

Effect of Slow Fading and Adaptive Modulation on TCP/UDP Performance of High-Speed Packet Wireless Networks

Xuanming Dong



Electrical Engineering and Computer Sciences
University of California at Berkeley

Technical Report No. UCB/EECS-2006-109

<http://www.eecs.berkeley.edu/Pubs/TechRpts/2006/EECS-2006-109.html>

August 25, 2006

Copyright © 2006, by the author(s).
All rights reserved.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission.

**Effect of Slow Fading and Adaptive Modulation on TCP/UDP
Performance of High-Speed Packet Wireless Networks**

Copyright 2006

by

Xuanming Dong

Abstract

Effect of Slow Fading and Adaptive Modulation on TCP/UDP Performance of High-Speed Packet Wireless Networks

by

Xuanming Dong

Doctor of Philosophy in Electrical Engineering and Computer Sciences

University of California, Berkeley

Professor Pravin Varaiya, Chair

High speed data wireless networks in multipath environments suffer channel impairment from many sources such as thermal noise, path loss, shadowing, and fading. In particular, short-term fading caused by mobility imposes irreducible error floor bounds on system performance. We study the effect of fading on the performance of the widely used TCP/UDP protocol, and investigate how to improve TCP performance over fading channels. Our solutions target upcoming mobile wireless systems such as IEEE 802.16e wireless MANs (Metropolitan Area Networks) where adaptive modulation is enabled and the underlying medium access scheme is On-Demand Time Division Multiple Access (On-Demand TDMA).

Adaptive modulation is used in the new generation of wireless systems to increase the system throughput and significantly improve spectral efficiency by matching parameters of the physical layer to the time-varying fading channels. Most high-rate applications for such wireless systems rely on the reliable service provided by TCP protocol. The effect of adaptive modulation on TCP throughput is investigated. A semi-Markov chain model for TCP congestion/flow control behavior

and a multi-state Markov chain model for Rayleigh fading channels are used together to derive the steady state throughput of TCP Tahoe and Reno. The theoretical prediction based on our analysis is consistent with simulation results using the network simulator NS2. The analytical and simulation results triggered the idea of cross-layer TCP protocol design for single-user scenarios. The fading parameters of wireless channels detected in the physical layer can be used to dynamically tune the parameters (such as packet length and advertised receiver window size) of the TCP protocol in the transport layer so that TCP throughput is improved.

For multi-user scenarios, we study how multi-user diversity can be used to improve the aggregate TCP throughput of base stations in fading channels. Since TCP performance involves complex interactions among layers of the networking protocol stack, the cross-layer design approach is adopted to tackle the problem. The performance improvement is achieved through channel-aware packet scheduling algorithms and active delay of TCP ACK packets in the buffer. Based on the adaptive modulation information from the physical layer, the advertised receive window size of TCP ACK packets is dynamically changed to accommodate the rate changes resulting from adaptive modulation. Our simulation results show that the new cross-layer approach increases TCP throughput.

Contents

List of Figures	iii
List of Tables	v
1 Introduction	1
1.1 Evolution of Wireless Networks	1
1.1.1 From 1G to 4G	1
1.1.2 From Circuit Switching to Packet Switching	4
1.2 Wireless Channels	7
1.2.1 Path Loss Model and Average Received Signal Power	8
1.2.2 Linear Time-Variant Channel Model	12
1.2.3 Wide-Sense Stationary Uncorrelated Scattering Channel Model	14
1.3 Main Ideas of OFDM	16
1.4 IEEE 802.11a, IEEE 802.11p, and IEEE 802.16e	18
1.4.1 IEEE 802.11a	18
1.4.2 IEEE 802.11p	20
1.4.3 IEEE 802.16e	22
1.5 Dissertation Organization	24
2 Effect of Doppler Spread and Delay Spread on Performance of High-Speed Packet Wireless Networks	26
2.1 Introduction	26
2.2 Multipath Fading Channel Model	28
2.3 IEEE 802.11a and IEEE 802.11b Simulators	31
2.4 Numerical Results	32
2.4.1 Effect of Doppler Spread	33
2.4.2 Effect of Delay Spread	34
2.5 Concluding Remark	36
3 Communication Range in Multipath Fading Channels	42
3.1 Introduction	42
3.2 System and Channel Model	43
3.3 Communication Range vs Speed	46
3.4 Numerical Results	48
3.5 Concluding Remark	50

4	Impact of Adaptive Modulation on TCP Throughput in Rayleigh Fading Channels	53
4.1	Introduction	53
4.2	TCP Tahoe and Reno	55
4.3	System Model and Basic Assumptions	60
4.3.1	System Model	60
4.3.2	Adaptive Modulation	61
4.3.3	Channel Modeling	65
4.4	Markov Model of TCP Evolution	68
4.4.1	Computation of Transition Probability	69
4.4.2	Discrete Semi-Markov Process and Throughput Calculation	70
4.4.3	Some Useful Functions	72
4.5	Analysis of TCP Tahoe and Reno	74
4.5.1	Computation of $\phi_1(z)$	74
4.5.2	Computation of $\phi_2(z)$	75
4.6	Simulation and Numerical Results	79
4.6.1	Simulation Setup and Configuration	79
4.6.2	Simulation Results and Analysis	82
4.6.3	Adaptive Configuration of TCP Parameters	88
4.7	Conclusion and Future Work	89
4.8	Proof of the Average Bit Error Probability in State i	90
4.9	A Recursive Approach to Derive the Average Delay	92
5	Exploiting Multi-user Diversity for Improvement of TCP Performance over Fading Channels	95
5.1	Introduction	95
5.2	System and Channel Model	97
5.2.1	System Model	97
5.2.2	Basic Mechanisms for Adaptive Modulation in IEEE 802.16e	100
5.2.3	Channel Model	101
5.3	A Channel-Dependent Scheduling Algorithm for TCP Flows	102
5.3.1	Buffer Management and Packet Scheduling	102
5.3.2	Proportional Fair Scheduling	103
5.3.3	A Channel-Dependent Scheduling Algorithm	107
5.4	Maximum Advertised Window Size and Rate Variation	108
5.5	Numerical Results	112
5.6	Concluding Remark	114
6	Conclusion and Future Work	115
6.1	Main contribution of this dissertation	115
6.2	Future work	116
	Bibliography	118

List of Figures

1.1	The evolution of wireless systems	3
1.2	Two-ray model	8
1.3	Path loss based on two-ray model	10
1.4	Received signal strength: path loss, shadowing, and fading	11
1.5	Fourier transform pairs of deterministic LTV wireless channels	13
1.6	Fourier transform pairs of WSSUS wireless channels	15
1.7	IEEE 802.11a PHY frame format	19
1.8	Pilot placement of IEEE 802.16e PHY layer	23
2.1	Multiplicative and additive model for fading channels	30
2.2	Tapped delay line channel model	31
2.3	IEEE 802.11a PHY simulator	37
2.4	IEEE 802.11b PHY simulator	38
2.5	Channel blocks for the effect of Doppler spread	39
2.6	Effect of Doppler spread on IEEE 802.11a	39
2.7	Effect of Doppler spread on IEEE 802.11b	40
2.8	Channel blocks for the effect of Delay spread	40
2.9	Effect of Delay spread on IEEE 802.11a	41
2.10	Effect of Delay spread on IEEE 802.11b	41
3.1	Basic elements of a mobile packet wireless system	44
3.2	Block structure of a multiplicative and additive channel model	44
3.3	SNR_{min} vs Rician coefficient K	49
3.4	Additional TX Power vs Rician coefficient K	50
3.5	Communication range vs Rician coefficient K	51
3.6	SNR_{min} vs Speed	51
3.7	Communication range vs Speed	52
3.8	Additional TX power vs Speed	52
4.1	The sliding window for TCP protocol	56
4.2	The block structure of the communication system	61
4.3	Maximum achievable data rate vs SNR: BPSK, QPSK, QAM-16, QAM-64	63
4.4	Average Bit Error Rate vs SNR: BPSK, QPSK, QAM-16, QAM-64	63
4.5	Average Packet Error Rate vs SNR: BPSK, QPSK, QAM-16, QAM-64	64
4.6	Finite state Markov chain model for a Rayleigh fading channel	66
4.7	Average throughput vs packet length ($f_d=10\text{Hz}$; $W_{max}=80\text{Kb}$; Dup=3)	83
4.8	Average throughput vs packet length ($f_d=20\text{Hz}$; $W_{max}=80\text{Kb}$; Dup=3)	83

4.9	Average throughput vs packet length ($f_d=30\text{Hz}$; $W_{max}=80\text{Kb}$; Dup=3)	84
4.10	Average throughput vs Doppler spread (SNR=25dB; $W_{max}=80\text{Kb}$; Dup=3)	85
4.11	Average throughput vs Doppler spread (SNR=30dB; $W_{max}=80\text{Kb}$; Dup=3)	86
4.12	Average throughput vs Doppler spread (SNR=35dB; $W_{max}=80\text{Kb}$; Dup=3)	86
4.13	Average throughput of Reno vs maximum received window size (SNR=35dB; DUP=3)	87
4.14	Average throughput of Tahoe vs maximum received window Size (SNR=35dB; DUP=3)	88
4.15	Adaptive configuration of TCP parameters	89
4.16	Integration along different direction	91
4.17	A recursive approach to calculate the average delay	93
5.1	Network model	97
5.2	Packet scheduling in the buffer of base station	98
5.3	Frame structure for time division duplexing	99
5.4	Burst profile threshold	101
5.5	Protocol header of TCP packets	110
5.6	Receive window of TCP protocol	111
5.7	Network topology	112
5.8	Aggregate TCP throughput vs number of users	113

List of Tables

1.1	IEEE 802.15.1, IEEE 802.11a, and IEEE 802.16a	6
1.2	Major OFDM parameters of 802.11a and 802.11p	22
3.1	Parameters for the computation of communication range	49
4.1	Level crossing rate N_i ($1 \leq i \leq M = 4$) (crossings/second)	80
4.2	Stationary probability distribution of channel state s_i ($1 \leq i \leq M = 4$)	80
4.3	Coefficients for Bit Error Probability $f_{e,i}(\alpha)$	81

Acknowledgments

The dissertation has been an arduous but enriching and growing experience for me. I would like to acknowledge many people without whom this dissertation would not have been possible.

First and foremost, I would like to express my sincere gratitude to my advisor, Professor Pravin Varaiya, for his technical guidance, patience, encouragement, and funding support. I am fortunate to have him as my advisor and am deeply appreciative of what makes him a top scientist and scholar. His critical thinking, fundamental understanding, breadth, precision, intuition, passion, easygoing personality and humor, set a perfect role model for me.

I am especially grateful to have Professor Jean Walrand and Professor Armen Der Kiureghian for being on my dissertation committee. Their remarks and comments have made significant contributions to the improvement of my work.

Dr Anuj Puri introduced me to this field. I would like to thank him for the valuable advice and intensive discussions in the early years of this research, and for his continuous friendship.

I would like to acknowledge the important feedback on my early research from Professor Ion Stoica, Professor Ahmad Bahai, Professor David Tse and Professor George Shanthikumar.

It has been my great pleasure to work with so many talented graduate students and engineers, especially, members of the Web over Wireless (WoW) group, members of EVII group in California PATH, and members of Distributed Sensor Network group. I would also like to thank them all for providing a pleasant working environment, for their friendship, and for numerous fruitful discussions.

The department administrative staff has been very helpful and encouraging. I would especially like to thank Ruth Gjerde and Mary Byrnes, for their kind assistance throughout my graduate studying.

Last, but certainly not least, I feel strongly indebted to my wife Hua for her love, understanding, patience, support, and constant encouragement, as well as her significant contributions to our family.

The research reported in this dissertation was supported in part by California Department of Transportation to the California PATH Program, National Science Foundation under Grant CMS-0408627, and by ARO-MURI UCSC-WN11NF-05-1-0246-VA-09/05.

Chapter 1

Introduction

1.1 Evolution of Wireless Networks

1.1.1 From 1G to 4G

The massive deployment of modern wireless systems began with the introduction of analog cellular service called Advanced Mobile Phone Service (AMPS). AMPS was invented at Bell Labs and first deployed in the U.S. in the early 1980's [57, 58]. The AMPS operation involves call setup/termination procedures, radio resource management (channel assignment, and power control, etc.), and mobility management (handoff and paging, etc.) functions. About the same time, similar wireless systems were introduced in other countries, e.g., Nordisk Mobiltelefon (NMT) in Scandinavia, Total Access Communication System (TACS, the European version of AMPS) and its extended version ETACS in Europe, and JTACS (Japan TACS) system in Japan [57, 75].

These wireless systems were connected to the Public Switched Telephone Network (PSTN). They were analog and circuit switched cellular systems that only carried voice traffic in wide areas. The medium access scheme of these systems is Frequency Division Multiple Access (FDMA), in which each conversation gets its own, unique, and relatively narrow radio channel. In addition,

these analog systems transmitted voice messages vaguely. As the very first generation of mobile telephony systems, they are usually referred to as 1G [57, 75]. Due to many technical and non-technical reasons (for example, wireless spectrum management policy, intellectual property rights, history, political, economics, and market), there are usually many standards that are adopted in different geographic areas and have different evolution.

Although 1G was a great success, it suffered from poor voice quality, low spectrum efficiency, short battery life, no security, etc. In the second generation of mobile telephony systems (2G), digital technology was introduced to improve performance [27, 57, 58, 75, 78, 92]. For example, digital vocoders, forward error correction (FEC), encryption, and high-order digital modulation schemes are used to improve spectrum efficiency, voice quality, and enhanced security; Very Large Scale Integration (VLSI) technology is used to dramatically reduce the phone size; Time Division Multiple Access (TDMA) and Code Division Multiple Access (CDMA) are used to further increase the spectrum efficiency. Other 2G features are expanded services such as short messaging, caller ID, and seamless roaming across multiple service providers.

Well-known 2G systems are Global System for Mobile communication (GSM) in Europe and Asia, Digital AMPS (D-AMPS, also called IS-54) in USA, CDMA(IS-95A or cdmaOne) in USA and Japan, and Personal Digital Cellular (PDC) in Japan. GSM, D-AMPS and PDC are TDMA-based technologies, which replace the analog TACS/ETACS systems in Europe, the AMPS systems in USA, and JTACS systems in Japan [27, 58, 69, 75, 92].

Nowadays 2G systems are widely deployed. However, motivated by user demand for higher data rates (up to 2 Mbps) and new data applications such as email and web browsing, the third generation of mobile telephony systems (3G) is being designed and developed [27, 36, 83, 67]. The services associated with 3G provide the ability to transfer both voice and data. Typical 3G systems or upcoming 3G systems include Universal Mobile Telecommunication System (UMTS) in Europe [94], Japan, and Asia, and CDMA2000 in USA [54]. Note that UMTS is based on Wideband CDMA (W-CDMA), using either Frequency Division Duplex (FDD) or Time Division Duplex (TDD). The

first 3G system (W-CDMA) started being by NTT DoCoMo in Japan at the end of the year 2001.

In the meantime, 2G wireless providers are upgrading their systems to 2.5G [27, 58, 78, 84], a technology that extends 2G networks with features such as packet-switched data service, enhanced data rates, and improved idle time charges. GSM networks are currently being upgraded with General Packet Radio System (GPRS), High-Speed Circuit Switched Data (HSCSD) and Enhanced Data rates for GSM Evolution (EDGE) in Europe and Asia. The 2.5G version of PDC networks is i-Mode in Japan, while the 2.5G version of cdmaOne (IS-95A) is IS-95B in USA [84].

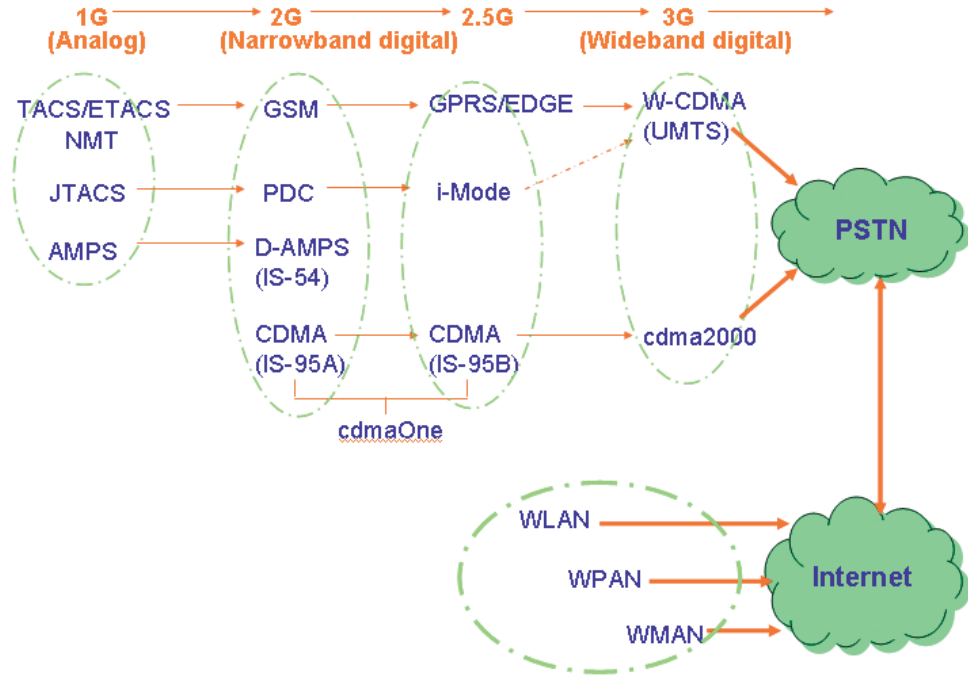


Figure 1.1: The evolution of wireless systems

While 3G deployment is promising, researchers and vendors are already thinking about the fourth generation of mobile telephony systems (4G) [86, 25, 104, 42]. Since 3G systems need to be backward compatible with 2G systems, they are a combination of existing and evolved equipments with data rate up to 2 Mbps. Therefore 3G systems may not be sufficient to meet needs of future high-performance applications such as full-motion video and wireless teleconferencing. 4G systems

are proposed to extend 3G capacity by an order of magnitude (up to 100 Mbps) and enable entirely packet-switched networks. In addition, as shown in Fig.1.1, 4G systems will be hybrid networks that integrate different kind of wireless networks such as IEEE 802.11 wireless LAN and wide-area cellular networks, etc. To achieve a 4G standard, one promising underlying technology, multicarrier modulation (MCM) [5, 12], will be adopted. MCM is a baseband process that uses parallel equal bandwidth subchannels to transmit information, as we discuss later.

1.1.2 From Circuit Switching to Packet Switching

Internet and World Wide Web (WWW) services have expanded dramatically over the past years, and have become important new communication media in daily life. The increased reliance on Internet and the strong desire of anywhere and anytime access to Internet have driven the need to design and develop data-oriented wireless systems with high capacity and reliability [69, 70, 93]. As we can see from the evolution of cellular systems, the voice-oriented application has driven the design of cellular systems in a way which is less than optimal for data applications. Although cellular systems use different combinations of digital technology, they are still mainly circuit-switched systems that provide narrow band voice and data services.

When the user is connected using a cellular system, the radio spectrum is dedicated to a single user over the entire conversation period. This is very similar to the dialup phone over a PSTN network. The local wired loop of phone service belongs to a single phone user, but the radio spectrum is shared by many users. If a user occupies the radio spectrum for a long time but uses it for a short time, the radio spectrum will not be fully utilized. The long delay to establish a wireless connection discourages short-lived data applications such as web browsing. In addition, users of cellular systems are usually charged based on the connection time regardless of the amount of through traffic. But users that need data services prefer to be charged based on the amount of through traffic instead of the connection time. Overall, due to the burstiness that data traffic usually exhibits, introducing data wireless networks based on packet switching technology can lead

to better use of the radio spectrum and attract more users [93].

Wireless local area networks (wireless LANs or WiFi) were developed in the 1990s as an extension of Ethernet, the dominant technology that enables today's Internet. Wireless LANs aim to transmit data and operate local networks without constraints of wires and associated infrastructure normally required by Ethernet. Originally, the IEEE 802.11 standard specified the MAC (Media Access Control) layer and PHY (Physical) layer in the 2.4 GHz band with data rates of 1 and 2 Mbps using either Direct Sequence Spread Spectrum (DSSS) or Frequency Hopping Spread Spectrum (FHSS) [43]. In 1999, the IEEE defined two high rate extensions: 802.11b that is based on DSSS and CCK (Complementary Code Keying) with data rates up to 11 Mbps in the 2.4 GHz band, and 802.11a that is based on OFDM (Orthogonal Frequency Division Multiplexing) technology with data rates up to 54 Mbps in the 5 GHz band [44, 45]. In 2003, the 802.11g standard that extends the 802.11b PHY layer to support data rates up to 54 Mbps in the 2.4 GHz band was finalized [46]. The IEEE 802.11 standard and its several important supplements (IEEE 802.11a/b/g) provided a basis for interoperability of different products, and triggered the explosive growth of the wireless LAN market. IEEE 802.11n, a new amendment to the 802.11 standard, is on the way [100]. IEEE 802.11n promises higher data rate in excess of 100 Mbps and longer operating distance than IEEE 802.11a/b/g wireless LANs using MIMO (Multiple Input Multiple Output) technology, and pre-standard 11n products are already widely available.

Wireless personal area networks (wireless PANs) are used for communications among devices that are very close to each other physically (typically within a few meters). The IEEE 802.15 standard series is proposed for wireless PANs. Of these, IEEE 802.15.1 is derived from the MAC layer and PHY layer of the Bluetooth specification, which is an industry standard for short-range RF-based connectivity for portable personal devices [17, 47, 105]. Advances in circuits and embedded systems have led to the development of small sensor nodes with low complexity and little power consumption. The complex protocol stack and power-consuming radios defined in IEEE 802.15.1 make it unsuitable for wireless PANs with sensor nodes. The IEEE 802.15.4 is proposed to enable sensor

networks which consist of very low-cost and battery-operated sensor nodes that can communicate with each other and a centralized device at low data rates.

Wireless metropolitan area networks (wireless MANs or WiMAX) provides the ‘last mile’ connection when DSL, cable and other broadband access methods are not available or too expensive [20]. The original IEEE 802.16 standard published in April 2002 defines a point-to-multipoint broadband wireless access standard for systems in the licensed frequency range 10-66 GHz. IEEE 802.16a published in January 2003 is designed for the licensed and unlicensed spectrum ranges of 2 GHz to 11 GHz with support for enhanced Quality of Service (QoS) features of MAC layer (for example, support multiple polling and piggyback polling requests). IEEE 802.16c deals with updates in the 10 GHz to 66 GHz range. However, it also addresses issues such as performance evaluation, testing and detailed system profiling. The 802.16 REVd published in June 2004 incorporates the original 802.16, 802.16a and 802.16c amendments [17, 48]. Among the changes is support of MIMO antennas, which will likely increase reliable range in multipath channels. IEEE 802.16 REVd makes wireless MAN ready for prime time.

Parameters	802.15.1	802.11a	802.16a
Frequency	Varies	5.4GHz	2-11GHz
Range	10Meters	100Meters	30Miles
Data Rate	20Kbps	54Mbps	70Mbps
Number of Nodes	Dozens	Dozens	Thousands

Table 1.1: IEEE 802.15.1, IEEE 802.11a, and IEEE 802.16a

The latest IEEE 802.16 update is IEEE 802.16e [49], which is in the process of being finalized. Whereas the 802.16 REVd mainly addresses fixed wireless applications, 802.16e can serve the dual purpose of adding extensions for mobility and including new enhancements to the Orthogonal Frequency Division Multiple Access (OFDMA) physical layer. This new enhanced 802.16e physical layer is now being referred to as Scalable OFDMA (SOFDMA) and includes a number of important features for fixed, nomadic, and mobile networks. With longer range (about 30 miles in a fixed

network or about 1-3 miles in a mobile (802.16e) scenario), wireless MANs are envisioned as a complementary technology to Wireless LANs and Wireless PANs. Table 1.1 provides a quick comparison of several typical networks.

1.2 Wireless Channels

In wired networks, signals are transmitted over guided medium, whose system impulse response is relatively stable and isolated from the environment. However, in wireless networks, signals are transmitted over wireless channels in open space, and suffer time-varying environmental interference. Compared to wired links, the wireless channels exhibit many different forms of channel impairments such as multipath, fading, and carrier frequency/phase noises [69, 70]. An accurate understanding of wireless channels, including approximate mathematical models and related important parameters, enables wireless engineers to predict signal transmissions and design wireless transceivers with better performance.

In early narrowband wireless communication systems with low data rate, the average signal strength along with fading description could be used to adequately predict system performance [51, 58]. Hence a path loss model meets many requirements of design tasks such as cell placement, configuration, and management of wireless systems. For today's high rate wireless systems, to achieve better spectrum/power efficiency and combat the new signal distortions and interferences introduced by high data rate, engineers need to know more about the wireless channels [69, 70].

This section introduces several important models and parameters that help to characterize a wireless channel and optimize the wireless transceiver design. These models and parameters include: path loss model, time-varying channel model, channel correlation functions or Wide-Sense Stationary Uncorrelated Scattering (WSSUS) model, distance-power gradient, Delay spread, coherence bandwidth, Doppler spread, and coherence time.

1.2.1 Path Loss Model and Average Received Signal Power

In any real channel, signals attenuate as they propagate. For single-path communications in free space, the mean signal power P_R decays as the square of the distance,

$$P_R = \frac{P_T G_T}{4\pi d^2},$$

where P_T is the transmitter power, G_T is the transmitting antenna gain in the direction of the receiver, and $4\pi d^2$ is the surface area of the sphere at radius d [69, 70].

The mean power the receive antenna intercepts depends on the antenna's effective aperture A_e . The effective aperture of the antenna is given as

$$A_e = \frac{G_R \lambda^2}{4\pi},$$

where G_R is the gain of the receive antenna in the direction of the transmitter and λ is the wavelength.

The free-space path loss or attenuation between the transmitter and the receiver is simply the ratio of the received power to the transmitted power as

$$\frac{P_R}{P_T} = G_T G_R \left(\frac{\lambda}{4\pi d} \right)^2.$$

Defining $P_0 = P_T G_T G_R \left(\frac{\lambda}{4\pi} \right)^2$ as the normalized received power at a distance of 1 m, the identity above reduces to $P_R = \frac{P_0}{d^2}$.

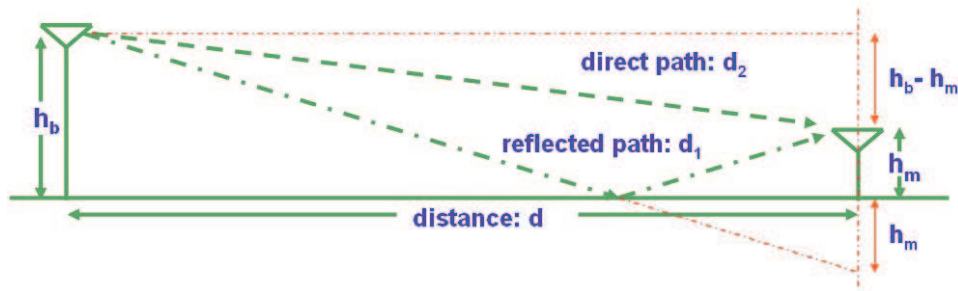


Figure 1.2: Two-ray model

Now we consider the wireless communication between a base station and a mobile terminal on the ground, as shown in Fig.1.2. Their antenna heights are h_b and h_m , respectively. we assume

that h_b is greater than or equal to h_m . The distance d is assumed to be much larger than either antenna height, i.e., $d \gg h_b$ and $d \gg h_m$. All reflection coefficients are -1, which means an ideal lossless reflector. With these assumptions, the mean received power is given by

$$P_r = \frac{P_0}{d^2} |1 - e^{j\Delta\phi}|^2, \quad (1.1)$$

where $\Delta\phi$ is the phase difference between the reflected path d_1 and the direct path d_2 . Defining the length difference Δd be $d_1 - d_2$, we can calculate $\Delta\phi$ by

$$\Delta\phi = 2\pi \cdot \frac{\Delta d}{\lambda},$$

where f is the carrier frequency, λ is the wavelength, and c is the speed of light.

Since $d \gg h_b$ and $d \gg h_m$ are assumed, we have $(h_b + h_m)^2 \ll 2d$ and $(h_b - h_m)^2 \ll 2d$. Therefore, we have $\left(\frac{(h_b + h_m)^2}{2d}\right)^2 \simeq 0$ and $\left(\frac{(h_b - h_m)^2}{2d}\right)^2 \simeq 0$. Based on these approximations, as given in [69, 70], we can calculate the lengths of the two paths by

$$\begin{aligned} d_1 &= \sqrt{(h_b + h_m)^2 + d^2} \simeq \sqrt{\left(\frac{(h_b + h_m)^2}{2d}\right)^2 + 2 \cdot d \cdot \frac{(h_b + h_m)^2}{2d} + d^2} = d + \frac{(h_b + h_m)^2}{2d}, \\ d_2 &= \sqrt{(h_b - h_m)^2 + d^2} \simeq \sqrt{\left(\frac{(h_b - h_m)^2}{2d}\right)^2 + 2 \cdot d \cdot \frac{(h_b - h_m)^2}{2d} + d^2} = d + \frac{(h_b - h_m)^2}{2d}. \end{aligned}$$

Thus we have

$$\Delta d = \frac{2h_b h_m}{d},$$

and

$$\Delta\phi = \frac{4\pi h_b h_m}{\lambda d}.$$

The other approximations are that, for very small values of $\Delta\phi$, $\cos(\Delta\phi) \simeq 1$ and $\sin(\Delta\phi) \simeq \Delta\phi$.

We have

$$|1 - e^{j\Delta\phi}| \simeq |\Delta\phi|.$$

Based on formula (1.1), we can obtain the approximate mean received power by

$$P_r = \frac{P_0}{d^4} \left(\frac{4\pi h_b h_m}{\lambda}\right)^2 = P_T G_T G_R \times \frac{h_b^2 h_m^2}{d^4}.$$

The conclusion is that, over two paths, if the distance between the transmitter and receiver is long enough, the path-loss exponent of the distance-power relationship is increased to four.

More generally, we can also find that the multiple paths change the distance-power relationship. In mobile wireless environments, similar situations are observed [57, 58, 70]. The mean signal strength decays as the α th power of the distance,

$$P_R = \frac{P_0}{r^\alpha},$$

where α typically ranges from 1.5 to 5. Actually α is called the path loss coefficient or distance-power gradient [57, 58, 31, 70, 89], which is a factor that limits the coverage of a transmitter. The value of α depends on the environment and the distance between the transmitter and receiver.

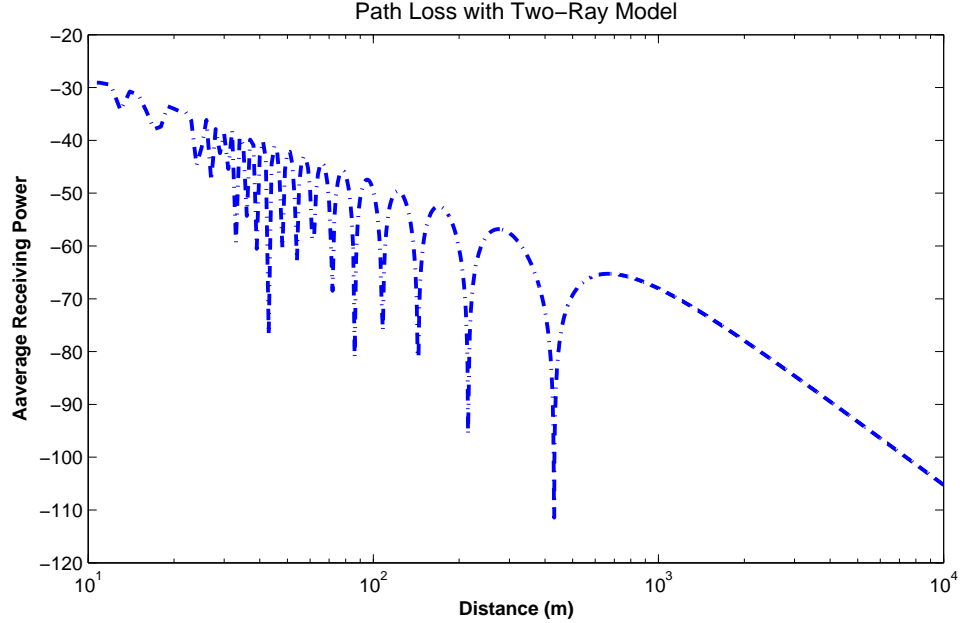


Figure 1.3: Path loss based on two-ray model

If $d \gg h_b$ and $d \gg h_m$ are not satisfied, the path loss based on two-ray model can be directly plotted using formula (1.1), as shown in Fig.1.3. we find that the mean power of received signals is not strictly inverse proportional to a fixed power of the distance d . This phenomenon is caused by the shadowing behavior. Since it is hard to predict the exact length of each arrival path

of signals, shadowing is modeled as a slowly time-varying random process. For a given observation interval, we assume that the mean power of received signals is a constant, which is usually modeled as a lognormal random variable. The shadowing behavior is also called long-term fading [31, 70, 89, 93].

In case of multipath wireless communications, due to the mobility of terminals and environment, the length of each arrival path will change even over a short distance. The waveforms used to carry signals in today's wireless data communications usually have short wavelength. As a result, the power of received signals may vary widely in amplitude and phase rapidly over a short period of time. Fading, or short-term fading, is the term used to describe the rapid fluctuations in power of received signals over a short period of time [31, 70, 89, 93]. System performance can be severely degraded by fading.

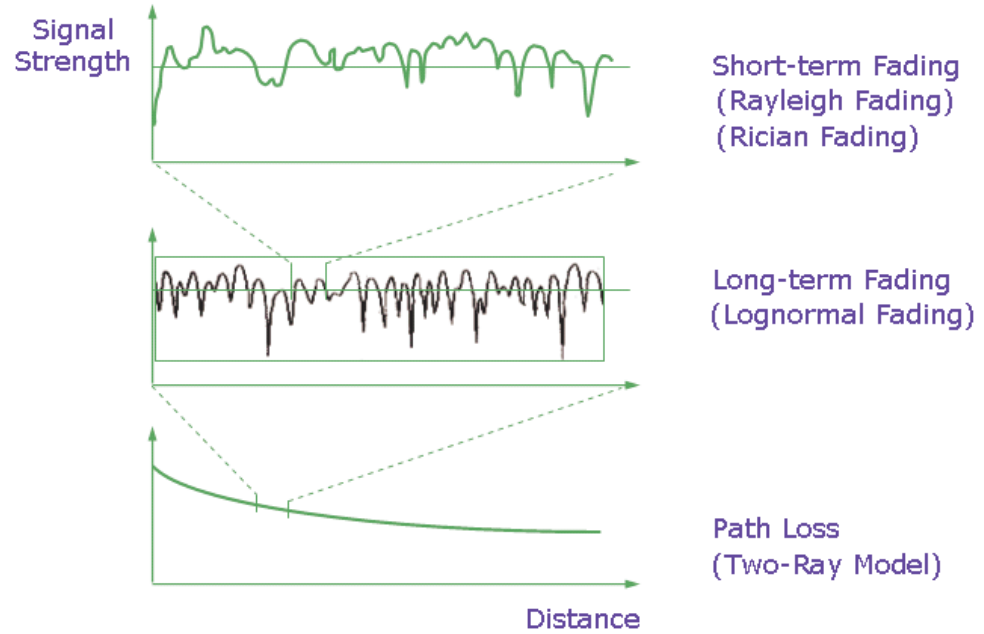


Figure 1.4: Received signal strength: path loss, shadowing, and fading

In summary, the received signal power is mainly affected by path loss, shadowing, and fading [31, 70, 89, 93], as shown in Fig.1.4. To further characterize the fading behavior, we'll discuss the stochastic channel model and several important parameters.

1.2.2 Linear Time-Variant Channel Model

Since the multipath signal propagation in wireless environment always changes as the terminal moves and/or any scattering sources in surrounding environment randomly move, the propagation channel is normally time-variant.

Define the impulse response of a Linear Time-Variant (LTV) channel to be $h(\tau, t)$, which is the channel output at time t in response to an impulse applied to the channel at $t - \tau$ [69, 70]. The variable τ represents the relative propagation delay. If the channel input is $x(t)$, then the channel output $r(t)$ can be represented by

$$r(t) = \int_{-\infty}^{\infty} h(\tau, t)x(t - \tau)d\tau.$$

If the signals arrive along N different paths in a channel, the channel impulse response has the general form

$$h(\tau, t) = \sum_{i=1}^N \alpha_i(t)e^{-j\phi_i(t)}\delta(\tau - \tau_i(t)),$$

where $\tau_i(t)$ is the relative signal delay along path i at time t , $\phi_i(t)$ is the phase of signals along path i , and $\alpha_i(t)$ is the amplitude along path i [63]. In addition, $\phi_i(t)$ and $\tau_i(t)$ are closely related by $\phi_i(t) = 2\pi f_c \tau_i(t)$.

Let $H(f, t)$ be the time-variant channel transfer function, then we have the Fourier transform pair:

$$H(f, t) = \mathcal{F}_{\tau} [h(\tau, t)] = \int_{-\infty}^{\infty} h(\tau, t)e^{-j2\pi f\tau}d\tau,$$

$$h(\tau, t) = \mathcal{F}_f^{-1} [H(f, t)] = \frac{1}{2\pi} \int_{-\infty}^{\infty} H(f, t)e^{+j2\pi f\tau}df.$$

We are interested in time-varying channels with user mobility or environment changes. Considering a scenario that narrow-band signals (sinc waves) are being transmitted at carrier frequency f_c , the receiver is moving at a constant velocity V , and ϕ is the angle of incoming signals with respect to the moving direction of the receiver, the wireless channel will introduce frequency shift $v(t)$ to transmitted signals at the receiver. The frequency shift $v(t)$ is called the Doppler shift

and is given by

$$v(t) = \frac{V \cdot f_c}{c} \cos \phi(t),$$

where c is the speed of light. For wideband signals, the wireless channel will introduce continuous Doppler shifts in a range including the zero Doppler shift for the original signal, called Doppler spread, but not a simple frequency shift. Doppler shift is a consequence of terminal motion, while Doppler spread is caused by mobile multipath propagation.

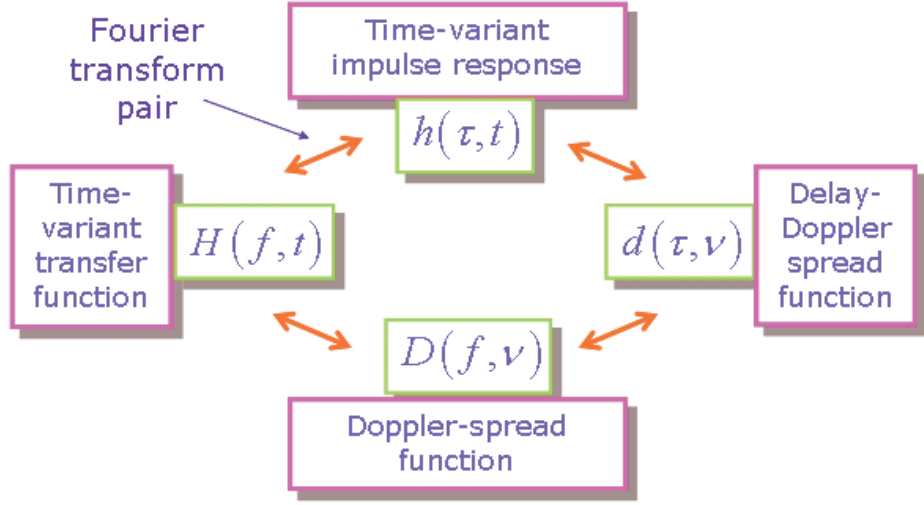


Figure 1.5: Fourier transform pairs of deterministic LTV wireless channels

Doppler spread is related to the changing rate of a Linear Time-Variant (LTV) wireless channels. To study the variation of LTV channels, the Doppler spread function $D(f, v)$ is defined as

$$D(f, v) = \mathcal{F}_t [H(f, t)] = \int_{-\infty}^{\infty} H(f, t) e^{-j2\pi vt} dt,$$

where $H(f, t)$ is the time-variant channel transfer function. It is easy to see that $H(f, t)$ and $D(f, v)$ are a Fourier transform pair [63, 69].

Similarly, the Delay-Doppler spread function is defined as

$$d(\tau, v) = \mathcal{F}_t [h(\tau, t)] = \int_{-\infty}^{\infty} h(\tau, t) e^{-j2\pi vt} dt.$$

It can be shown that $D(f, v)$ and $d(\tau, v)$ are a Fourier transform pair, as illustrated in Fig.1.5.

1.2.3 Wide-Sense Stationary Uncorrelated Scattering Channel Model

When the channel changes randomly with time, the channel impulse response $h(\tau, t)$, time-variant transfer function $H(f, t)$, Doppler spread function $D(f, v)$, and Delay-Doppler spread function $d(\tau, v)$ are random processes that are difficult to characterize. Under the assumption that the random processes have zero mean, we are interested in the correlation functions of the random processes. For the simplicity of analysis [6, 51, 70], we assume that

- The channel impulse response $h(\tau, t)$ is a wide-sense stationary (WSS) process of t
- The channel impulse response $h(\tau_1, t)$ and $h(\tau_2, t)$, are uncorrelated if $\tau_1 \neq \tau_2$ for any t

A channel under these two assumptions is said to be a WSSUS channel [6, 7, 51, 58, 63, 70]. In WSSUS channels, the statistical description of channels will be independent of the absolute time. Then the autocorrelation function of the channel impulse response $h(\tau, t)$ is denoted by $\phi_h(\tau, \Delta t)$ and is given by

$$\phi_h(\tau, \Delta t) = E [h^*(\tau, t) h(\tau, t + \Delta t)],$$

where the superscript $(*)$ denotes complex conjugation. By Fourier transforming the autocorrelation function we can obtain the scattering function of the channel,

$$S_h(\tau, v) = \mathcal{F}_{\Delta t} [\phi_h(\tau, \Delta t)].$$

The scattering function $S_h(\tau, v)$ provides a single measure of the average power output of the channel as a function of a time domain variable (delay τ) and a frequency domain (Doppler frequency v) variable. Similarly we can obtain the space-time, space-frequency correlation function

$$\phi_H(\Delta f, \Delta t) = \mathcal{F}_{\tau} [\phi_h(\tau, \Delta t)],$$

and

$$S_H(\Delta f, v) = \mathcal{F}_{\Delta t} [\phi_H(\Delta f, \Delta t)].$$

Suppose that we send a very narrow pulse over a fading channel. We can measure the received power as a function of time delay. The average received power or the average Power

Spectral Density (PSD) $\phi_h(\tau)$ as a function of delay τ is called the channel intensity profile or the delay power spectrum, which is denoted as

$$\phi_h(\tau) = \phi_h(\tau, 0) = E[|h(\tau, 0)|^2].$$

The range of delay τ over which $\phi_h(\tau)$ is essentially non-zero is called the Delay spread of the multipath channel, and is often denoted by T_m . It tells us the maximum delay between paths of significant power in the channel.

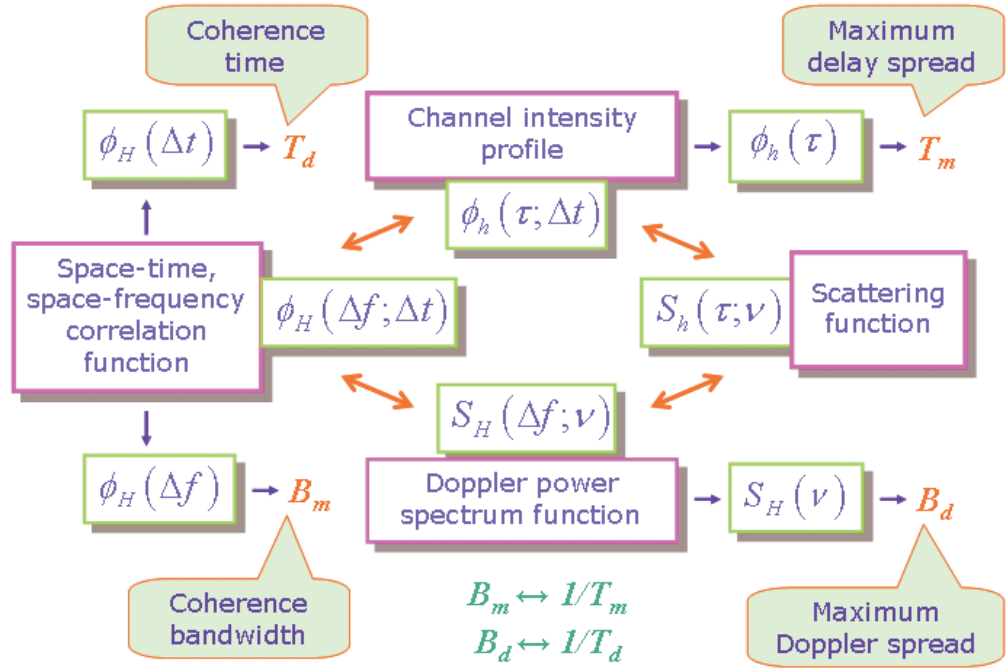


Figure 1.6: Fourier transform pairs of WSSUS wireless channels

The delay power spectrum $\phi_h(\tau)$ portrays the time domain behavior of the fading channel. If we need to study the frequency domain behavior, we need the frequency correlation function of the channel, which is defined as

$$\phi_H(\Delta f) = \phi_H(\Delta f, 0).$$

The range of Δf making $\phi_H(\Delta f)$ non-zero is called the coherence bandwidth of the multipath channel, denoted as B_m . In a fading channel, signals with different frequency contents can undergo

different degrees of fading. If two sinusoids are separated in frequency by more than B_m , they will undergo independent fading. It can be shown that B_m is related to T_m by $B_m \simeq \frac{1}{T_m}$ [6, 31, 89].

Due to the time-varying nature of the channel, a signal propagating in the channel may undergo Doppler shift. When a sinusoid is transmitted through the channel, the received power spectrum can be plotted against the Doppler shift. The average received power $S_H(v)$ as a function of Doppler shift v is called the Doppler power spectrum, which is denoted by

$$S_H(v) = S_H(0, v) = E [|D(0, v)|^2] .$$

The Doppler spread, denoted by B_d , is the range of v making the Doppler power spectrum non-zero. It gives the maximum range of Doppler shifts.

The time correlation function of the channel is defined as

$$\phi_H(\Delta t) = \phi_H(0, \Delta t).$$

The range of delay Δt over which $\phi_h(\Delta t)$ is non-zero is called the coherence time of the multipath channel, denoted as T_d . In a time-varying channel, the coherence time gives a measure of the time duration over which the channel impulse response is almost invariant or highly correlated. It can be shown that T_d is related to B_d by $B_d \simeq \frac{1}{T_d}$ [6, 31, 89].

The Fourier transform relation verifies that being time-variant in the time domain can be equivalently described by having Doppler shifts in the frequency domain, as shown in Fig.1.6.

1.3 Main Ideas of OFDM

Orthogonal Frequency Division Multiplexing (OFDM) has been adopted as the modulation and demodulation technique in the physical layer of several communication standards, including digital audio broadcast (DAB), digital video broadcast (DVB), IEEE 802.11a/g wireless LANs, and IEEE 802.16e wireless MANs [5, 37, 64].

The ideas behind OFDM have been around for a long time. Guard bands are usually required in conventional Frequency Division Multiplexing (FDM) systems to avoid sidebands of

adjacent subcarriers overlapping with each other creating Inter-Carrier Interference. However, if subcarriers are mathematically orthogonal, the spectral overlapping among subcarriers is allowed (orthogonality will ensure the subcarrier separation at the receiver), so that better spectral efficiency can be achieved. The idea of OFDM as a special case of FDM was originally published in mid 1960s and patented in 1970 [14]. In OFDM, the data are split and transmitted over a large number of subcarriers and modulated at a low rate. The subcarriers are made orthogonal to each other by appropriately choosing the frequency spacing between them. For example, the frequency of the modulating sinusoid in each subcarrier is an integer multiple of a base frequency $1/T_s$ where T_s is the symbol period.

In addition to ICI, OFDM also aims to address the Inter-Symbol Interference (ISI) problem. For a single carrier wireless system, its symbol rate R_s is inversely proportional to the symbol period T_s . Higher R_s means smaller T_s ($T_s = 1/R_s$), which might cause serious ISI problems in multipath channels with Delay spread larger than T_s . In OFDM systems with data rate R_s , since data are split to N subcarriers, the symbol rate in each subcarrier is R_s/N , and the symbol period in each subcarrier is $N \cdot T_s$. Therefore, in multipath channels, OFDM systems are more robust to ISI than single carrier systems with equivalent data rate.

The problem with the original OFDM idea is that it needs many local oscillators whose frequencies are the accurate multiples to maintain orthogonality, which is difficult and expensive to implement. The new major finding is that OFDM is closely related to Discrete Fourier Transform (DFT) [97]. As the introduction of DFT and Inverse Discrete Fourier Transform (IDFT) digital components in OFDM transceivers, the multiple carrier scheme using a bank of parallel subcarrier oscillators for modulation and coherent demodulation in analog hardware is abandoned. The idea of baseband modulation and demodulation based on DSP hardware and software enables more efficient and flexible OFDM implementation with lower complexity and cost. Moreover, the fast implementation technique of DFT and IDFT, Fast Fourier Transform (FFT) and Inverse FFT (IFFT) further improve the performance and reduce the cost.

Another major breakthrough is the introduction of a Cyclic Prefix (CP) as a Guard Interval (GI) to extend the OFDM symbol period. The CP is taken from samples in the end of an OFDM time-domain signal (an OFDM data symbol) and is appended to the front of OFDM data symbol. Due to the CP, the transmitted time domain signal becomes periodic, as long as the length of CP is larger than the Delay spread of multipath channel. Therefore, in the time domain, the effect of the multipath channel becomes a circular convolution with the channel impulse response function. In the frequency domain, the effect of the multipath channel is just a point-wise multiplication of the constellation symbols by the channel transfer function. Periodicity of an OFDM symbol also ensures that subcarriers after FFT are orthogonal to each other. Thus the ICI may be reduced. In addition, the periodicity reduces OFDM systems' sensitivity to timing synchronization [5, 37].

With all of these techniques, OFDM significantly changes the landscape of wireless systems. It will continue to play an important role in future high speed wireless systems.

1.4 IEEE 802.11a, IEEE 802.11p, and IEEE 802.16e

The convergence between different wireless technologies inevitably leads to the mobility support requirement for high-speed packet wireless systems. By far, IEEE 802.11p and IEEE 802.16e are the two major upcoming standards for mobile packet wireless networks. IEEE 802.11p is revised from IEEE 802.11a, while IEEE 802.16e extends IEEE 802.16a [49, 102, 106]. To combat multipath Delay spread, OFDM is adopted in these standards due to its inherent multipath resistance. In this section, we'll compare the PHY layer design of IEEE 802.11a, IEEE 802.11p, and IEEE 802.16e 256FFT mode.

1.4.1 IEEE 802.11a

The PHY layer of the 802.11a is divided into two entities: the Physical Layer Convergence Protocol (PLCP) and the Physical Medium Dependent (PMD) sublayers. The PLCP sublayer maps

MPDUs (MAC Protocol Data Units) from MAC layer into a frame format suitable for PMD layer and delivers incoming frames from PMD layer to the MAC layer. PMD layer deals with actual transmission and reception over the air medium [44].

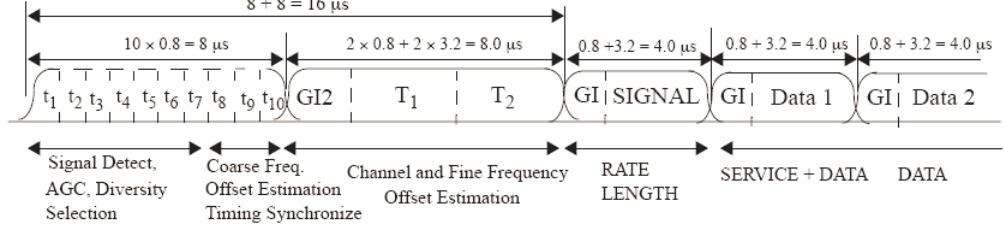


Figure 1.7: IEEE 802.11a PHY frame format

The PHY frame format for 802.11a is shown in Fig.1.7 (The figure is from [44]). The PLCP preamble field has 10 repetitions of a short training symbol, a Guide Interval, and two repetitions of a long training symbol. In the receiver, the short training symbols are used for Automatic Gain Control (AGC, to prevent signals from saturating the output of the A/D converter) convergence, timing acquisition, and coarse frequency acquisition, while the long training symbols are used for channel estimation and fine frequency acquisition. Here the frequency-domain channel estimation based on two long training symbols T_1 and T_2 assumes that the wireless channel remain the same till the end of packet transmission.

Binary input data comes from MAC layer is collected and coded with Forward Error Correction (FEC) schemes such as convolutional encoder. Then the coded bits are punctuated to achieve high rates and interleaved to combat bursty bit errors. Afterwards, based on the selected baseband modulation scheme and its number of bits per symbol, interleaved binary bits are grouped and mapped to corresponding constellation points or data symbols. These symbols are serial complex numbers, and divided into groups of 48 symbols. Each such group is associated with one OFDM symbol. Assuming that the active subcarriers are sequentially numbered from -26 to 26, in each group, the baseband symbols are mapped into 48 OFDM subcarriers numbered -26 to -22, -20 to -8, -6 to -1, 1 to 6, 8 to 20, and 22 to 26. At this point, pilot symbols with known modulation scheme

may be inserted into 4 subcarriers -21, -7, 7, and 21 with a known pattern. In addition, zero symbols for unloaded subcarriers are also inserted, so that the number of subcarriers is a power of 2 which is required by the IFFT/FFT operation (64 in 802.11a) [44].

The OFDM modulation in frequency domain, IFFT operation, is performed on the parallel symbols to generate parallel symbols. These parallel symbols are then serialized and inserted with the Cyclic Prefix to form an OFDM baseband symbol for one OFDM symbol period in time domain. The data portion of an OFDM symbol (as seen from Fig.1.7) comes from the time-domain signals. A Cyclic Prefix is introduced as Guard Interval (GI) by prepending to the IFFT waveform a circular extension of itself and truncating the resulting periodic waveform to a single OFDM symbol length. Append the OFDM symbols one after another, starting after the PLCP header, until the length is reached. The PLCP header field has information about the transmission RATE, the LENGTH of the payload, a parity bit, and six zero tail bits. The RATE field conveys information about the type of modulation and coding rate used in the rest of the packet. The LENGTH field specifies the number of bytes in the Physical Layer Service Data Unit (PSDU).

The operating frequencies of 802.11a in the U.S. fall into the National Information Structure (U-NII) bands: 5.15-5.25 GHz, 5.25-5.35 GHz, and 5.725-5.825 GHz. Within this spectrum, there are twelve 20 MHz channels, and each band has different output power limits. The complex baseband time domain waveform will be upconverted to an RF frequency according to the center frequency of the desired channel and transmitted.

1.4.2 IEEE 802.11p

Recently, upcoming IEEE 802.11p has been endorsed by ASTM as the platform for the PHY and MAC layers of Dedicated Short Range Communications (DSRC). Providing wireless communications between vehicles and the roadside, and between vehicles, DSRC enables a whole new class of applications that enhance the safety and productivity of the transportation system.

As a variant of IEEE 802.11a, 802.11p follows many design features of 802.11a, such as frame

structure, training sequences, scrambler, convolutional coding, interleaving, modulation schemes, pilot subcarriers, IFFT/FFT size, Cyclic Prefix, and pulse shaping [102, 106]. IEEE 802.11a is basically designed for low mobility indoor applications, where the wireless channel is assumed to be stationary for the frame duration. Therefore, all system parameters are chosen to achieve best performance in indoor propagation environments. However, IEEE 802.11p is proposed for high mobility outdoor environments. Thus, it has to deal with the impairments brought by high mobility. Some of these challenging requirements for PHY layer design include:

- The transceiver has to combat increased multipath Delay spread
- The transceiver has to combat increased Doppler spread, which means that the multipath channel varies more rapidly. It is no longer a realistic assumption that the channel estimation acquired at beginning of transmission will be valid till the end of transmission
- Longer communication range.

IEEE 802.11p works in the 5.850 to 5.925 GHz Intelligent Transportation Systems Radio Service (ITS-RS) Band that accommodates seven channels in a total spectrum of 75 MHz. In 802.11p, each mandatory channel operates with 10 MHz bandwidth, rather than 20 MHz bandwidth as used in 802.11a. Accordingly, rates of 802.11p are half of those 802.11a rates, and symbol period of 802.11p is twice the symbol period of 802.11a. The advantage with longer symbol period is that the longer Cyclic Prefix of each OFDM symbol enables 802.11p to combat possibly larger Delay spread introduced by outdoor channel environments. In addition, the transmission power limits designated by 802.11p are different from the power limits of 802.11a. The maximum antenna input power for some DSRC mandatory channels is 28.8 dBm (750 mW) that enables longer range.

Major OFDM parameters for 802.11a and 802.11p are given in table 1.2. IEEE 802.11p has the same pilot placement as 802.11a, e.g., block-type placement.

Parameter	802.11a	802.11p
Number of data subcarriers	48	48
Number of pilot subcarriers	4	4
Number of total subcarriers	52	$52 (N_{SD} + N_{SP})$
IFFT/FFT size	64	64
IFFT/FFT period	3.2	6.4
PLCP preamble duration	16	$32 (T_{SHORT} + T_{LONG})$
GI duration (μs)	0.8	$1.6 (T_{FFT}/4)$
Training symbol duration (μs)	1.6	$3.2(T_{FFT}/2)$
Symbol Period (μs)	4	$8 (T_{GI} + T_{FFT})$
Short training period (μs)	8	$16 (10T_{FFT}/4)$
Long training period (μs)	8	$32(T_{GI2} + 2T_{FFT})$
Data Rate (Mbps)	6, 9, 12, 18, 24, 36, 48, 54	3, 4.5, 6, 9, 12, 18, 24, 27
Modulation	BPSK/QPSK/16-QAM/64-QAM	BPSK/QPSK/16-QAM/64-QAM
Convolutional Code	$K = 7$ (64 states)	$K = 7$ (64 states)
Coding Rates	$1/2, 2/3, 3/4$	$1/2, 2/3, 3/4$
Channel Spacing (MHz)	20	10
Signal Bandwidth (MHz)	16.66	8.3
Subcarrier Spacing (kHz)	312.5	156.2

Table 1.2: Major OFDM parameters of 802.11a and 802.11p

1.4.3 IEEE 802.16e

Following the finalization of IEEE 802.16 REVd, the main task of 802.16 standards working group switches to 802.16e, which aims to add mobility to 802.16a amendment [49]. The multi-carrier technique of IEEE 802.16e is mainly based on two of 802.16a basic forms: OFDM with 256 point transform (OFDM-256) and Orthogonal Frequency Division Multiple Access (OFDMA) with 2048 point transform. OFDMA distributes subcarriers among users so all users can transmit and receive at the same time within a single channel. These groups of subcarriers for one user are then known as subchannels. The assignment of subcarriers is based on the location and channel condition of each user. Therefore, the carriers of one subchannel are spread along the channel spectrum to mitigate the frequency selective fading. The new feature of 802.16e is Scalable OFDMA (SOFDMA) that allows the deployment of variable FFT sizes such as 2048-FFT, 1024-FFT, 512-FFT and 128-FFT, etc., to further improve performance. In this section, we'll only introduce how the OFDM-256 for downlink transmissions are revised for mobile users.

The original OFDM-256 in 802.16a has 256 subcarriers with 192 data subcarriers, 8 pilot subcarriers, and 56 null subcarriers. Since 802.16a assumes a stationary scenario, the block-type pilot placement scheme is similar to 802.11a and 802.11p, where pilot subcarriers are evenly distributed among data subcarriers. As we discussed in previous sections, mobility brings significant changes to the channel characteristics. The Doppler spread and Delay spread of the time-varying channel are strongly related to users' mobile speed. However, the QAM modulation schemes used with OFDM require coherent demodulation. Therefore, to support mobility demands, the receiver needs to estimate its channel more frequently to keep track of the phase of received signals for correct demodulation and the amplitude distortion to combat fading. Actually two mechanisms, including the hopping pilot mechanism and the mid-amble insertion mechanism, have been proposed in 802.16e.

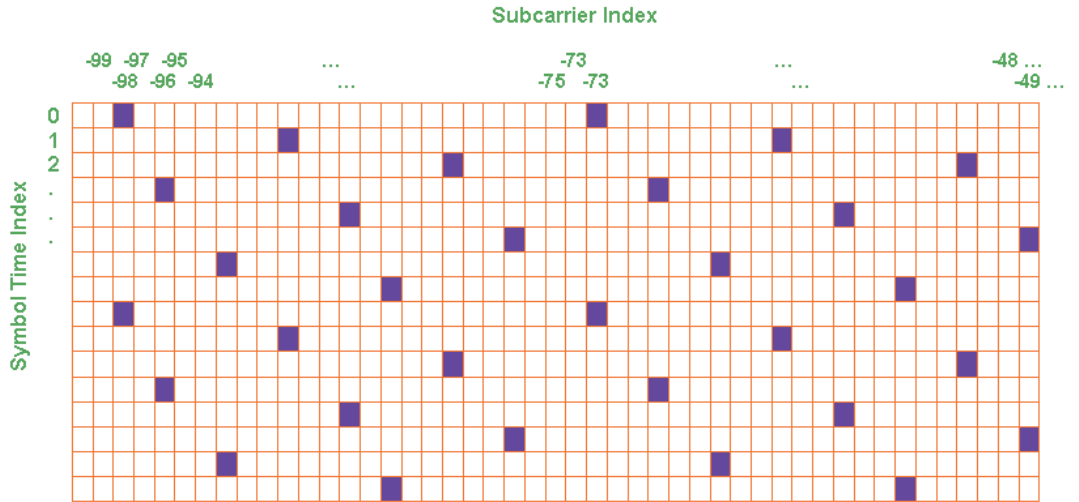


Figure 1.8: Pilot placement of IEEE 802.16e PHY layer

The hopping pilot mechanism doesn't use fixed pilot subcarriers (block-type placement) as used in 802.11a and 802.11p. The pilot subcarriers are regularly and cyclically changed symbol by symbol with a periodicity of 8 OFDM symbols. Thus, it is possible to interpolate the channel estimate for each subcarrier without pilot at a particular time. Assuming that the active subcarriers are sequentially numbered from -100 to +100 and that t ($t \geq 0$) is the number of symbols from

current symbol to the beginning symbol with pilots of a frame, the arrangement of pilot subcarriers follows

$$P_t = \{-98, -73, -48, -23, +2, +27, +52, +77\} + \text{mod}(9 \cdot t, 24),$$

where P_t is the index of pilot subcarriers at time t , as shown in Fig.1.8. With this new pilot placement scheme, 64 of 200 subcarriers are covered, while only 8 of 200 subcarriers are covered in 802.16a.

In 802.11a and 802.11p, training symbols are in the preamble of a MAC frame. However, with the mid-amble insertion mechanism, training symbols may be inserted in the middle of data symbols of a MAC frame, so that the channel estimate can be restored in the middle of a MAC frame in case that the channel changes too quickly.

Overall, hopping pilot mechanism and the mid-amble insertion mechanism help 802.16e combat the adverse effect of time-varying channel brought by user mobility.

1.5 Dissertation Organization

The main focus of this dissertation is an investigation of the effect of multipath fading channels and adaptive modulation on TCP/UDP performance of high-speed packet wireless networks such as wireless LANs and wireless MANs.

Following this introductory chapter, the simulation study of effect of Doppler spread and Delay spread on 802.11a/b/p is presented in chapter 2. Then chapter 3 uses 802.11a/p as an example to investigate the effect of mobile speed on the communication range of packet-based wireless systems that serve users with high mobility in Rayleigh and Rician fading channels. Both chapter 2 and chapter 3 assume a UDP traffic source.

Taking into account adaptive modulation and Rayleigh fading channel condition, chapter 4 shows how a semi-Markov chain model for TCP congestion/flow control behavior and a multi-state Markov chain model for Rayleigh fading channels are used together to derive the bulk throughput

of TCP Tahoe and Reno in steady state. In addition, chapter 4 introduces the idea of cross-layer TCP protocol design for single-user scenarios.

Considering the multi-user scenarios, chapter 5 explores how multi-user diversity can be used to improve the aggregate TCP throughput of base stations in fading channels.

Chapter 6 concludes this dissertation, summarizing the main contributions and suggesting areas for future research.

Chapter 2

Effect of Doppler Spread and Delay Spread on Performance of High-Speed Packet Wireless Networks

2.1 Introduction

Most high-speed packet wireless networks are not designed for mobile users. For example, IEEE 802.11a standard explicitly states that it is proposed for indoor applications with low mobility [44]. However, the ongoing evolution of wireless devices and service requirements makes inevitable the development of new high-speed packet wireless standards for mobile deployment. Upcoming wireless standards such as 802.11p and 802.16e illustrate this trend [49, 102].

The convenience brought by mobility support in an urban environment comes with the

harsh and challenging time-varying channel condition. In addition to the Line-Of-Sight (LOS) or direct-path component (if there is one), the radio signal emitted by the transmitter is reflected by many environment objects such as buildings, trees, and cars, to reach the receiver over many different paths. Since the wavelength of modern high-speed packet wireless networks is very short (for example, only about 5 mm at a carrier frequency of 5.9 GHz), the small distance difference between paths followed by the same signal means huge phase difference between received radio waves. These radio waves at the receiver may add (constructively) or cancel each other (destructively). Due to the relative movement between receiver and transmitter in the communications system or the random movement of environment reflecting objects, the travel paths of signals for mobile applications always change, and accordingly the constructive/destructive effect on signals keep changing too. As a result, the wireless channel for mobile users, called the multipath fading channel, is much more unpredictable than a static channel [31, 89].

Although multipath fading channel is unpredictable, it shows some statistical characteristics that can be modeled in terms of several important parameters. Among them, as we discussed in the previous chapter, Delay spread and Doppler spread are two fundamental parameters of a mobile multipath channel. Doppler spread indicates how fast the channel changes, while Delay spread shows the time dispersion of signal arrival along different paths.

We study how these two important parameters of multipath fading channels affect the performance of high-speed packet wireless networks. Many techniques may be used to perform the evaluation. Running experiments using real radio, hardware, and software is practical. However, since the behavior of mobile multipath channels strongly depends on environment such as terrain, foliage, buildings, and other moving objects like cars, it will be hard to reproduce experiment results. Therefore, a simulation-based approach is preferred to study the effect of time-varying channel conditions, because of its great advantages such as repeatability, controllability, and low cost.

Extensive simulations are reported in this chapter investigating how Doppler spread and

Delay spread affect the performance of several popular high-speed packet wireless networks including IEEE 802.11a and IEEE 802.11b. Since many other packet wireless networks use the same techniques such as OFDM and Spread Spectrum, our simulation results and conclusion may be used as a reference for the design of new high-speed packet wireless networks.

2.2 Multipath Fading Channel Model

A baseband multipath fading channel is usually modeled with a multiplicative fading component and an additive noise component [31, 89]. Two typical models for the multiplicative component are the Rayleigh and Rician distributions.

Assuming that the received signal is the sum of signals with different phases caused by different paths, the amplitude of the received signal, r , can be modeled as a random variable with a Rayleigh distribution, whose Probability Density Distribution $f(r)$ is given by

$$f(r) = \frac{r}{\sigma^2} \exp\left(-\frac{r^2}{2\sigma^2}\right), r \geq 0,$$

where $2\sigma^2$ is the pre-detection mean power of the received signal [51].

If $x(t)$ and $y(t)$ are two uncorrelated and zero-mean Gaussian processes with the same statistical properties, as if they are the in-phase and quadrature signal components of the received signal, we can calculate the joint PDF of $x(t)$ and $y(t)$ by

$$f(x, y) = \frac{1}{2\pi\sigma^2} \exp\left(-\frac{x^2 + y^2}{2\sigma^2}\right).$$

Let $r(t) = \sqrt{x^2(t) + y^2(t)}$ and $\theta(t) = \arctan\left(\frac{y(t)}{x(t)}\right)$. Then after transformation, the joint PDF of $r(t)$ and $\theta(t)$ becomes

$$f(r, \theta) = \frac{r}{2\pi\sigma^2} \exp\left(-\frac{r^2}{2\sigma^2}\right).$$

This is the PDF of the Rayleigh distribution.

It can be seen that the average power $p(t) = r^2(t)$ of the received signal is an exponential

random variable with mean $p_0 = 2\sigma^2$,

$$f(p) = \frac{1}{p_0} \exp\left(-\frac{p}{p_0}\right), p \geq 0.$$

Now let $s(t) = x(t) + iy(t)$ be a complex Gaussian random process. Although the amplitude $r(t)$ of $s(t)$ has the same form of PDF as a Rayleigh distribution, it is not the proper model for the fading of a wireless channel. Extensive measurements indicate that the complex fading envelope coefficient $r(t)$ is a random variable that changes slowly over time, which means it has correlation over time or it is a narrowband random process. In most cases, $r(t)$ can be modeled as the output of a low pass filter with white complex Gaussian random process $s(t)$ as input [51]. The low pass filter determines the power spectrum shape and the temporal correlation function of the random process $r(t)$, i.e., the narrowband fading channel envelope.

In Jakes' model [51], the most widely used Rayleigh simulation model, $r(t)$ is assumed to have following temporal correlation function,

$$R(\tau) = E[r(t) \cdot r^*(t - \tau)] \sim I_0(2\pi\tau f_{doppler}),$$

where $r^*(t)$ is the conjugate of $r(t)$, I_0 is the zero-order Bessel function, and $f_{doppler}$ is the Doppler frequency. The power spectrum of $r(t)$ is

$$R(f) = \frac{1}{\sqrt{1 - \frac{f^2}{f_{doppler}^2}}}, |f| \leq f_{doppler}.$$

The Doppler frequency $f_{doppler}$ can be approximately computed as

$$f_{doppler} = \frac{v}{c} \cdot f_{carrier},$$

where v is the velocity of the mobile terminal, c is the speed of light or electromagnetic wave in the air, and $f_{carrier}$ is the carrier frequency of the passband signal transmission. Typical $f_{doppler}$ in current wireless networks ranges from 5 Hz to 300 Hz, depending on the specific situation. For example, for a carrier frequency $f_{carrier}$ of 3 GHz and a mobile speed v of 5 m/second (11.3 mile/hour), the Doppler frequency is

$$f_{doppler} = \frac{5}{3 \cdot 10^8} \cdot 3 \cdot 10^9 = 50Hz.$$

If $f_{doppler}$ is below 100 Hz, it is usually called slow fading.

We give an example of the relationship of the fading channel changing rate and the transmission symbol rate. As stated in IEEE 802.16 wireless MANs, for a 25 MHz channel width, the symbol rate is 20 MBaud/Second. For the wireless mobility scenario of the previous paragraph, the fading rate normalized to symbol rate is $\frac{50}{20 \cdot 10^6} = \frac{1}{0.4 \cdot 10^6}$, which means that the channel does not change much over $0.4 \cdot 10^6$ symbols. Thus we can assume that the channel condition remains constant over a symbol in slow fading channels. That is the reason that Adaptive Modulation and Coding has a chance to play an important role in IEEE 802.16 wireless MANs.

Although Rayleigh fading is often a good approximation of realistic channel conditions, it is considered by many to be a worst-case scenario of signal fading. If a wireless receiver works in a Rayleigh fading channel, then it is likely to work in other types of channels.

In simulations, as shown in Fig.2.1, the Rayleigh fading channel is modeled as

$$r(t) = h(t) \cdot s(t) + n(t),$$

where $r(t)$ is the received signal, $n(t)$ is the noise, $s(t)$ is the transmitted signal, and $h(t)$ refers to a multiplicative distortion of the transmitted signal $s(t)$. The Rayleigh channel simulator generates the fading envelope coefficients $h(t)$ using a statistically accurate and computationally efficient approach.

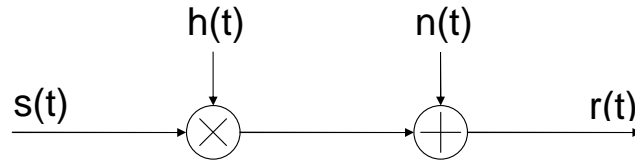


Figure 2.1: Multiplicative and additive model for fading channels

Rayleigh fading with a strong line of sight (LOS) Line of sight is called Rician fading with Probability Density Distribution $f(r)$

$$f(r) = \frac{r}{\sigma^2} \exp\left(-\frac{r^2 + K^2}{2\sigma^2}\right) I_0\left(\frac{rK}{\sigma^2}\right),$$

where K is the coefficient that indicates how strong the LOS component is compared to the rest of the received signal. I_0 is the zero-order modified Bessel function [31, 51, 58, 70].

Both Rayleigh fading and Rician fading are in the standard SIMULINK block library. Jakes' model is used to generate the Rayleigh fading coefficients in the block.

A general model, Tapped Delay Line Channel Model [31, 58, 70, 89], is used to model a time-varying multipath fading channel with L multipath signal components, as illustrated in Fig.2.2. The channel model consists of a tapped line with different delays. The tap coefficients, denoted by $\alpha_i(t)$, are usually modeled as complex-valued Rayleigh fading coefficients or Rician fading coefficients that are uncorrelated with each other. The tap delay τ_i corresponds to the amount of time dispersion in the multipath fading channel. In this model, Doppler spread is used to generate the tap coefficients, and Delay spread is used to configure tap delays.

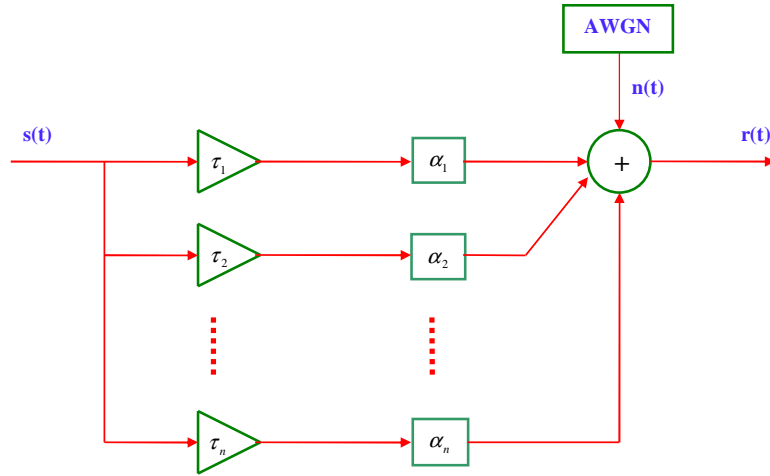


Figure 2.2: Tapped delay line channel model

2.3 IEEE 802.11a and IEEE 802.11b Simulators

The two simulators used for this chapter are modified from the existing IEEE 802.11a PHY SIMULINK model and IEEE 802.11b PHY SIMULINK model in MATLAB Central web site

(<http://www.mathworks.com/matlabcentral/>). To make a fair performance comparison, the two simulators share the same code of multipath fading channel model.

As shown in Fig.2.3, the IEEE 802.11a digital baseband transceiver can be used to perform a complete end-to-end simulation. The baseband simulator assumes that the underlying passband subsystem works perfectly without any ICI (InterCarrier Interference). Inside the Frequency Domain Equalizer, the first four training symbols in each of these 52 subcarriers are extracted and divided by the corresponding four known training symbols to calculate the channel estimate, which is used to correct the received data symbols within the same OFDM block. Each subcarrier has its own channel estimator. The Adaptive Modulation block is disabled in our simulations. The simulator works in (BPSK, 1/2) mode, e.g., at data rate of 6 Mbps.

Fig.2.4 shows the DSSS simulator for an IEEE 802.11b digital baseband transceiver. The simulator follows the IEEE 802.11b standard [45]. It supports 1 Mbps, 2 Mbps, 5.5 Mbps, and 11 Mbps rates. The basic components includes DBPSK and DQPSK modulation, Barker code spreading, and Complementary Code Keying (CCK), etc. However, only 1 Mbps rate, e.g., DPSK and Barker code spreading, is used in our simulations. Perfect synchronization is assumed by the simulator.

2.4 Numerical Results

We now present key numerical results to highlight the effects of Delay spread and Doppler spread. Simulations have been performed under various mobile multipath channel conditions (via different Doppler spread and Delay spread). Each run of simulation lasts 50 seconds in the simulated system time. In our simulations, PHY traffic is assumed to be in saturation state, i.e., the PHY layer buffer is never empty.

2.4.1 Effect of Doppler Spread

To concentrate on the effect of Doppler spread, only one tap of multipath fading channel is used in our simulations, as shown in Fig.2.5. The SNR for all simulations is set to 10dB. As the controllable parameter, the Doppler spread in Rayleigh multipath fading block or Rician multipath fading block is changed for each simulation run.

From Fig.2.6 and Fig.2.7, it is observed that both 802.11a and 802.11b perform worse when the mobile multipath channel changes rapidly (the speed and Doppler spread becomes larger). However, the performance of 802.11a degrades much more than 802.11b. The fundamental reason is that 802.11b has a much higher symbol rate or shorter symbol period than 802.11a.

In the 802.11a standard, channel estimation and fine frequency offset estimation are done via training symbols at the beginning of a PHY frame. The channel estimates obtained during this period are used to compensate for multipath effects for the entire frame. The implicit assumption is that the channel will remain stationary for the duration of the entire OFDM frame. Although the assumption may be valid for low mobility channel conditions, it doesn't hold for mobile multipath channels. The fast-varying channels are not equalized adequately based on the available training symbol placement scheme. In addition, Doppler spread may translate to carrier frequency offset, which causes subcarriers of OFDM to lose their orthogonality relative to each other (the perfect synchronization assumption in this chapter produces more optimistic simulation results). For 802.11a in mobile multipath channels, an equalizer is certainly required to adapt and compensate for the fast changing channel. To combat distortion introduced by large Doppler spread, OFDM-based wireless systems need more advanced channel estimation techniques than what is currently proposed for 802.11a. Indeed, 802.16e does propose a better pilot symbol placement schemes so that efficient channel estimation can be implemented [49].

Since 802.11b has a high symbol rate, well above usual Doppler spread (less than 5000 Hz), the channel is considered to be fairly constant during each symbol transmission. In addition, 802.11b's differential BPSK modulation scheme enhance its intrinsic ability to handle a larger

Doppler spread than 802.11a.

2.4.2 Effect of Delay Spread

Delay spread (or equivalently coherence bandwidth) describes the time dispersive nature of multipath fading channels. But it does not provide information about the time varying nature (caused by relative motion of transmitter and receiver). There are two ways to define delay spread: Maximum Excess Delay and Root Mean Square (RMS) Delay.

Maximum Excess Delay is defined as the overall time span from the earliest arrival to the latest arrival. As the simplest way to define Delay spread, maximum excess delay does not exhibit the relative amplitudes of multipath components (or intensity-delay profiles), which will strongly affect the system performance. Thus a better measure of Delay spread is the root mean square (RMS) Delay spread, which is given mathematically by

$$\tau_{rms} = \sqrt{\overline{\tau^2} - (\overline{\tau})^2},$$

where, given L propagation paths,

$$\overline{\tau^n} = \frac{\sum_{i=1}^L \tau_i^n |\alpha_i|^2}{\sum_{i=1}^L |\alpha_i|^2}, n = 1, 2.$$

The RMS Delay spread reveals the distribution of arriving signal power along different paths, while maximum excess Delay spread only tells the time difference between the first arrival signal and the last arrival signal. However, to simplify the parameter configuration in our simulations, we use maximum excess delay as the metric for Delay spread.

As shown in Fig.2.8, a two tap multipath channel model is used in the simulations. There is no fading for either of these two signal components. Adding fading for each component will lead to worse performance. Here we are trying to isolate the effect of Delay spread from Doppler spread. During the simulations, the maximum excess delay has been gradually increased by setting the

parameter of delay block in one of two signal components. The ratio of signal power on each tap is 0.73:0.27.

BER versus maximum excess delay for 802.11a and 802.11b are plotted in Fig.2.9 and Fig.2.10 respectively. It is easy to see that performance of both 802.11a and 802.11b is degraded, as the Delay spread increases. These figures show that 802.11a tolerates a larger Delay spread (around 1.6ms) than 802.11b (around 1ms), although 802.11a (6Mbps) has a higher data rate than 802.11b (1Mbps). In addition, even with the same Delay spread, 802.11a achieves better BER performance than 802.11b. For example, with Delay spread around 2ms, the BER scale of 802.11a is 10^{-5} , but the BER scale of 802.11b is only 10^{-3} . We have to emphasize that maximum excess delay has been used as Delay spread instead of RMS Delay spread. Otherwise it is easy to get confused by simulation results reported elsewhere.

IEEE 802.11a's relative immunity to Delay spread is due to a combination of the slower symbol rate (longer symbol period) and placement of significant guard time (Cyclical Prefix) around each symbol, which provide protection against ISI interference. To combat very larger Delay spread in an outdoor urban environment, the upcoming 802.11p standard enhances the robustness of 802.11a by reducing the symbol rate by half. By contrast, 802.11b is very sensitive to larger Delay spread because of its higher symbol rate or short symbol period. Reception of 802.11b in a multipath channel can be substantially improved by techniques such as the RAKE receiver principle and equalization [69, 70]. However, they come with complexity and cost.

One more thing worthy to point out is that the introduction of additional taps (paths) or fading effect on each path will degrade the system performance further. What we see in Fig.2.9 and Fig.2.10 is for the best multipath scenarios.

2.5 Concluding Remark

The signal impairment caused by Delay spread and Doppler spread in multipath fading wireless channels places a big challenge to the design of high-speed packet wireless networks. This chapter studied the impact of Delay spread and Doppler spread on performance of two typical high-speed packet wireless systems including IEEE 802.11a and IEEE 802.11b. Simulation results show that 802.11a is robust to Delay spread and susceptible to distortions introduced by Doppler spread, while 802.11b performs well with larger Doppler spread and is sensitive to Delay spread. The modulation technique selected in a high-speed packet wireless network, for example, OFDM or DSSS, should be based on the degree of mobility support, the typical channel environment, data rates, and the BER/PER requirement, etc. The simulation results presented in this chapter provide a useful reference for selecting the modulation technique.

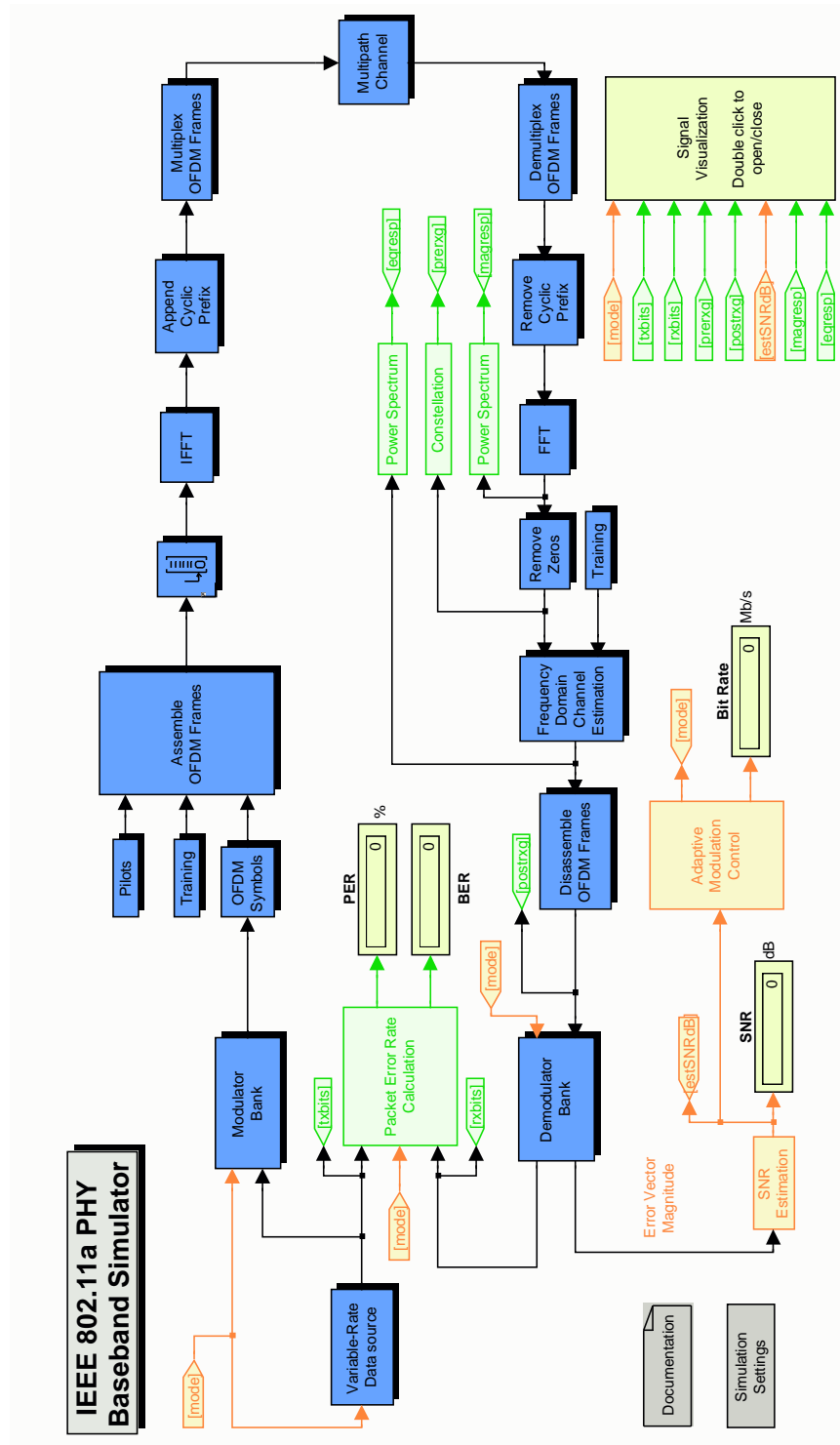


Figure 2.3: IEEE 802.11a PHY simulator

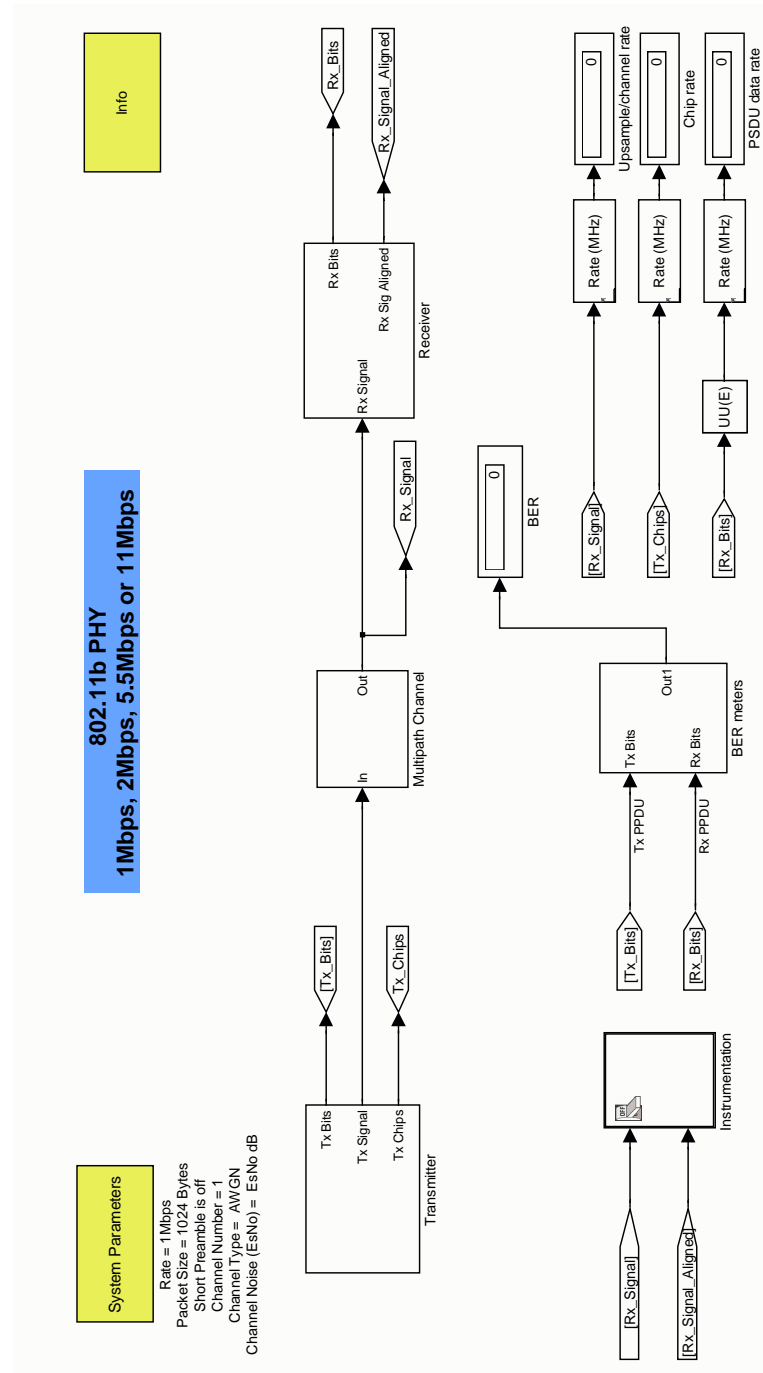


Figure 2.4: IEEE 802.11b PHY simulator

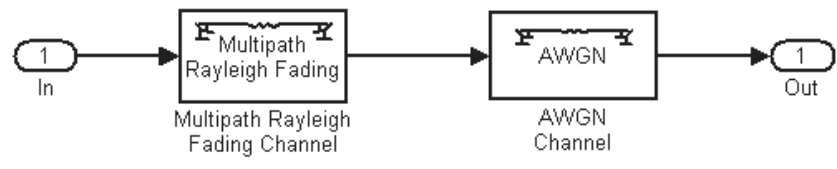


Figure 2.5: Channel blocks for the effect of Doppler spread

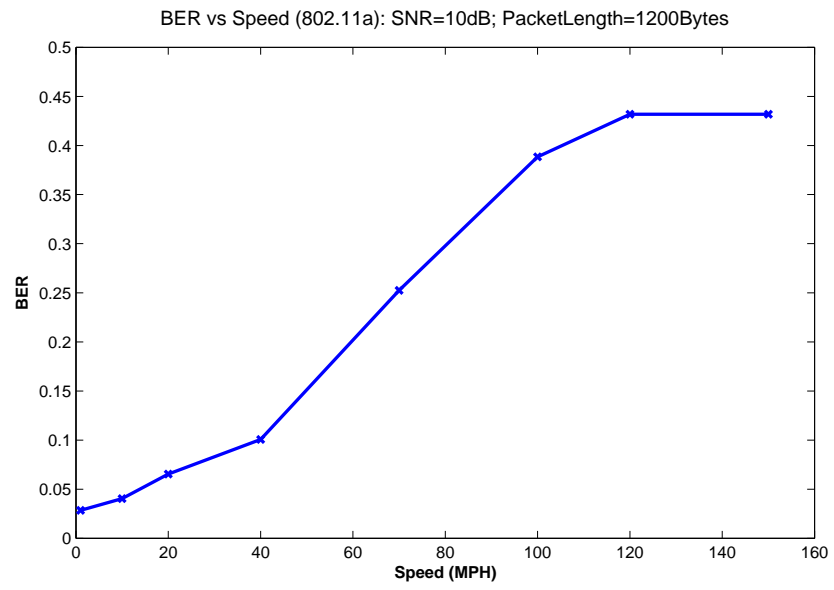


Figure 2.6: Effect of Doppler spread on IEEE 802.11a

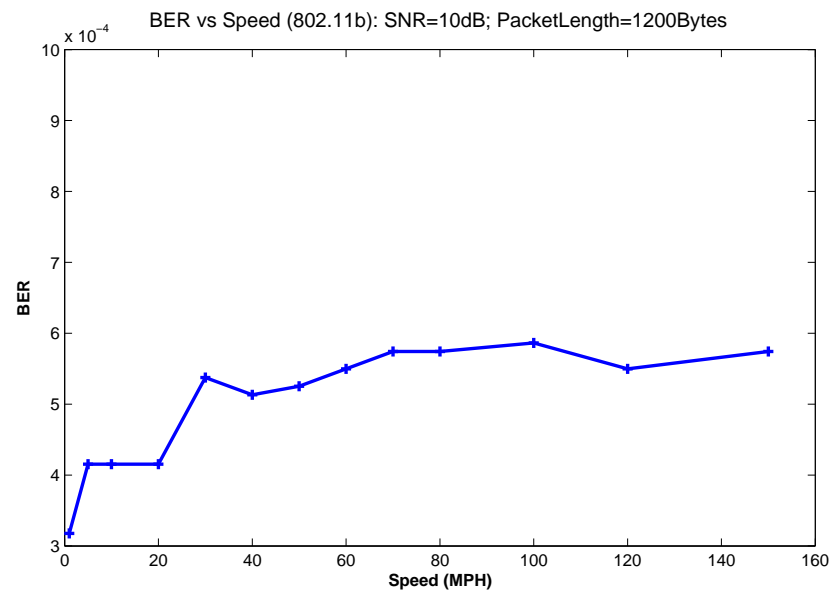


Figure 2.7: Effect of Doppler spread on IEEE 802.11b

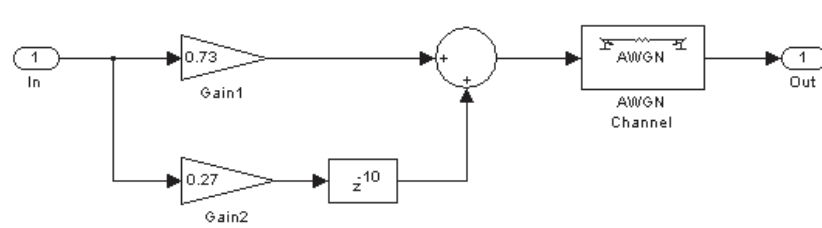


Figure 2.8: Channel blocks for the effect of Delay spread

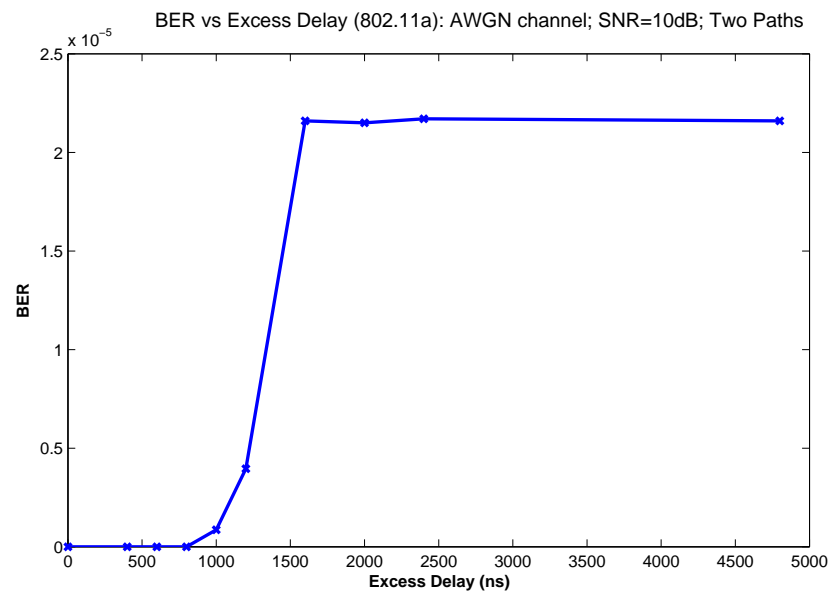


Figure 2.9: Effect of Delay spread on IEEE 802.11a

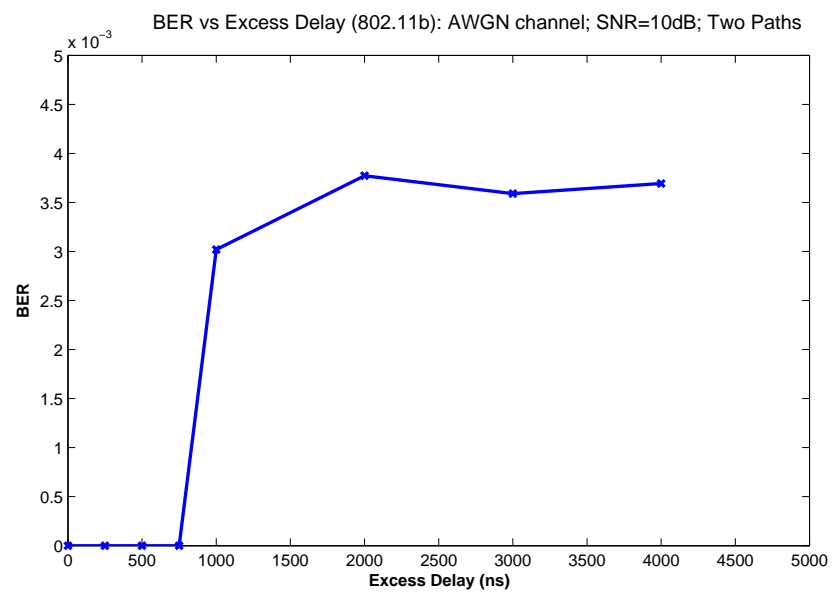


Figure 2.10: Effect of Delay spread on IEEE 802.11b

Chapter 3

Communication Range in Multipath Fading Channels

3.1 Introduction

Communication range is one of the most fundamental parameters characterizing packet-switched wireless systems including IEEE 802.11 wireless LANs (WiFi) and IEEE 802.16 wireless MANs (WiMAX). It does not only determine the coverage area of wireless systems, but also strongly affects capacity planning, network management, seamless handoff design, and admission control [31, 58, 70, 89, 103]. An in-depth understanding of communication range is critical to achieve desired overall system performance when deploying a high-speed packet wireless system.

In this chapter, communication range is defined as the maximum distance between the transmitter and the receiver to achieve given packet error rate (PER) at given packet size. For example, within the communication range of 802.11p, the PER should be less than 10% at a packet size of 1000 Bytes, based on the standard. Since packet is the basic unit of a successful transmission attempt in packet wireless systems, PER is a better metric for service quality than Bit Error Rate

(BER).

Traditionally range is determined from receiver sensitivity, which is ultimately related to the Signal-to-Noise Ratio (SNR), BER, and PER. Assuming an AWGN (Additive White Gaussian Noise) Channel and stationary scenario (for example, the positions of both the transmitter and the receiver do not change, the antenna orientation stays the same, and there is no other environment change), the minimum SNR to achieve given PER or BER determines the receiver sensitivity. Therefore, receiver sensitivity, BER, or PER are interchangeable to calculate the communication range for a given wireless system.

Recently, high-speed packet wireless networks are being designed for mobile users in urban environments. Examples include the upcoming IEEE 802.11p and IEEE 802.16e standards [49, 102]. In typical urban environments, due to the Delay spread introduced by multipath and the Doppler Spread introduced by the mobility of users and environment objects, the wireless channel varies rapidly, and becomes much more unpredictable. Mobility changes everything, including the relationship between receiver sensitivity, BER, and PER. Now PER also depends on speed of user terminal. As a result, communication range is a function of user speed.

Assuming a Rician fading channel model and a two-ray path loss model, this chapter shows by extensive simulations how communication range of 802.11a/p varies with user speed. It is a preliminary attempt to establish the quantitative relationship between range and speed for a particular packet wireless system.

3.2 System and Channel Model

The wireless system we consider consists of a Road Side Unit (RSU) setup along the highway and an On Board Unit (OBU) in a moving vehicle [80, 102]. The distance between the RSU and the OBU is d , and the vehicle is moving at speed v , as depicted in Fig.3.1.

We ignore packet errors caused by MAC protocols, for example, collisions due to multiple

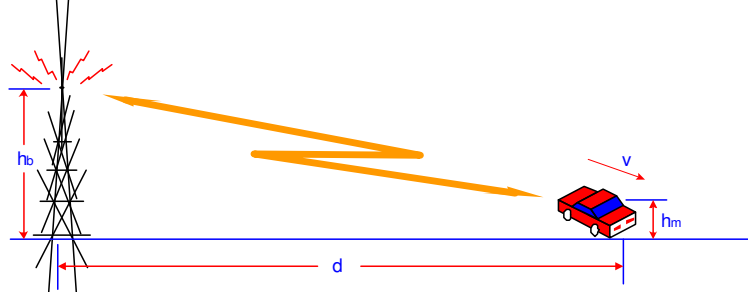


Figure 3.1: Basic elements of a mobile packet wireless system

access. It is assumed that the phase can be fully recovered at the receiver. Actually fast fading will cause phase estimation error, and degrade performance. So our results are optimistic performance bounds.

For simplicity, the modulation scheme for the wireless connection between the RSU and the OBU is assumed to be Binary Phase Shift Keying (BPSK). Coherent reception is considered. Let the center carrier frequency be f_c , and the speed of light be C . Then the wavelength λ is C/f_c , and the Doppler frequency f_D is v/λ .

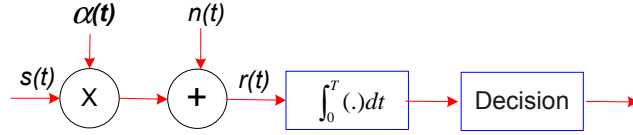


Figure 3.2: Block structure of a multiplicative and additive channel model

We assume an Additive White Gaussian Noise (AWGN) channel with Rician fading and Rayleigh fading between the RSU and OBU, as shown in Fig.3.2. For any transmitting signal $s(t)$, the input to the receiver is $r(t)$ given by

$$r(t) = \alpha(t) \cdot s(t) + n(t),$$

where $n(t)$ represents the AWGN noise process with zero mean and two-sided power spectral density $N_0/2$, and $\alpha(t)$ refers to a frequency non-selective and time-varying multiplicative envelope distortion of the transmitted signal $s(t)$. The probability density function $f(\alpha)$ of α in a Rayleigh fading channel

is expressed as

$$f(\alpha) = \frac{\alpha}{\sigma^2} \exp\left(-\frac{\alpha^2}{2\sigma^2}\right), \forall \alpha \geq 0, \quad (3.1)$$

where $E(\alpha) = \int_0^{+\infty} \alpha f(\alpha) d\alpha = \sqrt{\frac{\pi}{2}}\sigma$ and $E(\alpha^2) = 2\sigma^2$.

Rayleigh fading with a strong LOS component is called Rician fading with Probability Density Distribution $f(\alpha)$,

$$f(\alpha) = \frac{\alpha}{\sigma^2} \exp\left(-\frac{\alpha^2 + K^2}{2\sigma^2}\right) I_0\left(\frac{\alpha K}{\sigma^2}\right),$$

where K is the coefficient that indicates how strong the LOS component compared to rest of the received signal. I_0 is the zero-order modified Bessel function. I_0 is given by

$$I_0(s) = \frac{1}{2\pi} \int_0^{2\pi} \exp(s \cos \theta) d\theta,$$

and K is the ratio of the LOS component power over the total power of all other scattered power, i.e.,

$$K = \frac{\alpha^2}{2\sigma^2}.$$

As K goes to 0, the LOS component is so weak that Rician fading becomes Rayleigh fading. As K approaches infinity, the LOS component dominates, and there is no fading at all.

The calculation of communication range for a given wireless communications format depends on the RF propagation model employed. We use the two-ray path loss model for flat terrain presented in [58] to predict the average strength of receiving signals within a small interval. First, we obtain the Fresnel zone distance D_f by

$$D_f = \frac{4h_b h_m}{\lambda},$$

where h_b is the antenna height of the RSU and h_m is the antenna height of OBU. Before the Fresnel zone distance D_f , the average received power P_r at the OBU can be approximated by the free-space propagation-path loss formula

$$P_r(d) = P_t G_t G_r \left(\frac{\lambda}{4\pi d}\right)^2, \forall d \leq D_f,$$

where P_t is the transmitting power, G_t is the gain of transmitting antenna, and G_r is the gain of receiving antenna. After the Fresnel zone distance D_f , the average received power $P_r(d)$ at the OBU can be approximated by the two-way model's formula,

$$P_r(d) = P_t G_t G_r \left(\frac{h_b h_m}{d^2} \right)^2, \forall d > D_f.$$

Since the formula does not include wavelength as a variable, it is not an accurate prediction of path loss. However, it is useful in explaining the 40 dB/decade path loss as a function of distance d .

In summary, the average power of received signals $P_r(d)$ can be expressed in decibels as follows,

$$P_r(d) = \begin{cases} P_{TX} - L_1 - 20 \log(d), & 1 \leq d \leq D_f \\ P_{TX} - L_1 - 20 \log(D_f) - 40 \log\left(\frac{d}{D_f}\right), & d > D_f \end{cases} \quad (3.2)$$

where P_{TX} is the transmission power, and L_1 is the average power loss in decibels at 1 meter.

3.3 Communication Range vs Speed

Receiver sensitivity, denoted as R_s , is defined as the minimum receiving power to achieve given packet error rate (PER) for a given packet size. Theoretically R_s can be calculated as

$$R_s = N_t + N_s + SNR_{min}, \quad (3.3)$$

where N_t is the thermal noise (dBm) caused by electron activity in a resistive source, N_s is the system noise figure (dB), and SNR_{min} is the minimum theoretical SNR (dB) required for a given PER using the desired modulation in the absence of any channel interference.

Formula (3.4) illustrates how thermal noise N_t can be computed,

$$N_t = 10 \log(k \cdot T \cdot BW), \quad (3.4)$$

where k is the Boltzman's constant (1.38×10^{-23} Joules/K), T is the resistor's working temperature in Kelvin (K), and BW is the frequency spectrum used by the wireless system. Formula (3.4) can

be rewritten as

$$N_t = 10 \log(k \cdot T) + 10 \log(BW). \quad (3.5)$$

The first component is the thermal noise power in a 1 Hz bandwidth, called thermal noise floor. To simplify the comparison of receivers, a reference noise temperature is set to be room temperature 290K. At room temperature, the mean thermal noise floor is

$$\log(1.38 \cdot 10^{-23} \cdot 290 \cdot 1) = \log(4.002 \cdot 10^{-23}) = -174 \text{ dBm}.$$

The second component in (3.5) indicates how frequency contributes to the overall thermal noise. The interesting point is that a doubling of data rate increases the receive sensitivity by 3dB ($10 \log 2 = 3 \text{ dB}$) because of increased thermal noise.

The RF front-end system of receivers generates noise by itself, called system noise figure. The system noise figure N_s indicates the design quality of a radio receiver. It determines the sensitivity of the radio receiver and it is associated with elements such as the RF amplifier. Typical system noise figures for practical radio receivers are in the range of 2 to 20 dBm, depending on circuit, process, supply power, etc.

SNR_{min} is the acceptable minimum SNR for a receiver to achieve given PER at given packet size. For example, in IEEE 802.11a, at data rate of 6Mbps, SNR_{min} must be at least 5dB so that the PER is less than 10% at a PSDU length of 1000 bytes.

In general, SNR_{min} strongly depends on modulation/demodulation and coding schemes. In case of multipath fading wireless channels, SNR_{min} also depends on Doppler spread. Therefore, SNR_{min} is a function of speed v . As a result, in multipath fading channels, R_s also depends on speed v .

Using the equation (3.2) for the path loss model, and solving for the distance between transmitter and receiver, we have the following

$$d(v) = \begin{cases} 10^{[P_{TX} - L_1 - R_s(v)]}, & P_r(D_f) \leq R_s(v) \leq P_{TX} - L_1 \\ D_f \cdot 10^{[P_{TX} - L_1 - R_s(v) - 20 \log(D_f)]}, & R_s(v) \leq P_r(D_f), \end{cases} \quad (3.6)$$

where $P_r(D_f)$ is obtained by substituting $d = D_f$ into equation (3.2). Note here that communication range d is a function of speed v , because receiver sensitivity R_s depends on SNR_{min} , which is closely related to speed v .

3.4 Numerical Results

An end-to-end baseband PHY layer model of IEEE 802.11a from MATLAB file exchange center has been used for our extensive SIMULINK simulations. The model supports all mandatory/optional data rates with modulation/demodulation, coding, and interleaving schemes defined in 802.11a standard, although only 6 Mbps data rate with BPSK and 1/2 Convolutional coding and puncturing is used in our simulations. The Rayleigh fading block and Rician fading block from standard SMULINK library are used to configure the fading channels.

In most cases it is complex to get a closed-form SNR_{min} formula for different parameters of 802.11a. To simplify the problem, our SNR_{min} is obtained from extensive SIMULINK simulations.

Fig.3.3 plots the SNR_{min} as a function of Rician K -factor. As K becomes larger, a smaller SNR_{min} is required to achieve desired performance. Therefore, as seen in Fig.3.5, a larger LOS component means better channel condition and longer communication range, which is consistent with our intuitive understanding.

Fig.3.4 shows the additional TX power required to achieve the same communication range as that of LOS case with K -factor value of 10. In the design of a packet wireless system, a fade margin has been taken into account to combat multipath fading channels. Considering the worst case fading, Rayleigh fading, we usually use a fade margin of x dB, i.e., increase the receiver sensitivity by x dB. Similar results are found in Fig.3.4, if we compare the Rician fading case with $K = 10$ and Rayleigh fading case.

Fig.3.6 plots the SNR_{min} as a function of user speed for different Rician K -factor values. Substituting SNR_{min} into equation (3.6), we get Fig.3.7, which plots the range as a function of user

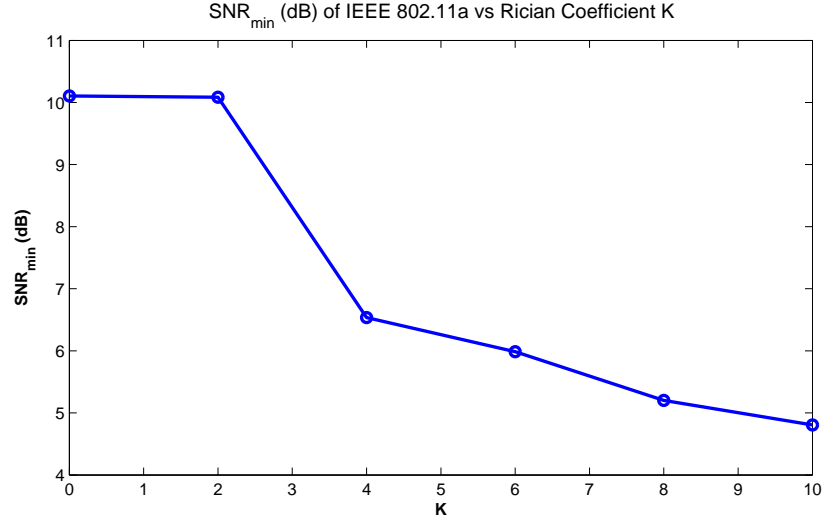


Figure 3.3: SNR_{min} vs Rician coefficient K

speed. Generally, as speed increases, SNR_{min} increases, and the communication range is shorter. It implies that as user speed increases, the cell coverage of packet wireless networks will be reduced, and denser deployment is needed to service the same area.

One way to combat the range reduction introduced by mobility is to increase the maximum power limit. For example, 802.11p increases the power limit from 20dBm of 802.11a to 28.8dBm. Additional transmission power is needed to keep the same range is shown in Fig.3.8. In a Rayleigh fading channel, the worst condition fading channel, as speed increases, it is not realistic to keep the same communication range by increasing the TX power.

Some parameters used for our calculation can be found in Table 3.1.

Parameter	Value	Parameter	Value
f_c	5.9GHz	h_b	3m
P_t	20dBm	h_m	1m
G_t	0dB	G_r	0dB

Table 3.1: Parameters for the computation of communication range

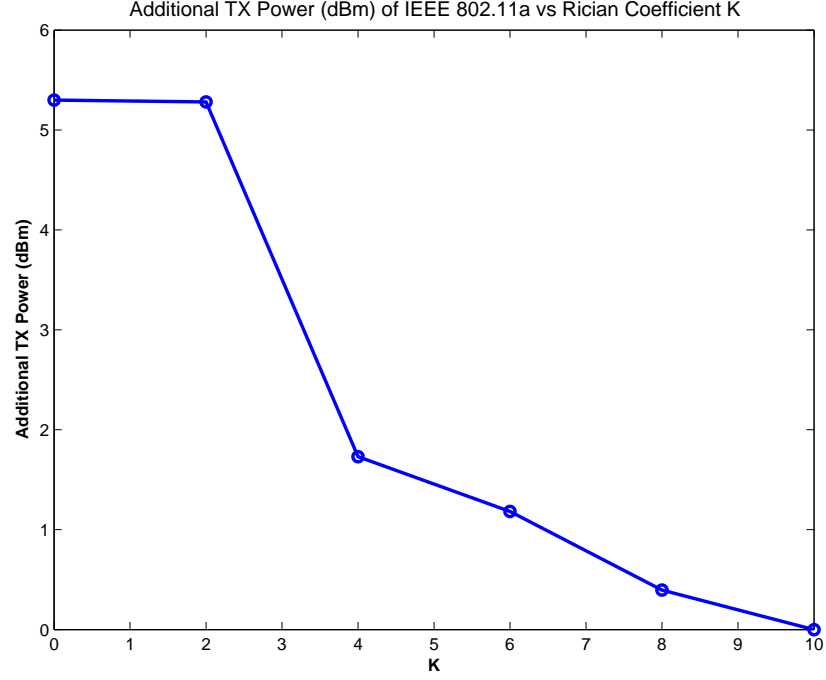


Figure 3.4: Additional TX Power vs Rician coefficient K

3.5 Concluding Remark

The increased mobility of users in packet wireless networks often results in fading, which changes some basic parameters of wireless systems. Extra efforts such as more transmission power, short cable between the radio and the transceiver, high-gain antennas, or denser deployment are needed to maintain the same service quality.

Analyzing, understanding and managing important parameters such as communication range in common network deployments can help alleviate the confusion about performance degradation and help wireless network designers make an informed choice. However, the calculation provided in this chapter only determines the theoretical maximum range. Due to many other factors such as antenna efficiency, cable loss, interference, and physical environment, etc., the actual communication range may be much smaller than the theoretical value.

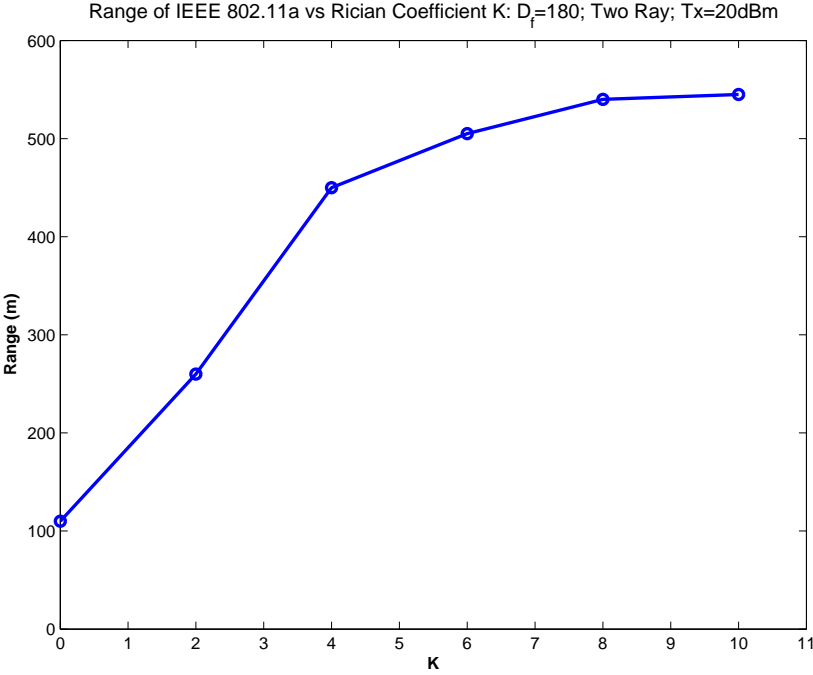


Figure 3.5: Communication range vs Rician coefficient K

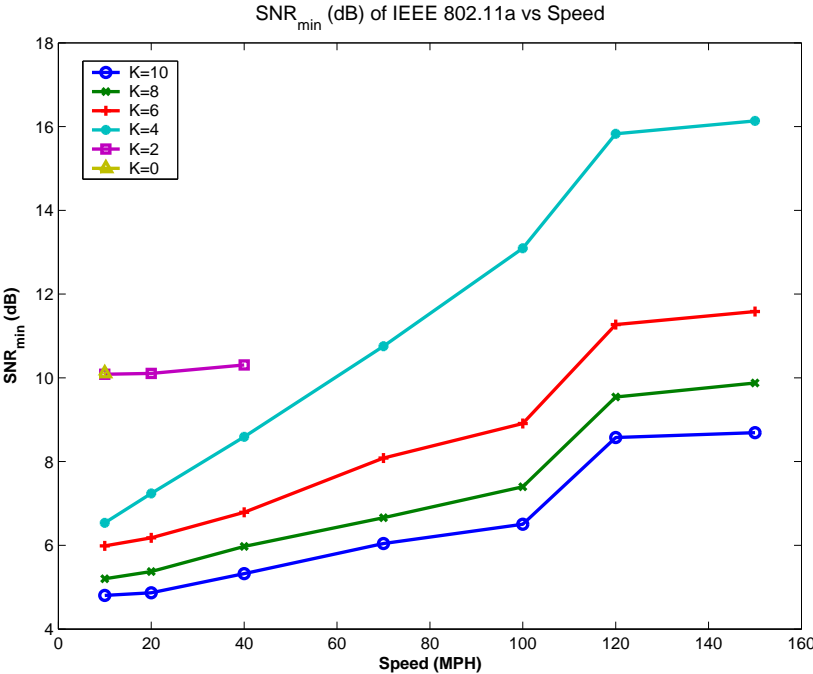


Figure 3.6: SNR_{min} vs Speed

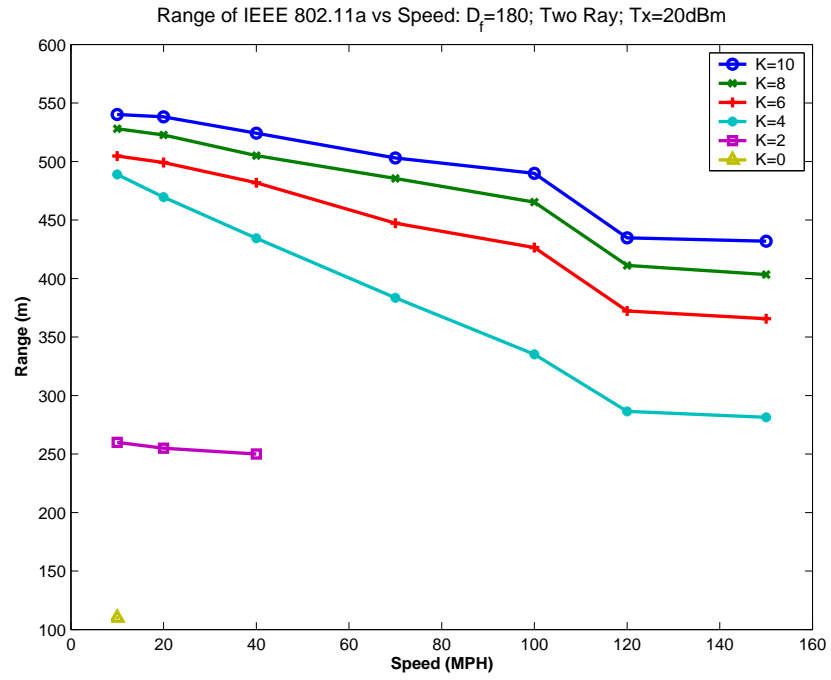


Figure 3.7: Communication range vs Speed

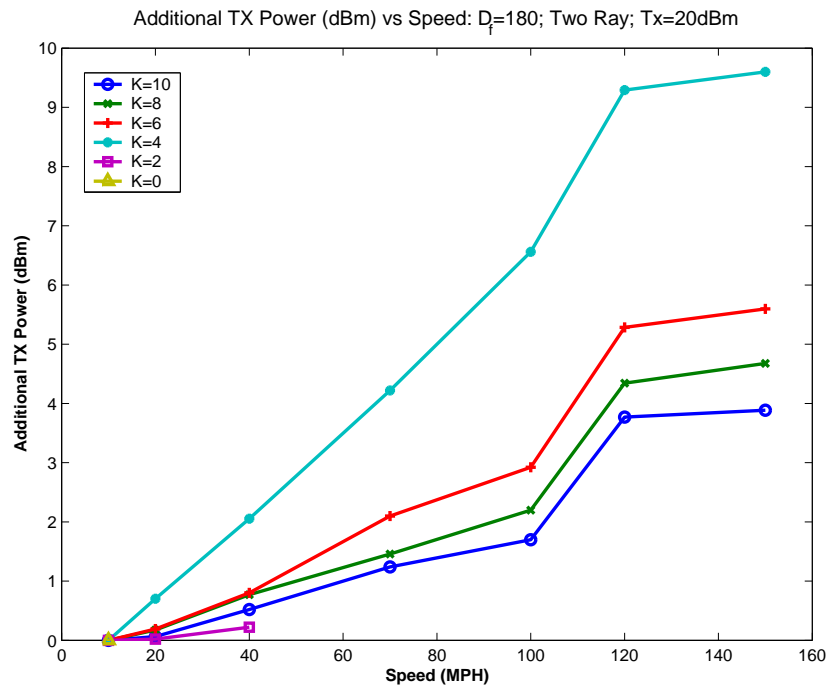


Figure 3.8: Additional TX power vs Speed

Chapter 4

Impact of Adaptive Modulation on TCP Throughput in Rayleigh Fading Channels

4.1 Introduction

Data-oriented high-speed wireless networks, such as wireless LANs and fixed broadband wireless networks, provide connectivity to a packet-based wired backbone network like Internet. Popular applications on wired networks including WWW, email, and file transfer applications, use TCP (Transmission Control Protocol) as the transport layer protocol for reliable data delivery across networks. To keep backward compatibility and allow seamless integration with wired networks, the use of TCP protocol over new generation high-speed wireless networks is inevitable.

The performance of all applications built on the TCP protocol is significantly affected by the TCP throughput in steady state. TCP throughput is determined mainly by the congestion/flow control mechanisms. In wired networks, TCP relies on packet drops as the indication of congestion.

It assumes a relatively reliable underlying network so that most packet losses are due to congestion that results in queue overflows in routers. The TCP source infers the presence of congestion either from the receipt of three duplicate acknowledgements (ACKs) or after the timeout of a retransmit timer, and then invokes congestion control mechanisms. However, in wireless networks, path loss, thermal noise, fading, and interference cause significant bit errors. Therefore, congestion is not the main reason of packet loss in wireless networks.

Recent research efforts have been made to model TCP congestion/flow control behavior and derive the steady TCP throughput. The analytical models of TCP over wired links (packet losses caused only by queue overflow in routers) have been studied in [55, 68, 81]. They all assume independent packet losses. But this assumption might not hold for wireless links with fading channels, where the channel has memory and packet losses are highly correlated. Assuming constant symbol transmission rate for all packets, [1, 108] present analytical models of TCP over fading wireless channels.

As wireless technology advances, the assumption of constant symbol rate transmission and fixed modulation scheme may not hold. Wireless communications occur in the public space, where signal transmissions suffer from many deterministic and nondeterministic factors such as path loss, shadowing, fading, etc. As a result, the wireless channel has a time-varying condition and capacity. Thus transmission techniques such as adaptive modulation, will play an important role in increasing the spectrum efficiency (throughput/channel bandwidth). For a given symbol transmission rate, the modulation scheme determines the data rate. Together with the modulation scheme, the symbol error rate introduced by channel condition will determine the efficient throughput. The best modulation scheme needs to be dynamically chosen for a time-varying channel condition to maximize the spectrum efficiency. Adaptive Modulation is one of the key enabling techniques in the new generation standards for wireless systems that have been developed to achieve high spectral efficiency on fading channels.

However, the impact of adaptive modulation on TCP throughput has not been well studied.

We present an analytical model of TCP congestion/control behavior with adaptive modulation enabled in Rayleigh fading wireless channels. The remainder of this chapter is organized as follows: section 4.2 briefly introduces the TCP protocol and its two variants, Tahoe and Reno. The system model, the threshold-based adaptive modulation, and the discrete Markov chain model of Rayleigh fading channel model are introduced in section 4.3. Section 4.4 models TCP congestion/flow control behavior using the discrete Markov chain. The analysis of TCP Tahoe and Reno is given in section 4.5. The model validation using NS2 simulations is presented in section 4.6, where we also propose a cross-layer approach to adaptively optimize TCP throughput based on the detection of physical layer parameters.

4.2 TCP Tahoe and Reno

TCP provides a connection-oriented and reliable end-to-end communication between pairs of processes in computers that are interconnected via data networks. As an important component in the layered TCP/IP protocol architecture, TCP obtains a simple and unreliable datagram service from the IP protocol [50, 76, 93]. TCP performs the following typical transport layer functions: connection setup/teardown; transferring a continuous byte stream; recovering from loss, corruption, duplication, delay, or reordering of packets that can occur in IP layer; flow control; congestion control; multiplexing/demultiplexing of multiple connections within a single computer. In this chapter, we are only interested in modeling TCP congestion control and flow control that play a dominant role in determining the steady state performance of a large file transfer. The algorithms for TCP congestion control and flow control contribute to the great success of today's Internet in spite of limited networking resource and diverse usage patterns.

Each byte of data sent over a TCP connection has a sequence number. The sequence number is monotonically increasing with wrapping back. Thus each byte that is received successfully can be acknowledged. TCP divides the contiguous byte stream from applications into TCP segments

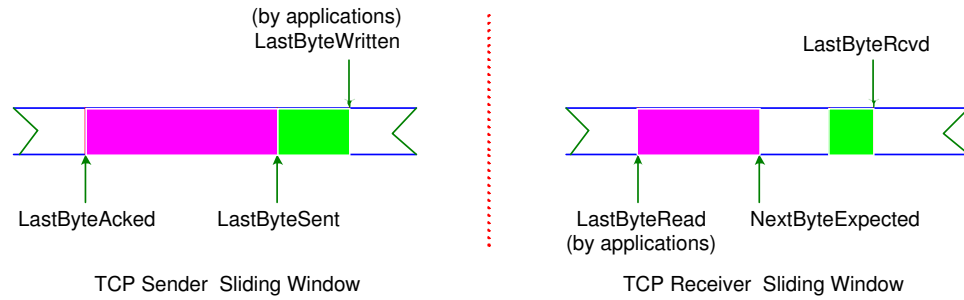


Figure 4.1: The sliding window for TCP protocol

with variable length to transmit them. Cumulative positive acknowledgments from the receiver to the sender specify the sequence number of the next byte that the receiver expects to receive, and confirm that bytes with previous sequence numbers have been received successfully.

A window-based flow control mechanism is deployed in TCP. As shown in Fig.4.1, both the TCP sender and the TCP receiver maintain a sliding window. The sender updates three variables: LastByteAcked, LastByteSent, LastByteWritten. The receiver also updates three variables: LastByteRead, NextByteExpected, LastByteRcvd. Sending Window Size W is the upper bound on the number of outstanding bytes (sent but not acknowledged) the sender can transmit, which means $\text{LastByteSent} - \text{LastByteAcked} \leq W$. Receiving window size, W_{max} , is the upper bound on the number of out-of-order bytes the receiver is willing to accept, which means $\text{LastByteRcvd} - \text{LastByteRead} \leq W_{max}$. Actually W_{max} is the maximum buffer size allocated for a TCP connection by the receiver. TCP flow control prevents sender from overflowing receivers buffer, and properly match the transmission rate of sender to that of the receiver. During the three-way handshaking phase of TCP connection setup, the receiver advertises W_{max} via AdvertisedWindow field in TCP header to the sender. During the data transfer, each ACK will carry variable AdvertisedWindows ($W_{max} - (\text{LastByteRcvd} - \text{LastByteRead})$) that is never greater than W_{max} .

The sender may send all packets within the sending window without receiving an ACK, but it must start a timeout timer (RTO) for each of them. The receiver must acknowledge each packet received, indicating the sequence number of the last well-received packet. Whenever an ACK

is received, the window variables of TCP sender is updated so that the sending window slides to the right.

Although the sequence number and window size are byte-oriented, we will switch to packet-oriented sliding window model to simplify analysis and reduce unnecessary computation complexity. With the assumption of constant-length TCP segment, the analysis in this chapter is equivalent to the byte-oriented model. In addition, TCP segment is called TCP packet in this chapter.

TCP Congestion control introduces congestion window W , a variable set by the TCP sender for each connection, to keep sender from overrunning buffers in routers. Congestion window W is also the upper bound on the number of outstanding packets. Based on the current traffic load condition in networks, a variety of sophisticated congestion control algorithms dynamically configure the window size W . Since the slowest part of the network and receiver should be accommodated, TCP flow control and congestion control work together to configure the TCP sending window to $\min\{W, W_{max}\}$. Since the sending window W is adaptively changed mostly by congestion control algorithms, it is also called the congestion window.

Since the Internet suffered from congestion collapse in the late 80's [50], several congestion control algorithms were proposed and implemented to prevent the TCP senders from overwhelming the resources of the network. In 1988, TCP Tahoe was introduced with three congestion control algorithms: slow start, additive increase/multiplicative decrease (AIMD), and fast retransmit. Since then, many modifications have been made to TCP and several different versions of TCP have been implemented [22, 72, 93]. TCP Reno revises TCP Tahoe by modifying the fast retransmit to include fast recovery. TCP Vegas implements a fundamentally different congestion avoidance algorithm than TCP Reno [13]. It uses the difference between expected and actual flow rates to estimate the available bandwidth in the network. Another conservative extension of TCP Reno is TCP SACK, which adds Selective Acknowledgment to TCP [22]. In this chapter we study the two most widely used variants of TCP protocol, Tahoe and Reno.

The slow start is used to probe the network to estimate the available bandwidth at the

beginning of a transfer or after loss recovery. A slow start threshold variable, W_{th} , is introduced to determine whether the slow start should be invoked. The minimum value of W_{th} is 2. When a TCP connection is established, W is first set to 1, and then on each received non-duplicate ACK, W is updated by

$$W = W + 1, \forall W < W_{th}$$

which implies doubling of W for each RTT.

Slow start is followed by congestion avoidance when $W \geq W_{th}$. During congestion avoidance, AIMD mechanism is invoked. The congestion window W is incremented by 1 per RTT (Additive Increase)

$$W = W + \frac{1}{W}, \forall W \geq W_{th}$$

until congestion or a packet loss is observed, e.g., either a retransmit timer expires or three duplicated ACK (four identical ACKs) are received. If a retransmit timer expires, W and W_{th} are updated by:

$$\begin{cases} W &= 1 \\ W_{th} &= \max \left\{ 2, \left\lceil \frac{W}{2} \right\rceil \right\} \end{cases}$$

and the sender retransmits the lost packet. If the retransmission fails, exponential backoff algorithm of the retransmit timer is triggered. In this chapter, we assume that multiple consecutive retransmit Timer expiration will not happen.

Duplicate ACKs can be either caused by lost packets or by reordered packets. If only one duplicate ACK is received, the sender may not know what really happened. However, if several (3 in RFC 793 [76]) duplicate ACKs are received, it is reasonable to infer that a packet loss has occurred. Then fast retransmit is used by both Tahoe and Reno to speed up the retransmission process. With fast retransmit, after receiving three duplicate ACKs, the sender will not wait for the retransmit timer to expire, and it will retransmit the lost packet immediately instead.

By and large, the TCP congestion window evolution is the same for TCP Tahoe and Reno. However, after the retransmit caused by three duplicate ACKs as the congestion indication, Tahoe

and Reno respond in different ways. TCP Tahoe simply updates W and W_{th} in the same way as in the case of a retransmit timer expiration. However, if W is large and available bandwidth is not small, TCP sender has to spend some time probing the network again for the already known W_{th} . It would be reasonable for the sender to continue from the start of the congestion avoidance phase. The fast recovery algorithm is proposed in Reno TCP to follow the fast retransmit phase until a non-duplicate ACK arrives.

At the beginning of fast retransmit and fast recovery phase, TCP Reno updates W and W_{th} by

$$\begin{cases} W_{th} &= \max \left\{ 2, \left\lceil \frac{W}{2} \right\rceil \right\} \\ W &= W_{th} + 3 \end{cases}$$

and retransmits the lost packet without waiting for the retransmit timer to expire. Note that three duplicate ACKs also means three packets have left the network and the pipe has empty slots for three packets. So the congestion window W is $W_{th} + 3$ instead of W_{th} . For each additional duplicate ACK received, W will be increased by 1. Thus the congestion window W will be artificially inflated by the number of packets that have left the network and been buffered by the receiver. To preserve the ACK-clocking property of TCP, each arrival of such additional duplicate ACKs at the sender is used to clock the transmission of a new data packet, if it is allowed by the congestion window. The arrival of a non-duplicate ACK (a recovery ACK) advances TCP sender from fast recovery phase to the congestion avoidance phase. At the end of fast recovery, TCP Reno updates W and W_{th} by

$$\begin{cases} W_{th} &= \max \left\{ 2, \left\lceil \frac{W}{2} \right\rceil \right\} \\ W &= W_{th} \end{cases}.$$

The congestion window W is deflated to W_{th} so that the Additive Increase starts again.

TCP Reno improves the performance of TCP Tahoe when a single packet from the outstanding packets within the same congestion window (loss window) is lost. However, it suffers from multiple packet losses in the same loss window. As pointed out in [22], the TCP Reno sender is

often forced to wait for a retransmit timeout when two or more packets in the same loss window are dropped. In this chapter, we are interested in two or more packet losses in the same loss window before the fast recovery for the first lost packet is initialized. To simplify the modeling, we adopt the same assumption as in [108], i.e., in small bandwidth \times delay links, two or more packet losses in the same congestion window before the fast recovery of the first lost packet is initialized, will finally cause the expiration of a retransmit timer.

4.3 System Model and Basic Assumptions

4.3.1 System Model

We consider a TCP connection between a server in a wired network and a mobile client with a dedicated wireless link to the wired network. A large file is being transferred from the server to the mobile client via FTP. The FTP transmission time is long enough for the connection to reach the steady TCP state. Ignoring TCP connection setup/teardown procedure, we are only interested in the bulk data transfer performance in steady state. The rate in wired links is considered much higher than the rate in wireless links. To ignore the effect of queuing caused by the intermediate system, the Drop-Tail buffer in the wireless base station is assumed nonempty at all time.

The wireless channel between the base station and the mobile client is assumed to be a single user discrete-time channel degraded by multipath fading, as shown in Fig.4.2. The channel has stationary and ergodic time-varying random gain and additive white Gaussian noise. We assume that the channel is slowly varying when compared to the symbol or frame transmission time. Therefore the channel between the transmitter and receiver is considered to remain approximately the same during a packet transmission. Define the signal transmission bandwidth to be B , the two-sided power spectral density of noise to be $N_0/2$, the average bit energy to be E_{avg} , and the Rayleigh channel fading amplitude to be α , then the Signal-Noise-Ratio (SNR) becomes $\gamma = \frac{\alpha^2 \cdot E_{avg}}{B \cdot N_0}$. In

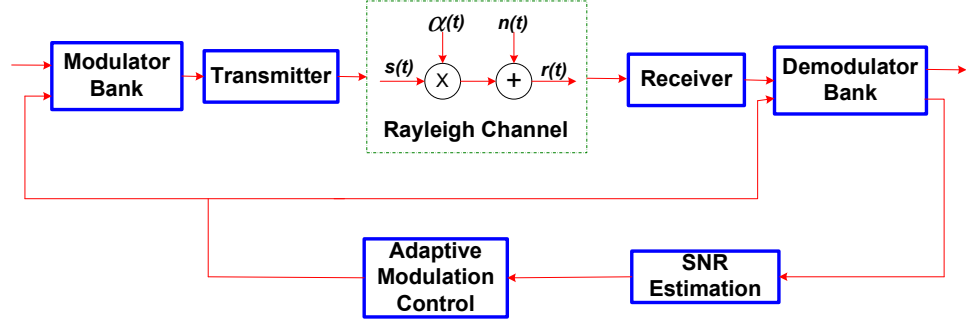


Figure 4.2: The block structure of the communication system

Rayleigh fading channel, the Probability Density Function (PDF) is given by:

$$f(\gamma) = \frac{\gamma}{\delta^2} \cdot \exp\left(-\frac{\gamma^2}{2\delta^2}\right), \gamma > 0$$

where δ is the Root Mean Square (RMS) value of received voltage before envelope detection, called the fading envelope of the Rayleigh channel.

The block diagram for both the transmitter and the receiver that we are working with is shown in Fig.4.2. The input to this system is a bit sequence. The modulation bank accepts the bit sequence and produces the modulated symbol sequence. Only BPSK, QPSK, QAM-16, and QAM-64 are implemented in the modulation bank. The change in modulation scheme follows the decision made by Adaptive Modulation and Coding Control block. It can switch among these four modulation schemes before the transmission of a packet. The decision of modulation scheme selection is sent via feedback channel or command to the transmitter based on the continuous monitoring of Signal Noise Ratio at the receiver. Ideal detection conditions are assumed (matched filter, no Inter-Symbol Interference). The symbol rate is the same for all modulation schemes, and the different constellation size for the modulation scheme causes a different data rate.

4.3.2 Adaptive Modulation

Let E_s be the average transmission energy per symbol (Watts) at the receiver input, E_b be the average transmission energy per bit (Watts, over one-bit interval) at the receiver input, $N_0/2$ be

the single-sided noise power density (Watts/Hz), and M be the modulated bits per symbol. Define

$$Q(x) = \frac{1}{2} \operatorname{erfc}\left(\frac{x}{\sqrt{2}}\right) = \frac{1}{\sqrt{2\pi}} \int_x^\infty \exp\left(-\frac{t^2}{2}\right) dt.$$

Then the symbol error probabilities for BPSK ($M=2$) and QPSK ($M=4$) [74] are

$$P_{s,BPSK} = Q\left(\sqrt{\frac{2E_b}{N_0}}\right),$$

$$P_{s,QPSK} = 2Q\left(\sqrt{\frac{2E_b}{N_0}}\right).$$

The approximate symbol error probabilities for QAM-16 and QAM-64 [74] are

$$P_{s,QAM-M} = 2 \cdot \left(1 - \frac{1}{\sqrt{M}}\right) \cdot Q\left(\sqrt{\frac{3}{(M-1)} \cdot \log_2^M \cdot \frac{2E_b}{N_0}}\right),$$

where $M = 16$ and 64 respectively. Let P_b be the bit error probability and P_s be the symbol error probability, then we have

$$P_b = \frac{P_s}{\log_2^M}.$$

The Signal-Noise-Ratio $\alpha = \frac{E_s}{N_0}$ for BPSK, QPSK, QAM-16, and QAM-64 is $\frac{E_b}{N_0}$, $\frac{2E_b}{N_0}$, $\log_2^{16} \frac{E_b}{N_0}$, and $\log_2^{64} \frac{E_b}{N_0}$, respectively [74].

We further assume that the constant frame length is L bits and bits of the same frame are corrupted uniformly with probability P_b for given SNR α . Let P_f be the frame error probability for given SNR α , R_s be the symbol rate, and R_d be the efficient average data rate, then we have

$$P_f = 1 - (1 - P_b)^L,$$

and

$$R_d = R_s \cdot (1 - P_f) \cdot \log_2^M.$$

Fig.4.3 shows the efficient Data Rates as a function SNR α for different modulation schemes. It is found that no modulation scheme achieves best data rates over a wide range of signal noise ratio. The range of received SNR α can be partitioned into a finite number of non-overlapping intervals, such as $[A_1, A_2)$, $[A_2, A_3)$, ..., $[A_{k-1}, A_k)$, ..., $[A_M, \infty)$ (In this chapter $M = 4$ and $0 = A_1 < A_2 <$

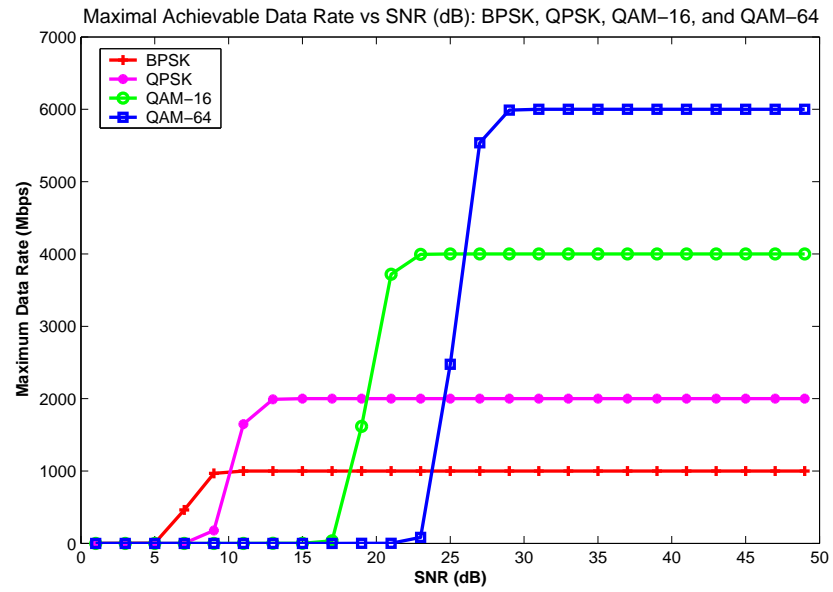


Figure 4.3: Maximum achievable data rate vs SNR: BPSK, QPSK, QAM-16, QAM-64

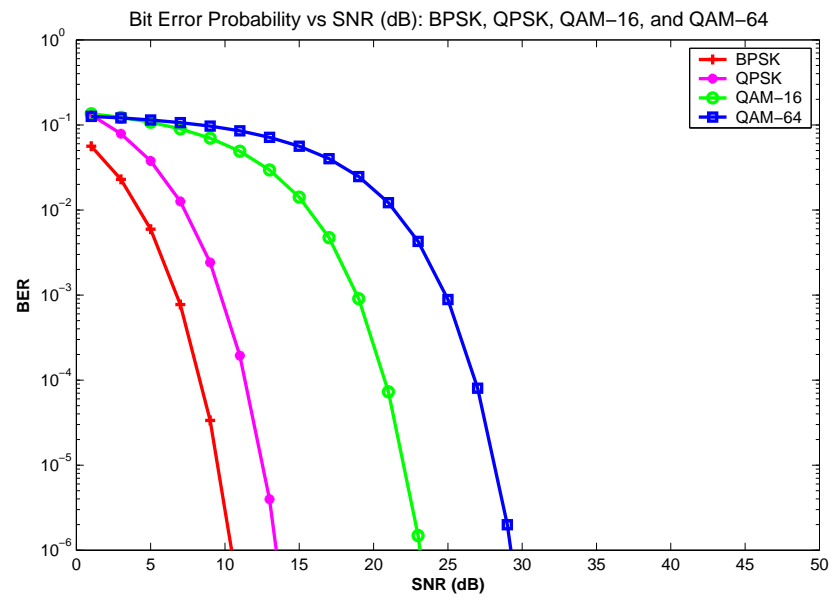


Figure 4.4: Average Bit Error Rate vs SNR: BPSK, QPSK, QAM-16, QAM-64

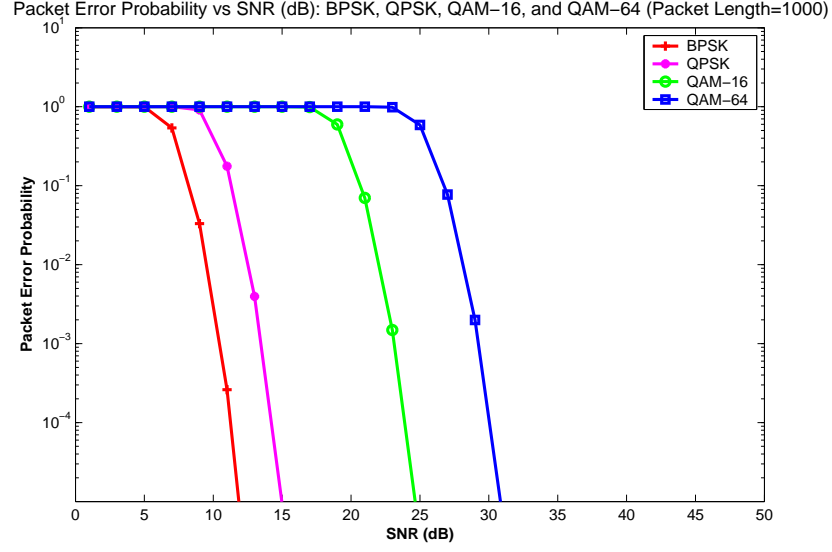


Figure 4.5: Average Packet Error Rate vs SNR: BPSK, QPSK, QAM-16, QAM-64

... $A_M < A_{M+1} = \infty$). For each SNR within interval $[A_i, A_{i+1})$ ($1 \leq i \leq M$), a particular modulation scheme provides the best data rate R_i . Since the wireless channel condition is always changing and the signal quality is unpredictable, adaptive modulation always chooses the best modulation scheme based on the Channel Side Information (CSI) obtained from the receiver.

In this chapter, we assume that the threshold-based adaptive modulation is implemented. Ideal conditions are assumed for adaptive modulation: no channel estimation error, no feedback transmission error, no channel state feedback delay, and no peak power constraint, etc. It is also assumed that all management and control frames are transmitted correctly and in time. The basic procedure for adaptive modulation follows:

- The receiver constantly measures and estimates current SNR α .
- The SNR α is converted into the packet error probability for each modulation scheme.
- Based on a target packet error probability, select a modulation scheme that yields the highest data rate while remaining within the BER target bounds.
- The receiver sends the decision of modulation scheme back to the transmitter via control

channel or management frames. Both the receiver and the transmitter switch to the new modulation scheme.

4.3.3 Channel Modeling

We are interested in closed-form analytical formulas for the channel. As discussed in [95], any partition of the received SNRs into a finite number of ranges leads to a finite-state Markov chain (FSMC) model, which we can use to approximate the Rayleigh fading channel. In our study, the finite ranges come from the intervals defined for the adaptive modulation in section 4.3.2, i.e., $[A_i, A_{i+1}), i = 1, 2, \dots, M$. Let $S = \{s_1, s_2, \dots, s_M\}$ be the set of channel states in the FSMC model, and let $\alpha (\alpha \geq 0)$ be the received SNR. If $A_k \leq \alpha < A_{k+1} (k = 1, 2, \dots, M)$, the channel is said to be in state s_k . Let $S_n (S_n \in S; n = 1, 2, \dots)$ be the discrete Markov stochastic process for the channel evolution. Then the fading channel can be fully defined by the $M \times M$ state transition probability matrix T , the $M \times 1$ steady state probability vector $\vec{\pi} = [\pi_1, \pi_2, \dots, \pi_M]$, and the $M \times 1$ crossover packet error probability vector $\vec{e} = [e_1, e_2, \dots, e_M]$.

Let ρ be the expectation of α over all ranges ($\rho = E[\alpha]$). The PDF function of the Rayleigh fading is

$$f_A(\alpha) = \frac{1}{\rho} \cdot \exp \left\{ -\frac{\alpha}{\rho} \right\}.$$

Then the steady state probability of state i ($i = 1, 2, \dots, M$) is

$$\begin{aligned} \pi_i &= \int_{A_i}^{A_{i+1}} f_A(\alpha) d\alpha = \int_{A_i}^{A_{i+1}} \frac{1}{\rho} \cdot \exp \left\{ -\frac{\alpha}{\rho} \right\} d\alpha \\ &= \exp \left\{ -\frac{A_i}{\rho} \right\} - \exp \left\{ -\frac{A_{i+1}}{\rho} \right\}. \end{aligned}$$

Each element $t_{i,j}$ ($i, j \in 1, 2, \dots, M$) of the state transition probability matrix T is defined as

$$t_{i,j} = \Pr(S_{n+1} = s_j | S_n = s_i),$$

for $n = 0, 1, 2, \dots$. Assuming that the one-step state transition only occurs between neighboring

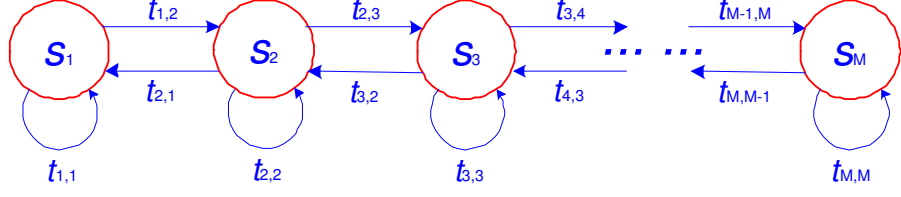


Figure 4.6: Finite state Markov chain model for a Rayleigh fading channel

states as shown in Fig.4.6, we have

$$t_{i,j} = 0, \forall |i - j| \geq 2. \quad (4.1)$$

Other non-zero state transition probabilities include

$$\begin{aligned} t_{i,i+1} &= \frac{N_{i+1}}{R_t^{(i)} \cdot \pi_i}, i = 1, 2, \dots, M-1 \\ t_{i,i-1} &= \frac{N_i}{R_t^{(i)} \cdot \pi_i}, i = 2, 3, \dots, M \end{aligned} \quad (4.2)$$

$$t_{i,i} = \begin{cases} 1 - t_{i,i+1} - t_{i,i-1}, & \text{if } 1 < i < M \\ 1 - t_{1,2}, & \text{if } i = 1 \\ 1 - t_{M,M-1}, & \text{if } i = M, \end{cases}$$

where $N_i (i = 1, 2, \dots, M)$ is the level crossing rate function of state i for a Rayleigh fading channel [95] and $R_t^{(i)} (i = 1, 2, \dots, M)$ is the long-term average transmitted symbols per second in channel state i ($i = 1, 2, \dots, M$). With modern high-speed wireless packet network systems, we assume that $R_t^{(i)} \gg N_i$ or N_{i+1} ($i = 1, 2, \dots, M-1$) hold. From [95], $N_i (i = 0, 1, 2, \dots, M)$ is given by

$$N_i = \sqrt{\frac{2\pi A_i}{\rho}} \cdot f_d \cdot \exp \left\{ -\frac{A_i}{\rho} \right\},$$

where $\rho = E[\alpha]$ and $f_d = v/\lambda$ (v is the terminal moving speed, λ is the wavelength, and f_d is called

the maximum Doppler frequency). Then the state transition matrix looks like

$$T = \begin{bmatrix} t_{1,1} & t_{1,2} & \cdots & \cdots & 0 \\ t_{2,1} & t_{2,2} & t_{2,3} & \cdots & \vdots \\ 0 & \ddots & \ddots & \ddots & 0 \\ \vdots & 0 & t_{M-1,M-2} & t_{M-1,M-1} & t_{M-1,M} \\ 0 & \cdots & 0 & t_{M,M-1} & t_{M,M} \end{bmatrix}. \quad (4.3)$$

Let the average bit error probability or crossover probability vector \vec{P}_e of state i ($i = 1, 2, \dots, M$) be $P_{e,i}$ and the modulation scheme for state i be scheme i . Then $P_{e,i}$ can be approximately calculated based on the bit error probability function $f_{e,i}(\alpha)$ for any given SNR α under modulation scheme i . Following the derivation in [95] (see section 4.8), we have

$$P_{e,i} = \frac{\int_{A_i}^{A_{i+1}} \frac{1}{\rho} \cdot \exp\left\{-\frac{\alpha}{\rho}\right\} \cdot f_{e,i}(\alpha) \cdot d\alpha}{\int_{A_i}^{A_{i+1}} \frac{1}{\rho} \cdot \exp\left\{-\frac{\alpha}{\rho}\right\} \cdot d\alpha}.$$

If the bit error probability function has the general form like

$$f_{e,i}(\alpha) = \eta \cdot Q\left(\sqrt{\xi \cdot \alpha}\right), (0 < \eta \leq 1, 0 < \xi < 1), \quad (4.4)$$

the closed-form average bit error probability in state i can be derived by a similar manipulation in [95],

$$P_{e,i} = \frac{\eta \cdot (\gamma_i - \gamma_{i+1})}{\pi_i}, \quad (4.5)$$

where

$$\gamma_i = \sqrt{\frac{\xi \cdot \rho}{\xi \cdot \rho + 2}} \cdot F\left(\sqrt{\frac{2 \cdot A_i \cdot (\xi \cdot \rho + 2)}{\xi \cdot \rho}}\right) + \exp\left\{-\frac{A_i}{\rho}\right\} \cdot F\left(\sqrt{\xi \cdot A_i}\right),$$

and π_i is from equation (4.1).

Given the constant packet length L and the average bit error probability $P_{e,i}$ for modulation scheme i , the packet error probability e_i is

$$e_i = 1 - (1 - P_{e,i})^L. \quad (4.6)$$

4.4 Markov Model of TCP Evolution

The sustainable TCP throughput is vital to a large file transfer, in which the TCP connection setup and termination only take a small fraction of overall transmission time. Therefore this chapter ignores TCP connection setup/termination procedures, and concentrates on the stable TCP stages. So TCP flow control and congestion control mechanisms play the most important role in the TCP throughput. From the discussion in Section 4.2, it is found that the combined behavior of TCP flow control and congestion control can be fully characterized by the following variables: slow start threshold, current congestion window, maximum receiver window, current outstanding packets in the sending window, and the timer status of each outstanding packet. If fast recovery is used, its status also affects the TCP behavior. However, it is not practical to construct Markov states with reasonable size based on all of these status variables. Moreover, the congestion window size and the timer status of each outstanding packet may depend on past state transitions, which doesn't satisfy the Markov property.

To make the problem tractable, we only consider three state variables and sample these states at some particular points $k = t_k$. In the sampling time slots, one of following events occurs:

- A timeout timer for an outstanding packet expires
- A fast recovery phase is successfully finished in case of TCP Reno or a fast retransmit is triggered in case of TCP Tahoe.

At these sampled time slots, all previous outstanding packets are acknowledged. So TCP's memory of outstanding packets and related timers are cleared. In addition, its congestion window only depends on the TCP states at previous sampled time slots. Assume that the state is (A_{k-1}, W_{th}, W) at sampled time slot t_{k-1} . If a timeout timer expires at sampled time slot t_k , the state goes to $(A_k, \lceil W/2 \rceil, 1)$; if a fast recovery is successfully completed at sampled time slot t_{k-1} , the state goes to $(A_k, \lceil W/2 \rceil, \lceil W/2 \rceil)$.

The state space Ω_x is

$$\Omega_x = \{(m, W_{th}, 1)\} \cup \{(m, W_{th}, W_{th})\},$$

where $m \in \{1, 2, \dots, M\}$ and $2 \leq W_{th} \leq \lceil \frac{W_{max}}{2} \rceil$. The first set includes all states caused by timeout events of outstanding packets. States in the second set are caused by a successful fast recovery. Here the congestion window size is not less than 2, which is consistent with related TCP RFCs and is different from that in [108]. The simulator ns2 we used for this chapter is consistent with the TCP RFCs in terms of the evolution of slow start threshold and congestion window size. The congestion window size is always reset to be an integer after timeout expiration or successful fast recovery. So the total number of states is $M \cdot 2 \left(\lceil \frac{W_{max}}{2} \rceil - 1 \right)$.

For this Markov chain model, we need to calculate its stationary probability $\pi_{t,i}$ ($i \in \Omega_x$), and the state transition probability $P_{i,j}$ ($i \in \Omega_x, j \in \Omega_x$).

4.4.1 Computation of Transition Probability

With the assumption of small bandwidth \times delay product and negligible ACK transmission time, there is at most one packet being transmitted at any instance.

The exploration of only the sample space Ω_x is not enough to analyze the transition probability $P_{i,j}$ of the Markov Chain Model. The TCP dynamic behavior has to be incorporated into the model for the performance analysis. It directly decides the transition probability, the delay, and successful transmissions between two states in Ω_x . To fully model the intermediate TCP evolution procedure, we only need to memorize the correct window size, denoted by Y . All possible values of Y can be calculated if the maximal window size W_{max} and slow start threshold W_{th} are given. But it will lead to a large state space. To reduce the state space, we define the intermediate TCP state space to be Ω_Y .

Given the current window size Y , the protocol only cares how many more packets can be transmitted and the values of the next window size and slow start threshold in case of packet losses.

Assuming W_{max} to be an even number, Y can be mapped into one element of set Ω_Y : $[1, 2)$, 2 , $[2, 3)$, $[3, 4)$, ..., $[W_{max} - 1, W_{max})$, and W_{max} . If Y is mapped to the same element in Ω_Y , they will lead to the same TCP behavior.

During the derivation of transition probability at beginning state i , we track all possible different transition paths from state i to any destination state j via one or more intermediate states k in Ω_Y , and record all associated transform variables. If there are multiple paths from state i to state j , these paths must constitute mutually exclusive transmission events. Therefore, we add all probabilities for these paths from state i to state j together to get the transition probability $P_{i,j}$.

4.4.2 Discrete Semi-Markov Process and Throughput Calculation

In the previous section, to simplify the Markov chain formulation, we assume that the holding time in each state out of the state space Ω_x is the same. However, to calculate the average throughput and delay, we have to take into account the random holding time in each state. This section uses the semi-Markov process introduced in [108] to model the dynamic behavior of TCP.

Consider the TCP stochastic process $X_k(t)$ with finite state space Ω_x . Let $D_{i,j}$ be the average amount of holding time, $E_{i,j}$ be the average number of transmission attempts, and $S_{i,j}$ be the average number of successful transmissions from TCP state $i \in \Omega_x$ to TCP state $j \in \Omega_x$. The corresponding matrices for them are denoted as D , E , and S , respectively.

To keep track of variables such as average delay, number of transmission attempts, and number of successful transmissions that help to derive throughput, we define a vector of transform variables $z = (z_d, z_t, z_s)$. Let $\xi_{i,j}(N_d, N_t, N_s)$ be the probability that the system switches from TCP state i to TCP state j after N_t transmission attempts with N_s successful transmissions and average delay N_d . Since N_d , N_t , and N_s take non-negative discrete values, we have

$$\Phi_{i,j}(z_d, z_t, z_s) = \sum_{N_d, N_t, N_s} \xi_{i,j}(N_d, N_t, N_s) z_d^{N_d} z_t^{N_t} z_s^{N_s}.$$

Based on the time of first packet loss, the evolution between two TCP states can be

divided into two phases: before loss and after loss. In the phase before loss, TCP Tahoe and TCP Reno have the same behavior and lead to the same intermediate TCP state. In the phase after loss, due to different window adaptation mechanism, TCP Tahoe and TCP Reno demonstrate different behavior and finally transfer to different TCP state, even though they have the same starting intermediate TCP state and channel state. Since these two phases are independent, they can be analyzed separately and then combined together to get the complete model of evolution between two TCP states. Therefore, we introduce function $\Phi_{i,k}^{(1)}(z)$ for the phase before loss from TCP state i to intermediate TCP state k , and $\Phi_{k,j}^{(2)}(z)$ for the phase after loss from intermediate TCP state k to intermediate TCP state j . Then we have

$$\Phi(z) = \Phi^{(1)}(z)\Phi^{(2)}(z),$$

where $\Phi(1, 1, 1)$, $\Phi^{(1)}(1, 1, 1)$, and $\Phi^{(2)}(1, 1, 1)$ represents the transition probability matrix between two TCP states, from a TCP state to an intermediate TCP state, and from an intermediate TCP state to a TCP state, respectively.

Then the matrix of average delay between any two TCP states can be calculated by

$$D = \left. \frac{\partial \Phi(z_d, z_t, z_s)}{\partial z_d} \right|_{z_d, z_t, z_s=1} = D_1 \Phi^{(2)}(1, 1, 1) + \Phi^{(1)}(1, 1, 1) D_2,$$

where

$$\begin{aligned} D_1 &= \left. \frac{\partial \Phi^{(1)}(z_d, z_t, z_s)}{\partial z_d} \right|_{z_d, z_t, z_s=1}, \\ D_2 &= \left. \frac{\partial \Phi^{(2)}(z_d, z_t, z_s)}{\partial z_d} \right|_{z_d, z_t, z_s=1}. \end{aligned} \quad (4.7)$$

Similarly, we can calculate the matrix of the number of transmission attempts, E , and the matrix of the number of successful transmission, S .

Note that S can not be accurately estimated, because this will involve some packet transmissions in at least two cycles, which is against the memoryless assumption of each cycle. For example, a particular outstanding packet in current cycle may be transmitted successfully without

being acknowledged. Then possibly it will be retransmitted again in future cycles. If the first successful transmission attempt is counted, we'll get redundant counting for the same packets in future cycles. However, it is also possible that the packet will be acknowledged cumulatively without retransmission in future cycles, then we do need to count it in the initial cycle to get an accurate value. To avoid the computation complexity introduced by enumerating the transmission states of all outstanding packets in the phase after loss, we adopt the bounding strategy presented in [108] to handle this. If further information from future transmissions is needed to decide whether a successful transmission should be counted, we get an optimistic throughput prediction (upper bound) by counting it and get a pessimistic throughput estimation (lower bound) by ignoring it. In this chapter we'll use the optimistic approximation approach.

With the assumption of constant frame length L , since $\Phi^{(1)}(1, 1, 1)$, $\Phi^{(2)}(1, 1, 1)$, D , E , and S are known, the steady average TCP throughput can be computed by

$$\text{Throughput} = \frac{\sum_{i \in \Omega_X} \pi_i \sum_{j \in \Omega_X} P_{i,j} S_{i,j} L}{\sum_{i \in \Omega_X} \pi_i \sum_{j \in \Omega_X} P_{i,j} D_{i,j}}.$$

4.4.3 Some Useful Functions

Given the channel state transition matrix T , we define $d(c, n)$ to be the average delay function as of the beginning channel state c , and the number of transmission attempts n . Define

$$Z = \begin{bmatrix} z^{d_1} & z^{d_2} & \dots & z^{d_M} \\ z^{d_1} & z^{d_2} & \dots & z^{d_M} \\ \vdots & \vdots & \ddots & \vdots \\ z^{d_1} & z^{d_2} & \dots & z^{d_M} \end{bmatrix}.$$

We get a new matrix B by inner-product of T and Z (e.g., $B_{i,j} = T_{i,j} \cdot Z_{i,j}$). The average delay can be calculated by using an intermediate transformation function $\beta_{c,n}(z)$ defined by

$$\beta_{c,n}(z) = b_c \cdot B^n, n \geq 1, 1 \leq c \leq M,$$

where b_c is the c th row vector of the M -dimensional unit matrix $I_{M \times M}$. Then we have

$$d(c, n) = \left. \frac{\partial \beta_{c,n}(z)}{\partial z} \right|_{z=1}.$$

Denote the packet success probability matrix T_S and the packet error probability matrix T_F by

$$\begin{cases} T_S = T \cdot B \\ T_F = T \cdot F \end{cases},$$

where

$$F = \begin{bmatrix} e_1 & e_2 & \dots & e_M \\ e_1 & e_2 & \dots & e_M \\ \vdots & \vdots & \ddots & \vdots \\ e_1 & e_2 & \dots & e_M \end{bmatrix},$$

(e_i is the cross-over packet error probability in modulation scheme i) and

$$F = 1 - B.$$

Let $G(n)$ be the event that $n - 1$ consecutive successful transmissions followed by a failure transmission, the probability matrix for this event is

$$P[G(n)] = T_S^{n-1} \cdot T_F, n \geq 1.$$

If the starting TCP state is i and the TCP state after n such transmissions is j , the corresponding probability is $\alpha_{i,j}(n) = P[G(n)]_{i,j}$.

Let $H(k, n)$ be the event that there are k success transmissions in n consecutive transmission attempts. If the starting TCP state is i and the TCP state after n such transmissions is j , the probability is $P[H(k, n)]_{i,j}$. The probability matrix $P[H(k, n)]$ is:

$$P[H(k, n)] = \sum_{i_1, i_2, \dots, i_K} A_1 \cdots A_{i_1} \cdots A_{i_2} \cdots A_{i_k} \cdots A_n,$$

where

- $1 \leq i_1 < i_2 \leq \dots \leq i_K \leq n$
- $A_{i_1} = A_{i_2} = \dots = A_{i_K} = T_S$
- $A_j = T_F, \forall j \neq i_1, i_2, \dots, i_K$.

Since the k for this chapter is usually small, we can use exhaustive list to calculate it.

Let $B(k)$ be the event that a packet failure is immediately followed by k consecutive successful transmissions. The probability matrix for this event is

$$P[B(k)] = T_F \cdot T_S^k, n \geq 1.$$

Let $B(k, l, i), 1 < l \leq k < i$ be the event that a packet loss is followed by k successful transmissions within i consecutive transmissions and a second packet loss occurs at transmission attempt l . The probability matrix for this event is

$$P[B(k, l, i)] = T_F \cdot T_S^{l-1} \cdot T_F \cdot P[H(k-l+1, i-l)].$$

Note that after the first loss, there are $l-1$ consecutive successful transmissions.

Define the channel state function $\varphi(x) = c$ ($1 \leq c \leq M$) for any $x = (c, W_{th}, W \in \Omega_X$, and $\varphi(x) = c$ ($1 \leq c \leq M$) for any $y = (c, W) \in \Omega_Y$.

4.5 Analysis of TCP Tahoe and Reno

4.5.1 Computation of $\phi_1(z)$

For each of these TCP variants, the evolution of TCP states from $X \in \Omega_X$ to $Y \in \Omega_Y$ is determined and can be tabulated before a packet failure is detected.

As to the transform variables, the number of transmission attempts N_t is known to be n , the average delay $N_d = d(\varphi(X), N_t)$, and the number of successful transmissions N_s is $n-1$. Therefore, the transform function $\Phi_{XY}^{(1)}(z_d, z_t, z_s)$ is:

$$\sum_{n \in C(X, Y)} \alpha(n) z_d^{d(\varphi(X), N_t)} z_t^n z_s^{n-1}, n \geq 1,$$

where $n = C(X, Y)$ is the number of transmission attempts so that the TCP state goes from X to intermediate TCP state Y .

4.5.2 Computation of $\phi_2(z)$

TCP variants differ in the way to detect and respond to packet losses, and in the packet loss recovery mechanism. These difference will affect the model for each TCP variant and thus the performance evaluation. Two TCP variants, TCP Tahoe and TCP Reno, will be discussed in this chapter.

TCP Reno

Case 1: Fast retransmission is not triggered

If the TCP sender can't receive K duplicate ACKs or less than $K - 1$ packets successfully arrive at the TCP receiver, the retransmission will not be triggered. The reason for this can be further divided into two cases:

Case 1.1: The case of $\lfloor Y \rfloor \leq K$

In this case, the number of allowable outstanding packets that triggers the duplicate ACK packet is less than $K - 1$. Even if all outstanding packets and corresponding duplicate ACK packets can be transmitted successfully, they are still not enough to start the fast retransmit caused by K duplicate ACK packets. The TCP sender has to wait for the expiration of timeout timer. After the timeout timer expiration, the slow start threshold will be halved, and current congestion window will be reset to be 1. Therefore, the TCP state goes to $X(k + 1) = (C, \lceil \frac{Y}{2} \rceil, 1), 1 \leq C \leq M$.

Since there is $N_t = \lfloor Y \rfloor - 1$ transmission attempts within $N_d = d_{RTO}$ after the first packet failure, the transform function $\Phi_{YX}^{(2)}(z_d, z_t, z_s)$ is

$$W_{XY}(d_{RTO}) z_d^{d_{RTO}} z_t^{\lfloor Y \rfloor - 1} z_s^{N_s},$$

where $W_{XY}(d_{RTO}) = T^{d_{RTO}}$ (here there is an approximation). Since we don't enumerate all transmission possibilities exhaustively in this case, N_s is uncertain. Thus the bounding approach is used

here. In case that the value of a transform variable is unknown for the transition from TCP state space to Intermediate TCP state space, we are only interested in the mean of that transform variable, i.e., $E[N_s]$ for this case. Using the union bound, we have:

$$0 \leq E[N_s] \leq \left[\sum_{l=1}^{\lfloor Y \rfloor - 1} \sum_{m=1}^M \vec{b}_j \cdot T_s \right]_{i,j}.$$

Case 1.2: The case of $\lfloor Y \rfloor > K$

In this case, less than K of $\lfloor Y \rfloor$ outstanding packets are successful transmitted so that the fast retransmit can be triggered. After timeout timer expiration, the slow start threshold will be halved, and current congestion window will be reset to be 1. Therefore, the TCP state goes to $X(k+1) = (C, \lceil \frac{Y}{2} \rceil, 1)$, $1 \leq C \leq M$.

The transform function $\Phi_{YX}^{(2)}(z_d, z_t, z_s)$ should be

$$W_{XY} z_d^{d_{RTO}} z_t^{\lfloor Y \rfloor - 1} z_s^{N_s},$$

where the probability W_{XY} is

$$W_{XY} = T_F^{\lfloor Y \rfloor} + P[H(1, \lfloor Y \rfloor)] + P[H(2, \lfloor Y \rfloor)] + \dots + P[H(K-1, \lfloor Y \rfloor)].$$

As to N_s , it should be less than $K-1$.

Case 2: Fast retransmit is triggered

Case 2.1: Single Loss within the same congestion window

Retransmit Success

If the retransmission of lost packet is successful, K transmissions (including the retransmission) in all will be attempted. After receiving the ACK for the retransmitted packet, TCP sender will change both its slow start threshold and congestion window to be half of current congestion window. Then TCP state will go to $X(k+1) = (C, \lceil \frac{Y}{2} \rceil, \lceil \frac{Y}{2} \rceil)$.

The transform function $\Phi_{YX}^{(2)}(z_d, z_t, z_s)$ should be

$$T_S^K \cdot z_d^{d(\varphi(X), K)} z_t^K z_s^K,$$

where the delay can be equivalently obtained from this intermediate transform function as presented in section 4.9.

Retransmit Failure

If the retransmission fails, then the TCP sender has to wait for the timeout timer expiration of the lost packet. After receiving the K th duplicate ACK, the congestion window becomes $W' = \min \left\{ \lceil \frac{Y}{2} \rceil + K, W_{max} \right\}$. After the timer expires, Then TCP state will go to $X(k+1) = \left(C, \lfloor \frac{W'}{2} \rfloor, 1 \right)$.

The transform function $\Phi_{YX}^{(2)}(z_d, z_t, z_s)$ should be

$$T_S^{K-1} \cdot T_F z_d^{d_{RTO}} z_t^K z_s^{K-1}.$$

Case 2.2: Multiple Loss before Fast Retransmit

Suppose that the first packet loss occurs at transmission attempt 0. Before the K th duplicate ACK is received and the fast retransmit is triggered at transmission attempt i , a second packet loss occurs at transmission attempt l . Before the second packet loss, $l - 1$ consecutive packets are successfully transmitted. From transmission attempt $l + 1$ to transmission attempt i , exactly $K - l + 1$ packets (not necessarily consecutive) need to be transmitted successfully to trigger the fast retransmit. The relationship $0 < l < K < i$ holds in this case. Note here the multiple loss before fast retransmit is a more strict condition than multiple loss within the same congestion window.

Based on the result of fast retransmit at transmission attempt $i + 1$, we can discuss two further cases:

Retransmit Success

If the fast retransmit is successful at transmission attempt $i + 1$, the TCP fast recovery is getting started, and the congestion window becomes $W'' = \min \left\{ \lceil \frac{Y}{2} \rceil + K + 1, W_{max} \right\}$. (Since one more ACK caused by successful retransmit is received.) This is the result of inflation for fast recovery. It will stop until a new ACK is received. Due to lack of more ACKs to trigger the second retransmit, the retransmit timer for the second lost packet will finally expire at $d_{RTO} + l - t_i$. The transform

function $\Phi_{YX}^{(2)}(z_d, z_t, z_s)$ is

$$P[B(i, l)]T^{d_{RTO}+l-i}z_d^{d_{RTO}+l-i}z_t^{i+1}z_s^{N_s},$$

where $i \leq N_s \leq K + 1$. Here we count the successful retransmission at time $i + 1$.

Retransmit Failure

If the fast retransmit is not successful at transmission attempt $i + 1$, the TCP sender will wait for the ACK until the expiration of retransmit timer for the lost packet. Since no more ACKs come back, the retransmit timer will expire at time $d_{RTO} - t_i$. The TCP state will go to $X(k+1) = (C, \lceil \frac{W'}{2} \rceil, 1)$.

The transform function $\Phi_{YX}^{(2)}(z_d, z_t, z_s)$ is

$$P[B(K, l, i)]T^{d_{RTO}}z_d^{d_{RTO}}z_t^{i+1}z_s^{N_s},$$

where $i - 1 \leq N_s \leq K$.

TCP Tahoe

Without fast recovery, the sampling time slots of TCP Tahoe for the Markov chain model is different from that of TCP Reno. The TCP state space Ω_X of Tahoe is a subset of the TCP state space of Reno, i.e.,

$$\Omega_X = \{(m, W_{th}, 1)\},$$

where $m \in \{0, 1, 2, \dots, M - 1\}$ and $2 \leq W_{th} \leq \lceil W_{th} \rceil$. So the total number of states is $M \cdot 3(\lceil W_{max} \rceil / 2 - 1)$. The intermediate TCP state space Ω_Y of Tahoe is the same as that of Reno.

Regarding the computation of $\Phi_{YX}^{(2)}$, Case 1 of Tahoe is the same as Reno, while Case 2 of Tahoe becomes much easier than Reno. In Case 2, whenever the fast retransmit is triggered by K duplicate ACKs in Tahoe, the slow start threshold will be halved and the congestion window will be set to 1. The success or failure of the retransmission doesn't need to be memorized by current TCP state. Therefore, the TCP state will go to $X(k+1) = (C, \lceil \frac{W}{2} \rceil, 1)$, and the transform function

$\Phi_{YX}^{(2)}(z_d, z_t, z_s)$ is

$$P[B(K, i, l)]z_d^{d_{c,i}}z_t^i z_s^{N_s},$$

where $0 \leq N_s \leq K$.

4.6 Simulation and Numerical Results

4.6.1 Simulation Setup and Configuration

We use the network simulation tool NS2 to perform the simulations. The NS2 simulator is developed by LBNL, UC Berkeley, and USC VINT project, with wireless extensions from the CMU Monarch project. NS2 is an event-driven network simulator embedded into the OTCL (object oriented version of TCL) Language. An extensible simulation engine is implemented in C++ and is configured and controlled via an OTCL interface. A simulation is defined by an OTCL script. NS2 simulator is not used to reproduce the accurate behavior of a specific TCP variant. It is only used to study the effect of adaptive modulation and congestion/flow control mechanisms on steady TCP performance.

The NS2 simulator supports a multi-state Markov chain (MultiState) error model derived from the ErrorModel class whose parent class is the Connector base class. The Connector base class is for a link, so the MultiState loss module does not have the ability to attach the loss process to an individual TCP flow. If we want to support multiple TCP flows over the same link, we need to extend the original MultiState error model. With the MultiState error model enabled, the link will be in one of a set of channel states. The MultiState module supports a separate loss module for each channel state, with fixed state sojourn times. The unit of error can be specified in term of packet, bits, or time-based. Packet is used as the unit of error in our simulations. The transition state model matrix and an initial state have to be defined to control the activation of these loss modules or channel states. To support our simulations, three files, errmodel.cc, errmodel.h, and

ns-errmodel.tcl, have been changed to simulate the packet loss in a Rayleigh fading wireless channel.

The Rayleigh fading channel can be fully characterized by the maximum Doppler shift, f_d , and the average signal power, ρ . The Rayleigh fading envelope coefficients for the simulated link are sampled every frame. The center carrier frequency is 2.4 GHz. We also do these simulations in different scenarios of average Signal Noise Ratio ρ and Doppler shift f_d . The threshold values for the adaptive modulation are $A_1 = -\infty, A_2 = 12\text{dB}, A_3 = 22\text{dB}, A_4 = 28\text{dB}, A_5 = +\infty$. Three different average Signal Noise Ratios, 25dB, 30dB, and 35dB, are simulated. The values of f_d are taken as 10 Hz, 20 Hz, and 30 Hz (At a carrier frequency of 2.4 GHz, the moving speed of user can be 2.8125 mph, 5.625 mph, and 8.4375 mph respectively). Four modulation schemes are used: BPSK, QPSK, QAM-16, and QAM-64. No coding scheme is discussed in this chapter.

The Level Crossing Rates for different channels are shown in Table 4.1.

Speed(MPH)	ρ (dB)	N_1	N_2	N_3	N_4
2.8125	25	0	9.17142	7.98630	0.37365
2.8125	30	0	5.91137	10.6217	6.46814
2.8125	35	0	3.47076	7.84180	10.7505
5.625	25	0	18.3429	15.9726	0.74729
5.625	30	0	11.8227	21.2434	12.9363
5.625	35	0	6.94152	15.6835	21.5009
8.4375	25	0	27.5143	23.9589	1.12094
8.4375	30	0	17.7340	31.8651	19.4044
8.4375	35	0	10.4123	23.5254	32.2514

Table 4.1: Level crossing rate N_i ($1 \leq i \leq M = 4$) (crossings/second)

The stationary probability distribution of channel state s_i ($1 \leq i \leq M$) used in our simulations can be found in Table 4.2.

ρ (dB)	s_1	s_2	s_3	s_4
25	0.1808813	0.5351597	0.2773006	0.0066584
30	0.0611464	0.2672635	0.4666204	0.2049697
35	0.0197549	0.0985355	0.2758986	0.6058110

Table 4.2: Stationary probability distribution of channel state s_i ($1 \leq i \leq M = 4$)

Supposing that the packet length is 1000 Bytes, the Doppler frequency is 10 Hz, and the average received SNR is 30 dB, the channel state transition probability matrix T then is

$$T = \begin{bmatrix} 0.922660 & 0.077340 & 0 & 0 \\ 0.017695 & 0.950512 & 0.031793 & 0 \\ 0 & 0.018210 & 0.970701 & 0.011089 \\ 0 & 0 & 0.025245 & 0.974755 \end{bmatrix}.$$

If only the Doppler frequency becomes 20 Hz, then the matrix T changes to

$$T = \begin{bmatrix} 0.845319 & 0.154681 & 0 & 0 \\ 0.035389 & 0.901023 & 0.063588 & 0 \\ 0 & 0.036421 & 0.941400 & 0.022179 \\ 0 & 0 & 0.050491 & 0.949509 \end{bmatrix}.$$

If the Doppler frequency changes to be 30 Hz, then the matrix T is

$$T = \begin{bmatrix} 0.767979 & 0.232021 & 0 & 0 \\ 0.053084 & 0.851535 & 0.095381 & 0 \\ 0 & 0.054631 & 0.912101 & 0.033268 \\ 0 & 0 & 0.075736 & 0.924264 \end{bmatrix}.$$

Table 4.3 shows all coefficients η and ξ for modulation schemes used in this chapter to calculate the average Bit Error Probability based on formula in equation (4.5). The related average Packet Error Probability is based on formula given in equation (4.5).

	BPSK	QPSK	QAM-16	QAM-64
η	1	2	3/2	7/4
ξ	2	2	8/5	32/63

Table 4.3: Coefficients for Bit Error Probability $f_{e,i}(\alpha)$

In the simulation, a simple scenario with a sending node and a receiving node connected by a wireless link is simulated. an FTP application is attached with the TCP agent at the sending

node. The timeout timer is 500ms. Other TCP parameters are set to be their default values in NS2. We ran the simulations for time long enough so that the TCP can reach steady state. Some parameters for TCP protocol, for example, the packet size L , the advertised window size W_{max} , and the fast retransmit threshold K , are changed to observe how TCP throughput will be affected by them under different link condition.

4.6.2 Simulation Results and Analysis

In all figures of this section, each symbol represents the average of many simulation results with different random seeds. Some symbols are superposed because the simulation results are very close to each other.

Effect of Packet Length

Fig.4.7, Fig.4.8, and Fig.4.9 show the long-term average throughput through long time simulations and the theoretical average throughput with different packet length. The maximum Doppler shifts for these three figures are 10 Hz, 20 Hz, and 30 Hz, respectively. Most simulation results (lines) are close to the theoretical results (symbols), which validate our analytical model.

As you can find from these figures, the packet length is not directly proportional to the throughput. Before a threshold of packet length that achieves the best throughput, TCP throughput increases as the packet length becomes larger. However, after the threshold, the larger the packet length, the lower the TCP throughput.

If the packet length is smaller, the packet loss probability will be lower, and the chance to incur retransmit and timeout will be less frequent. Therefore, it is more likely that the TCP maintains big congestion window. However, this doesn't necessarily mean that TCP throughput will increase as the TCP packet length decrease. Since the TCP throughput of payload portion also depends on extra overhead introduced by protocol (TCP, IP, and MAC) headers. Small packet length means large extra header overhead. There is a tradeoff between packet length and throughput.

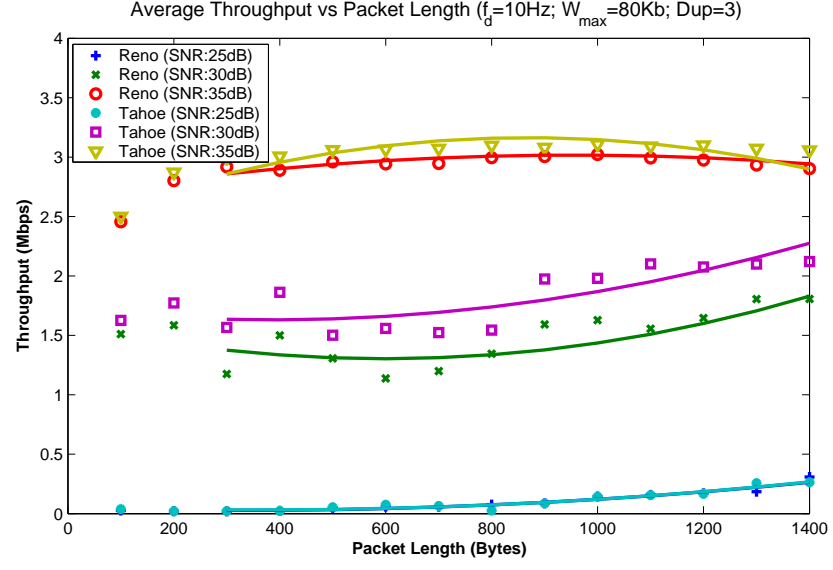


Figure 4.7: Average throughput vs packet length ($f_d=10\text{Hz}$; $W_{\max}=80\text{Kb}$; Dup=3)

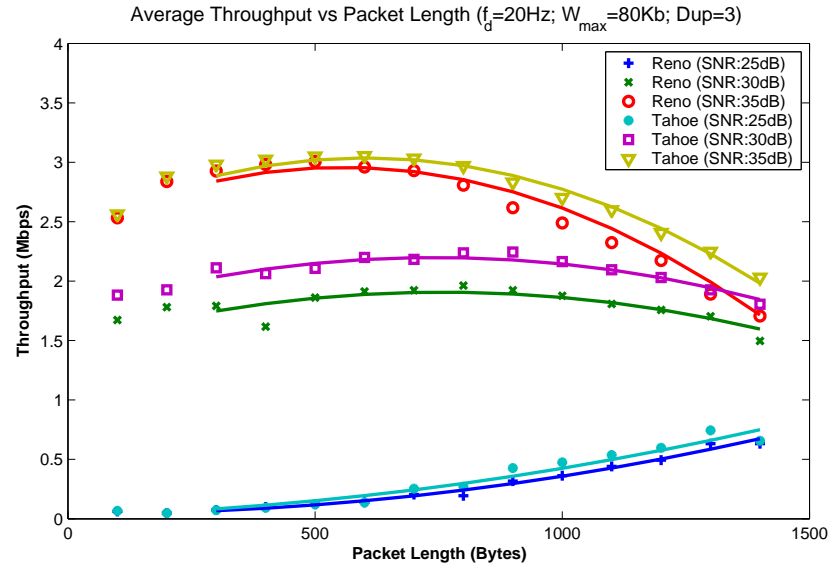


Figure 4.8: Average throughput vs packet length ($f_d=20\text{Hz}$; $W_{\max}=80\text{Kb}$; Dup=3)

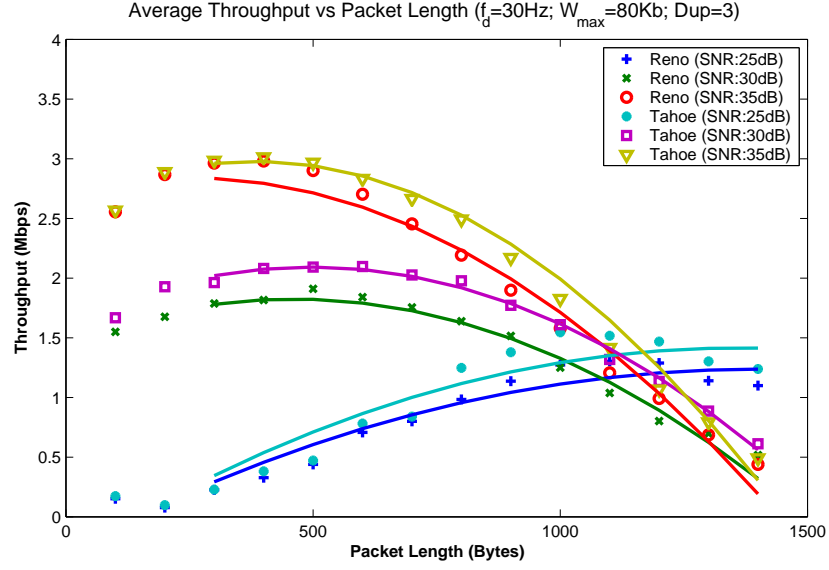


Figure 4.9: Average throughput vs packet length ($f_d=30\text{Hz}$; $W_{max}=80\text{Kb}$; $\text{Dup}=3$)

If the packet length is larger, the packet loss probability will be higher. The resulting frequent retransmit and timeout will keep the average congestion window of TCP being small. In addition, large packet size makes extra protocol header overhead negligible. Thus, it is more likely the TCP throughput will decrease as the packet length increases after a threshold.

Effect of SNR α

It is not a surprise that TCP in fading channels with higher Signal-Noise-Ratio (SNR) performs better, as shown in Fig.4.7, Fig.4.8, and Fig.4.9. With higher SNR, the average time that channel spends in high-rate state is longer, and there will be fewer packet losses due to fading. As a result, the average window size of TCP will be larger, and more packets will be transmitted successfully.

Effect of Doppler Spread f_d

It is found that Doppler spread has a big impact on TCP performance. From Fig.4.10, Fig.4.11, and Fig.4.12, we see that a larger Doppler spread leads to a lower TCP throughput if

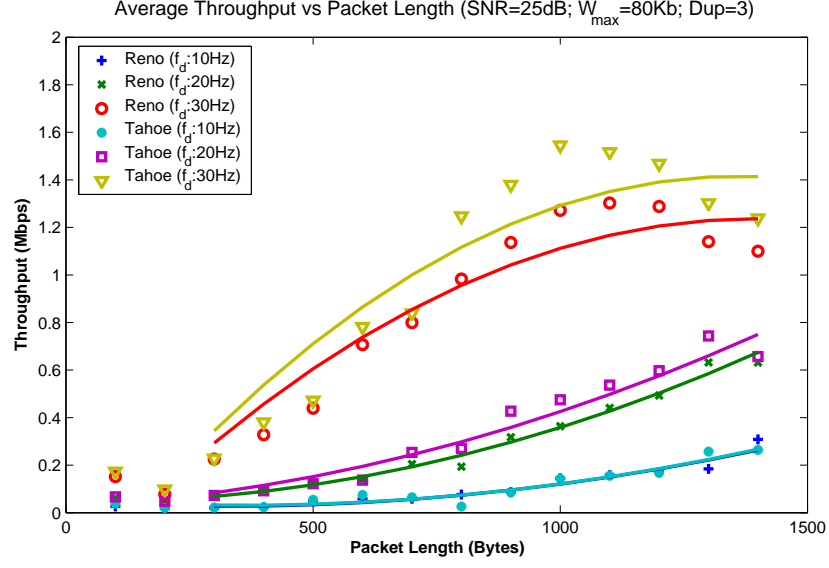


Figure 4.10: Average throughput vs Doppler spread (SNR=25dB; W_{max} =80Kb; Dup=3)

other TCP parameters are the same. Since larger Doppler spread causes more frequent transitions of channel states and modulation schemes, it is hard for TCP congestion control mechanisms to accommodate the dynamic link characteristics. In addition, large Doppler spread increases the transition probability between different channel states, and thus reduces the average time between two corrupted packets (in bad channel states). Accordingly, as the Doppler spread becomes larger, TCP congestion control mechanisms will be invoked more frequently, which will causes more serious TCP performance impairment, as shown in Fig. Fig.4.10, Fig.4.11, and Fig.4.12.

One more interesting finding is that, as the Doppler spread becomes larger, the packet length that achieves the best throughput becomes smaller, if other TCP parameters and the average SNR are the same. The longer the packet length, the higher the packet error rate is. Fast fading with large Doppler spread as well as high packet error rate pushes the packet length for the best TCP throughput to the left side, as shown in Fig.4.10, Fig.4.11, and Fig.4.12. Therefore, in practical applications, the selection of a good packet length should take the Doppler spread into account.

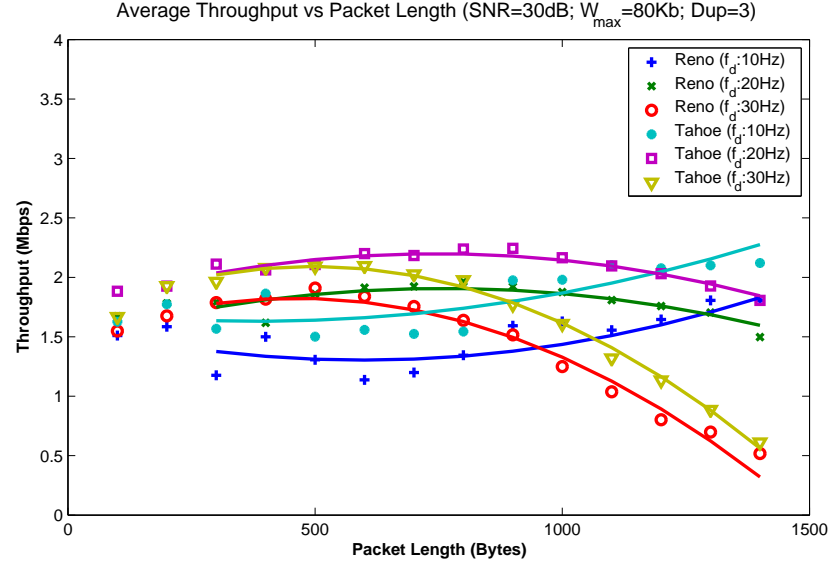


Figure 4.11: Average throughput vs Doppler spread (SNR=30dB; W_{max} =80Kb; Dup=3)

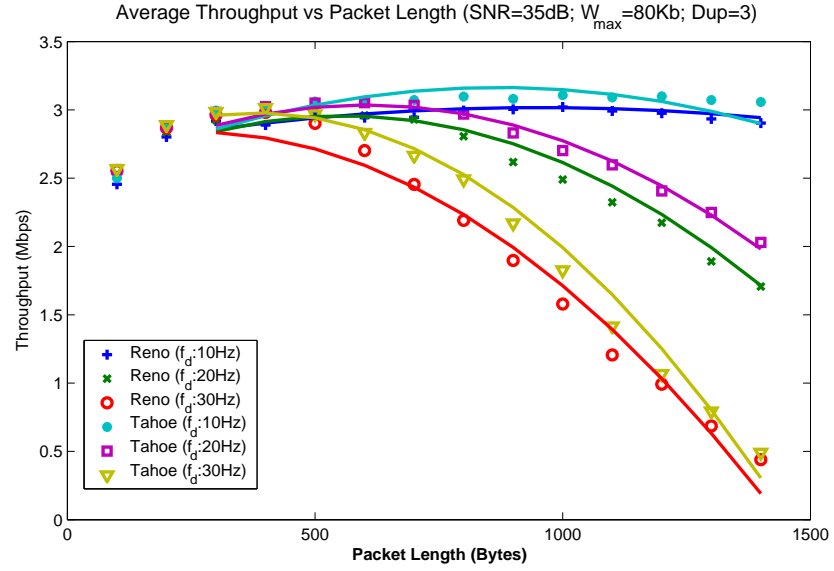


Figure 4.12: Average throughput vs Doppler spread (SNR=35dB; W_{max} =80Kb; Dup=3)

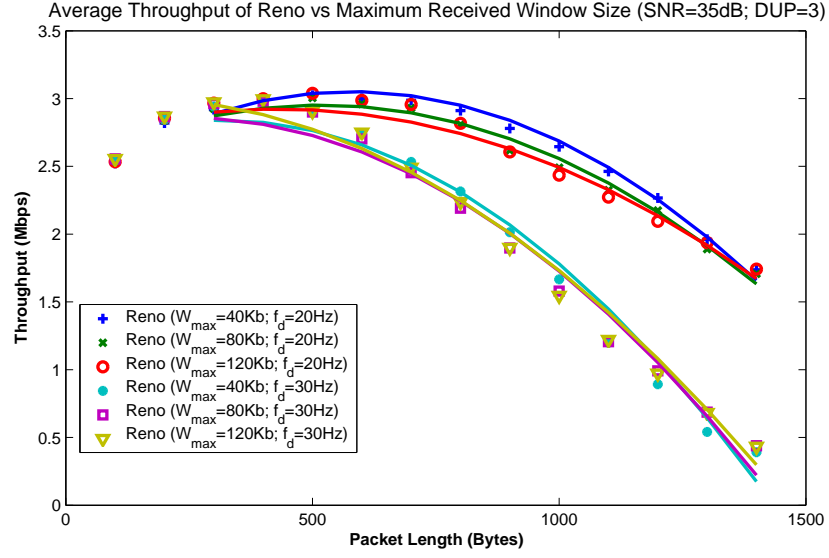


Figure 4.13: Average throughput of Reno vs maximum received window size (SNR=35dB; DUP=3)

Effect of Maximum Advertised Window Size W_{max}

From Fig.4.13 and Fig.4.14, we find that TCP performance in Rayleigh fading channels depends weakly on the value of maximum advertised window size W_{max} . The TCP throughputs in simulations with $W_{max} = 40Kb$, $W_{max} = 80Kb$, and $W_{max} = 120Kb$ are very close to each other. In wired networks, since congestion is the only reason for packet losses, the maximum window size usually should be greater than the bandwidth-delay product of the link to achieve best performance. However, in fading wireless channels, congestion control will be mis-invoked by packet losses in bad channel states. The highly correlated packet errors usually will trigger RTO timeout more frequently and more retransmissions in case of larger W_{max} . In addition, the rate adaptation introduced by adaptive modulation also makes TCP performance less sensitive to the choice of W_{max} . Therefore, in fading wireless channels, with adaptive modulation enabled, the maximum advertised window size W_{max} is not determinant to improve performance.

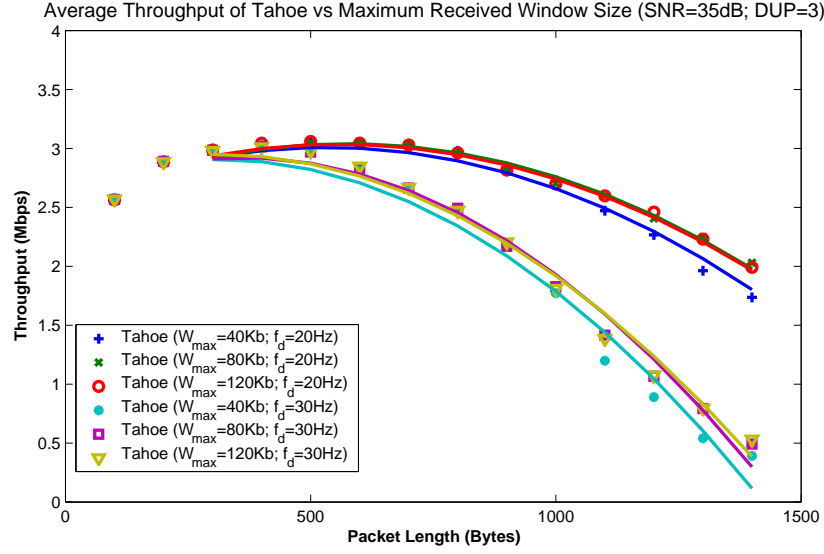


Figure 4.14: Average throughput of Tahoe vs maximum received window Size (SNR=35dB; DUP=3)

4.6.3 Adaptive Configuration of TCP Parameters

As discussed before, in slowly Rayleigh fading channels, average TCP throughput F_t is a function of average SNR α , maximum Doppler shift f_d , packet length l , and maximum advertised window size W_{max} . Since α and f_d come from system inputs, an adaptive control system can be setup to dynamically configure packet length l , and maximum advertised window size W_{max} . The target of such an adaptive control system is to maximize the average TCP throughput for large file transfer.

Numerous techniques have been proposed and used to estimate the Doppler spread. Most of the prior techniques can be categorized into two classes: the techniques based on the level crossing rate (LCR) [88] and the ones based on the covariance of the channel estimates [77]. The performance of these techniques has been extensively studied and discussed in [40]. All these prior techniques have proven to be efficient and robust to the variation of propagation medium provided that the signal-to-ratio (SNR) is high enough.

As shown in Fig.4.15, two input state variables α and f_d will be used. These two variables can be detected or estimated by physical layer, and then be stored in MIB (Management Information

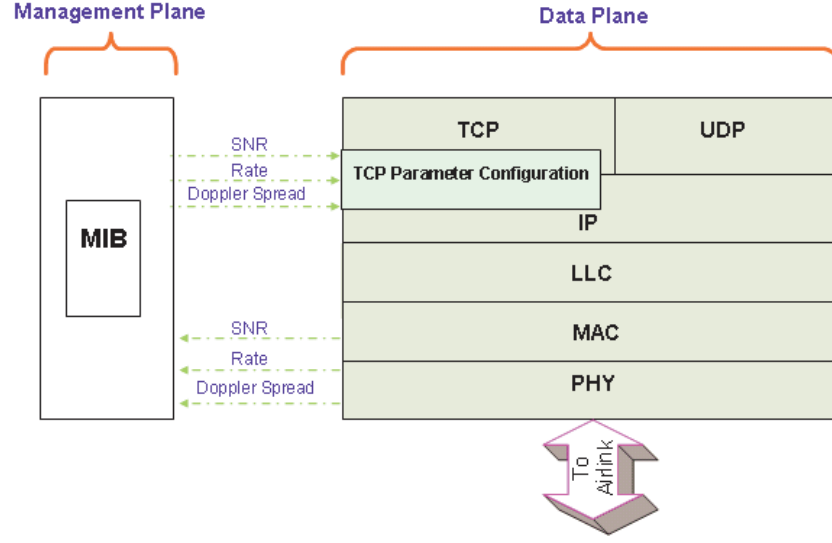


Figure 4.15: Adaptive configuration of TCP parameters

Base) that is accessible to IP layer and transport layer. The controllable variables are l and W_{max} . When TCP service is requested in transport layer, transport layer and IP layer can exert appropriate control actions by setting the TCP header or doing fragmentation using packet length that achieves the best TCP throughput. Previous analytical and simulation results indicate that IP fragmentation will be an efficient approach to improve TCP throughput.

The adaptive configuration of TCP parameters does not change the TCP module and functionality that are defined in the standard or implemented in the operating system. Applications have the same TCP interface as before. Only a piece of glue software between TCP layer and IP layer is required. Our solution maintains backward compatibility with both TCP protocol and TCP applications.

4.7 Conclusion and Future Work

We extended the analysis in [108] to study the TCP throughput with adaptive modulation enabled over slowly Rayleigh fading channels. Although many approximations and bounding tech-

niques are introduced to make the problem tractable, NS2 simulation results show that our model is good enough to accurately predict steady throughput of TCP Tahoe and Reno. The analysis presented in this chapter also provides solid foundation for dynamically adjusting TCP parameters such as maximum advertised window and packet length to accommodate fading wireless links with different condition.

4.8 Proof of the Average Bit Error Probability in State i

The proof is given here: Define

$$F(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{t^2}{2}} dt = 1 - Q(x)$$

to be the cumulative probability distribution function of a standard Gaussian random variable.

$$\begin{aligned} \int_{A_i}^{A_{i+1}} \frac{1}{\rho} e^{-\frac{\alpha}{\rho}} \cdot f_{e,i}(\alpha) d\alpha &= \int_{A_i}^{A_{i+1}} \frac{1}{\rho} e^{-\frac{\alpha}{\rho}} \cdot \eta \cdot Q(\sqrt{\xi \cdot \alpha}) d\alpha \\ &= \eta \cdot \int_{A_i}^{A_{i+1}} \frac{1}{\rho} e^{-\frac{\alpha}{\rho}} \cdot (1 - F(\sqrt{\xi \cdot \alpha})) d\alpha \\ &= \eta \cdot \left[\int_{A_i}^{A_{i+1}} \frac{1}{\rho} e^{-\frac{\alpha}{\rho}} d\alpha - \int_{A_i}^{A_{i+1}} \frac{1}{\rho} e^{-\frac{\alpha}{\rho}} \cdot F(\sqrt{\xi \cdot \alpha}) d\alpha \right]. \end{aligned} \quad (4.8)$$

Now we take the second term

$$\int_{A_i}^{A_{i+1}} \frac{1}{\rho} e^{-\frac{\alpha}{\rho}} \cdot F(\sqrt{\xi \cdot \alpha}) d\alpha = \int_{A_i}^{A_{i+1}} \frac{1}{\rho} e^{-\frac{\alpha}{\rho}} \cdot \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\sqrt{\xi \cdot \alpha}} e^{-\frac{\beta^2}{2}} d\beta d\alpha$$

(If we change the integration order of variable α and β , the integration area is divided into

two parts, as shown in Fig.4.16)

$$\begin{aligned} &= \int_{\sqrt{\xi \cdot A_i}}^{\sqrt{\xi \cdot A_{i+1}}} \left[\int_{\beta^2/\xi}^{A_{i+1}} \frac{1}{\rho} e^{-\frac{\alpha}{\rho}} d\alpha \right] \cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{\beta^2}{2}} d\beta + \int_{-\infty}^{\sqrt{\xi \cdot A_i}} \left[\int_{A_k}^{A_{k+1}} \frac{1}{\rho} e^{-\frac{\alpha}{\rho}} d\alpha \right] \cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{\beta^2}{2}} d\beta \\ &= \int_{\sqrt{\xi \cdot A_i}}^{\sqrt{\xi \cdot A_{i+1}}} \left[e^{-\frac{\beta^2}{\xi \cdot \rho}} - e^{-\frac{A_{k+1}}{\rho}} \right] \cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{\beta^2}{2}} d\beta + \int_{-\infty}^{\sqrt{\xi \cdot A_i}} \left[e^{-\frac{A_i}{\rho}} - e^{-\frac{A_{i+1}}{\rho}} \right] \cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{\beta^2}{2}} d\beta \\ &= \int_{\sqrt{\xi \cdot A_i}}^{\sqrt{\xi \cdot A_{i+1}}} e^{-\frac{\beta^2}{2} \cdot \frac{\xi \cdot \rho + 2}{\xi \cdot \rho}} d\beta - e^{-\frac{A_{k+1}}{\rho}} \cdot \int_{\sqrt{\xi \cdot A_i}}^{\sqrt{\xi \cdot A_{i+1}}} e^{-\frac{\beta^2}{2}} d\beta + \left[e^{-\frac{A_i}{\rho}} - e^{-\frac{A_{i+1}}{\rho}} \right] \cdot \int_{-\infty}^{\sqrt{\xi \cdot A_i}} e^{-\frac{\beta^2}{2}} d\beta \end{aligned}$$

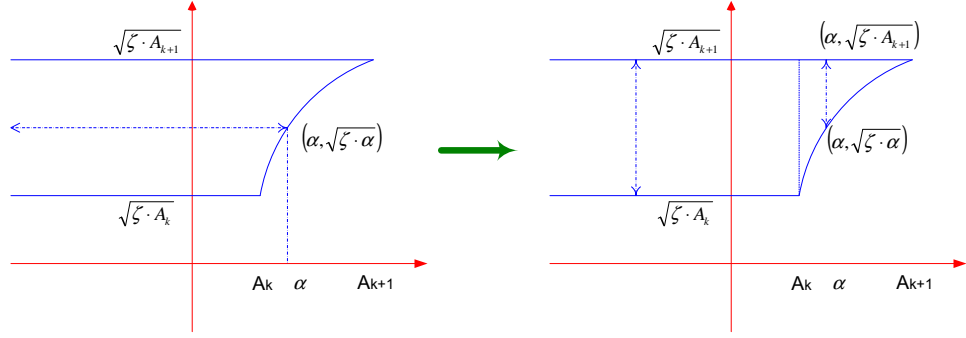


Figure 4.16: Integration along different direction

$$\begin{aligned}
&= \sqrt{\frac{\xi \cdot \rho}{\xi \cdot \rho + 2}} \cdot \int_{\sqrt{\xi \cdot A_i}}^{\sqrt{\xi \cdot A_{i+1}}} e^{-\frac{(\sqrt{\frac{\xi \cdot \rho + 2}{\xi \cdot \rho}} \cdot \beta)^2}{2}} d\left(\sqrt{\frac{\xi \cdot \rho + 2}{\xi \cdot \rho}} \cdot \beta\right) - e^{-\frac{A_{i+1}}{\rho}} \cdot \int_{\sqrt{\xi \cdot A_i}}^{\sqrt{\xi \cdot A_{i+1}}} e^{-\frac{\beta^2}{2}} d\beta \\
&\quad + \left[e^{-\frac{A_i}{\rho}} - e^{-\frac{A_{i+1}}{\rho}}\right] \cdot \int_{-\infty}^{\sqrt{\xi \cdot A_i}} e^{-\frac{\beta^2}{2}} d\beta \\
&= \sqrt{\frac{\xi \cdot \rho}{\xi \cdot \rho + 2}} \cdot \left[F\left(\sqrt{\frac{(\xi \cdot \rho + 2) \cdot A_{k+1}}{\xi \cdot \rho}}\right) - F\left(\sqrt{\frac{(\xi \cdot \rho + 2) \cdot A_k}{\xi \cdot \rho}}\right)\right] - e^{-\frac{A_{k+1}}{\rho}} \cdot F\left(\sqrt{\xi \cdot A_{i+1}}\right) \\
&\quad - F\left(\sqrt{\xi \cdot A_i}\right) + \left[e^{-\frac{A_i}{\rho}} - e^{-\frac{A_{i+1}}{\rho}}\right] \cdot F\left(\sqrt{\xi \cdot A_i}\right) \\
&= \sqrt{\frac{\xi \cdot \rho}{\xi \cdot \rho + 2}} \cdot \left[F\left(\sqrt{\frac{(\xi \cdot \rho + 2) \cdot A_{k+1}}{\xi \cdot \rho}}\right) - F\left(\sqrt{\frac{(\xi \cdot \rho + 2) \cdot A_k}{\xi \cdot \rho}}\right)\right] \\
&\quad - e^{-\frac{A_{k+1}}{\rho}} \cdot F\left(\sqrt{\xi \cdot A_{i+1}}\right) + e^{-\frac{A_i}{\rho}} \cdot F\left(\sqrt{\xi \cdot A_i}\right) \\
&= \left[\sqrt{\frac{\xi \cdot \rho}{\xi \cdot \rho + 2}} \cdot F\left(\sqrt{\frac{(\xi \cdot \rho + 2) \cdot A_{k+1}}{\xi \cdot \rho}}\right) - e^{-\frac{A_{k+1}}{\rho}} \cdot F\left(\sqrt{\xi \cdot A_{i+1}}\right)\right] \\
&\quad - \left[\sqrt{\frac{\xi \cdot \rho}{\xi \cdot \rho + 2}} \cdot F\left(\sqrt{\frac{(\xi \cdot \rho + 2) \cdot A_k}{\xi \cdot \rho}}\right) - e^{-\frac{A_i}{\rho}} \cdot F\left(\sqrt{\xi \cdot A_i}\right)\right] \tag{4.9}
\end{aligned}$$

Substitute the result from equation (4.9) into equation (4.8), and then we have:

$$\int_{A_i}^{A_{i+1}} \frac{1}{\rho} e^{-\frac{\alpha}{\rho}} \cdot f_{e,i}(\alpha) d\alpha = \eta \cdot (\gamma_i - \gamma_{i+1})$$

if we define

$$\gamma_i = \sqrt{\frac{\xi \cdot \rho}{\xi \cdot \rho + 2}} \cdot F\left(\sqrt{\frac{2 \cdot A_i \cdot (\xi \cdot \rho + 2)}{\xi \cdot \rho}}\right) + \exp\left\{-\frac{A_i}{\rho}\right\} \cdot \left[1 - F\left(\sqrt{\xi \cdot A_i}\right)\right]$$

4.9 A Recursive Approach to Derive the Average Delay

Assume an initial source channel state i ($1 \leq i \leq M$), a destination channel state j ($1 \leq j \leq M$), and n channel state transitions between i and j . The channel state transition probability matrix is T , as discussed in section 4.3.3, and the delay for each channel state i ($1 \leq i \leq M$) is d_i , as discussed in section 4.3.3. We are interested in the average delay $d(i, j, n)$ from channel state i to channel state j after exact n state transitions. In section 4.4.3, a formula based on transform function is given. However, the actual calculation of $d(i, j, n)$ can be simplified by using the recursive approach presented here.

Suppose the intermediate channel states from state i to state j are c_2, c_3, \dots , and c_{n-1} . By definition, the average delay is

$$\begin{aligned} d(i, j, n) &= \sum_{i, c_2, \dots, c_{n-1}, j} [T_{i, c_2} T_{c_2, c_3} \cdots T_{c_{n-2}, c_{n-1}} T_{c_{n-1}, j} \cdot (d_i + d_{c_2} + d_{c_3} + \dots + d_{c_{n-1}} + d_j)] \\ &= \sum_{l=1}^M \sum_{i, c_2, \dots, l, j} [T_{i, c_2} T_{c_2, c_3} \cdots T_{c_{n-2}, l} T_{l, j} \cdot (d_i + d_{c_2} + d_{c_3} + \dots + d_l + d_j)] \\ &= \sum_{l=1}^M T_{l, j} \left\{ \sum_{i, c_2, \dots, c_{n-2}, l} [T_{i, c_2} T_{c_2, c_3} \cdots T_{c_{n-2}, l} \cdot (d_i + d_{c_2} + d_{c_3} + \dots + d_l)] \right. \\ &\quad \left. + d_l \cdot \sum_{i, c_2, \dots, l} [T_{i, c_2} T_{c_2, c_3} \cdots T_{c_{n-2}, l}] \right\}. \end{aligned}$$

Now we define the cumulative probability $X(i, j, n)$ from initial channel state i to final channel state j after n transmission attempts as

$$X(i, j, n) = \sum_{i, c_2, \dots, c_{n-1}, j} [T_{i, c_2} T_{c_2, c_3} \cdots T_{c_{n-1}, j}].$$

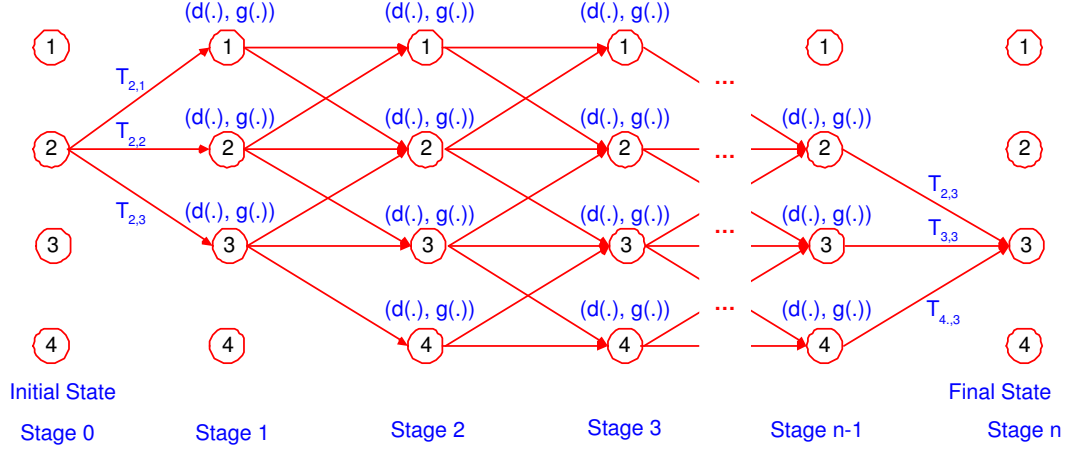


Figure 4.17: A recursive approach to calculate the average delay

Then the average delay $d(i, j, n)$ can be rewritten as

$$d(i, j, n) = \sum_{l=1}^M [T_{l,j} \cdot d(i, l, n-1) + d_l \cdot X(i, l, n-1)]. \quad (4.10)$$

Similar to dynamic programming or Viterbi decoding algorithm, equation (4.10) may be used to recursively calculate the average delay.

The computation can be divided into stages based on the number of transmission attempts, as shown in Fig.4.17. Each stage has a number of channel states associated with it. For any channel state at each stage, we store two metrics: the average delay $d(i, j, n)$ and the average probability $X(i, j, n)$. The transition between channel states at adjacent stages follows the finite state Markov channel model and its associated transition probability matrix. For stage n , we only need two metrics of all channel states at previous stage $n-1$, based on equation 4.10. We begin by computing all $d(i, ., n)$ and $X(i, ., n)$. Then we need to compute all $d(i, ., n-1)$ and $X(i, ., n-1)$, continuing to work backward in this fashion until all $d(i, ., 1)$ and $X(i, ., 1)$ have been computed based on the channel state transition probability matrix.

However, forward recursions are used to calculate the average delays in practical case. We begin by calculating $d(i, ., 1)$ and $X(i, ., 1)$ for all possible states at stage 1. Then we calculate $d(i, ., 2)$ and $X(i, ., 2)$ for all possible states at stage 2, and so on. Note that the computation needs

to be done only once and stored in a table. Therefore, explicit exhausted enumeration of all paths from the initial state to the final state can be avoided. Thus the computation effort for the average delay is significantly reduced by the forward recursive technique.

Chapter 5

Exploiting Multi-user Diversity for Improvement of TCP Performance over Fading Channels

5.1 Introduction

Until recently, wireless networks were built in an environment where voice traffic dominated. However, data traffic has grown dramatically over the last decade. The wireless industry is undergoing an unprecedented wave of innovation. As a result of this rapid evolution, many wireless standards are currently in deployment and in various stages of ongoing development, including IEEE 802.11 wireless LANs (WiFi), IEEE 802.16 wireless MANs (WiMAX) and lots of their variants [17]. The transmission rates of new wireless networks based on these standards will be significantly higher than those derived from voice-oriented and circuit-switching wireless networks such as GPRS and CDMA2000. For example, the highest rate of IEEE 802.11a/g is 54 Mbps, IEEE 802.11n may work at rate of more than 100 Mbps, and IEEE 801.16 is intended to support data rate of 70 Mbps. The

new wireless networks help to alleviate the limits on many bandwidth-intensive applications imposed by current low-rate local access infrastructure.

The convenience of wireless networks such as untethered connection, mobility support, low cost, fast setup, and flexible configuration, is built on a scarce resource, the radio spectrum. Since wireless communications occur in the public space, it is also subject to distortion and attenuation by various factors such as terrain, buildings, moving objects, thermal noise, interference from other channels, etc. With today's huge demand for wireless connections, there is a seriously shortage of radio spectrum. The promise of high data rates for new wireless data services and applications can be realized only if radio spectrum is used in an efficient way.

We are interested in service quality of data wireless communications to mobile users, where the time-varying nature of multipath fading causes amplitude and phase fluctuations, and time delay in the received signals. To mitigate the effect of multipath fading caused by user's mobility in high-speed wireless networks, many forms of diversity that benefit from independent fading channels have been utilized. Successful examples include spatial diversity and multi-user diversity [2, 31, 52, 56, 89, 90]. Spatial diversity increases the performance of a single wireless link by use of multiple antennas between the transmitter and the receiver, while multi-user diversity increases the aggregate performance of many independently fading links.

Most multi-user diversity research seeks to increase the UDP throughput or PHY throughput in saturation state. However, UDP is much simpler than TCP, another basic component of TCP/IP protocol suite, which has been widely used by many classical TCP/IP applications such as WWW, FTP, and Email, etc. In this chapter, Multi-user diversity is used to combat multipath fading to achieve better TCP performance in high-speed wireless networks. We propose a simple modification of the widely used Proportional Fair scheduling algorithm, and then we make the TCP receive window aware of channel states. Simulations demonstrate our revisions lead to performance improvement. The main target system is IEEE 802.16e wireless MANs, which provide efficient basic mechanisms to support our solution.

5.2 System and Channel Model

5.2.1 System Model

As shown in Fig.5.1, a wireless-cum-wired network scenario is considered. We focus on wireless subsystem within the dotted frame. The subsystem is an IEEE 802.16e wireless MAN with a Base Station (BS) and N active Subscriber Stations (SS). BS is connected to the Internet via a high speed wired link. Uplink (from SS to BS) transmission is addressed. However, since transmission scheduling in base station of IEEE 802.16e is bidirectional, the results may apply to downlink transmission as well.

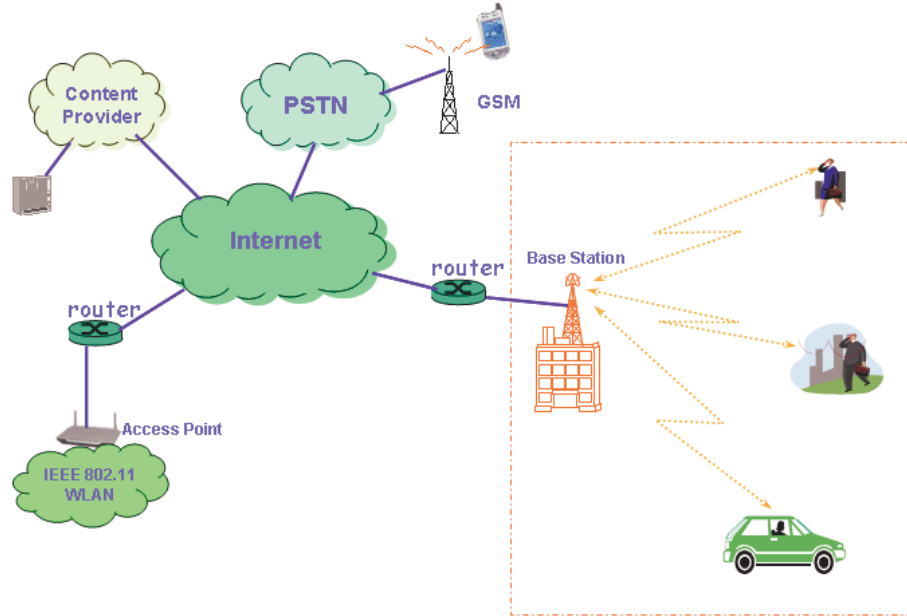


Figure 5.1: Network model

In an On-Demand TDMA system like IEEE 802.16e, all users operate on the same channel at any instance, and centralized packet scheduling is performed in the base station. The structure of packet scheduling and buffer management in base station is shown in Fig.5.2. It is assumed that buffering is usually performed on a per-flow basis (each user may have multiple flows). The base station has knowledge of channel conditions and buffer status of all active connections. All packets

in the buffer are classified into multiple flows. Each flow forms a queue.

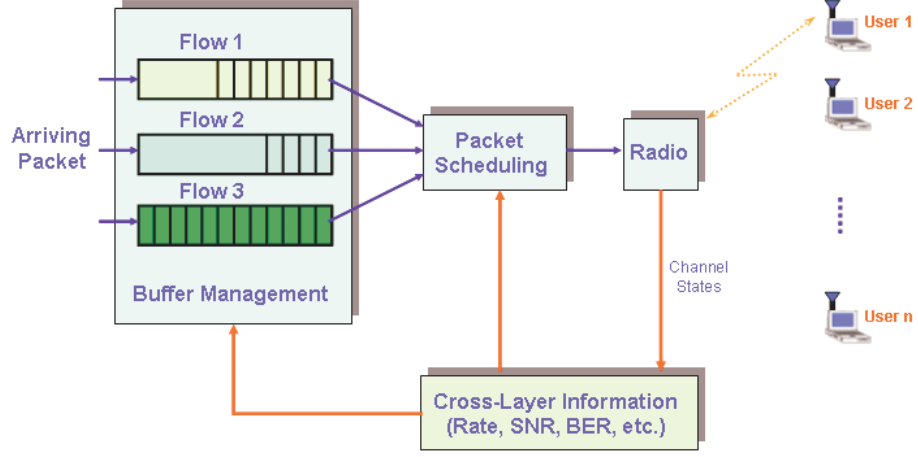


Figure 5.2: Packet scheduling in the buffer of base station

Data transmission in IEEE 802.16e networks occurs in the unit of Physical slot (PS). The length of PS depends on the physical layer specification, and usually is the duration of 4 modulation symbols at the symbol rate of the downlink transmission. The basic unit of an uplink bandwidth allocation is minislot. The length of minislot is equivalent to n physical slots, where $n = 2^m$ and m is an integer ranging from 0 through 7. m is defined by Uplink Channel Descriptor (UCD) message that is transmitted by the base station periodically [20, 48].

In this chapter, we are mainly interested in frame-based Time Division Duplexing (TDD) medium access. A TDD frame has a fixed duration and contains one downlink subframe and one uplink subframe, as shown in Fig.5.3. The downlink subframe comes first, followed by the uplink subframe. The downlink subframe begins with information necessary for frame synchronization and control. Each SS attempts to receive all downlink traffic, except in cases that the burst profile is either not implemented by the SS or is less robust than the SS's current downlink burst profile. The uplink subframe consists of minislots. The partition between uplink and downlink is controlled by MAC layer for bandwidth efficiency, and is vendor-specific.

The usage of the downlink intervals is given in the Downlink Map (DL-MAP) message.

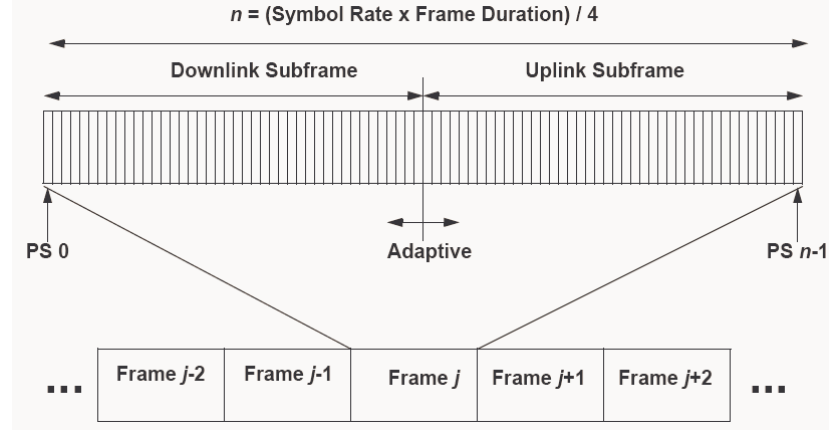


Figure 5.3: Frame structure for time division duplexing

The uplink map (UL-MAP) message defines the usage for the uplink minislots using a series of Information Elements (IEs). Each IE consists of at least three of the following fields: Connection Identifier (CID), Uplink Interval Usage Code (UIUC), Offset from the previous IE start (the length) in numbers of minislots. The IEs are in strict chronological order within UL-MAP messages [20, 48].

The transmission properties of these uplink or downlink subframes are defined in the Burst Profiles with an interval usage code, which defines a set of parameters such as modulation type, forward error correction type, preamble length, guard times, etc. Information about Burst Profiles are conveyed via Downlink Channel Descriptor (DCD) messages and Uplink Channel Descriptor (UCD) messages, which are transmitted by base station regularly. Changing of Burst Profiles can be also explicitly requested by other messages. In this chapter, we only consider a very simple case with four different modulation schemes: BPSK, QPSK, QAM-16, and QAM-64.

In IEEE 802.16e MAC layer, data transmission is connection-oriented. The concept of connection is a unidirectional mapping between BS and SS for the purpose of transporting a service flows traffic. Connections are identified by a connection identifier (CID), which maps to a Service Flow Identifier (SFID). SFID defines the quality of service (QoS) parameters of the service flow provided by upper layers. Note here upper layer protocols such as TCP can be directly mapped to a MAC connection.

5.2.2 Basic Mechanisms for Adaptive Modulation in IEEE 802.16e

Channel quality measurement is mandatory for IEEE 802.16e wireless MANs that operate in bands below 11 GHz (or license-exempt bands). Two metrics, Receive Signal Strength Indicator (RSSI) and Carrier-to-Interference-and-Noise Ratio (CINR), need to be reported. Since the measurement of RSSI doesn't require receiver demodulation lock, it is a reliable metric for channel strength even at low signal level. CINR shows the operating condition of receiver/channel, i.e., the effects of signal strength, interference, and noise level on the receiver. RSSI and CINR are measured message by message. The mean and the standard deviation statistics of RSSI and CINR should be derived and recorded in units of dB (in case of signal strength it is dBm) after every measurement. An SS needs to report the mean and standard deviation of RSSI and CINR only when solicited via REP-RSP messages [20, 48].

Most previous work on Adaptive Modulation focuses on a point-to-point link [4, 16, 28, 29]. Thus the channel feedback information can be obtained on time. However, IEEE 802.16e defines a broadband wireless point-to-multipoint (LMDS) communication system. Previous assumptions and approaches for Adaptive Modulation need to be revised.

The channel status feedback approach in IEEE 802.16e is different from that in a point-to-point wireless link. In 802.16e, the channel status information is transmitted via MAC layer messages but not PHY layer messages in a point-to-point wireless link. The Adaptive Modulation of downlinks is different from that of uplinks. No matter what kind of Adaptive Modulation and Coding occurs, it is announced by BS via Interval Usage Code for downlink or uplink at the beginning of a frame. Since BS makes final decision on burst profile transition and BS is also the uplink receiver, the Adaptive Modulation in uplinks takes fewer steps. However, the Adaptive Modulation in downlinks needs information exchange between BS and SS. At least three ways can be used as feedback transmission from SS to BS: Report-Request/Response (RR-REQ/RR-RSP) messages, Ranging-Request/Ranging-Response (RNG-REQ/RNG-RSP) messages, or Downlink Burst Profile Change Request/Response (DBPC-REQ/RES) messages.

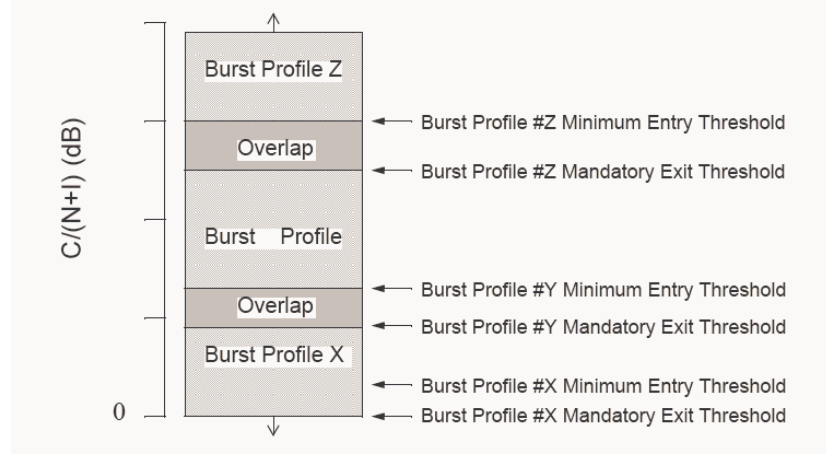


Figure 5.4: Burst profile threshold

Basically the Adaptive Modulation can be implemented in MAC layer using related control frames defined in IEEE 802.16e standard. However, the standard only suggests a framework. It is assumed that there are N Burst Profiles that represent combinations of different Modulation schemes, as shown in Fig. 5.4. For the i th Burst Profile, a Minimum Entry Threshold E_i and a Mandatory Exit Threshold X_i of CINR are defined. It is required that $X_i \leq E_i$. If CINR of received signals is below X_i , it is mandatory to switch the burst profile i to a more robust one. If CINR of received signals is above E_i , switching to burst profile i is recommended.

How to calculate E_i and X_i is not given in the standard. They are vendor-specific or can be configured by ISPs based on their application scenarios. In addition, when and how frequently the AMC should be done is not defined in the standard. These issues are left to vendors for their own optimization.

5.2.3 Channel Model

We use the same channel model as presented in subsection 4.3.3.

5.3 A Channel-Dependent Scheduling Algorithm for TCP Flows

5.3.1 Buffer Management and Packet Scheduling

As discussed in section 5.2, the buffer in the base station consists of many per-user queues. The manipulation of buffer space mainly includes enqueueing and dequeuing activities, which are called buffer management and packet scheduling respectively. Buffer management determines how to allocate buffer space for each queue and how to drop some of received packets based on a pre-determined policy if the space for a queue is not large enough to accommodate all received packets. Packet scheduling determines how to select a packet from the buffer to transmit based on certain policy.

A large amount of research has been performed on buffer management. Typical algorithms for buffer management include First-In-First-Out (FIFO) with Tail-Drop and Random Early Detection (RED) [23]. FIFO drops packets only if buffer overflows. RED tries to control the queue length by randomly dropping packets with a small probability as buffers grow beyond a low threshold. If the buffer continues to grow, packets will be randomly dropped with higher and higher probability. If buffers grow beyond a high threshold, all incoming packets will be dropped.

Since packet scheduling algorithms can significantly improve the Quality of Service (QoS) performance and link utilization, it attracted extensive research in the past. The simplest packet scheduling algorithm is First-In-First-Out. To achieve per-flow max-min fairness, many other algorithms have been proposed. The typical examples are Round Robin and its variants Weighted Round Robin (WRR) and Deficit Round-Robin (DRR) [24, 35, 79, 85]. Based on the idea of Generalized Processor Sharing (GPS) that emulates a fluid-flow system within the constraints of a packet system [71], some other canonical packet scheduling algorithms have been proposed. They are Weighted Fair Queueing (WFQ) [19] and its variants such as Self-Clocked Fair Queueing [32] and Worst-case

Fair Weighted Fair Queueing (WF2Q) [9], etc.

Buffer management has significant impact on packet loss and bandwidth allocation, while packet scheduling algorithm affects bandwidth allocation, delay, and jitter. Although their effects are not the same, they have to be jointly addressed. The packet scheduling algorithm may give transmission opportunity to a particular queue. However, it helps only if the queue is not empty, which is influenced by buffer management. Otherwise, the transmission chance will be wasted.

5.3.2 Proportional Fair Scheduling

The basic assumption of packet scheduling algorithms in wired networks is that the packet transmission from current node to the receiving node will be successful. With this assumption, the packet scheduling algorithms know that the bandwidth resource will be fully utilized once it is assigned to a flow. However, this assumption cannot be directly applied to wireless networks with Adaptive Modulation enabled. Unlike in wired networks, wireless networks over fading channels suffer highly correlated packet errors. Thus, in the short term, the available bandwidth for each flow is not necessarily proportional to its assigned transmission opportunity. When the channel of a flow is in a bad state, the throughput will still be very low even if it gets most of the transmission chances. A proportional fair algorithm that utilize cross-layer parameters provided by physical layer and MAC layer has been proposed to combat the channel-dependent performance degradation [8, 41, 52, 98, 89].

Before the proportional fair algorithm, the base station used two strategies to schedule the transmissions: round robin and best channel condition (SNR). Round robin provides fair transmission opportunity to each user, but it ignores the time-varying channel condition and fails to fully exploit the multi-user diversity among links with independent channel conditions. Thus the overall system throughput is not maximized. If the link with best channel condition is always selected, the system throughput will be maximized. However, the best channel condition strategy favors users close to the base station. As we know, the average signal power is mainly determined by the power

loss model, which means that longer distance causes more power loss and worse channel condition. Users close to the base station will be more likely to reach a particular link rate or higher rates than users far away from the base station at any time slot. Therefore, with the best channel condition strategy, users close to the base station will get more transmission opportunities, and fairness among users is hard to achieve. Tradeoff needs to be made to obtain reasonable system throughput with acceptable fairness among users in practical wireless systems.

To understand more about the proportional fair algorithm, we explain some basic concepts and terms [52]:

- N : The active number of users
- M : The number of data rates supported on each link
- $r_i(t) (1 \leq i \leq M)$: data rate supported on each link. Each data rate is determined by a particular modulation scheme
- t : the index of scheduled transmission time slot
- $\alpha_i(t) (1 \leq i \leq N)$: data rate of user i at time t $\alpha_i(t) \in \{r_1, r_2, \dots, r_M\}$
- $d_i(t) (1 \leq i \leq N)$: service indicator of user i at time t . It depends on the scheduling algorithm. Only one user is selected by the scheduler to transmit. At any time slot t , if user i is chosen to receive service, then $d_i(t) = 1$; otherwise $d_i(t) = 0$
- $\beta_i(t) (1 \leq i \leq N)$: estimated average throughput of user i at time t .
- T_c : the constant time window size used by an exponentially weighted low pass filter to update the average data rate. Only after the time threshold T_c , will the scheduler be able to detect the abrupt transition from a good channel condition to a bad channel condition. Within time T_c , the scheduler will understand any drop in channel condition as temporary deterioration caused by short-term fadings. Since the PF scheduler aims to serve each user at the peak of its channel condition, T_c is related to the maximum period that a user will not be served. A

higher T_c enables the scheduler to wait longer for the improvement of a user's channel condition and improve the overall throughput, but the maximum packet delay time for that user will be longer too. A lower T_c makes the scheduler too sensitive to the channel variation and hence unstable. There is a tradeoff between overall system throughput and packet delay constraint when choosing an appropriate T_c .

The criterion of proportional fairness that considers the variance of long-term channel condition is: if the throughput of a specific user is increased by $x\%$ over what that user receives under the proportional fairness, the aggregate throughput of all other users will be decreased by more than $x\%$. The proportional fair scheduling algorithm aims to maximize

$$\sum_{i=1}^N \log \beta_i(t) \quad (5.1)$$

The strategy adopted by the PF algorithm to achieve this goal is to schedule the user whose data rate $\alpha_i(t)$ is closer to the peak compared to its recent average throughput $\beta_i(t)$, i.e., always serve user i so that

$$i = \arg \max_i \left\{ \frac{\alpha_i(t)}{\beta_i(t)} \right\} \quad (5.2)$$

where we can find that the numerator is related to the throughput improvement and the denominator is relevant to the proportional fairness.

The PF algorithm works as follows:

- initialization:

$N, M, r_i (1 \leq i \leq M), T_c$ are from system information;

$\alpha_i(0)$ is determined by channel condition;

$\beta_i(0) = \alpha_i(0)$;

time slot $j = 0$.

- scheduling:

$d_i(t) = 0, 1 \leq i \leq N$;

for any user who has traffic, select i so that

$$i = \arg \max_i \left\{ \frac{\alpha_i(t)}{\beta_i(t)} \right\}$$

assign time slot j to user i ;

$d_i(t) = 1$ if i is not empty.

- updating:

for all users

$$\beta_i(t+1) = \left(1 - \frac{1}{T_c}\right) \cdot \beta_i(t) + \frac{d_i(t)}{T_c}$$

$j = j + 1$.

The PF algorithm has been reported to achieve excellent simulation performance in MAC layer. The assumption is that the scheduling algorithm will not affect the MAC traffic pattern or the scheduling is independent from MAC traffic pattern [8, 52, 98]. However, for TCP data traffic, the PF algorithm might not work well due to following reasons

- To provide reliable end-to-end connection-oriented transmission service in transport layer, TCP uses mechanisms such as congestion control that requires closed-loop feedback from the other side. The scheduling of an individual packet and resulting transmission sequence of packets will affect the behavior of TCP sender and the timely delivery of remaining packets.
- Since TCP sender sets a timeout timer for each outstanding packet and keeps estimating the round-trip time over all links for a connection, the large variations of intervals between two consecutive transmissions for a particular TCP connection caused by the scheduling algorithm will adversely affect TCP's loss detection mechanism and result in degraded performance.
- TCP is originally proposed for wired networks. It is assumed that the link bandwidth is stable so that TCP sender can explore the available link bandwidth by congestion window size. However, in wireless networks with Adaptive Modulation enabled, the link bandwidth changes as the channel condition changes. However, TCP sender will not understand what happens at

the underlying links. Therefore, the variation of link bandwidths will also negatively impact on TCP performance.

In addition, the TCP is not the only traffic source in transport layer. The UDP traffic is connectionless and independent from the scheduling algorithm. However, if the scheduler (either IP layer or MAC layer) treats them the same, it is likely that UDP traffic will squeeze TCP traffic in the buffer and slow down the TCP. Therefore, intelligent scheduling algorithms need to be developed to achieve satisfactory TCP performance.

5.3.3 A Channel-Dependent Scheduling Algorithm

To address the delay variation caused by bandwidth oscillation and the packet scheduling algorithm, we introduce a delay component for the PF algorithm. It is used to interleave packets as much as possible to avoid the dramatic changes of TCP congestion window. First, several new terms are given:

- $C_i^d(t)$: Current service interval for user i at time t , i.e., the time between two transmission slots for the same user
- $A_i^d(t)$: Average service interval for user i at time t

The objective function becomes:

$$\arg \max_i \left\{ \frac{\alpha_i(t)}{\beta_i(t)} + \mu \cdot \frac{C_i^d(t)}{A_i^d(t)} \right\} \quad (5.3)$$

where μ is a weight coefficient that decides how significantly the packet delay is going to affect the packet scheduling algorithm. The revision of PF algorithm follows

- initialization:

$N, M, r_i (1 \leq i \leq M), T_c$ are from system information;

$\alpha_i(0)$ is determined by channel condition;

$\beta_i(0) = \alpha_i(0)$;

$$C_i(0) = 0;$$

$$A_i(0) = C_i(0);$$

time slot $j = 0$.

- scheduling:

$$d_i(t) = 0, 1 \leq i \leq N;$$

for any user who has traffic, select i so that

$$i = \arg \max_i \left\{ \frac{\alpha_i(t)}{\beta_i(t)} \right\} + \mu \cdot \frac{C_i^d(t)}{A_i^d(t)}$$

assign time slot j to user i ;

$$d_i(t) = 1 \text{ if } i \text{ is not empty.}$$

- updating:

for all users

$$\beta_i(t+1) = \left(1 - \frac{1}{T_c}\right) \cdot \beta_i(t) + \frac{d_i(t)}{T_c}$$

$$C_i(t+1) = (C_i(t) + 1) \cdot (1 - d_i(t))$$

$$j = j + 1.$$

5.4 Maximum Advertised Window Size and Rate Variation

It is assumed that the last-mile wireless link over multipath fading channel is the bottleneck of the communication system for a TCP flow. Multipath fading wireless channel is well known for its time-varying characteristics. As a result, bandwidth variation or rate oscillation occurs in wireless systems with Adaptive Modulation and coding enabled, for example, 802.16e wireless MAN. To fully utilize the available channel capacity without traffic congestion or radio starvation, TCP has to probe current available bandwidth quickly and accommodate rate variation.

Conventional techniques used by TCP to sense available bandwidth are congestion control mechanisms such as slow start, AIMD, fast retransmission, and fast recovery, etc. If current available bandwidth for a TCP connection is stable for enough time, at steady state, TCP congestion window

will finally get to an equilibrium point and present small oscillations like saw tooth around the equilibrium point. These TCP congestion control mechanisms help to alleviate Internet congestion collapse problem. However, they might not be enough to deal with rate variation caused by multipath fading wireless channels and Adaptive Modulation. The main reasons include:

- High rate variation is inherent to wireless links over multipath fading channels. Wireless link rates usually change within seconds, but it takes TCP congestion control at least several RTT (Round-Trip Time, more than ten ms in Internet) to probe the available bandwidth. Actually TCP suffers from slow convergence time. However, fast responsiveness is the dominant factor to keep track of the rate variation of wireless links and efficiently utilize the radio spectrum. Therefore, TCP congestion control is simply too slow to handle the wireless rate changes.
- TCP congestion control takes packet loss as implicit feedback to variation of available bandwidth (congestion or rate changes). Therefore, at least one packet has to be lost before TCP is able to detect and react to bandwidth variation. In other words, TCP congestion window increases steadily until packet loss is detected, which necessitates the subsequent error recovery and adaptation of available bandwidth. Since most Internet congestion is unpredictable and there is no explicit notification of congestion, it is reasonable to use drops as a sign of variation of available bandwidth. In mobile wireless systems over multipath fading channels, the rate variation and wireless channel state are known by systems. We believe that the introduction of intelligent TCP congestion/flow control based on channel state or link rate will reduce the impact of losses on TCP throughput and improve the radio spectrum utilization.

A preliminary attempt to improve TCP throughput using both channel state and multi-user diversity is presented in this section.

The maximum receive window size, W_{max} , is an important tunable parameter in protocol header of TCP packets, as shown in Fig.5.5. TCP window size W , the maximum number of bytes

that can be in the network at any time for a single connection, is determined by

$$W = \min \{W_{cong}, W_{max}\},$$

where W_{cong} is the congestion window dynamically determined by TCP protocol based on update packet transmission.

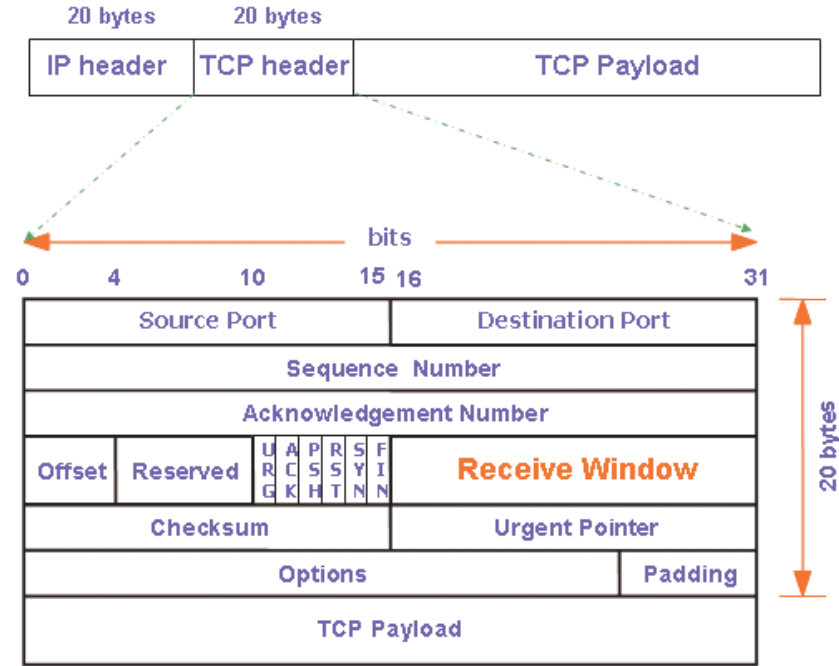


Figure 5.5: Protocol header of TCP packets

Applications determine the initial receive window size W_{max} for a TCP connection at the initial synchronization (the three-way handshake). Afterwards, TCP sender needs to update W_{max} continuously based on the new value in the field of TCP acknowledgement packets sent by the TCP receiver, as shown in Fig.5.6, until the TCP connection is closed.

Originally W_{max} was introduced to avoid the buffer overflow in TCP receiver during a TCP connection. It is the maximum amount of receive data (in bytes) that can be buffered in TCP receiver. Besides the buffer limit and processing speed in TCP receiver, the link between TCP sender and TCP receiver also has a big impact on the selection of appropriate W_{max} . It is preferable

to have a W_{max} so that available bandwidth will be fully utilized and there will be less buffering. For a TCP flow, given receive window size W_{max} and round trip time RTT , its maximal achievable throughput T is constrained by

$$T \leq \frac{W_{max}}{RTT}. \quad (5.4)$$

Assuming that the delay variation of wireless bottleneck link do not significantly change the RTT, the maximum sending rate of TCP source can be controlled by appropriately configuring W_{max} . In mobile wireless networks such as 802.16e, the base station knows the channel state of each link and determines the active modulation scheme. Therefore, the base station can incorporate channel state and link rate into the configuration of W_{max} with little delay.

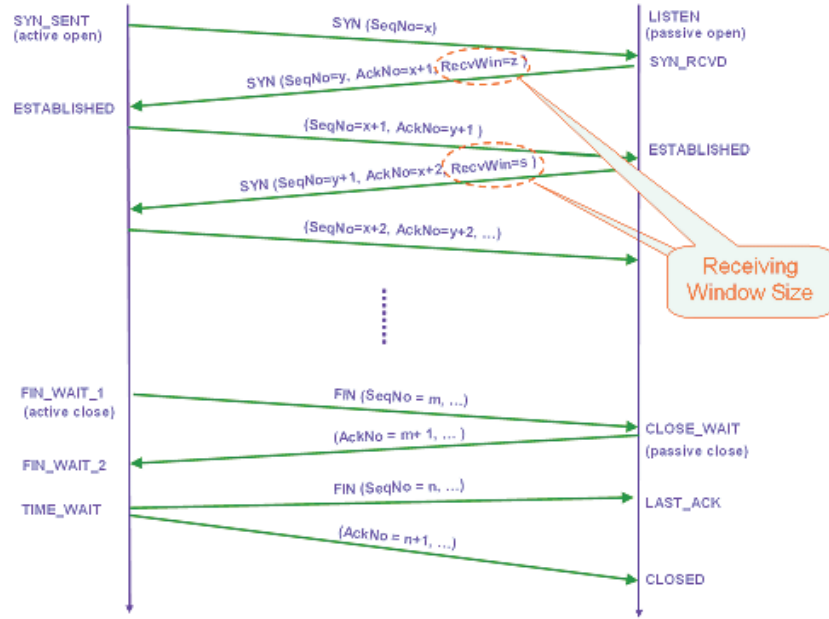


Figure 5.6: Receive window of TCP protocol

At time t , let the number of active connections from mobile user j be $N_j(t)$ ($1 \leq j \leq N$) and the data rate be $\alpha_i(t)$ ($1 \leq i \leq M$). Note that $\alpha_i(t)$ must be one of the set $\{r_1, r_2, \dots, r_M\}$. For each link rate r_k ($1 \leq k \leq M$), we map it to a maximum receive window size of W_{max}^k . Since TCP connections from mobile user j should have the same channel state, the base station can change the

value of W_{max} in TCP header to be $\frac{W_{max}^i}{N_j(t)}$. Here W_{max}^k ($1 \leq k \leq M$) is just an experimental value that depends on the estimation of RTT . However, W_{max}^k is used as an efficient mechanism to force TCP sender to slow down or speed up, based on current channel conditions. Thus the selection of appropriate W_{max} will alleviate the misbehavior of TCP congestion control mechanisms caused by high bit error probability in multipath fading wireless channels.

Sometimes in case of deep fade, it is not good to transmit any TCP packet even if it deserves the radio resource. So we introduce the third revision: if the channel condition is below a threshold, set W_{max} to be 1 in size to freeze the TCP source. When channel condition becomes good, reset W_{max} to recover the normal TCP operation.

5.5 Numerical Results

Since the interaction between multiple TCP flows from more than one users is complicated, it is hard to predict the base station performance mathematically using closed-form analytical expressions. Instead, NS2 simulations have been conducted to investigate the performance of algorithms proposed in this chapter. Despite a wide variety of TCP implementations in NS2, TCP Reno may be the most widely used version and the current de facto. Therefore, it is used in our simulations.

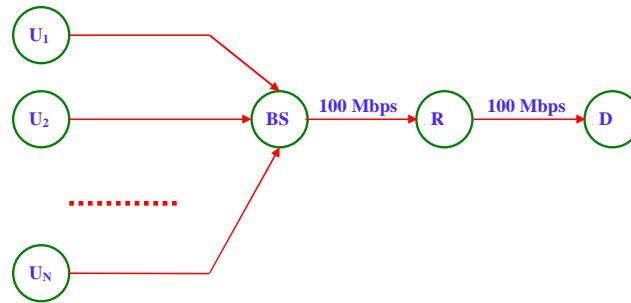


Figure 5.7: Network topology

The general network topology is shown in Fig.5.7. Many users, U_1 , U_2 , ..., and U_N , subscribe to base station BS , which is connected to the backbone network via node R . Each

wireless link from $U_i (1 \leq i \leq N)$ to node BS is exactly like the one described in the last chapter, including the threshold-based policy of Adaptive Modulation. The wireless links are independent of each other. Based on the condition of multipath fading wireless channel, the rate of link from U_i to BS switches between 1Mbps, 2Mbps, 4Mbps, and 6Mbps. The link rate from node BS to node R is 100Mbps, which is large enough to handle all incoming traffic from mobile users. The link rate from node R to traffic destination node D is 100Mbps.

Traffic begins from mobile users $U_i (1 \leq i \leq N)$, and enters the base station BS . Finally all packets come to the sink node D via intermediate router node R . Since it is assumed that each base station can reserve particular bandwidth for TCP traffic, we only consider upstream TCP traffic in our simulations. Actually, an FTP application is attached with the TCP agent at node U_i , and goes to node D .

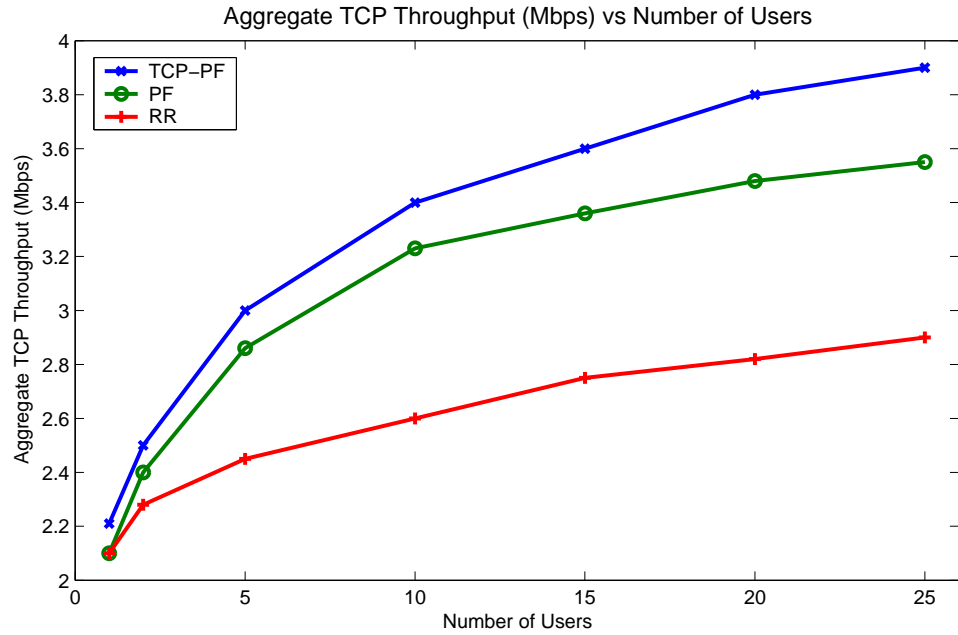


Figure 5.8: Aggregate TCP throughput vs number of users

To make a comparison, Round Robin (RR) algorithm and Proportional Fair (PF) algorithm are taken as the reference. The new algorithm proposed in this chapter is called TCP-PF algorithm.

Fig.5.8 plots the aggregate TCP throughput of upstream traffic in the base station versus number of active users. From Fig.5.8, we observe that the aggregate TCP throughput becomes larger, as the number of active user increases. When there is only one user, PF algorithm achieves the same performance as RR algorithm, while TCP-PF gets better performance than both PF and RR (because of the link-aware W_{max} revision). As the number of users increases, PF algorithm obtains better performance than RR algorithm by utilizing multi-user diversity. Our simulation results also illustrate that the new TCP-PF algorithm generally demonstrates better TCP performance than PF algorithm. Because TCP-PF algorithm is aware of the special congestion/flow control mechanisms of TCP protocol and tries to avoid the mis-invoking of TCP congestion control mechanisms. As shown in Fig.5.8, the larger the number of active users, the bigger gap between RR, PF, and TCP-PF algorithms. More users yield better multi-user gain.

5.6 Concluding Remark

The traditional Proportional Fair algorithm doesn't take into account the interaction between TCP protocol and its lower layer scheduling algorithms. Several revisions of proportional fair algorithm have been proposed in this chapter to improve the aggregate TCP throughput in the base station. Our strategy is to put a delay penalty function in the objective function of PF algorithm and make the TCP parameter W_{max} be aware of wireless channel state and link rate. Simulation results show that our cross-layer design does improve the aggregate TCP throughput and obtains better multi-user diversity.

Chapter 6

Conclusion and Future Work

6.1 Main contribution of this dissertation

The primary goal of this research work is to understand the basic behavior, explore fundamental concepts, and push the state of the art in the direction of high-speed mobile packet wireless networks over multipath fading channels. We concentrate our efforts on some concepts and techniques that promise to provide improvement of spectrum efficiency and system capacity. The main contributions of this research work are summarized below

- Effects of Doppler spread and Delay spread on the UDP performance of several popular high-speed packet wireless networks including IEEE 802.11a and IEEE 802.11b have been presented and analyzed. Since many other packet wireless networks use the same techniques such as OFDM and Spread Spectrum, our simulation results and conclusion may be used as reference for the design of new high-speed packet wireless networks.
- We investigated the effect of mobile speed on the communication range of new generation of packet-based wireless systems that serve users with high mobility in fading channels. Our study shows that high mobile speed will reduce the valid communication range. We need to

pay attention to the mobile speed when we plan such kind of wireless systems or design upper layer protocols that can combat mobile speed and fading.

- We developed a new performance model that is accurate in predicting single TCP flow throughput over Rayleigh fading channels with Adaptive Modulation enabled in the wireless systems. The model captures the key aspects of TCP congestion/flow control, Rayleigh fading channels, and Adaptive Modulation. The analytical and simulation results triggered the idea of cross-layer TCP protocol design for single-user scenarios. The fading parameters of wireless channels detected in physical layer can be used to dynamically tune the parameters of TCP protocol in transportation layer (such as packet length and advertised window size) so that TCP throughput can be improved.
- Considering the multi-user scenarios, we study how multi-user diversity can be used to improve the aggregate TCP throughput of base stations in fading channels. The multi-user diversity gain is achieved through channel-aware packet scheduling algorithms and active delay of TCP ACK packets in the downstream buffer. Based on the Adaptive Modulation information from physical layer, the receive window size of TCP ACK packets will be dynamically changed to accommodate the rate changes resulted from Adaptive Modulation.

6.2 Future work

Based on the experience and knowledge acquired during this research work, several areas for future research can be suggested. They are:

- Effect of fading channel prediction algorithms on these cross-layer protocols over multipath fading channels.
- How will multiple-input multiple-output (MIMO) technology affect the cross-layer protocol design in high-speed packet wireless systems [2, 30, 38]?

- How can multi-user diversity be used by OFDMA (Orthogonal Frequency Division Multiple Access) to allocate subcarriers and power to each user for given Quality of Service requirement [53, 99, 101]? OFDMA is currently the modulation of choice for high speed mobile packet wireless systems such as IEEE 802.16e.

Bibliography

- [1] A. Abouzeid, S. Roy, and M. Azizoglu, Stochastic Modeling of TCP over Lossy Links, Proceedings of IEEE INFOCOM 2000, vol. 3, pp1724-33, March 2000.
- [2] S. A. Alamouti, Simple Transmit Diversity Technique for Wireless Communications, IEEE Journal on Selected Areas in Communications, vol. 16, No. 8, October 1998.
- [3] J. B. Anderson, T. S. Rappaport, and S. Yoshida, Propagation Measurements and Models for Wireless Communications, IEEE Communications Magazine, vol. 33, no. 1, January 1995, pp42-49.
- [4] F. Babich, Considerations on Adaptive Techniques for Time-division Multiplexing Radio Systems, IEEE Transaction on Vehicular Technology, vol. 48, no. 6, pp1862-1873, November 1999.
- [5] A. R. S. Bahai and B. R. Saltzberg, Multi-Carrier Digital Communications - Theory and Applications of OFDM, Kluwer Academic/Plenum, 1999.
- [6] P. A. Bello, Characterization of Randomly Time-variant Linear Channels, IEEE Transactions on Communication Systems, vol. 11, December 1963.
- [7] P. A. Bello, Measurement of Random Time-Variant Linear Channels, IEEE Transactions on Information Theory, vol. IT-15, pp469-475, July 1969.

- [8] P. Bender, et al., CDMA/HDR: A Bandwidth-Efficient High-Speed Wireless Data Service for Nomadic Users, *IEEE Communications Magazine*, vol. 38, no. 7, pp70-77, July 2000.
- [9] Jon C.R. Bennett and Hui Zhang, WF2Q: Worst-case Fair Weighted Fair Queueing, *Proceedings of IEEE INFOCOM 1996*, pp120-128, San Francisco, CA, 1996.
- [10] G. Bianchi, Performance Analysis of the IEEE 802.11 Distributed Coordination Function, *IEEE Journal on Selected Areas in Communications*, vol. 18, no. 3, March 2000.
- [11] S. Borst, User-Level Performance of Channel-Aware Scheduling Algorithms in Wireless Data Networks, *Proceedings of INFOCOM 2003*, San Francisco, CA, March 2003.
- [12] John A. C. Bingham, ADSL, VDSL, and Multicarrier Modulation, John Wiley & Sons, Inc., 2002.
- [13] L. Brakmo, S. O'Malley, and L. Peterson. TCP Vegas: New Techniques for Congestion Detection and Avoidance, *Proceedings of ACM SIGCOMM94*, pp24-35, August 1994.
- [14] Robert W. Chang, Synthesis of Band-limited Orthogonal Signals for Multichannel Data Transmission, *The Bell Systems Technical Journal*, December 1966.
- [15] D. Chiu, and R. Jain, Analysis of Increase and Decrease Algorithms for Congestion Avoidance in Computer Networks, *Computer Networks and ISDN Systems*, vol. 17, pp1-14, 1989.
- [16] S. T. Chung and A. J. Goldsmith, Degrees of Freedom in Adaptive Modulation: A Unified View, *IEEE Transactions on Communications*, vol. 49, no. 9, September 2001.
- [17] T. Cooklev, *IEEE Wireless Communication Standards*, IEEE Press, New York, NY, 2004.
- [18] J. Schwarz daSilva, S. Mahmoud, Capacity Degradation of Packet Radio Fading Channels, *Proceedings of the sixth symposium on Data communications*, Pacific Grove, California, United States, Nov 27-29, pp96-101, 1979.

- [19] A. Demers, S. Keshav, and S. Shenker. Design and Analysis of a Fair Queuing Algorithm. Proceedings of ACM SIGCOMM, pp1-12, Austin, September 1989.
- [20] C. Eklund, R. B. Marks, K. L. Stanwood, S. Wang, IEEE Standard 802.16: A Technical Overview of the WirelessMAN Air Interface for Broadband Wireless Access, IEEE Communications Magazine, vol. 40, no. 6, June 2002.
- [21] V. Erceg, L. Greenstein, S. Tjandra, S. Parkoff, A. Gupta, B. Kulic, A. Julius, and R. Bianchi, An Empirically Based Path Loss Model for Wireless Channels in Suburban Environments, IEEE Journal of Selected Areas in Communications, vol. 17, no. 7, 1999.
- [22] Kevin Fall, Sally Floyd, Simulation-based Comparisons of Tahoe, Reno and SACK TCP, ACM SIGCOMM Computer Communication Review, vol. 26, no. 3, pp5-21, July 1996.
- [23] S. Floyd and V. Jacobson, Random Early Detection Gateways for Congestion Avoidance, IEEE Transactions on Networking, vol. 1, no. 4, pp397-413, August 1993.
- [24] S. Floyd, V. Jacobson, Link-sharing and Resource Management Models for Packet Networks, IEEE/ACM Transactions on Networking, vol. 3, no. 4, pp365-386, August 1995.
- [25] S. Frattasi, H. Fathi, F. H. P. Fitzek, R. Prasad, M. D. Katz, Defining 4G technology from the users perspective, IEEE Network, vol. 20, no. 1, pp 35-41, Jan.-Feb. 2006.
- [26] R. Ganesh and K. Phalavan, Statistical Modeling and Computer Simulation of Indoor Radio Channel, Proceedings of the IEEE, vol. 138, no. 2, pp153-161, 1991.
- [27] V. K. Garg, Wireless Network Evolution: 2G to 3G, Prentice Hall, Upper Saddle River, NJ, 2002.
- [28] A. J. Goldsmith and S. G. Chua, Variable-Rate Variable-Power MQAM for Fading Channels, IEEE Transactions on Communications, vol. 45, no. 10, pp1218-1230, October 1997.

- [29] A. J. Goldsmith and P. P. Varaiya, Capacity of Fading Channels with Channel Side Information, *IEEE Transactions on Information Theory*, vol. 43, pp1986-1992, Nov. 1997.
- [30] A. Goldsmith, Capacity Limits of MIMO Channels, *IEEE Journal on Selected Areas in Communications*, vol. 21, no. 5, pp684-704, June 2003.
- [31] Andrea Goldsmith, *Wireless Communications*, Cambridge University Press, 2005.
- [32] S.J. Golestani, A Self-clocked Fair Queuing Scheme for Broadband Applications, *Proceedings of IEEE INFOCOM 1994*, Toronto, Canada, pp636-646, 1994.
- [33] J. Greenstein, et al., eds., Channel and Propagation Models for Wireless System Design I, *IEEE Journal on Selected Areas in Communications*, vol. 20, no. 3, April 2002.
- [34] J. Greenstein et al., eds., Channel and Propagation Models for Wireless System Design II, *IEEE Journal on Selected Areas in Communications*, vol. 20, no. 6, August 2002.
- [35] E. L. Hahne, Round-Robin Scheduling for Max-Min Fairness in Data Networks, *IEEE Journal on Selected Areas in Communications*, vol. 9, no. 7, pp1024-1039, September 1991.
- [36] L. Harte, et al., *3G Wireless Demystified*, McGraw-Hill, New York, 2002.
- [37] J. Heiskala, and J. Terry, *OFDM Wireless Lans: a Theoretical and Practical Guide*, Sams Publishing, Indianapolis, IN, USA, 2001.
- [38] Ilan Hen, MIMO Architecture for Wireless Communication, *Intel Technology Journal*, vol. 10, no. 2, May 2006.
- [39] G. Holland and N. H. Vaidya, Analysis of TCP performance over Mobile Ad Hoc Networks, *Proceedings of IEEE/ACM MOBICOM 99*, pp219-230, August 1999.
- [40] J. M. Holtzman and A. Sampath, Adaptive Averaging Methodology for Handoffs in Cellular Systems, *IEEE Transaction on Vehicular Technology*, vol. 44, no. 1, pp59-66, 1995.

- [41] J. M. Holtzman, Asymptotic Analysis of Proportional Fair Algorithm, 12th IEEE International Symposium on Personal, Indoor and Mobile Radio Communications, vol. 2, 30 Sept.-3 Oct., pp33-37, 2001.
- [42] H. Honkasalo, K. Pehkonen, M. T. Niemi, A. T. Leino, WCDMA and WLAN for 3G and beyond, IEEE Wireless Communications, vol. 9, no. 2, pp14-18, April 2002.
- [43] IEEE 802.11 Standard Part II: Wireless LAN Medium Access Control (MAC) and Physical Layer (PHY) Specifications, August, 1999.
- [44] IEEE 802.11a (Supplement to IEEE 802.11 Standard Part II): High-Speed Physical Layer Extension in the 5 GHz Band, September, 1999.
- [45] IEEE 802.11b (Supplement to IEEE 802.11 Standard Part II): High-Speed Physical Layer Extension in the 2.4 GHz Band, September, 1999.
- [46] IEEE 802.11g (Supplement to IEEE 802.11 Standard Part II): Further Higher Data Rate Extension in the 2.4 GHz Band, 2003.
- [47] IEEE 802.15.1 standard Part 15.1: Wireless Medium Access Control (MAC) and Physical Layer (PHY) Specifications for Wireless Personal Area Networks (WPANs), 2005.
- [48] IEEE 802.16, IEEE Standard for Local and Metropolitan Area Networks - Part 16: Air Interface for Fixed Broadband Wireless Access Systems, Oct. 2004.
- [49] IEEE 802.16e (Supplement to IEEE 802.16 Standard Part 16): Physical and Medium Access Control Layers for Combined Fixed and Mobile Operation in Licensed Bands and Corrigendum 1, 2006.
- [50] V. Jacobson, M. J. Karels, Congestion Avoidance and Control, ACM SIGCOMM Computer Communication Review, vol. 18, no. 4, pp314-329, August 1988.
- [51] W. C. Jakes, Jr. Microwave Mobile Communications, John Wiley & Sons, N.Y., 1974.

- [52] A. Jalali, R. Padovani, and R. Pankaj, Data Throughput of CDMA-HDR a High-Efficiency - High Data Rate Personal Communication Wireless System, Proceedings of IEEE Vehicular Technology Conference, vol. 3, pp1854-1858, Tokyo, Japan, May 2000.
- [53] D. Kivanc and H. Liu, Subcarrier Allocation and Power Control for OFDMA, IEEE Conference Record of the Thirty-Fourth Asilomar Conference on Signals, Systems and Computers, pp147-151, 2000.
- [54] D. N. Knisely, S. Kumar, S. Laha, S. Nanda, Evolution of Wireless Data Services: IS-95 to cdma2000, IEEE Communications Magazine, vol. 36, no. 10, pp140-149, 1998.
- [55] A. Kumar, Comparative Performance Analysis of Versions of TCP in a Local Network with a Lossy Link, IEEE/ACM Transactions on Networking, vol. 6, pp485-498, August 1998.
- [56] E. G. Larsson, On the Combination of Spatial Diversity and Multiuser Diversity, IEEE Communications Letters, vol. 8, no. 8, pp517-519, 2004.
- [57] William C. Y. Lee, Mobile Cellular Telecommunications Systems, McGraw-Hill, Inc., New York, NY, 1990.
- [58] William C. Y. Lee, Mobile Communications Engineering, McGraw-Hill, 1998.
- [59] S. Lu, V. Bharghavan, and R. Srikant, Fair Scheduling in Wireless Packet Networks, Proceedings of ACM SIGCOMM'97, Cannes, France, Sept. 1997.
- [60] Stefan Mangold, Sunghyun Choi, Guido R. Hiertz, Ole Klein, and Bernhard Walke, Analysis of IEEE 802.11e for QoS Support in Wireless LANs, IEEE Wireless Communications, Special Issue on Evolution of Wireless LANs and PANs, vol. 10, no. 6, December 2003.
- [61] J. W. McJown and R. L. Hamilton, Jr., Ray Tracing as a Design Tool for Radio Networks, IEEE Network Magazine, vol. 5, no. 6, 1991.

- [62] S. Nanda, K. Balachandran, and S. Kumar, Adaptation Techniques in Wireless Packet Data Services, *IEEE Communications Magazine*, vol. 38, no. 1, pp54-64, January 2000.
- [63] J. W. Mark, W. Zhuang, *Wireless Communications and Networking*, Prentice Hall, Upper Saddle River, New Jersey, 2003.
- [64] R. V. Nee and R. Prasad, *OFDM for Wireless Multimedia Communications*, Artech House Publishers, 2000.
- [65] The ns Manual, <http://www.isi.edu/nsnam/ns/doc/index.html>.
- [66] G. E. Oien, H. Holm, and K. J. Hole, Impact of Channel Prediction on Adaptive Coded Modulation Performance in Rayleigh Fading, *IEEE Transaction on Vehicular Technology*, vol. 53, no. 3, pp759-769, May 2004.
- [67] T. Ojanpera and R. Prasad, An Overview of Third-generation Wireless Personal Communications: A European Perspective, *IEEE Personal Communication Magazine*, vol. 5, no. 6, December 1998.
- [68] J. Padhye, V. Firoiu, D. Towsley, and J. Kurose, Modeling TCP Reno Performance: a Simple Model and Its Empirical Validation, *IEEE/ACM Transactions on Networking*, vol. 8, pp133-145, April 2000.
- [69] K. Pahlavan and A. Levesque, *Wireless Information Networks*, John Wiley and Sons, 1995.
- [70] K. Pahlavan and P. Krishnamurthy, *Principles of Wireless Networks - A Unified Approach*, Prentice Hall, 2002.
- [71] A. K. Parekh and R. G. Gallager, A Generalized Processor Sharing Approach to Flow Control in Integrated Services Networks: The Single-Node Case, *IEEE/ACM Transactions on Networking*, vol. 1, no. 3, pp344-357, 1993.

- [72] Larry Peterson and Bruce Davie, Computer Networks: A Systems Approach, Third Edition, Morgan Kauffmann, 2003
- [73] S. Pilosof, R. Ramjee, D. Raz, Y. Shavitt, and P. Sinha, Understanding TCP fairness over Wireless LAN, Proceedings of INFOCOM 2003, San Francisco, CA, April 2003.
- [74] John Proakis, Digital Communications (4 edition), McGraw-Hill, August, 2000.
- [75] Theodore S. Rappaport and Theodore Rappaport, Wireless Communications: Principles and Practice (2nd Edition), Prentice Hall, December, 2001.
- [76] RFC 793, Transmission Control Protocol, September 1981 (<http://www.faqs.org/rfcs/rfc793.html>).
- [77] A. Sampath and J. Holtzman, Estimation of Maximum Doppler Frequency for Handoff Decisions, Proceedings of IEEE Vehicular Technology Conference, pp859-862, Secaucus, NJ, May 1993.
- [78] Mischa Schwartz, Mobile Wireless Communications, Cambridge University Press, 2005.
- [79] M. Shreedhar and George Varghese, Efficient Fair Queuing Using Deficit Round Robin, IEEE/ACM Transactions on Networking, vol. 4, no. 3, pp375-385, June 1996.
- [80] S. Sibecas, C. A. Corral, S. Emami and G. Stratis, On the Suitability of 802.11a/RA for High-mobility DSRC, Proceedings of VTC 2002, vol. 1, pp229-234, 2002.
- [81] B. Sikdar, S. Kalyanaraman and K. S. Vastola, An Integrated Model for the Latency and Steady-State Throughput of TCP Connections, Performance Evaluation, vol. 46, no. 2-3, pp139-154, September 2001.
- [82] B. Sklar, Rayleigh Fading Channels in Mobile Digital Communication Systems Part I: Characterization, IEEE Communications Magazine, vol. 35, no. 9, pp90-100, September 1997.

- [83] C. Smith and D. Collins, 3G Wireless Networks, McGraw-Hill, 2002.
- [84] Raymond Steele, Chin-Chun Lee, Peter Gould, GSM, cdmaOne and 3G Systems, John Wiley & Sons, 2001.
- [85] D. Stiliadis and A. Verma, Efficient Fair Queuing Algorithms for Packet-switched Networks, IEEE/ACM Transactions on Networking, vol. 6, no. 2, pp175-185, April 1998.
- [86] Jun-Zhao Sun, J. Sauvola, D. Howie, Features in Future: 4G Visions from a Technical Perspective, Proceedings of IEEE GLOBECOM 2001, vol. 6, pp3533-3537, San Antonio, TX, 25-29 Nov. 2001.
- [87] E. Teletar, Capacity of Multi-antenna Gaussian Channels, AT&T-Bell Labs, Technical Report, June 1995.
- [88] C. Tipedelenlioglu, A. Abdi, G. B. Giannakis, and M. Kaveh, Estimation of Doppler Spread and Signal Strength in Mobile Communications with Applications to Handoff and Adaptive Transmission, Wireless Communications and Mobile Computing, vol. 1, pp221-242, August 2001.
- [89] David Tse and Pramod Viswanath, Fundamentals of Wireless Communication, Cambridge University Press, 2005.
- [90] Z. Tu and R. S. Blum, Multiuser Diversity for a Dirty Paper Approach, IEEE Communication Letters, vol. 7, pp370-372, August 2003.
- [91] N. H. Vaidya, P. Bahl, and S. Gupta, Distributed Fair Scheduling in a Wireless LAN, Proceedings of MobiCom 2000, Boston, MA, USA, August 2000.
- [92] Andrew J. Viterbi, CDMA: Principles of Spread Spectrum Communication, Addison Wesley Longman Publishing Inc., Redwood City, CA, 1995.

- [93] J. Walrand and P. Varaiya, High Performance Communication Networks, Morgan Kaufmann, San Francisco, 2000.
- [94] B. Walke, P. Seidenberg, M. Althoff, UMTS: The Fundamentals, John Wiley, 2004
- [95] H. S. Wang and N. Moayeri, Finite-state Markov Channel-A Useful Model for Radio Communication Channels, IEEE Transactions on Vehicular Technology, February 1995.
- [96] W. T. Webb and R. Steele, Variable Rate QAM for Mobile Radio, IEEE Transactions on Communications, vol. 43, no. 7, pp2223-2230, July 1995.
- [97] S. B. Weinstein, P. M. Ebert, Data Transmission of Frequency Division Multiplexing Using The Discrete Frequency Transform, IEEE Transactions on Communications, COM-19(5), pp623-634, October 1971.
- [98] C. Westphal, Monitoring Proportional Fairness in cdma2000 High Data Rate Networks, Proceedings of Globecom 2004, vol. 6, pp3866-3871, November 2004.
- [99] C. Y. Wong, R. S. Cheng, K. B. Letaief, and R. D. Murch, Multiuser OFDM with Adaptive Subcarrier, Bit and Power Allocation, IEEE Journal on Selected Areas in Communications, vol. 17, no. 10, pp1747-1758, Oct. 1999.
- [100] Yang Xiao, IEEE 802.11n: Enhancements for Higher Throughput in Wireless LANs, IEEE Wireless Communications, vol. 12, no. 6, pp82-91, 2005.
- [101] Yaghoobi, H., Scalable OFDMA Physical Layer in IEEE 802.16 Wireless MAN, Intel Technology Journal, vol. 8, no. 3, August 2004.
- [102] J. Yin, T. ElBatt, G. Yeung, B. Ryu, S. Habermas, H. Krishnan, T. Talty, Performance Evaluation of Safety Applications over DSRC Vehicular Ad Hoc Networks, Proceedings of the first ACM workshop on Vehicular ad hoc networks, October 01-01, 2004, Philadelphia, PA, USA.

- [103] Jens Zander, Seong-Lyun Kim, Radio Resource Management for Wireless Networks, Artech House, Boston, MA, 2001.
- [104] K. Zheng, L. Huang, W. Wang, G. Yang, TD-CDM-OFDM: Evolution of TD-SCDMA toward 4G, IEEE Communications Magazine, vol. 43, no. 1, pp45-52, Jan. 2005.
- [105] Jianliang Zheng, M. J. Lee, Will IEEE 802.15.4 Make Ubiquitous Networking a Reality?: a Discussion on a Potential Low Power, Low Bit Rate Standard, IEEE Communications Magazine, vol. 42, no. 6, pp140-146, 2004.
- [106] J. Zhu and S. Roy, MAC for Dedicated Short Range Communications in Intelligent Transport System, IEEE Communications Magazine, vol. 41, no. 12, 2003.
- [107] M. Zorzi and R. Rao, On the Statistics of Block Errors in Bursty Channels, IEEE Transactions on Communications, vol. 45, pp. 660-667, June 1997.
- [108] M. Zorzi, A. Chockalingam, and R. Rao, Throughput Analysis of TCP on Channels with Memory, IEEE Journal on Selected Areas in Communications, vol. 18, pp1289-1300, July 2000.