through beam switching also minimizes interference.

Adaptive antennas increase data transfer rates, reduce interference, and reduce the amount of transmit power required.

### 5.11.5 Microwave Frequency Operation

Low frequency operation refers to 5G operating below 6 GHz. The bands being adopted for early stage deployment are the 700 MHz, 3300–4200 MHz and 4400–5000 MHz bands but the full coverage is not available in all countries. Many are predicting that the bulk of 5G communication will be at these frequencies, particularly at 3.5 GHz, up until the mid or late 2020s. It is important to 5G that there be globally available spectrum. The 700 MHz band is of particular interest as it will provide wide area coverage and deep penetration into buildings and so is a candidate for providing high reliability links important for the 5G mMTC applications. At this frequency 5G is compatible with 4G but with evolved characteristics such as supporting up to 32 carrier aggregation and  $8 \times 8$  MIMO instead of the five carrier aggregation and (usually) maximum  $3 \times 3$  MIMO of 4G. Also 5G in the low frequency range supports 256 QAM modulation but then the order of QAM modulation supported with 4G has evolved with latter 3GPP releases.

### 5.11.6 Millimeter-Wave Operation

The very high data rates of 5G are obtained by operating at millimeter-wave (mm-wave) frequencies where much more bandwidth is available. The bands being focused on are the 24.25–29.5 GHz and 37–40.5 GHz. (Yes, technically 24.25–29.5 GHz is not millimeter-waves which requires a wavelength to be 10 mm or less, but that is being called mm-waves in the 5G community.) At mm-wave frequencies very tight beamforming can be achieved if the same overall antenna size is used. However there is a penalty, signal attenuation is high and it is not possible to send signals through walls and into buildings. As such window mounted units are required to receive signals. Millimeter-wave operation is envisioned to use  $2 \times 2$  MIMO only, but still use massive MIMO, be useful for low speed mobility, but provide very high data rates.

### 5.11.7 Non Orthogonal Multiple Access

One of the distinguishing features of 4G and 5G is OFDMA, an **orthogonal multiple access (OMA)** scheme. OFDMA in (both both 4G and) 5G allocates a dedicated resource block to each user and orthogonality ensures that there is almost no inter-user interference. The system utilization is optimum when each resource block is fully utilized. However this is not the case with IoT devices, the support for which is one of the main features of 5G. When IoT devices communicate very little information is usually exchanged and so most of the resource block used by an IoT device is unused. The expected growth of 5G will see a great proliferation of IoT devices with these devices having diverse data rate and latency requirements. Communicating with each IoT device by allocating a dedicated resource block does not use the capacity effectively. In addition, using the same basestation transmit power for each user is said to be ‘unfair’ in that each user typically has a different channel quality with a user having a poorer channel quality needing to operate with lower-order modulation. The ‘unfairness’ arises because the channel capacity is not equally shared but it would be if the the basestation transmitted higher powers to the user with a poorer channel. Non orthogonal multiple access (**NOMA**) is designed to addresses this imbalance and to

ensure optimum use of each resource block and hence of the communication system.

The key concept of NOMA is that low-rate devices will share a resource block, they will use differing symbol rates, and the transmit power level from the basestation to each user will be adjusted according to the quality of the communication channel [36, 37]. In NOMA, and in contrast to OFDMA, there will be inter-user interference and schemes have been developed that will enable individual users to be separated. One possible NOMA scheme to separate users by using code domain multiplexing with each user sharing a resource block being assigned a unique code. The concept is very similar to that used in WCDMA but applied at the level of a resource block.

The dominant NOMA scheme in 5G and the one to be first deployed (as of the time of writing this book) in a future 3GPP release of the 5G standard is **multi-user superposition transmission (MUST)**. **Successive interference cancellation (SIC)** is the scheme used to separate users.

In MUST users are allocated different power levels and will use different symbol rates. For example, a distant user with a poor quality channel will receive more power in NOMA and a user with a good channel will receive less since the maximum power available from the basestation transmitter is fixed. The result is that the throughput of the close and distant users will be similar. This is called 'fairness' and is particularly important when the system is operating at or near capacity. A further advantage of NOMA is that the demands on the basestation transmitter hardware is reduced. One of the consequences of having a common symbol rate in OFDMA is that the PMEPR is higher than it would be if there was a range of symbol rates. Thus NOMA will result in the total RF signal being transmitted from the basestation having a lower PMEPR.

It is expected that there will be a significant increase in throughput as many users, think of IoT devices, only need to communicate at very low data rates. Current cellular systems reach an abrupt limit on their throughput. In NOMA overload is supported at the cost of inter-user interference. Studies have shown that a threefold increase incapacity is possible with NOMA schemes that tolerate overload and inter-user interference [37].

### 5.11.8 Summary

The 5G system introduces many system optimization strategies and a high-level of system optimization across perhaps hundreds of basestations is required to achieve full potential. The 5G systems will inevitably be followed by 6G and this is the subject of current research. The defining characteristic is that 6G is being targeted to operate above 100 GHz and provide enormous bandwidths since very wide bandwidths are available.

## 5.12 6G, Sixth Generation Radio

Sixth generation radio will operate in the 95 GHz to 275 GHz range, and will provide data rates of 100 Gbit/s [38]. While 5G hopes to provide these data rates at 60 GHz this may not come to pass as the required spectral efficiency of 14 bits/s/Hz and requires 60 GHz hardware with very high dynamic range and this is unlikely to be available for some time. With 6G these data rates will be possible with much lower-order modulation and hence reduced

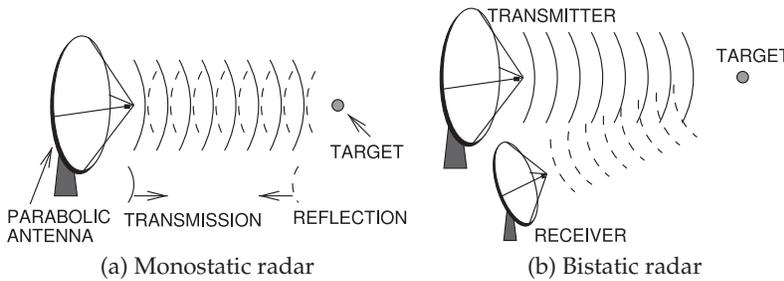
demands on hardware Regulatory bodies have allocated several bands in this range. In the USA these bands are 116–123 GHz, 174.8–182 GHz, 185–190 GHz, and 244–246 GHz for a total of 21 GHz of spectrum. Atmospheric losses can be large at these frequencies but, except for a peak of absorption around 140 GHz, this is more than compensated by the pencil-like beams generated by many-element phased array antennas possible at these high frequencies. The 100+ GHz frequencies have wavelengths of 3 mm or less enabling functionality not available at 60 GHz and below. These include precise positioning with millimeter resolution, near visual-quality imaging through fog and clouds, wireless cognition (off-loading large data sets from a mobile units for fixed computation), and unique sensing modalities by exploiting the many molecular resonances that occur above 100 GHz. The overwhelming challenge is the development of physical hardware and the development of array processing technologies [38].

### 5.13 Radar Systems

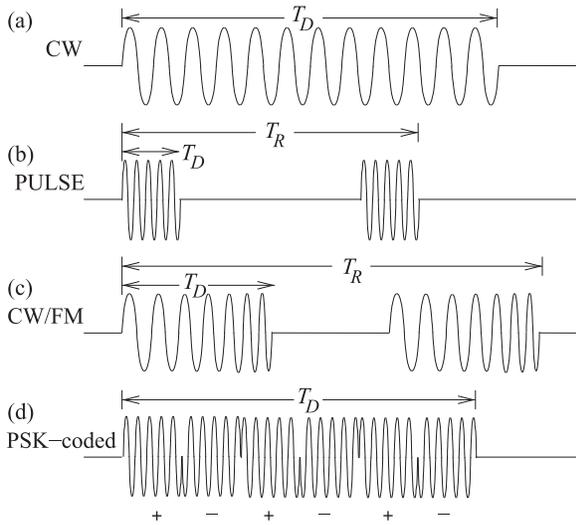
Radar uses EM signals to determine the range, altitude, direction, and speed of objects called targets by looking at the signals received from transmitted signals called radar waveforms. Common radar bands and applications are given in Table 5-16. The earliest use of EM signals to detect targets was demonstrated in 1904 by Christian Hülsmeyer using a spark gap generator [39]. This system was promoted as a system to avoid collisions of ships and detected the direction of targets only. Research contributed to further developments, with a significant acceleration during World War II. Radar is now a word in its own right, but in 1941 the term RADAR was created as an acronym for radio detection and ranging.

**Table 5-16:** IEEE radar bands and applications.

Band	Frequency	Wavelength	Application
HF	3–30 MHz	10–100 m	Over-the-horizon radar, oceanographic mapping
VHF	30–300 MHz	1–10 m	Oceanographic mapping, atmospheric monitoring, long-range search
UHF	0.3–1 GHz	1 m–30 cm	Long-range surveillance, foliage penetration, ground penetration, atmospheric monitoring
L	1–2 GHz	15–30 cm	Satellite imagery, mapping, long-range surveillance, environmental monitoring
S	2–4 GHz	7.5–15 cm	Weather radar, air traffic control, surveillance, search, IFF (identify, friend or foe)
C	4–8 GHz	3.75–7.5 cm	Hydrological radar, topography, fire control, weather
X	8–12 GHz	2.5–3.75 cm	Cloud radar, air-to-air missile seeker, maritime, air turbulence, police radar, high-resolution imaging, perimeter surveillance
Ku	12–18 GHz	1.7–2.5 cm	Remote sensing, short-range fire control, perimeter surveillance; pronounced “kay-you”
K	12–8–27 GHz	1.2–1.7 cm	Police radar, remote sensing, perimeter surveillance
Ka	27–40 GHz	7.5–12 mm	Police radar, weapon guidance, remote sensing, perimeter surveillance, weapon guidance; pronounced “kay-a”
V	40–75 GHz	4–7.5 mm	Perimeter surveillance, remote sensing, weapon guidance
W	75–110 GHz	2.7–4 mm	Perimeter surveillance, remote sensing, weapon guidance



**Figure 5-25:** Radar system: (a) monostatic radar with the same site used for transmission of the radar signal and receipt of the reflection from the target; and (b) bistatic radar with different transmit and receive sites.



**Figure 5-26:** Radar waveforms: (a) continuous wave; (b) pulsed wave; (c) frequency modulated continuous wave; and (d) phase-encoded (PSK-coded) waveform.

In a radar system, typically a high-gain antenna such as a parabolic antenna is used to transmit a radar signal, but always a high-gain antenna is used to receive the signal. If the same antenna is used for transmit and receive (possibly two similar antennas at the same site) the system is called a **monostatic radar** (see Figure 5-25(a)). Radar with transmit and receive antennas at different sites is called **bistatic radar** (shown in Figure 5-25(b)).

In a monostatic radar using the same antenna for transmit and receive, the space is painted with a radar signal and the received signal is captured after a propagation delay from the antenna to the target and back again. A radar image can then be developed. In many radars the receive antenna is mechanically steered and often a regular rotation is used. With so-called synthetic aperture radars, a platform such as an aircraft moves the radar in one direction and a one-dimensional mechanical or electrical scan enables a two-dimensional image to be developed.

The categories of radar waveforms are shown in Figure 5-26. The continuous wave (CW) waveform shown in Figure 5-26(a) is on all or most of the time and is used to detect a reflection from a target. This reflected signal is much smaller than the transmitted signal and it can be difficult to separate the transmitted and received signals. A monostatic CW radar architecture is shown Figure 5-27(a), where a circulator<sup>6</sup> is used to separate the transmitted and received signals. The received signal is converted to digital form using

<sup>6</sup> For the circulator shown, power entering port 1 of a circulator leaves at port 2 and power entering port 2 is delivered to port 3. So ports 1 and 3 are isolated.

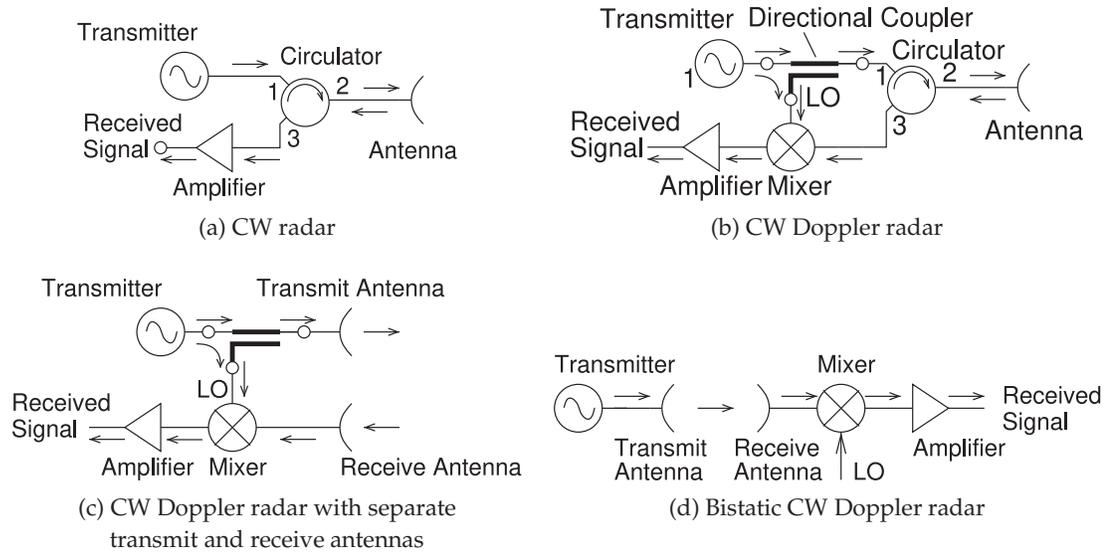
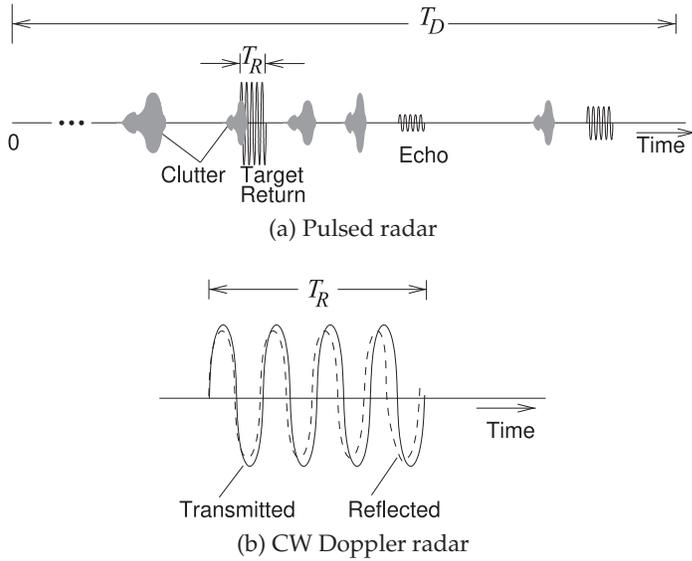


Figure 5-27: Radar architectures.

an ADC and the bandwidth of the ADC with the required dynamic range determines the limit on the bandwidth of the radar signal. Generally the broader the bandwidth, the better the radar system is at identifying objects. A CW signal can be used to develop an image, but it is not good at determining the range of a target. For this, a pulsed radar waveform, as shown in Figure 5-26(b), is better. The repetition period,  $T_R$ , is more than the round-trip time to the target, thus the time interval between the transmitted signal and the returned signal can be used to estimate range. Direction is determined by the orientation of the antenna.

The CW architecture can also be used with **pulsed radar**. In pulsed radar, the signal received contains the desired target signal, multipath and echoes, and clutter, as shown in Figure 5-28(a). These effects also appear in the signal received in CW radar, but it is much easier to see in pulsed radar. Identifying clutter and multipath effects is a major topic in radar processing. Alternative waveforms, especially digitally modulated waveforms, aid in extracting the desired information. The frequency modulated waveform, or chirp waveform, in Figure 5-26(c), will have a reflection that will also be chirped, and the difference between the frequency being transmitted and that received indicates the range of the target, provided that the target is not moving. The radar architecture that can be used to extract this information is shown in Figure 5-27(b). The directional couplers tap off a small part of the transmit signal and use it as the LO of a mixer with the received signal as input. A similar architecture is shown in Figure 5-27(c), but now separate transmit and receive antennas are used to separate the transmitted and received signals rather than using a circulator. Better separation of the transmitted and received signals can be obtained this way. The frequency of the IF that results is proportional to the range of the target.

An important concept in radar is that of **radar cross section (RCS)** denoted  $\sigma$  (with SI units of  $\text{m}^2$ ). The RCS of a target is the equivalent area that



**Figure 5-28:** Radar returns. The strength of the reflected signal depends on the radar cross section (RCS),  $\sigma$ , of the target.  $\sigma = 0.01 \text{ m}^2$  for a bird,  $< 0.1 \text{ m}^2$  for a stealth aircraft,  $1 \text{ m}^2$  for a human,  $2\text{--}6 \text{ m}^2$  for a conventional fighter plane,  $100 \text{ m}^2$  for a large commercial aircraft,  $200 \text{ m}^2$  for a truck,  $10,000\text{--}100,000 \text{ m}^2$  for a container ship.

intercepts the power of a transmitted signal and re-radiates all of that power isotropically to produce the observed power density at a receiver [40]. The RCS therefore depends on the frequency of the transmitted signal, the size of the target, the incident and reflected angle of the signal reflected by the target, and the reflectivity of the target.

The power of the signal reflected by the target and captured by the receive antenna is given by the **radar equation** [40]:

$$P_R = \frac{P_T G_T A_R \sigma F^4}{(4\pi)^2 R_T^2 R_R^2}, \tag{5.32}$$

where  $P_T$  is the transmit power delivered to the transmit antenna,  $G_T$  is the antenna gain of the transmitter,  $A_R$  is the effective aperture area of the receive antenna,  $F$  is the pattern propagation factor,  $R_T$  is the distance from the transmitter to the target, and  $R_R$  is the distance from the target to the receiver.  $F$  captures the effective loss due to multipath, which was captured in communications by introducing a signal dependence of  $1/d^n$ , where  $n$  ranges from 2 to 4 depending on the environment. In free space,  $F = 1$ . If the same antenna is used for transmit and receive, and multipath is not important (so  $F = 1$ ), then Equation (5.32) becomes

$$P_R = \frac{P_T G_T A_R \sigma}{(4\pi)^2 R^4}, \tag{5.33}$$

where  $R = R_T = R_R$ . In Section 4.5.4 the antenna effective area was related to the antenna gain. From Equation (4.27),

$$A_R = \frac{G_R \lambda^2}{4\pi}, \tag{5.34}$$

and so, with  $G_R = G_T$ , Equation (5.33) becomes

$$P_R = \frac{P_T G_T^2 \lambda^2 \sigma}{(4\pi)^3 R^4}. \tag{5.35}$$

**Table 5-17:** Doppler frequency shifts for targets moving toward a radar at speed  $v_R$ .

Radar frequency, $f_T$	Relative speed, $v_R$		
	1 m/s	100 km/hr	1000 km/hr
500 MHz	3.3 Hz	92.6 Hz	925.9 kHz
2 GHz	13.3 Hz	370.4 Hz	3.704 kHz
10 GHz	66.7 Hz	1.852 kHz	18.519 kHz
40 GHz	266.7 Hz	7.407 kHz	74.074 kHz

If a signal of frequency  $f_T$  is transmitted and reflected from a moving target there is a **Doppler shift**,  $f_D$ , and the received signal will be at frequency

$$f_R = f_T + f_D \quad (\text{target moving toward the radar}) \quad (5.36)$$

$$f_R = f_T - f_D \quad (\text{target moving away from the radar}), \quad (5.37)$$

$$\text{where } f_D = 2v_R f_T / c. \quad (5.38)$$

In Equation (5.38)  $v_R$  is the radial component of the speed of the target relative to the radar, and  $c$  is the speed of light. Typical Doppler shifts are shown in Table 5-17.

If the target is moving, there will be a Doppler shift. If the target is moving toward a CW radar, then the frequency of the returned signal will be higher, as shown in Figure 5-28(b). A similar architecture to that used with chirp radar can be used (see Figures 5-27(bc)). The concept can be extended to bistatic radars, but now the local oscillator reference must be generated. As was seen previously, the frequency of the transmit carrier can be recovered for digitally modulated signals such as the PSK-encoded signals shown in Figure 5-26(d). Advanced high-end radars use digital modulation and CDMA-like waveforms and exploit space-time coding. Typically the radar waveforms to be transmitted are chirped, which is a technique that takes the desired transmit waveform and stretches it out in time so that it can be more efficiently amplified and more power can be transmitted. At the receiver, the radar signal is compressed in time so that it corresponds to the original transmit signal prior to chirping.

It should be apparent that radars and radar waveforms can be optimized for imaging or for exploiting Doppler shifts to track moving targets. Imaging is suitable when there is little clutter, such as looking into the air. However, it is difficult to detect targets such as cars that are moving on the ground. So-called **ground moving target indication (GMTI) radar** relies on Doppler shifts to discriminate moving targets, and then the ability to image accurately is compromised. Considerable effort is devoted to developing waveforms that are difficult to detect (stealthy) and are optimized for imaging or GMTI.

## 5.14 Summary

Radio frequency and microwave design continues to rapidly evolve, responding to new communication, radar, and sensor architectures. The long-term evolution also exploits opportunities made available with larger-scale monolithic integration and by advances in high-performance, low-power digital signal processing.

If filters and other hardware in a communication receiver are ideal, the final  $E_b/N_o$  and BER achieved are directly related to the RF SIR as described in this chapter. With practical filters and analog hardware there is a performance degradation and an SIR higher than the theoretical SIR is required to achieve a particular BER, or  $E_b/N_o$ . The difference is

Modulation	Implementation margin
BPSK	0.5 dB
QPSK	0.8 dB
8-PSK	1–1.6 dB
16-QAM	1.5–2.1 dB
CDMAOne	0.5 – 1 dB
WCDMA	2 dB

**Table 5-18:** Implementation margins for modern communication receivers [41, p. 328], [25, 42]. These implementation margins are what can be achieved by good system designs.

captured by the **implementation margin**,  $k$ , usually specified in decibels. To achieve a specific BER, SIR must be greater than the theoretical SIR by  $k$ . The implementation margin is therefore a measure of the performance of RF hardware and is an important metric in design and in planning design. The implementation margin captures many imperfections. A design group and a company learn what it can achieve, e.g. see Table 5-18, and use this in budgeting design costs and planning design effort. The design cost of RF systems is considerable and the ability to manage the design process and be able to estimate design effort is critical to timely success. Higher implementation margins result from the choice of lower-performing technologies, perhaps resulting from a compromise of performance, cost, and design effort.

## 5.15 References

- [1] "Radio regulations," International Telecommunications Union, 2004.
- [2] "Recommendation ITU-R V.662-3 (1986-1990-1993-2000), updated in 2005 for editorial reasons only, 2005." International Telecommunications Union.
- [3] "American National Standard T1.523-2001, Telecom Glossary 2011," available on-line with revisions at <http://glossary.atis.org>, 2011, sponsored by Alliance for Telecommunications Industry Solutions.
- [4] D. Ring, "Mobile telephony—wide area coverage," *Bell Laboratories Technical Memorandum*, Dec. 1947.
- [5] J. Schulte, H.J. and W. Cornell, "Multi-area mobile telephone system," *IRE Trans. on Vehicular Communications*, vol. 9, no. 1, pp. 49–53, Jun. 1957.
- [6] W. Lewis, "Coordinated broadband mobile telephone system," *IRE Trans. on Vehicular Communications*, vol. 9, no. 1, pp. 43–48, Jun. 1957.
- [7] A. Joel, "Mobile communication system," US Patent 3 663 762, May 16, 1972.
- [8] F. Ikegami, "Mobile radio communications in Japan," *IEEE Trans. on Communications*, vol. 20, no. 4, pp. 738–746, Aug. 1972.
- [9] "Cost final report," <http://www.lx.it.pt/cost231>.
- [10] F. Ikegami, T. Takeuchi, and S. Yoshida, "Theoretical prediction of mean field strength for urban mobile radio," *IEEE Trans. on Antennas and Propagation*, vol. 39, no. 3, pp. 299–302, 1991.
- [11] J. Doble, *Introduction to Radio Propagation for Fixed and Mobile Communications*. Norwood, MA, USA: Artech House, Inc., 1996.
- [12] H. Markey and G. Antheil, "Secret communication system," US Patent US Patent 397 412, 04 11, 1942.
- [13] K. Gilhousen, I. Jacobs, L. Weaver Jr., and E. Armstrong, "Spread spectrum multiple access communication system using satellite or terrestrial repeaters," US Patent 4 901 307, Feb 13, 1990.
- [14] A. Viterbi, *CDMA Principles of Spread Spectrum Communications*. Addison-Wesley, 1995.
- [15] V. Garg, K. Smolik, and J. Wilkes, *Applications of CDMA in Wireless/Personal Communications*. Prentice Hall, 1996.
- [16] J. Barry, E. Edwards, and D. Messerschmitt, *Digital Communication*, 3rd ed. Kluwer Academic Publishers, 2004.
- [17] D. Smith, *Digital Transmission Systems*, 3rd ed. Kluwer Academic Publishers, 2004.
- [18] G. Lewis, *Communications Systems: Engineers' Choices*, 3rd ed. Focal Press, 1999.
- [19] W. Lee, "Spectrum efficiency in cellular [radio]," *IEEE Trans. on Vehicular Technology*, vol. 38, no. 2, pp. 69–75, May 1989.
- [20] R. Price and P. E. Green Jr., "Anti-multipath receiving system," US Patent 2 982 853, May 02, 1961.

- [21] R. Ziemer and W. Tranter, *Principles of Communications: Systems, Modulation, and Noise*, 5th ed. Wiley, 2001.
- [22] T. Rappaport, *Wireless Communications Principles and Practice*. Prentice Hall, 1996.
- [23] K. Cheun, "Performance of direct-sequence spread-spectrum rake receivers with random spreading sequences," *IEEE Trans. on Communications*, vol. 45, no. 9, pp. 1130–1143, Sep. 1997.
- [24] International Telecommunications Union, at <http://www.itu.int>.
- [25] "3rd generation partnership project (3gpp)," <http://www.3gpp.org>.
- [26] Third Generation Partnership Project specifications, at <http://www.3gpp.org/specifications>.
- [27] G. Foschini, "Layered space-time architecture for wireless communication in a environment when using multi-element antennas," *Bell Labs Technical J.*, vol. 1, no. 2, pp. 41–59, Autumn 1996.
- [28] G. Raleigh and J. Cioffi, "Spatio-temporal coding for wireless communications," in *Global Telecommunications Conf., 1996*), vol. 3, Nov. 1996, pp. 1809–1814.
- [29] A. Goldsmith, S. Jafar, N. Jindal, and S. Vishwanath, "Capacity limits of mimo channels," *IEEE J. on Selected Areas in Communications*, vol. 21, no. 5, pp. 684–702, Jun. 2003.
- [30] D. Gesbert, M. Shafi, D. shan Shiu, P. Smith, and A. Naguib, "From theory to practice: an overview of mimo space-time coded wireless systems," *IEEE J. on Selected Areas in Communications*, vol. 21, no. 3, pp. 281–302, Apr. 2003.
- [31] W. He and C. Georghiades, "Computing the capacity of a MIMO fading channel under psk signaling," *IEEE Trans. on Information Theory*, vol. 51, no. 5, pp. 1794–1803, May 2005.
- [32] D. Benyamina, A. Hafid, and M. Gendreau, "Wireless mesh networks design—a survey," *IEEE Communications Surveys Tutorials*, vol. 14, no. 2, pp. 299–310, quarter 2012.
- [33] B. Wang and K. Liu, "Advances in cognitive radio networks: A survey," *IEEE J. of Selected Topics in Signal Processing*, vol. 5, no. 1, pp. 5–23, Feb. 2011.
- [34] J. Wang, M. Ghosh, and K. Challapali, "Emerging cognitive radio applications: A survey," *IEEE Communications Magazine*, vol. 49, no. 3, pp. 74–81, Mar. 2011.
- [35] D.-C. Chang and C.-N. Hu, "Smart antennas for advanced communication systems," *Proc. IEEE*, vol. 100, no. 7, pp. 2233–2249, Jul. 2012.
- [36] Y. Saito, Y. Kishiyama, A. Benjebbour, T. Nakamura, A. Li, and K. Higuchi, "Non-orthogonal multiple access (noma) for cellular future radio access," in *2013 IEEE 77th vehicular technology conference (VTC Spring)*. IEEE, 2013, pp. 1–5.
- [37] L. Dai, B. Wang, Y. Yuan, S. Han, I. Chih-Lin, and Z. Wang, "Non-orthogonal multiple access for 5g: solutions, challenges, opportunities, and future research trends," *IEEE Communications Magazine*, vol. 53, no. 9, pp. 74–81, 2015.
- [38] T. Rappaport, Y. Xing, O. Kanhere, S. Ju, A. Alkateeb, and G. Trichopoulos, "Wireless communications and applications above 100 GHz: opportunities and challenges for 6g and beyond," *IEEE Access*, vol. 7, 2019.
- [39] Hulsmeier, "Wireless transmitting and receiving mechanism for electric waves," germany Patent US Patent 810 150, Jan. 16, 1906.
- [40] G. W. Stimson, *Introduction to Airborne Radar*, 3rd ed. Scitech Publ., 2014.
- [41] E. Tozer, *Broadcast Engineer's Reference Book*. Focal Press, 2004.
- [42] A. Richardson, *WCDMA Design Handbook*. Cambridge University Press, 2005.
- [43] K. Gard, H. Gutierrez, and M. Steer, "Characterization of spectral regrowth in microwave amplifiers based on the nonlinear transformation of a complex gaussian process," *IEEE Trans. on Microwave Theory and Techniques*, vol. 47, no. 7, pp. 1059–1069, Jul. 1999.

## 5.16 Exercises

- Research at Bell Labs in the 1960s showed that the minimum acceptable SIR for voice communications is 17 dB. This applies to analog modulated signals, but not digitally modulated signals, where BER is important. Consider a seven-cell cluster. If the power falls off as  $1/d^3$ , where  $d$  is distance, determine the worst possible SIR considering only interference from other radios. The worst situation will be when a mobile handset is at the edge of its cell. To do this you need to estimate the distance from the handset to the other basestations (in neighboring clusters) that are operating at the same power levels. Consider the cells to be hexagons. Develop a symbolic expression for the total interference signal level at the handset, assuming that all basestations are radiating at the same power level,  $P$ . You can use approximate distances. For example, each distance can be expressed in terms of integer multiples of cell radii,  $R$ . Is the 17 dB SIR achieved using a 7-cell cluster?
- Describe the following concepts.
  - Clusters in a cellular phone system.
  - Multipath effects in a central city area compared to multipath effects in a desert.
- Describe the concept of clusters in a cellular phone system in four lines.
- Short answer questions on modulation and spectral efficiency.
  - What is the PMEPR of a phase modulated signal?
  - In five lines explain your understanding of spectral efficiency as it relates to bits per hertz. That is, how can you have a spectral efficiency of  $n$  bit/s/Hz where  $n$  is more than 1? [Note that sometimes this is expressed as bit/s/Hz as well as bit/s/Hz.]
  - What is the modulation efficiency of a QPSK-modulated signal? Ignore the impact of the number of cells in a cluster.
- A cellular communication system uses a frequency reuse plan with seven cells per cluster to obtain the required minimum SIR. If a QPSK system is used, what is the radio spectrum efficiency in terms of bit/s/Hz/cell if all transitions on the constellation diagram are allowable? Assume that there is no coding.
- A cellular communication system uses GMSK modulation and a frequency reuse plan with three cells per cluster. What is the radio spectrum efficiency in terms of bit/s/Hz/cell? Assume that there is no coding.
- A cellular communication system uses  $\pi/4$ -DQPSK modulation and a frequency reuse plan with five cells per cluster. What is the radio spectrum efficiency in terms of bit/s/Hz/cell? Assume that there is no coding.
- A frequency reuse plan has three cells per cluster. If ideal 16-QAM modulation is used, what is the system spectral efficiency (in bits/s/Hz/cell)?
- A 2G cellular communication system uses a frequency reuse plan with seven cells per cluster. If ideal QPSK modulation is used, what is the system spectral efficiency (in bits/s/Hz/cell)?
- A 2G system has three cells per cluster. If  $3\pi/8$ -8PSK modulation is used, what is the system spectral efficiency (in bits/s/Hz/cell)?
- A cellular communication system uses  $3\pi/8$ -8PSK modulation and a frequency reuse plan with three cells per cluster. What is the radio spectrum efficiency in terms of bit/s/Hz/cell? Assume that there is no coding.
- A communication system uses a modulation with a modulation efficiency of 5 bit/s/Hz. Ignore coding so that  $R_b = R_c$ . What is the radio spectral efficiency in terms of bit/s/Hz/cell if there are three cells per cluster?
- A proposed modulation format has a modulation efficiency of 3.5 bit/s/Hz. Antenna sectoring and required SNR lead to a system with seven cells per cluster. You can ignore the impact of coding so you can assume that  $R_b = R_c$ . What is the radio spectral efficiency in terms of bit/s/Hz/cell modulated signal?
- A cellular radio system uses a frequency reuse plan with 12 cells per cluster. If ideal 8-PSK modulation is used, what is the system spectral efficiency in terms of bit/s/Hz/cell?
- Consider a cellular system having a frequency reuse plan with seven cells per cluster to obtain the minimum signal-to-interference ratio. If ideal QPSK modulation is used, what is the system spectral efficiency in terms of bits per second per hertz per cell?
- A monostatic free-space 10 GHz pulsed radar system is used to detect a fighter plane having a radar cross section,  $\sigma$ , of  $5 \text{ m}^2$ . The antenna gain is 30 dB and the transmitted power is 1 kW. If the minimum detectable received signal is  $-120 \text{ dBm}$ , what is the detection range?

17. An antenna with a gain of 10 dB presents an RF signal with a power of 5 dBm to a low-noise amplifier along with noise of 1 mW and an interfering signal of 2 mW.
- What is the RF SIR? Include both noise and the interfering signal in your calculation. Express your answer in decibels.
  - The modulation format and coding scheme used have a processing gain,  $G_P$ , of 7 dB. The modulation scheme has four states. What is the ratio of the energy per bit to the noise per bit, that is, what is the effective  $E_b/N_o$  after despreading?
18. The receiver in a digital radio system receives a 100 pW signal and the interference from other radios at the input of the receiver is 20 pW. The receiver has an overall gain of 40 dB and the noise added by the receiver, referred to the output of the receiver, is 100 nW.
- What is the RF SIR at the output of the receiver?
  - If 16-QAM modulation with a modulation efficiency of 2.98 bit/s/Hz is used and the processing gain is 30 dB, what is the effective SIR after despreading, i.e. what is  $E_{b,\text{eff}}/N_{o,b}$ ?
19. A new communication system is being investigated for sending data to a printer. The system will use GMSK modulation and a channel with 25 MHz bandwidth and an information bit rate of 10 Mbit/s. The modulation format will result in a spectrum that distributes power almost uniformly over the 25 MHz bandwidth. [Parallels Example 5.3]
- What is the processing gain?
  - If the received RF SIR is 6 dB, what is the effective system SIR (or  $E_{b,i}/N_{o,i}$ ) after DSP? Express your answer in decibels.
20. A 4 kHz bandwidth voice signal is coded by a vocoder as an 8 kbit/s data stream. Coding increases the data stream to 64 kbit/s. What is the processing gain that can be achieved at the receiver if QPSK modulation is used with a modulation efficiency of 1.4 bit/s/Hz?
21. A receive antenna with a gain of 10 dB presents a signal with a power of 5 dBm to a low-noise amplifier along with noise of 1 mW and an interfering signal of 2 mW.
- What is the SIR at the input to the amplifier? Express your answer in decibels.
  - BPSK modulation is used and coding results in a processing gain,  $G_P$ , of 7 dB. What is the ratio of the energy per bit to the noise power per bit (i.e., what is  $E_b/N_o$ ) after despreading?
22. If a received RF signal has an SIR of  $-5$  dB and the processing gain (calculated from bit rates) that can be achieved for the modulation and coding used is 15 dB, what is the  $E_b/N_o$  after processing? There are 4 bits per symbol.
23. A signal with a power of 13 dBm is input to a low-noise amplifier along with noise of 1 mW and an interfering signal of 2 mW.
- What is the SIR at the input to the amplifier? Express your answer in decibels.
  - QPSK modulation is used and coding results in a processing gain,  $G_P$ , of 13 dB. What is the ratio of the effective energy per bit to the noise power per bit (i.e., what is  $E_b/N_o$ ) after despreading?
24. Short answer questions. Each part requires a short paragraph of five lines and a figure.
- Explain how OFDM reduces the impact of multiple paths in a wireless communication system.
  - Explain how MIMO exploits multipath to enhance the capacity of a digital communication system.
25. A coding rate of  $2/3$  is required to manage transmission errors in a 54 Mbit/s data link. That is, the information bit rate is 54 Mbit/s. What is the total bit rate required (including data and coding bits)?
26. The channel bandwidth in the GSM cellular phone system is 200 kHz and the GMSK modulation scheme used has a spectral efficiency of 1.354 bit/s/Hz.
- What is the data rate of one frequency channel?
  - A time slot is 577  $\mu$ s long. How many bits are there in one (i.e. a duration of 8.25 bits). How many data bits are there in a GSM time slot?
  - A GSM frame duration is 4.615 ms long and has eight time slots and a voice user has one time slot every frame. How many data bits per second are available to a single user?
27. Consider an OFDM system with 48 subcarriers carrying data and which uses 16-QAM modulation of each subcarrier and a coding rate of  $2/3$ . There are also four pilot subcarriers that are used for frequency and phase reference, to ensure that spectral lines are not created, and to facilitate carrier recovery. The pilot carriers can be ignored in this problem so consider 48 subcarriers. The modulation efficiency achieved

- for this particular implementation of 16-QAM is 2.98 bit/s/Hz.
- (a) How many symbols are there for each subcarrier? That is, how many points are there in the constellation diagram for one subcarrier?
  - (b) How many coded bits are there on each subcarrier? That is, how many bits per symbol are there for each subcarrier?
  - (c) Considering all of the data subcarriers, how many coded bits are there per OFDM symbol? [Hint, there are 16 subcarriers, so for each OFDM symbol there will be 16 subcarrier symbols.]
  - (d) Considering the coding rate, determine the number of data bits per OFDM symbol. That is, ignore coding bits.
28. Consider an OFDM system with 12 subcarriers carrying data and which uses 8-PSK modulation of each subcarrier and a coding rate of 3/4. Pilot subcarriers they can be ignored in this problem so consider all 12 subcarriers.
- (a) How many symbols are there for each subcarrier? That is, how many points are there in the constellation diagram for one subcarrier?
  - (b) How many coded bits (code + data) are there on each subcarrier? That is, how many bits per symbol are there for each subcarrier?
  - (c) Considering all of the data subcarriers, how many coded bits are there per OFDM symbol? [Hint, there are 12 subcarriers, so for each OFDM symbol there will be 12 subcarrier symbols.]
  - (d) Considering the coding rate, determine the number of data bits per OFDM symbol. That is, ignore coding bits.
29. A digital radio system transmits a baseband digital signal of 100 Mbit/s over a channel that is 300 MHz wide. The digital modulation scheme effectively fills the 300 MHz channel with uniform power.
- (a) What is the processing gain that can be achieved with this system?
  - (b) Consider that the signal received and delivered to the input of the receiver front end is 100 pW and the interference from other radios delivered to the receiver front-end is 20 pW. What is the SIR at the input to the receiver electronics?
30. The L1 band of the global positioning system (GPS), is centered at 1.57542 GHz and has two overlapping spread-spectrum encoded signals. The stronger of these is the coarse acquisition (C/A) signal with an information bit rate of 50 bits/s and a transmission rate of 1.023 million chips per second using BPSK modulation with an RF bandwidth of 2.046 MHz. In ideal conditions the C/A signal received has a power of  $-130$  dBm. A GPS receiver has an antenna noise temperature is of 290 K.
- (a) What is the processing gain?
  - (b) What is the noise in dBm received in the 2.046 MHz bandwidth?
  - (c) What is the SNR in decibels?
  - (d) If a C/A signal is received from each of 10 satellites (so there are 9 interfering signals), what is the total interference power for one satellite's C/A signal?
  - (e) With respect to just one of the C/A signals, what is the SINR (signal to interference plus noise ratio) at the receiver?
  - (f) If the receiver does not contribute noise, what is the effective SNR of the despread bit-stream from each satellite?
  - (g) If the required minimum effective SNR is 6 dB, what is the minimum acceptable power, in dBm, of the GPS signal received from one satellite?
31. GLONASS is the Russian satellite navigation system with one of two open signals called the L1OF band at 1600.995 MHz. The system uses DSSS encoding and BPSK modulation and each GLONASS satellite transmits on a different frequency. The symbol rate is 511,000 chips/s, the bandwidth of the transmitted signal is approximately 540 kHz, and there are 50 information bits per second.
- (a) What is the system's processing gain?
  - (b) What is the noise in dBm received in the 540 kHz bandwidth?
  - (c) If the required system minimum effective SNR is 6 dB, what is the minimum acceptable power, in dBm, of the received signal? Assume that the receiver is noiseless.
32. A new communication system uses a channel that is 100 MHz wide and uses direct sequence CDMA to efficiently spread a 5 Mbit/s digital signal over the full channel.
- (a) What is the processing gain that can be achieved with the received signal?
  - (b) The analog signal at the output of the receive antenna has an SIR of 1 dB, what is energy per bit divided by the noise per bit?
33. A deep-space communication system will use direct sequence spread spectrum to code a data stream of 10 kbit/s then modulate the signal to transmit over a 5 GHz link to a ground sta-

- tion. Since the propagation loss is very high it has been determined that the processing gain must be 50 dB. If the link has a bandwidth of 1 MHz, what is the maximum baseband bit rate (in bit/s) that can be supported?
34. A direct sequence spread spectrum code of 10 Mbit/s is used to code a 4 kbit/s data stream that is modulated using  $3\pi/8$ -8PSK modulation to produce an RF signal at 1900 MHz. The modulation efficiency of  $3\pi/8$ -8PSK modulation is 2.7 bit/s/Hz.
    - (a) What is the bandwidth of the RF signal?
    - (b) What processing gain can be achieved in the receiver?
  35. An OFDM system with 12 data subcarriers, uses a coding rate of  $3/4$ , and each subcarrier uses 16-QAM modulation (with a modulation efficiency of 2.7 bit/s/Hz) with a bandwidth of 250 kHz. What is the maximum data rate supported?
  36. An OFDM system with 48 subcarriers carrying data uses 16-QAM modulation of each subcarrier and a coding rate of  $2/3$ . The actual modulation efficiency for the 16-QAM system here is 2.7 bit/s/Hz. What is the maximum data rate supported in Mbit/s when the bandwidth of each modulated subcarrier is 312 kHz?
  37. Explain using sentences and a diagram how OFDM reduces the impact of multiple paths in a wireless communication system.
  38. Explain using sentences and a diagram how MIMO exploits multipath to enhance the capacity of a digital communication system.
  39. A free-space 2 GHz pulsed monostatic radar system transmits a 2 kW pulse and has a minimum detectable received signal power of  $-90$  dBm. What is the antenna gain required to be able to detect a target with a radar cross section of  $10 \text{ m}^2$  at 10 km?
  40. A 10 GHz bistatic radar has a minimum detectable received signal power of  $-150$  dBm, an antenna gain of 26 dB, and a required range of 100 km. What is the transmitted pulse power in dBm needed to detect a
    - (a) conventional fighter aircraft having an RCS of  $5 \text{ m}^2$ ?
    - (b) a stealth aircraft with an RCS of  $0.05 \text{ m}^2$ ?
  41. The L5 band is a new public band of the GPS system and is centered at 1.176.5 GHz. The coarse acquisition (C/A) signal has an information bit rate of 50 bits/s and a spread-spectrum transmission rate of  $10.23 \cdot 10^6$  chips per second using BPSK modulation and occupying an RF bandwidth of 20.46 MHz. The noiseless GPS receiver has an omnidirectional antenna with a noise temperature is 290 K.
    - (a) What is the system's processing gain?
    - (b) What is the noise in dBm in the 20.46 MHz bandwidth?
    - (c) If overlapping C/A signals are received from 10 satellites, what is the total interference power for the signal from one satellite? The power of a C/A signal is  $S$ .
    - (d) If the required system minimum effective SNR is 6 dB, what is the minimum acceptable received power, in dBm, of the signal from one satellite?

### 5.16.1 Exercises By Section

†challenging, ‡very challenging

- |   |  |                      |
|---|--|----------------------|
| §5.3 1 <sup>†</sup> , 2 <sup>†</sup> , 3 <sup>†</sup> , 4 <sup>‡</sup> , 5 <sup>†</sup> , 6 <sup>†</sup> , 7 <sup>‡</sup> | §5.6 22, 23, 24, 25 <sup>‡</sup> , 26 <sup>†</sup> , 27 <sup>†</sup> ,   | §5.10 35, 36, 37, 38 |
| §5.5 8 <sup>†</sup> , 9, 10, 11 <sup>†</sup> , 12, 13, 14 <sup>†</sup> , 15 <sup>†</sup> ,                                | 28 <sup>†</sup> , 29 <sup>†</sup>  | §5.13 39, 40         |
| 16 <sup>†</sup> , 17 <sup>†</sup> , 18 <sup>†</sup> , 19, 20, 21  | §5.8 30 <sup>†</sup> , 31 <sup>†</sup> , 32 <sup>†</sup> , 33 <sup>‡</sup> , 34 <sup>‡</sup> , 41 <sup>†</sup> |                      |

### 5.16.2 Answers to Selected Exercises

- |                 |                     |               |
|-----------------|---------------------|---------------|
| 1 12.5 dB       | 6 0.45bit/s/Hz/cell | 27(d) 128     |
| 8(d) 3 bit/s/Hz | 32(b) 14 dB         | 28(d) 27      |
| 19 0.1          | 34(a) 3.704 MHz     | 29(b) 6.99 dB |
| 21(b) 7.228 dB  | 25 81 Mbit/s        |               |

# Appendix

## 5.A Mathematics of Random Processes

This appendix presents the essential mathematics required to describe the statistical properties of a random process such as noise. Also, it is difficult to analyze circuits with digitally modulated signals unless the signals are treated as being random with high-order statistics. Several statistical terms are introduced to describe the properties of a random variable  $X$  and how its value at one time is related to its value at other times. So  $X$  will, in general, vary with time (i.e., it can be written as  $X(t)$ ) and its value at a particular time will be random.

This section presents all of the probability metrics required to describe noise, interference, and digitally modulated signals. A random process in time  $t$  is a family of random variables  $\{X(t), t \in T\}$ , where  $t$  is somewhere in the time interval  $T$ . The probability that  $X(t)$  has a value less than  $x_1$  is denoted by  $P\{X(t) \leq x_1\}$ . This is also called the **cumulative distribution function (CDF)**, and sometimes called just the **distribution function (DF)**:

$$F_X(x_1) = P\{X \leq x_1\}. \quad (5.39)$$

In general, since  $X(t)$  varies with time, the CDF required to handle noise, that is, random voltages and currents, will depend on time. So the CDF used with noise and interference in communications will have two arguments, and the multivariate CDF is

$$F_X(x_1; t_1) = P\{X(t_1) \leq x_1\}. \quad (5.40)$$

That is,  $F_X(\infty, t) = 1$  and  $F_X(-\infty, t) = 0$ .  $F_X$  is being used with two arguments, as it is necessary to indicate the value of  $X$  at a particular time.  $F_X(x_1; t_1)$  is said to be the first-order distribution of  $X(t)$ .

Another probability metric often used is the **probability density function (PDF)**,  $f$ , which is related to the CDF as

$$F_X(x_1; t_1) = \int_{-\infty}^{x_1} f(x, t_1).dx. \quad (5.41)$$

Another property that needs to be captured is the relationship between the value of the random variable at one time,  $t_1$ , to its value at another time,  $t_2$ . Clearly, if the variables are completely random there would be no relationship. The statistical relationship of  $x$  at  $t_1$  and at  $t_2$  is described by the **joint CDF**,  $F_X(x_1, x_2; t_1, t_2)$ , which is the second-order distribution of the random process:

$$F_X(x_1, x_2; t_1, t_2) = P\{X(t_1) \leq x_1, X(t_2) \leq x_2\}. \quad (5.42)$$

This is the joint CDF, or probability, that  $X(t_1)$  will be less than  $x_1$  at time  $t_1$  and also that  $X(t_2)$  will be less than  $x_2$  at time  $t_2$ . So in the case of noise, the joint CDF describes the correlation of noise at two different times. In general, the  $n$ th order distribution is

$$F_X(x_1, \dots, x_n; t_1, \dots, t_n) = P\{X(t_1) \leq x_1, \dots, X(t_n) \leq x_n\}. \quad (5.43)$$

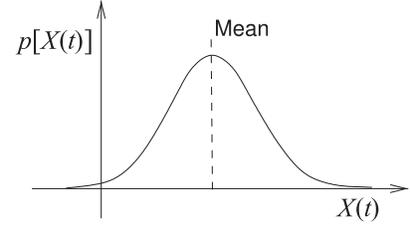
If the random process (i.e.,  $X$ ) is discrete, then the **probability mass function (PMF)** is used for the probability that  $x$  has a particular value. The general form of the PMF is

$$p_X(x_1, \dots, x_n; t_1, \dots, t_n) = P\{X(t_1) = x_1, \dots, X(t_n) = x_n\}, \quad (5.44)$$

and the general form of the PDF, used with continuous random variables, is (from Equation (5.41))

$$f_X(x_1, \dots, x_n; t_1, \dots, t_n) = \frac{\partial^n F_X(x_1, \dots, x_n; t_1, \dots, t_n)}{\partial x_1 \dots \partial x_n}. \quad (5.45)$$

The statistical measures above describe the properties of random variables and capture the way a variable is related to itself at different times, and how two different random processes are related to each other at the same time and at different times. Such characterizations are based on the **expected value** of a random variable. The expectation of a random variable  $X(t)$  is the weighted average of all possible



**Figure 5-29:** Gaussian distribution.

values of the random variable. The weighting is the probability for a discrete random variable, and is the probability density for a continuous random variable. So the expected value of the random variable  $X(t)$  is

$$E[X(t)] = \begin{cases} \sum_{-\infty}^{\infty} x(t)p_X(x, t) & \text{for a discrete random variable} \\ \int_{-\infty}^{\infty} x(t)p_X(x, t)dx & \text{for a continuous random variable.} \end{cases} \quad (5.46)$$

This is just the mean of  $X(t)$  defined as

$$\mu_X(t) = \bar{X}(t) = \langle X(t) \rangle = E[X(t)]. \quad (5.47)$$

The mean is also called the first-order moment of  $X(t)$ .  $E[\ ]$  is called the expected value of a random variable and the term is synonymous with the **expectation**, **mathematical expectation**, **mean**, and **first moment** of a random variable. The symbols  $\langle \rangle$  are a clean way of specifying the expectation. In general a computer program would need to be used to calculate the expected value. However, for some assumed probability distributions there are analytic solutions for  $E[\ ]$ .

The  $n$ th-order moment of  $X(t)$  is just the expected value of the  $n$ th power of  $X(t)$ :

$$\mu'_n = E[X^n(t)] = \langle X^n(t) \rangle = \begin{cases} \sum_{-\infty}^{\infty} x^n(t)p_X(x, t) & \text{for a discrete} \\ & \text{random variable} \\ \int_{-\infty}^{\infty} x^n(t)p_X(x, t)dx & \text{for a continuous} \\ & \text{random variable.} \end{cases} \quad (5.48)$$

Thus the **second moment**, sometimes called the **second raw moment**, is  $\mu'_2 = \langle X^2(t) \rangle = E[X^2(t)]$ . A more useful quantity for characterizing the statistics of a signal is the **second central moment**, which is the second moment about the mean. The second central moment of a random variable is also called its **variance**, written as  $\sigma^2$  or as  $\mu_2$ :

$$\begin{aligned} \sigma^2 = \mu_2 &= E[(X - \mu)^2] = E[X^2 - 2\mu X + \mu^2] = E[X^2] - 2\mu E[X] + \mu^2 \\ &= E[X^2] - 2\mu^2 + \mu^2 = E[X^2] - \mu^2 = E[X^2] - (E[X])^2 \\ &= \langle X^2(t) \rangle - \mu^2 = \langle X^2(t) \rangle - \langle X(t) \rangle^2. \end{aligned} \quad (5.49)$$

The variance is a measure of the **dispersion** of a random variable (i.e., how much the random variable is spread out). Variance is one of several measures of dispersion, but it is the preferred measure when working with noise and digitally modulated signals. The **standard deviation**,  $\sigma$ , is the square root of the variance  $\sigma^2$ . It is also common to denote the variance of  $X$  as  $\sigma_X^2$ , and in general, the variance can be a function of time,  $\sigma^2(t)$ .

It is common to approximate the statistical distribution of a digitally modulated signal as a Gaussian or normal distribution. This distribution is shown in Figure 5-29 and is mathematically described by its probability distribution

$$p[X(t)] = \frac{1}{\sqrt{2\pi}\sigma} e^{-(X(t)-\mu)/(2\sigma^2)}, \quad (5.50)$$

where the mean of the distribution is  $\mu$  and the variance is  $\sigma^2$ . The third and higher moments of the Gaussian distribution are zero. That is one of the reasons why this distribution is so commonly used as an approximation. Analysis using the Gaussian distribution is much simpler than it would be for other distributions. More realistically, a digitally modulated signal has  $I$  and  $Q$  components and the

distribution of each of these should be approximated as a Gaussian distribution. Such a distribution is called a **complex Gaussian distribution** and the analysis of the distortion produced by an amplifier using the complex Gaussian distribution is more accurate. Using a more sophisticated distribution, e.g. using the moments calculated from the actual digitally modulated signal, provides even greater accuracy in analysis [43] but now the complexity is beyond manual calculation.

If a digitally modulated signal is completely random, then there would be no correlation between the value of the signal at one time and its value at another time. However, there is a relationship and the relationship is described by the signal's **autocorrelation function**. The autocorrelation function is used to describe the relationship between values of a function separated by different instants of time. For a random variable, it is given by

$$R_X(t_1, t_2) = E[X(t_1)X(t_2)]. \quad (5.51)$$

When the random variable is discrete, the one-dimensional autocorrelation function of a random sequence of length  $N$  is expressed as

$$R_X(i) = \sum_{j=0}^{N-1} x_j x_{j+i}, \quad (5.52)$$

where  $i$  is the so-called **lag parameter**. The autocovariance function of  $X(t)$  is given by

$$\begin{aligned} K_X(t_1, t_2) &= E[\{X(t_1) - \mu_X(t_1)\}E[\{X(t_1) - \mu_X(t_1)\}]] \\ &= R_X(t_1, t_2) - \mu_X(t_1)\mu_X(t_2) \end{aligned} \quad (5.53)$$

while the variance is given by

$$\sigma_X(t) = E[\{X(t) - \mu_X(t)\}^2] = K_X(t, t). \quad (5.54)$$

A random process  $X(t)$  is stationary in the strict sense if

$$F_X(x_1, \dots, x_n; t_1, \dots, t_n) = F_X(x_1, \dots, x_n; t_1 + \tau, \dots, t_n + \tau) \quad (5.55)$$

$\forall t_i \in T, i \in \mathbf{N}$ . If the random process is **wide-sense stationary (WSS)**, then it is stationary with order 2. This means that in most cases only its first and second moments (i.e., the mean and autocorrelation) are independent of time  $t$  and only dependent on the time interval  $\tau$ . More precisely, if a random process is WSS, then

$$\begin{aligned} E[X(t)] &= \mu \quad (\text{i.e., its mean is a constant}) \\ \text{and } R_X(t, s) &= E[X(t)X(s)] = R_X(|s - t|) = R_X(\tau). \end{aligned} \quad (5.56)$$

Note that the autocorrelation of a WSS process is dependent only the time difference  $\tau$ . A random process that is not stationary to any order is **nonstationary** (i.e., its moments are explicitly dependent on time).

The discussion now returns to the properties of a Gaussian random process. A Gaussian random process is a continuous random process with PDF of the form

$$f_X(x, t) = \frac{1}{\sqrt{2\pi\sigma(t)^2}} \exp\left(\frac{-\mu(x, t)^2}{2\sigma(t)^2}\right), \quad (5.57)$$

where  $\mu(x, t)$  represents the mean of the random process and  $\sigma(t)$  represents the variance. A normal random process is a special case of a Gaussian random process in that it has a mean of zero and a variance of unity. A **Poisson random process** is a discrete random process with parameter  $\lambda(t) > 0$  and has a PDF given by

$$p_X(k) = P(X(t) = k) = e^{-\lambda t} \frac{(\lambda t)^k}{k!}, \quad (5.58)$$

where  $\lambda(t)$  is generally time dependent. The mean and variance of a Poisson random process is  $\lambda(t)$ . So, for a Poisson random process,

$$\mu_X = E[X(t)] = \lambda(t) \quad (5.59)$$

$$\sigma_X^2 = \text{Var}(X(t)) = \lambda(t). \quad (5.60)$$

The statistical measures above are those necessary to statistically describe digitally modulated signals and also describe most noise processes.



# Index

- $G_P$ , 187
  - $G_{PC}$ , 189
  - $G_{PS}$ , 190
  - $\nabla \cdot$ , 11
  - $\nabla \times$ , 11
  - MCS, 216
  - 0G, 16, 197
  - 16-QAM, 62, 64, 114, 115
  - 1G, 196, 197
  - 2.5G, 201
  - 2G, 197, 199
  - $3\pi/8$ -8PSK, 63
  - 32-QAM, 64
  - 3G, 196, 198, 201, 204
    - beyond 3G, 208
    - evolved, 218
    - radio interfaces, 205
  - 3GPP, 204, 205, 207
    - timeline, 207
  - 4G, 198, 208
    - FDD, 213
    - TDD, 213
  - 5G, 219
  - 64-QAM, 64
  - 6G, 223
  - 8-PSK, 63
  - 802.11n, 215
  - 802.16 (WiMAX), 185
  
  - $\epsilon_0$ , 7
  - $\mu_0$ , 6
  - $\pi/4$ -DQPSK, 59
  - $\pi/4$  Quadrature Phase Shift Keying, 57
  - $\pi/4$ -QPSK, 57
  - $\rho_V$ , 11
  - $\rho_{mV}$ , 11
  
  - AAS, 221
  - access
    - scheme, *see* multiple access scheme
  - ACI, 69
  - ACPR, 67, 69
  - active
    - device, 86
    - element, 86
  - active antenna systems, 221
  - adaptive antenna, 221
  - ADC, 85, 86, 226
    - subsampling, 95
  - adjacent channel
    - interference, 67, 69
    - power ratio, 69
  - Advanced Mobile Phone System, 178
  - AM, 13, 15, 32, 37, 44, 81
    - PMEPR, 34
    - radio, 44
  - Ampere's law, 7
  - amplifier, 86
    - low noise, 86
    - saturating, 47, 90
  - amplitude
    - modulation, 13
    - shift keying, 13
  - AMPS, 178, 180, 182, 197
    - attributes, 197
    - DAMPS, 197
    - SIR minimum, 180
  - AMTS, 197
  - analog
    - to-digital converter, 85
    - modulation, 28, 197
  - AM, 37, 44, 81
  - FM, 37, 41
  - FM bandwidth, 43
  - FM narrowband, 43
  - FM suppressed carrier, 90
  - FM wideband, 43
  - PM, 37, 41
  - PM bandwidth, 43
  - radio, 47, 197
- antenna, 129
  - adaptive, 221
  - aperture, 131, 141
  - directivity, 138
  - diversity, 148
  - efficiency, 138
  - gain, 138, 141
  - isotropic, 137
  - loss, 138, 140
  - main lobe, 137
  - major lobe, 137
  - resonant, 130
  - smart, 221
  - stacked dipole, 135
  - standing-wave, 130
  - traveling-wave, 130, 136
  - trisection, 160
  - Vivaldi, 136
  - wire, 133
- antenna array, 163
- aperture, 141
- aperture antenna, 131
- Armstrong, 44, 80, 82
- ARP, 197
- array gain, 215
- ASK, 13, 45
- audion, 80
- autocorrelation function, 237
- autodyne, 80, 81
- Autotel, 197
- average power, 33
- 
- B, 6, 18
- background noise, 180
- Bahrain Telephone Company, 178
- bandwidth, octave, 5
- base station, 16
- baseband, 13–15, 37, 51, 183
  - bitstream, 188
- beam bending, 153
- Bel, 18
- bel, 18
- Bell Laboratories, 178
- System, 178
- BER, 69, 174, 192, 193
- beyond 3G, 208
- BFSK, 65
- binary
  - phase shift keying, 53
  - PSK, 51
- Biot-Savart law, 7
- bistatic radar, 225
- bit error rate, 69, 174, 192
- bit rate, 46, 183, 186, 195
  - channel, 186
- bits per second, 10
- bitstream, 45, 51, 56
  - baseband, 188
  - channel, 188
  - information, 188
- block code, 189
- Bluetooth, 53
- BPF, 86
- bps, 10, 46
- BPSK, 51, 53
  - carrier recovery, 55
  - SBPSK, 62
  - shaped, 62
- broadcast
  - definition, 175
  - operation, 175
  - radio, 15
- BS, 16
- bursty RF, 200
- 
- $\hat{C}$ , 186, 218
- CA, 215
- call flow, 181
- canyon effect, 163
- capacity factor
  - MIMO, 218
- capacity limit, 185, 186
- car radio phone, 197
- carrier, 38, 45, 51
  - pseudo, 33
  - recovery, 55, 57
  - sense multiple access, 185
- carrier aggregation, 215
- Carson's rule, 43
- cathode, 81
- CDF, 235
  - joint, 235
- CDMA, 181, 183, 184, 197, 201, 205, 218
  - CDMA2000, 198, 205
  - CDMAOne, 184, 197, 201
  - direct sequence, 183
  - FH, 184
  - frequency hopping, 184
  - wideband, 184
- cell, 178
  - definition, 179
- cellular, 197
  - communications, 178
  - concept, 178
  - phone systems, 196
  - radio, 178
    - call flow, 181
    - definition, 178
    - hand-off, 181
- central moment, second, 236
- CF, 29
- channel
  - bit rate, 186
  - bitstream, 188
  - circuit-switched, 205
  - code, 189
  - dispersive multipath, 202
  - management, 182
  - spectrum efficiency, 186

- charge
  - electric, 11
  - magnetic, 11
- chip rate, 183
- circuit switched, 205
- circuit-switched
  - data, 197
  - radio, 205
- circulator, 225
- cluster, 180, 181
  - definition, 179
- CMOS, 95
- cochannel interference, 66, 182, 215
- code
  - block, 189
  - channel, 189
  - convolution, 189
  - division multiple access, 201
  - FEC, 189
  - forward error correction, 189
  - Gray, 66, 193
- coding, 183
  - bits, 46
  - gain, 189
  - space-time, 214
- cognitive radio, 220
- Colebrook, 81
- complex
  - Gaussian distribution, 237
- constellation, 56
  - diagram, 56, 60, 61
  - $\pi/4$ -DQPSK, 60
  - 16-QAM, 62, 64, 114, 115
  - 32-QAM, 64
  - 64-QAM, 64
  - and distortion, 68
  - and noise, 68
  - and phasor diagram, 48
  - FOQPSK, 62
  - GMSK, 62
  - OQPSK, 62
  - QAM, 64
  - RMS, 53
  - SOQPSK, 62
- continuous
  - phase modulation, 62
- conventional radio, 15
- conversion
  - frequency, 93
  - low IF, 95
  - subsampling, 95
  - zero IF, 94
- converter
  - digital-to-analog, 84
- convolution code, 189
- correlated signals, 35
- correlation, 35
- cosinusoid, 32
- Costas loop, 55
- CP, 212
- CPM, 62
- crest factor, 29
- crystal detector, 80
- CSD, 197
- CSMA, 183, 185
- cumulative
  - distribution function, 235
  - joint, 235
- curl, 11
- current
  - electric, 11
  - filament, 131
  - magnetic, 11
- cyclic prefix, 212
- D, 6
- D layer, 3
- DAC, 84, 86
- DAMPS, 199
- dB, 18
- dBm, 18
- dBV, 21
- dBW, 18
- decibel, 18
- delay spread, 148
- demodulator, 77
- DF, 235
- differential
  - coding, 60
  - encoding, 59
- differential quadra phase shift keying, 59
- diffraction, 144
- digital
  - advanced mobile phone system, 199
  - AMPS, *see* DAMPS
  - modulation, 28, 197
  - summary, 65
  - radio, 47, 199
  - signal processing, 77
  - signal processor, 14, 86
  - to analog converter, 84, 86
- diplex, 174
  - operation, 177
- diplexer, 86, 177
- direct
  - conversion, 94, 95
  - sequence, *see* CDMA
- directional
  - antenna, 160
- directivity, 137
  - antenna, 138
  - gain, 138
- dispersion, 236
- dispersive multipath
  - channel, 202
- distribution
  - complex Gaussian, 237
  - function, 235
  - cumulative, 235
  - cumulative joint, 235
- Gaussian, 236
  - normal, 236
- div, 11
- diversity
  - gain, 215
  - receiver, 204
- DL, 177
- Doppler shift, 228
- double
  - conversion receiver, 96
- downlink, 177, 182, 200
- DQPSK, 59
  - GMSK comparison, 68
  - MSK comparison, 68
- DS, 183
- DS-CDMA, *see* CDMA
- DSCDMA, 183
- DSP, 14, 77, 86
- DSSS, 183
- DTMF, 197
- ducting, 146, 149
- duplex, 174, 175, 177
  - half
    - ITU definition, 175
    - operation, 175, 177
    - TDD, 176
- duplexer, 86
- duplexing, 177
- dynamic
  - spectrum management, 220
- $E_b$ , 69
- E layer, 3
- e, natural number, 18
- early radio, 13
- EBNO, 189
- EDGE, 63, 198
- effective
  - aperture, 141
  - size, 151
  - isotropic radiated power, 141
  - radiated
    - isotropic power, 141
    - power, 141
- efficiency
  - antenna, 138
  - link spectrum, 46
  - modulation, 46
  - spectrum, 46, 185
- EHF, 3
- EIRP, 141
- electric
  - charge, 11
  - current, 11
  - field, 6
  - flux, 6
- emBB, 219
- enhanced data rates for
  - GSM evolution, 63
- enhanced mobile
  - broadband, 219
- envelope, 38, 44, 200
- GMSK, 48
- equivalent radiated power, 141
- Ericsson, 178
- ERP, 141
- error
  - rate
    - bit, 193
    - symbol, 192, 193
    - vector magnitude, 69
- EV-DO, 198, 205
- EVM, 69
- evolution-data optimized, 205
- evolved 3G, 218
- expectation, 236
- expected value, 235
- extremely high frequency, 3
- F layer, 3
- fading, 145, 147
  - ducting, 149
  - fast, 147
  - flat, 146
  - multipath, 147
  - rain, 148
  - Rayleigh, 144, 147
  - Rician, 147
  - Rician, 143
  - shadow, 146
  - slow, 146
  - thermal, 146
- family radio service, 175
- Faraday's law, 7
- fast
  - fading, 147
  - Fourier transform, 210
- fast fading, 147
- FCC, 178
- FDD, 176
  - 4G, 213
- FDMA, 182, 183, 197
- Federal Communications Commission, 178
- Feher
  - QPSK, 62
- FFT, 210
- FH, 184

- FH-CDMA, 184  
 FHSS, 184  
 field  
   electric, 6  
   magnetic, 6  
 fifth-generation radio, 219  
 filament, 131  
 filter  
   bandpass, 86  
 finite impulse response  
   filter, 115  
 FIR filter, 115  
 first  
   generation radio, 196, 197  
   moment, 236  
 FL, 177  
 flat fading, 146  
 flux  
   electric, 6  
   magnetic, 6  
 FM, 37, 41, 44  
   bandwidth, 43  
   narrowband, 43  
   suppressed carrier, 90  
   wideband, 43  
 FOMA, 198, 205  
 FOQPSK, 62  
 formulas  
   logarithm, 18  
 forward  
   link, 177  
   path, 177  
 fourth-generation radio,  
   198, 208  
 FQPSK, 62  
 freedom of mobile  
   multimedia access, 205  
 frequency  
   conversion, 93  
   division duplex, *see* FDD  
   hopping  
     multiple access, 184  
     spread spectrum, 184  
   reuse, 160, 178, 179  
   reuse plan, 160  
 Fresnel zone, 153  
 Friis transmission  
   equation, 142  
   formula, 142  
 front end, 86  
 FRS radio, 175  
 FSK, 45, 47, 48, 197  
   compared to QPSK, 69  
   constellation, 48  
  
 gain  
   antenna, 138  
   array, 215  
   coding, 189  
   diversity, 215  
   multiplexing, 215  
   processing, 187  
   spreading, 190  
 GAN, 198  
 Gauss's law, 8  
 Gaussian  
   distribution, 236  
   complex, 237  
   minimum shift keying,  
     48  
   random process, 237  
 Gbit/s, 208  
 generation  
   0G, 16  
   1G, 196, 197, 199  
   2G, 196, 197, 199  
   3G, 196, 198, 205  
   4G, 196, 198  
   5G, 196, 198  
   6G, 223  
 gigabit per second, 208  
 Global System for Mobile  
   Communications, *see*  
   GSM  
 GMSK, 48, 62, 63, 68, 197  
   compared to QPSK, 69  
   DQPSK comparison, 68  
   example, 195  
   QPSK comparison, 69  
 GMTI  
   radar, 228  
 GMTI radar, 228  
 GPRS, 198  
 grating lobe, 165  
 Gray code, 66, 193  
   mapping, 193  
 ground  
   reflection, 143  
   wave, 3  
 GSM, 48, 182, 199, 218  
   bands, 200  
   GMSK, 63  
 Groupe Spécial Mobile, *see*  
   GSM  
 guard band, 210  
  
*H*, 218  
 halfduplex  
   ITU definition, 175  
 hand-off, 181  
 Hartley  
   modulator, 39, 78, 87  
   receiver, 92  
 heterodyne, 15, 80, 82  
   frequency conversion, 94  
   mixing, 94  
   receiver, 80  
 HF, 3  
 Hicap, 197  
 high frequency, 3  
   homodyne, 80, 81, 93  
   HSCSD, 198  
   HSDPA, 198  
   HSUPA, 198  
  
 IEEE  
   802.11n, 215  
   802.16 (WiMAX), 185  
 IF, 14, 15, 86  
 IFFT, 210  
 IMD, 174  
 implementation margin,  
   71, 228, 229  
 IMT-2000, 205  
 IMT-DS, 205  
 IMT-FT, 205  
 IMT-MC, 205  
 IMT-SC, 205  
 IMTS, 197  
 information  
   bit rate, 46  
   bitstream, 188  
   rate, 46  
 instrumentation, scientific,  
   and medical, 14  
 Integrated Services Digital  
   Network, 199  
 inter-symbol interference,  
   210  
 interference, 15, 180  
   adjacent channel, 67  
   cochannel, 66, 182  
   inter-symbol, 210  
   intersymbol, 158  
   ISI, 158  
   radio link, 160  
 intermodulation  
   distortion, 174  
 International  
   Mobile  
     Telecommunications,  
     205  
   Telecommunications  
     Union, 175, 205  
 internet of things, 219  
 intersymbol interference,  
   158  
 ionosphere, 3  
   D layer, 3  
   E layer, 3  
   F layer, 3  
 IoT, 219  
 IRP, 141  
 IS-54, 199  
 IS-95, 197, 205  
 ISDN, 199  
 ISI, 158  
 ISM band, 14  
 isotropic  
   antenna, 137  
   power, 141  
   radiated power, 141  
 ITU, 175, 205  
  
**J**, 11  
 joint CDF, 235  
  
*k*, 71, 132  
  
 lag parameter, 237  
 LF, 3  
 line of sight, 143, 153  
 link, 129, 143  
   interference, 160  
   loss, 150  
   spectrum efficiency, 46  
 ln, 18  
 LNA, 86  
 LO, 84, 86  
 local oscillator, 84, 86  
 log, 18  
 log<sub>10</sub>, 18  
 logarithm, 18  
   formulas, 18  
 long  
   term evolution, 198  
 long term evolution, 208  
 long-term evolution, 205  
 LOS, 143, 153  
 low  
   frequency, 3  
   IF conversion, 95  
   noise amplifier, 86  
 LTE, 198, 205, 208  
   FDD, 213  
   TDD, 213  
   timeline, 207  
  
**M**, 11  
 magnetic  
   charge, 11  
   current, 11  
   energy, 6  
   field, 6  
   flux, 6  
   main lobe, 137  
   major lobe, 137  
   mapping, Gray code, 193  
 massive  
   machine-to-machine  
     communication, 219  
 mathematical  
   expectation, 236  
 Matsushita, 178  
 maximal-ratio combining,  
   204  
 Maxwell, 10, 12  
 Maxwell's equations, 10  
   point form, 10

- MCS, 216  
 mean, 236  
 medium  
   frequency, 3  
 MER, 71  
 mesh radio, 219  
 MF, 3  
 microcell, 161  
 MIMO, 198, 213, 215  
   capacity, 215, 218  
   capacity factor, 218  
 mixer, 86  
 mixing, 93  
 mMTC, 219  
 modulation, 36  
    $\pi/4$  Quadrature Phase Shift Keying, 57  
    $\pi/4$ -QPSK, 57  
   8PSK, 63  
   AM, 37, 44, 81  
   analog, 28  
   and coding scheme, 216  
   binary phase shift keying, 53  
   BPSK, 53  
   differential quadrature phase shift keying, 59  
   digital, 28  
   DQPSK, 59  
   efficiency, 46, 65  
   GMSK, 48  
   table, 65  
   error ratio, 71  
   FM, 37, 41  
   bandwidth, 43  
   narrowband, 43  
   suppressed carrier, 90  
   wideband, 43  
   GMSK- $\pi/4$ DQPSK  
     comparison, 68  
   Hartley, 78, 87  
   offset quadrature phase shift keying, 61  
   OQPSK, 61  
   phase, 40  
   Phase shift keying, 50  
   phase shift keying  
     8 state, 63  
   PM, 37, 41  
   bandwidth, 43  
   polar, 90  
   PSK, 50  
   QPSK, 56  
   quad-phase shift keying, 56  
   quadrature, 89  
   single-sideband, 79, 87  
   SSB, 79, 87  
   SSB-SC, 39, 79, 87  
   suppressed carrier, 79, 87  
   Weaver, 87  
   modulation and coding scheme, 216  
   modulator, 77  
     Hartley, 39  
   moment  
     second central, 236  
     second raw, 236  
   monostatic radar, 225  
   Morse, 9, 10  
     code, 8, 9, 13  
   moving target indicator, *see* radar  
   MRC, 204  
   MSK, 47, 48  
     DQPSK comparison, 68  
     envelope, 48  
   MTI, 228  
   MTS, 197  
   multi-user superposition transmission, 223  
   multipath, 59, 150, 202  
     dispersive, 202  
     fading, 147  
     knife-edge diffraction, 144  
     OFDM, 210  
     radar, 227  
   multiple access  
     carrier-sense, 185  
     code division, 201  
     orthogonal frequency division, 212  
     scheme, 182  
     CDMA, 182  
     CSMA, 185  
     FDMA, 182  
     OFDMA, 182  
     TDMA, 182  
   multiple input  
     multiple output, 213  
   multiple input multiple output, 213  
   multiplex, 174  
   multiplexing, 177  
     gain, 215  
   MUST, 223  
   MUXing, 177  
    $N_o$ , 69  
   NADC, *see* DAMPS, 199  
   natural  
     logarithm, 18  
     number, *e*, 18  
   next radio, 219  
   NextGen, 208  
   nibble, 60  
   NLOS, 145  
   NMT, 197  
   noise  
     background, 180  
     threshold, 15  
   Nokia, 178  
   NOMA, 222  
   non orthogonal multiple access, 222  
   nonstationary, 237  
   Nordic Mobile Telephone, *see* NMT  
   normal distribution, 236  
   North American digital cellular, 199  
   NR, 219  
   NTT, 178  
   Nyquist, 52, 95  
     signaling theorem, 52  
   octave bandwidth, 5  
   OFDM, 198, 209  
     frame, 211  
     multipath, 210  
     multiuser, 212  
     PMEPR, 210  
     resource block, 211  
   OFDMA, 185, 205, 212  
   offset quadrature phase shift keying, 61  
   Okumura-Hata model, 155  
   OMA, 222  
   OQPSK, 61  
   origins of radio, 10  
   orthogonal frequency division multiple access, 205, 212  
   orthogonal multiple access, 222  
   oscillator, 86  
     local, 86  
    $P_{avg}$ , 33  
   PA, 90  
   packet switched, 205  
   packet-switched radio, 205  
   PALM, 197  
   PAM, 204  
   PAPR, 29  
     compared to PMEPR, 36  
   PAR, 29  
     compared to PMEPR, 36  
   path loss, 150  
   PCS, 181  
   PDC, 197  
   PDF, 235  
   peak  
     -to-average power ratio, 29  
     -to-average power ratio, 29  
     -to-average ratio, 29  
     -to-mean envelope power ratio, 32  
     -to-mean envelope power ratio, 31  
     envelope power, 32  
   PEP, 32  
   PEPR  
     OFDM, 210  
   permeability, 6  
   relative, 7  
   permittivity, 7  
   relative, 7  
   personal  
     communication service, 181  
     digital cellular, 197  
     handyphone system, 197  
   phase  
     locked loop, *see* PLL  
     modulation, 40  
     modulator, *see* PLL  
     shift keying, 50  
       8 state, 63  
       binary, 51  
   phasor, 42  
   diagram  
     and constellation diagram, 48  
   PHS, 197  
   picocell, 161  
   pine needle, 145  
   PLL, 47, 90  
   PM, 37, 40, 41  
     bandwidth, 43  
   PMEPR, 31–35  
     AM, 34  
     compared to PAPR, 36  
     compared to PAR, 36  
     FM, 42  
     OFDM, 210  
     PM, 42  
   PMF, 235  
   point-to-point link, 153  
   Poisson  
     random process, 237  
   polar modulation, 90  
   power  
     amplifier, 90  
     average, 33  
     ratio  
       peak-to-average, 29  
       peak-to-mean envelope, 31  
   Poynting vector, 133  
   predetection combining, 204  
   probability  
     density function, 235  
     mass function, 235  
   processing gain, 187  
   propagation  
     loss, 151  
     model, 155

- Okumura–Hata, 155
- pseudo
  - carrier, 33
- PSK, 45, 50, 51
  - carrier recovery, 55, 57
- PTT, 176
- pulse radio, 36
- pulsed
  - amplitude modulation, 204
  - radar, 226
- push-to-talk, 197
  
- Q, 80
- QAM, 64
- QoS, 215
- QPSK, 52, 56
  - carrier recovery, 57
  - compared to FSK, 69
  - compared to GMSK, 69
  - constellation, 58
  - Feher offset, 62
  - GMSK comparison, 69
  - offset, 61
  - shaped offset, 62
  - staggered, 62
- quad-phase shift keying, 56
- quadra phase, 56, 61
- quadrature
  - amplitude modulation, 64
  - modulation, 89
  - phase shift keying, 52
- quality of service, 215
  
- $R_b$ , 186
- $R_c$ , 186
- radar, 224–228
  - band, 224
  - bistatic, 225
  - chirp, 228
  - cross section, 226, 227
  - CW, 228
  - Doppler, 227
    - shift, 228
  - equation, 227
  - GMTI, 228
  - ground moving target
    - indication, 228
  - monostatic, 225
  - moving target, 228
  - multipath, 227
  - PSK, 228
  - pulsed, 226, 227
  - RCS, 226, 227
  - synthetic aperture, 225
  - waveform, 225
    - CW, 225
    - FM/CW, 225
- phase-encoded, 225
  - pulse, 225
- radiation
  - density, 137
  - efficiency, 138
  - intensity, 137
- radio, 15
  - analog, 47
  - broadcast, 15
  - circuit-switched, 205
  - cognitive, 220
  - conventional, 15
  - digital, 47
  - early radio, 13
  - era, 13
  - frequency, *see* RF
  - generations, 197
  - link
    - interference, 160
    - reciprocity, 151
  - link interference, 160
  - mesh, 219
  - origins, 10
  - software-defined, 77
  - spectrum efficiency, 186
  - systems, 197
- Radiocommunication
  - Assembly, 205
- rain
  - fading, 148
  - scattering, 148
- rake receiver, 202, 203
- random process, 235
  - Gaussian, 237
  - mean, 236
  - nonstationary, 237
  - Poisson, 237
  - stationary, 237
  - wide-sense stationary, 237
- ratio-squared combining, 204
- Rayleigh fading, 144, 147
- RCS, 226, 227
- received signal strength
  - indicator, *see* RSSI
- receiver, 42
  - architecture, 92
  - autodyne, 80
  - direct conversion, 94
  - diversity, 204
  - double conversion, 96
  - early technology, 80
  - homodyne, 80
  - low IF, 95
  - quadrature, 95
  - superheterodyne, 82
  - syncrodyne, 80
- reflection
  - ground, 143
- regenerative circuit, 80
- resonant antenna, 130
- resource block, 211
- reverse path, 176
- RF, 1, 14, 15
  - bursts, 200
  - front end, 86
  - link, 129, 143
  - interference, 160
- Rician fading, 143, 147
- RL, 176
- RSSI, 181
- Rx, 96
  
- saturating
  - amplifier, 90
  - power amplifier, 91
- SBPSK, 62
- SC, 39
- SC-FDMA, 212
- SC-OFDMA, 212
- scattering
  - rain, 148
  - resonant, 145
- SCDMA, 205
- SCSS, 90
- SDR, 77, 97
- second
  - central moment, 236
  - generation radio, 197, 199
  - moment, raw, 236
  - raw moment, 236
- SER, 192
- shadow fading, 146
- Shannon's
  - capacity limit, 185
  - theorem, 185, 186
- shaped
  - binary phase sift keying, 62
  - offset QPSK, 62
- SHF, 3
- SIC, 223
- side lobe, 165
- signal
  - to-interference ratio, 66, 185
  - signal blocking, 154
- SIM, 181
- simplex, 174, 175
  - definition, 175
  - operation, 175
  - definition, 175
- single-carrier frequency
  - division multiple
    - access, 212
- single-carrier orthogonal
  - frequency division
    - multiple access, 212
- single-sideband, 39
- SIR, 66, 180, 185, 214
  - example, 130, 162, 172
  - sixth generation radio, 223
  - sky wave, 3
  - slow fading, 146
  - smart antenna, 221
  - SNR, 69, 185
  - software defined radio, 77, 97
  - SOQPSK, 62
  - space
    - diversity, 148
    - time coding, 214
  - spark gap, 12
  - spatial signature, 221
  - spatio-temporal coding, 214
  - spectrum
    - efficiency, 46, 185, 186
    - channel, 186
    - radio, 186
    - management, 220
  - spread spectrum, 183, 184, 198
  - spreading
    - gain, 190
  - SQPSK, 62
  - sr, 137
  - SS, 183
  - SSB, 39
  - SSB-SC, 39, 90
  - staggered quadrature
    - phase shift keying, 62
  - standard
    - deviation, 236
  - standards, wireless, 173
  - standing-wave
    - antenna, 130
  - stationary
    - random process, 237
  - steradian, 137
  - subcarrier, 209
  - subsampling receiver, 95
  - subscriber identification
    - module, 181
  - successive interference
    - cancellation, 223
  - super high frequency, 3
  - superheterodyne, 82
  - supersonic, 15, 82
  - suppressed-carrier, 39
  - surface
    - wave, 3
  - switch, 86
  - symbol
    - error, 192
    - error rate, 192
    - rate, 52
  - syncrodyne, 80
  - synthetic aperture radar, 225
  - system, 1

- TACS, 197
- TD-CDMA, 205
- TD-SCDMA, 198, 205
- TDD, 176
  - 4G, 213
- TDMA, 182, 183, 197, 198
- telegraph, 8, 12
- telephone, 12
- television, 15
- temperature inversion, 146, 150
- Tesla coil, 12
- thermal fading, 146
- thermionic
  - emission, 81
- THF, 3
- third
  - generation radio, 196, 198, 201, 204
- Third Generation
  - Partnership Project, 204
- three
  - tone signal, 34
- time division
  - duplex, 176
  - multiple access, 182
- timeline
  - 3GPP, 207
  - LTE, 207
- Titanic, 14
- tone
  - three-tone, 34
  - two-tone, 34
- transceiver, 86, 96
- transmitter
  - quadrature, 95
- traveling
  - wave antenna, 130, 136
- tremendously high
  - frequency, 3
- TRF, 80
- triode, 81
- tuned radio frequency, 80
- two-tone signal, 34
- Tx, 96
- UHF, 3
- UL, 176
- ultra
  - high frequency, 3
- ultra-wideband, 36
- UMA, 198
- UMTS, 198, 205
- uncorrelated signals, 35
- United Nations, 175
- uplink, 176, 182, 200
- urban
  - canyon effect, 163, 179
  - waveguide effect, 179
- UWB, 36
- variance, 236
- VCO, 81, 89, 91
  - Hartley, 81
- very
  - high frequency, 3
  - low frequency, 3
- VHF, 3
- VLF, 3
- VoIP, 198, 207
- voltage
  - controlled oscillator, 81, 89
  - gain and power gain, 20
- W, 18
- WARC, 178
- watt, 18
- wave-
  - channeling effect, 163
  - guiding effect, 163
- waveguide
  - bands, 6
- wavelength, 11, 15
- wavenumber, 132
- WBB, 198
- WCDMA, 184, 204
- Weaver modulator, 87
- wide-sense
  - stationary, 237
- WiDEN, 198
- WiFi, 215
- WiMAX, 185, 205
- wireless
  - standards, 173
- WLAN, 213
- World Administrative
  - Radio Conference, 178
- WSS, 237
- zero
  - generation radio, 16, 197
  - IF conversion, 94

*Microwave and RF Design: Radio Systems* is a circuits- and systems-oriented approach to modern microwave and RF systems. Sufficient details at the circuits and sub-system levels are provided to understand how modern radios are implemented. Design is emphasized throughout. The evolution of radio from what is now known as 0G, for early radio, through to 6G, for sixth generation cellular radio, is used to present modern microwave and RF engineering concepts. Two key themes unify the text: 1) how system-level decisions affect component, circuit and subsystem design; and 2) how the capabilities of technologies, components, and subsystems impact system design. This book is suitable as both an undergraduate and graduate textbook, as well as a career-long reference book.

## KEY FEATURES

- The first volume of a comprehensive series on microwave and RF design
- Open access ebook editions are hosted by NC State University Libraries at: <https://repository.lib.ncsu.edu/handle/1840.20/36776>
- 31 worked examples
- An average of 38 exercises per chapter
- Answers to selected exercises
- Coverage of cellular radio from 1G through 6G
- Case study of a software defined radio illustrating how modern radios partition functionality between analog and digital domains
- A companion book, *Fundamentals of Microwave and RF Design*, is suitable as a comprehensive undergraduate textbook on microwave engineering

## ABOUT THE AUTHOR

**Michael Steer** is the Lampe Distinguished Professor of Electrical and Computer Engineering at North Carolina State University. He received his B.E. and Ph.D. degrees in Electrical Engineering from the University of Queensland. He is a Fellow of the IEEE and is a former editor-in-chief of *IEEE Transactions on Microwave Theory and Techniques*. He has authored more than 500 publications including twelve books. In 2009 he received a US Army Medal, "The Commander's Award for Public Service." He received the 2010 Microwave Prize and the 2011 Distinguished Educator Award, both from the IEEE Microwave Theory and Techniques Society.

## OTHER VOLUMES

Microwave and RF Design  
Transmission Lines  
Volume 2  
ISBN 978-1-4696-5692-2

Microwave and RF Design  
Networks  
Volume 3  
ISBN 978-1-4696-5694-6

Microwave and RF Design  
Modules  
Volume 4  
ISBN 978-1-4696-5696-0

Microwave and RF Design  
Amplifiers and Oscillators  
Volume 5  
ISBN 978-1-4696-5698-4

## ALSO BY THE AUTHOR

Fundamentals of Microwave  
and RF Design  
ISBN 978-1-4696-5688-5