

**Delay Analysis for Wireless
Applications using a
Multiservice Multiqueue
Processor Sharing Model**

Yan Wang

Doctor of Philosophy

2008

RMIT

Delay Analysis for Wireless Applications using a Multiservice Multiqueue Processor Sharing Model

A thesis submitted in fulfilment of the requirements for
the degree of Doctor of Philosophy

Yan Wang

M. Eng.

School of Electrical and Computer Engineering
Science, Engineering and Technology Portfolio

RMIT University

Melbourne, Victoria, Australia

November 2008

Declaration

I certify that except where due acknowledgement has been made, the work is that of the author alone; the work has not been submitted previously, in whole or in part, to qualify for any other academic award; the content of the thesis is the result of work which has been carried out since the official commencement date of the approved research program; and, any editorial work, paid or unpaid, carried out by a third party is acknowledged.

Yan Wang

Melbourne, November 1, 2009

**Delay Analysis for Wireless Applications using a Multiservice
Multiqueue Processor Sharing Model**

To my faithful Lord and my family

Acknowledgement

Beginning in early 2004, I started this long journey for my PhD, which was filled with challenges right up until the present time. It was a period with many ups and downs, uncountable levels of help and support from many people has accompanied me along the way. Therefore, I would like to express my deepest honor and appreciation to all of them.

First, my deepest thanks go to Prof. Richard J. Harris and Prof. Moshe Zukerman for supervising me. During this time, they have been working in different universities. Currently, Prof. Richard Harris has a Chair of Telecommunications and Network Engineering at Massey University, New Zealand and Prof. Moshe Zukerman is currently working at the ARC Special Research Centre for Ultra-Broadband Information Networks (CUBIN) at The University of Melbourne, Australia. I am grateful for being their second collaboratively supervised student and a wonderful thing is that their first collaboratively supervised student, Assoc/Prof. Sammy Chan, and who is currently working at The City University of Hong Kong, Hong Kong, joined this work around 2007 as my third supervisor! Working together in different universities, even in different countries, it has been a big challenge but it has also been a big blessing for me. Distances can be shortened by modern communication technologies that our research work has focused on, most importantly, it has been shortened by our close inter-working spirit. I have learnt a great

deal from this collaborative spirit, and it is something that all research students dream of. A way of problem solving through comprehensive discussions, writing skills for presenting ideas in a research paper, language skills on presentation, publication revisions, and all aspects of the research skills that I required, they provided significant efforts in guiding me. I cannot begin to thank all of them for their assistance in my progress along the way.

Other people that have been supporting me throughout this candidature include Dr. Robert Suryasaputra and Dr. Suyong Eum, who have been part of the CATT Centre that operated in RMIT during much of my time, and Assoc/Prof. Ron Addie who works at the University of Southern Queensland, Australia. Robert and Suyong helped me a lot with programming skills and also have been a great help in networking. Ron has been working with us on the mathematical modelling aspects of the project and supported me with his deep knowledge and answers to my many questions. They all contributed greatly to my understanding of problem formulation and solution methodology with many useful discussion sessions and helpful advice.

The previous CATT Centre family were always a great source of inspiration and invaluable help. In the early days of my candidature, Dr. John Murphy, Dr. Alexander Kist, Dr. Irena Atov and Dr. Deddy Chandra when they worked in CATT used to explain to me about what was going on. I also would like to thank Dr. Suyong Eum, Dr. Robert Suryasaputra, Dr. Suresh Venkatachalaiah, Dr. Kevin Lin and Dr. Jimmy Lau for their friendship and for making my time so much more enjoyable. Being a part of CATT Centre family has also taught me that hard work and dedication is the *only* way to achieve your dreams.

I also wish to thank my parents, my grandmother's big family and my partner Rick for providing me with such a wonderful environment over the past

few years. Especially, my grandmother, who helped me grow up, but could not see the completion of this work before she passed away, which is one of my big regrets. My partner Rick, being a companion, has always been eager to listen to my story, he has supported me in many ways, and at the same time, encouraged me as I made many decisions. To my Christian family, brothers and sisters from different churches with their warm assistance, concerns and prayers, I would like to express my thanks for making me part of your wonderful family. Also, thanks go to all my teachers and friends during my time of school and university, colleagues who worked with me in China, who have contributed so much to my life. My final and greatest thanks go to my Lord, in whom I have great faith and who has influenced so many parts of my life and finally led me to make the decision to go this way.

Without these wonderful people, I am sure that you would be doing something else rather than reading this thesis, because this work would never exist without the continuous support from all these people.

Contents

Acknowledgement	vii
Abstract	xiii
1 Introduction	1
1.1 Focus of this Thesis	2
1.1.1 Scheduling for QoS	2
1.1.2 Performance Analysis by Queueing Models	3
1.1.3 PS Scheduling Strategy and Models	4
1.1.4 Wireless Applications using a Multiservice Multiqueue PS Model	6
1.2 Contributions of this Thesis	6
1.3 Organisation	10
2 Background	13
2.1 Wireless MAC Protocols and Related Scheduling Issues	14
2.1.1 QoS Metrics	14
2.1.2 Role of Wireless MAC Protocols	16
2.1.3 IEEE Standards 802.11 and 802.16	17
2.2 Scheduling Algorithms for Wireless Networks	19
2.2.1 Non-opportunistic Scheduling	20
2.2.2 Opportunistic Scheduling	23
2.2.3 Discussion	23
2.3 Processor Sharing Scheduling and its Queueing Models	25
2.3.1 Egalitarian Processor Sharing (EPS)	26
2.3.2 Generalized Processor Sharing (GPS)	30
2.3.3 Discriminatory Processor Sharing (DPS)	35
2.3.4 Multilevel Processor Sharing (MLPS)	39
2.4 Wireless Applications: PTT and WiMAX	42
2.4.1 PTT over GSM/GPRS	42

<i>CONTENTS</i>	xi
2.4.2 WiMAX	44
2.5 Conclusions	46
3 Multiservice Multiqueue PS Models	47
3.1 Introduction	47
3.1.1 Multiservice PS Scheduling	47
3.1.2 Multiservice PS Scheduling Models	48
3.2 Multiqueue PS (MPS) Model	50
3.3 Notation	51
3.4 Correction of a Subtle Incongruity in the LQ and the PS Mean Delays of the MPS Model	54
3.4.1 The Subtle Incongruity for L_p and $S_p(n)$	54
3.4.2 Derivation of LQ Waiting Time	57
3.5 Priority-Based Service Quanta (PBSQ) Model	61
3.5.1 Analytical Model	61
3.5.2 Mean Message Delay	63
3.5.3 Relevant Implementation Issues	72
3.5.4 Model Evaluation and Numerical Results	73
3.6 Conclusions	76
4 Modelling and Analysis of PTT delay	78
4.1 Introduction	78
4.2 GPRS LLC/RLC/MAC Layers	79
4.3 Analytical Model for PTT Uplink Delay	81
4.3.1 Partial Sharing Channel Allocation Scheme	81
4.3.2 Quasi-stationary Assumptions	82
4.3.3 Modelling GPRS Delay	83
4.4 Simulation Study	86
4.4.1 Scenario 1: Fixed Channel for GPRS Traffic	86
4.4.2 Scenario 2: Partial Sharing Channel Scheme for GPRS and GSM Traffic	88
4.5 PTT/GPRS Retransmissions	91
4.6 Effects of Traffic Load for GSM and GPRS	92
4.7 PTT/GPRS Priority Assignment	93
4.8 Conclusions	94
5 Modelling WiMAX SS Uplink Delay	96
5.1 Introduction	96
5.2 MAC Protocol of IEEE 802.16	99

<i>CONTENTS</i>	xii
5.3 Priority-based Fair Scheduling	102
5.4 Delay Analysis using the MPS Model	103
5.5 Model Evaluation	106
5.5.1 Traffic Models	107
5.5.2 Simulation Model	107
5.5.3 Simulation and Numerical Results	109
5.6 Conclusions	111
6 Conclusions	113
6.1 Summary of Contributions	116
6.2 Future Research	117
6.2.1 Enhanced Models	117
6.2.2 More Applications	118
6.2.3 Effects of Further MAC Scheduling and Modelling Issues	118
Bibliography	120
Abbreviations	144
List of Figures	147
List of Tables	149
Appendix	150

Abstract

The ongoing development of wireless networks supporting multimedia applications requires service providers to efficiently deliver complex Quality of Service (QoS) requirements. Compared to earlier simple service networks which only provided best effort service, the wide range of new applications in these networks significantly increases the difficulty of network design and dimensioning to meet QoS requirements. Traffic engineering that involves consideration of multimedia applications is targeted at meeting QoS requirements based on existing network resources. Improved service performance including guaranteed delay, jitter and bandwidth, cannot be easily maintained by legacy wireless networks.

Medium Access Control (MAC) protocols affect QoS achieved by wireless networks. Research on analysis and performance evaluation of different proposed MAC protocols is important for the efficient protocol design. As wireless networks feature scarce resources that are simultaneously shared by all users, processor sharing (PS) models were proposed for modelling resource sharing mechanisms in such systems. In this thesis, multi-priority MAC protocols are proposed for handling the various service traffic types. Then, an investigation of multiservice multiqueue PS models is undertaken to analyse the delay for some recently proposed wireless applications.

We start with an introduction of MAC protocols for wireless networks

which are specified in IEEE standards and then review scheduling algorithms which were proposed to work with the underlying MAC protocols to cooperatively achieve QoS goals. An overview of the relevant literature is given on PS models for performance analysis and evaluation of scheduling algorithms. We introduce four types of basic PS model, viz: egalitarian processor sharing (EPS), generalized processor sharing (GPS), discriminatory processor sharing (DPS) and multilevel processor sharing (MLPS) models.

We propose a multiservice multiqueue PS model using a scheduling scheme in multimedia wireless networks with a comprehensive description of the analytical solution. Firstly, we describe the existing multiqueue processor sharing (MPS) model, which uses a fixed service quantum at each queue, and correct a subtle incongruity in previous solutions presented in the literature. Secondly, a new scheduling framework is proposed to extend the previous MPS model to a general case. This newly proposed analytical approach is based on the idea that the service quantum arranged by a MAC scheduling controller to service data units (SDUs) can be priority-based. In this model, the arrival process and the service process of specified networks can be described by a specific mathematical distribution. We obtain a closed-form expression for the mean delay of each service class in this model. The proposed simple yet efficient MAC protocol for wireless multimedia networks can support a wide range of services. Our MAC protocol is designed to meet QoS requirements in terms of the average delay. Some implementation issues of the new model are addressed. In summary, the advantages of this model over the traditional PS approach for performance analysis of delay are: (i) it simplifies MAC protocols for multimedia applications into an analytical model that includes more complex and realistic traffic models without compromising details of the protocol; (ii) it significantly reduces the number of MAC headers, thus the overall

average delay will be decreased.

In response to using the studied multiservice multiqueue PS models, we apply the MPS model to two wireless applications: Push to Talk (PTT) service over GPRS/GSM networks and the Worldwide Interoperability for Microwave Access (WiMAX) networks, where delay performance is one of our primary concerns. We investigate the uplink delay of PTT over traditional GPRS/GSM networks and the uplink delay for WiMAX Subscriber Station (SS) scheduler under a priority-based fair scheduling. MAC structures capable of supporting dynamically varying traffic are studied for the networks, especially, with the consideration of implementation issues, such as channel sharing schemes in GPRS/GSM networks, retransmissions for the PTT service and traffic classes specified in WiMAX standard. The model provides useful insights into the dynamic performance behaviours of GPRS/GSM and WiMAX networks with respect to various system parameters and comprehensive traffic conditions. We then evaluate the model under some different practical traffic scenarios. The performance analysis based on the model under realistic traffic conditions shows the benefit of our model for determining the delay associated with real-time services.

As the layering model simplifies telecommunications processes, through modelling of the operation of wireless access systems, under a variety of multimedia traffic, our analytical approaches provide practical analysis guidelines for wireless network dimensioning.

Chapter 1

Introduction

Since wireless networks for commercial purposes first emerged in the 1970s, they have dramatically increased in size and scope in the networking industry. Many mobile users can access ubiquitous communication regardless of their location. With the broad range of services provided in wired networks, and as conventional wireless networks are progressively linked into wired networks, it is natural to see wired network services being extended into wireless networks.

However, service performance through wireless networks cannot simply rely on physical layer technologies, such as multiple antennas (e.g., multiple-input multiple-output (MIMO) [1]), advanced modulation techniques (e.g., orthogonal frequency-division multiplexing (OFDM) [2]), highly efficient coding schemes (e.g., Turbo codes [3]), and so on. The upper layer communication protocols, such as Medium Access Control (MAC), also play important roles in delivering network performance by controlling the number of users to access resources and effectively scheduling the different sizes and kinds of demands. The MAC layer is actually a sublayer of the data link layer, which is located between the physical and network layers in the Open Systems Interconnec-

tion (OSI) reference model [4]. As wireless bandwidth is a scarce resource shared by a number of individual users for multiple services, and these services can be real-time or non real-time services with specific quality of service (QoS) requirements, MAC protocols are required to effectively cope with growing bandwidth demands as well as featured service differentiation. So the design of MAC protocols for wireless networks is much more challenging than for wired networks.

1.1 Focus of this Thesis

1.1.1 Scheduling for QoS

Connection Admission Control (CAC) is implemented in order to achieve a range of QoS targets at the MAC layer. The decisions made by the CAC are highly dependent upon the performance of scheduling algorithms because these algorithms provide a mechanism for bandwidth allocation to different users and services. As a result, in the process of MAC design, choosing an efficient scheduling policy and evaluating the performance of the scheduling policy are vitally important. The ongoing development of scheduling schemes faces the challenge of providing end-to-end performance guarantees to heterogeneous and bursty traffic; at the same time, the service disciplines must be simple enough to be implemented in a high speed environment. Due to the complexity of wireless link conditions and the bursty nature of the traffic, it has been considered difficult to design an effective scheduler as well as to analyse its performance with a set of parameters under realistic traffic conditions.

The service discipline used at the system scheduler decides the principle

that the server capacity is allocated to the customers. The common disciplines are:

1. First Come First Served (FCFS) that serves jobs in the arriving order;
2. Last Come First Served (LCFS) that serves the job arrived most recently;
3. Processor Sharing (PS) service discipline that serves all jobs concurrently.

1.1.2 Performance Analysis by Queueing Models

Queueing theory provides fundamental methods to be applied in the performance analysis. Through modelling different telecommunication systems, we can obtain descriptions of networks or local access systems, and performance of traffic schedulers and access protocols by means of queueing analysis. A queueing system is mainly characterized by arrival process of the customers (that normally represent other items such as data-packets) and service discipline to allocate the resources to the customers. The performance of a system is decided by the complex interaction between these two characteristics. Therefore, essential analytical methods and approaches of queueing theory have been broadly applied to the evaluation and design of communication networks.

The PS model is one of the important queueing models. PS models are typically applied to the analysis of resource sharing systems due to the main property of the PS service discipline that all jobs present in the system share the common resources, such as data communication networks. The key performance issues are the sojourn time and the throughput of jobs.

1.1.3 PS Scheduling Strategy and Models

Since an ever-increasing number and diversity of services are required and various QoS requirements need to be supported in the Internet, the use of a processor sharing [5] scheduling strategy and its associated scheduling algorithms have been proposed to achieve realizable, efficient, flexible and fair service disciplines. As the original PS principle is simple to apply in situations where different users receive their shares of scarce common system resource, it benefits those applications with delay concerns by allowing fast scheduling decisions.

The performance of PS scheduling algorithms has been studied in a wide range of PS models, where the analysis usually relies on simplifying assumptions, such as assuming specific traffic arrival and service models. The standard PS model consists of a single server assigning each customer an equal fraction of the service rate, similar to the time-sharing model used in computer operating systems. More important extended PS models were then proposed, such as the generalized processor sharing (GPS) model [6, 7, 8] and the discriminatory processor sharing (DPS) model [5, 6]. Allowing the service rate to be class-dependent, these models can be applied in packet-switched networks supporting service differentiation. The packets of high quality traffic flows are served with strict priority over the packets of low quality traffic flows in the network nodes, such as routers in an Internet Protocol (IP) network or Base Stations (BS) in a cellular wireless network. For instance, if data traffic flows are subject to a flow control mechanism such as the Transport Control Protocol (TCP) and the available bandwidth of TCP connections is temporarily limited due to overloads, the flow control mechanism decreases the transmission rate of each flow according to the assigned rate to fairly

share the available bandwidth. Thus, these models significantly enhance the modelling capabilities of the PS model.

Over the past few decades, PS models have been widely studied in the queueing literature for performance evaluation of computer and communication systems. PS-type models specifically cover the main factors which determine system performance; on the other hand, they are still simple enough to allow for the determination of an exact or approximate analytical solution. Thus, in many cases, PS models are preferred over some other models that may involve more details but are too complicated to be tractable for use in practical situations. However, compared to some other simple queueing systems such as FCFS, PS models are still quite hard to analyse. Intractable problems have been investigated for many years. Classical queueing analysis usually only studies average performance for aggregate traffic under realistic traffic models. In a PS queue with Poisson arrivals, though it is well-known that the queue length distribution has a simple geometric distribution, regardless of the service requirement distribution, the sojourn time distribution is far less tractable, even for exponential service cases. Therefore, there is a need to develop an analytical model that provides performance analysis of a versatile multiservice multiqueue PS model and covers a wide range of practical scheduling mechanisms. Such analysis can then play a critical role in network design and dimensioning.

This thesis considers a PS discipline with multiservice multiqueue models for the analysis of resource sharing in communication networks, especially, in wireless networks. The multiservice multiqueue PS models with a fixed service quantum and with priority-based service quanta are presented respectively. A multiqueue processor sharing (MPS) model [9] is introduced that includes the correction of a subtle incongruity. We also extend the MPS model

to a priority-based service quanta (PBSQ) model to fit the general requirements of wireless multimedia applications. Extensive mathematical analyses are provided to quantitatively characterise the performance of the proposed PBSQ model. Combined with the MAC protocol operation in detail, under general and realistic arrival assumptions, it will be shown that our analytical approach using the PBSQ model leads to an accurate approximation for the delay in multimedia networks. The extended model takes into account features of current wireless networks that can be used to define multimedia traffic.

1.1.4 Wireless Applications using a Multiservice Multi-queue PS Model

We propose a priority-based scheduling algorithm for two wireless applications: Push to Talk (PTT) [10] over General Packet Radio Service/Global System for Mobile Communications (GPRS/GSM) and the Worldwide Interoperability for Microwave Access (WiMAX) [11], and apply the multiservice multi-queue PS model to analyse their delay performance. According to the MAC structures of these networks, our model involves different effects of implementation issues, such as channel sharing schemes in GPRS/GSM networks, retransmissions for the PTT service and traffic classes specified in WiMAX standard.

1.2 Contributions of this Thesis

The contribution of this thesis is mainly threefold. Firstly, we introduce a new multiservice multi-queue PS model, namely PBSQ, that enables a delay anal-

ysis of multiservice networks. Unlike other PS models, it extends the delay analysis to the case of priority-based multiservice PS scheduling where different priorities receive different service quanta at a time. As the priority-based multiservice PS scheduling inherits nice properties of traditional PS service disciplines such as ease of implementation and relative simplicity of analysis, it benefits systems such as wireless MAC, where the scheduler needs to make fast decisions. An accurate approximation for the mean message delay is derived for the PBSQ model. Delay analyses for the multi-class multi-connection systems are made possible by the analytical model in a close-formed solution. With regard to the MAC implementation issues, we demonstrate the comparison of the PBSQ model and the developed MPS model, where the throughput of services with large size requests can be improved by setting priority-based service quanta in the PBSQ model to reduce MAC headers. Our study offers a practical model for delay analysis of multimedia applications with heterogeneous traffic characteristics.

Secondly, we apply our analytical priority-based PS scheduling discipline model to two multiservice wireless networks applications: PTT/GPRS/GSM and WiMAX, for delay performance evaluation. More implementation issues of MAC layer specified in these two applications are taken into consideration. The delay performance is analysed under some realistic traffic conditions. This scheme can effectively protect real-time services on the delay requirement by setting them higher priorities. As a result of applying the model, the produced analytical results of delay can be used by operators in network dimensioning and management.

Thirdly, we discover some new insights into the delay analysis of the MPS model that better explains its components. Although the overall delay in the MPS model is correct, we demonstrate that one component in it, namely wait-

ing time in a local queue, may obtain negative value in a very particular situation. We provide a modified analysis that fixes this subtle inconsistency with the physical meaning of local queue waiting time.

The contributions led to the following publications: [12, 13, 14, 15, 16] and are summarised below:

The work described in Chapter 2 led to the publication of [12]. The contributions of this chapter can be summarised as follows:

- Overview and comparison of various options for scheduling algorithms in wireless networks. Both non-opportunistic and opportunistic scheduling algorithms are considered.

The work described in Chapter 3 led to the publication of [16]. The contributions of this chapter can be summarised as follows:

- Demonstration and correction of a subtle incongruity for the waiting time in a local queue in the original Potter and Zukerman's work [9].
- Development of a priority-based service quanta model that allows for different kinds of services to obtain different service quanta each time rather than obtaining a fixed service quantum in the MPS model.
- Study of the mean delay of the PBSQ model under specified assumptions for the given parameters.
- A C++ simulation study for model evaluation and delay comparisons with the MPS model.
- An implementation issue for the PBSQ model involving MAC headers.

The work in Chapter 4 led to the publication of [13, 14]. The contributions of this chapter can be summarised as follows:

- Description of PTT over the GPRS/GSM uplink architecture.
- Modelling and analysis of the uplink delay using the MPS model.
- Discussion of the quasi-stationary assumptions in the analytical model.
- Estimating the effects of GPRS retransmission on PTT delay.
- An ns-2 simulation study to evaluate accuracy of the analytic model.
- A C++ simulation to study the effects of GSM voice traffic.
- Numerical results that illustrate the benefit for PTT service by using a priority arrangement.

The work in Chapter 5 led to the publication of [15]. The contributions of this chapter can be summarised as follows:

- Description of the WiMAX MAC architecture and identification of QoS issues.
- Development of a priority-based scheduling scheme for the Subscriber Station (SS) uplink scheduler.
- The delay analysis of the scheduling algorithm using the MPS model.
- Discussion of WiMAX traffic models.
- An ns-2 simulation study to evaluate the model.
- A numerical study for the effects of traffic models and scheduling policy.

Publications by the author

1. S. Chan, Y. Wang, A. Haider, R. J. Harris, and M. Zukerman, "Algorithms for WiMAX Scheduling," in *Proceedings of APMC 2008*, pp.1-4, Hong Kong, December 2008.
2. Y. Wang, M. Zukerman, and R. J. Harris, "PTT packet delay analysis for GPRS/GSM links," *IEEE Communications Letters*, Vol.10, No.6, pp.456-458, June 2006.
3. Y. Wang, M. Zukerman, and R. J. Harris, "Modelling PTT uplink in GPRS/GSM networks," in *Proceedings of IEEE VTC 2006-Spring*, Vol.1, pp.420-424, Melbourne, Australia, May 2006.
4. Y. Wang, S. Chan, M. Zukerman, and R. J. Harris, "Priority-based fair scheduling for multimedia WiMAX uplink traffic," in *Proceedings of IEEE ICC 2008*, pp.301-305, Beijing, China, May 2008.
5. Y. Wang, M. Zukerman, R. Addie, S. Chan, and R. J. Harris, "A priority-based Processor Sharing Model for TDM Passive Optical Networks," *Submitted to IEEE Journal on Selected Areas in Communications*, August 2009 .

1.3 Organisation

Chapter 2: In this chapter, we start with the context for QoS requirements in wireless networks and tasks performed by MAC protocols. In the MAC layer, the scheduling mechanism is introduced for the delivery of QoS. Scheduling issues in MAC protocols become more complicated due to the growth of bandwidth demands and the service diversification. Herein, we outline the pro-

posed wireless scheduling algorithms. Afterwards, we focus on the PS service discipline and existing PS models used for handling heterogeneous traffic in multiple classes to provide an overview of the characteristics and results of applying the models. Two wireless applications: PTT and WiMAX are introduced at the end of this chapter, which are used as application examples for the multiservice multiqueue PS models in later chapters.

Chapter 3: We review the results of the existing MPS queueing model. We demonstrate a subtle incongruity in the original model and modify the analysis to correct it. Furthermore, we extend this model from the assumption of a fixed service quantum at each priority to using a priority-based service quantum at each priority. We provide an accurate approximation of the average delay with derivation of a closed-form expression for the new model. The effects of different parameters are demonstrated. Furthermore, through our approach, we show that having implemented PBSQ model reduces message delay significantly as it is possible to reduce the number of MAC headers by setting a bigger service quantum to cope with large service demands.

Chapter 4: We consider the case of PTT service delivered over GPRS/GSM networks. The reason for analysing PTT delay in this context is based on the fact that GPRS is initially designed for a normal data service (non real-time), but not for a voice (real-time) service. Firstly, in this chapter, the existing GSM standard is described in some detail. We then propose the use of the MPS model to analyse the delay performance of PTT over GPRS/GSM networks. Different types of traffic mix – including the original GSM voice traffic, GPRS data and PTT voice over GPRS are blended into consideration using our analytical model. The accuracy of our model is evaluated using simulation. Moreover, the effect of retransmissions over GPRS on PTT delay is analysed. Through this numerical study, we show that a PTT voice service can

be strictly protected by using priority assignment. Our model allows service providers to do comprehensive network dimensioning.

Chapter 5: In this chapter, we introduce WiMAX technology according to the IEEE 802.16 standards and the multiple services supported by WiMAX, including Voice over IP (VoIP), streaming video, data etc.. We describe the MAC architecture of WiMAX for the delivery of different QoS requirements. Then, we propose a priority-based scheduling policy for the SS scheduler. The delay of real-time traffic is protected by a priority arrangement. All scheduling services are considered in the model. We validate the model using simulation and demonstrate the effects of traffic characteristics on scheduling policy using our numerical and simulation results.

Chapter 6: This chapter is devoted to providing some final conclusions and discussing future research directions.

Chapter 2

Background

Nowadays, following the integration of wired and wireless networks in our communication systems, more multimedia applications have been launched into wireless networks. The traffic types from these applications include non real-time and real-time traffic which require performance guarantees. Such performance requirements need to be technically expressed using a range of QoS requirements such as delay, throughput, delay-jitter and loss rates. Such like in wireless local area networks (WLANs) [17], MAC protocols play a critical role to provide QoS guarantee by carrying out two essential functions: connection admission control (CAC) and service differentiation [18].

Due to limitations in bandwidth resources, CAC policies are very dependent on the specific scheduling disciplines used to handle different services. Scheduling algorithms which provide mechanisms for bandwidth allocation are considered to be a key issue for MAC layer performance. Once scheduling algorithms have been chosen, modelling the algorithms and getting analytical performance results involving specific parameters become important research issues, as such performance results can be used by the admission control module.

In this chapter, we provide a general description of a wireless MAC layer and explore difficulties in achieving QoS requirements due to features of the wireless link. After a brief introduction to wireless scheduling algorithms, we focus on current work involving PS disciplines to give a comprehensive review of a range of PS models which have been proposed for modelling various scheduling algorithms. Issues of packet scheduling associated with service differentiation for two wireless applications: PTT and WiMAX, are discussed at the conclusion of this chapter.

2.1 Wireless MAC Protocols and Related Scheduling Issues

We begin by briefly discussing QoS metrics and addressing the functions performed by MAC protocols in wireless networks. We then introduce two common sets of wireless standards with a MAC specification and discuss scheduling issues in these wireless networks.

2.1.1 QoS Metrics

Although there are many alternative definitions of QoS, in this thesis, QoS is under the field of packet-switched networks as a traffic engineering term. We focus on “objective” QoS which means the performance that can be measured rather than “subjective” QoS which corresponds to the quality from a user’s perspective and is now more commonly referred to as “Quality of Experience” and denoted by QoE.

QoS requirements vary for different service applications. Since real-time applications are delay sensitive and may also require certain data rates, the

QoS issues are critical. For example, streaming multimedia requires guaranteed throughput; IP telephony or VoIP require strict limits on jitter and delay; video teleconferencing requires low jitter and delay. For a comprehensive discussion of QoS for multimedia applications, readers can refer to [19].

Here, we review four general performance metrics which are also relevant to QoS of multimedia applications [20].

- Throughput (or bandwidth) refers to the data rate (also called the bit rate) generated by the application. The required throughput depends on application characteristics. Throughput is called bandwidth when it is considered to be the network resource allocated to applications.
- Delay (or latency) directly affects satisfaction of a user. This information is valuable QoS metric and gives direct insight of the total performance of the system. We shall focus on probabilistic delays by referring to the expected value of delay in this thesis.
- Jitter (or delay variation) as a QoS term refers to delay variation. Several definitions of jitter have been used in the literature, such as the maximum variation of the delay and the standard deviation of the delay. Jitter can be measured from the exact delay difference between sequential packets, but cannot be calculated from the mean delay. Large jitter can be smoothed out using a large buffer at the expense of longer delays.
- Packet loss has a direct impact on user's perceived quality. It can be caused by network congestion or by communication channel errors. The lost packets or the packets with errors might be recovered by several techniques, such as packet retransmission or error correction. In a wireless network, lost packets are usually detected using sequence numbers

and automatic repeat request methods.

Broad-band wireless networks need to provide services for heterogeneous traffic with different QoS requirements. In wireless networks, because of the limited capacity, the QoS issue is more challenging than in wired networks. More complex MAC Layer implementations have to be considered for the QoS solutions in wireless networks. Such as the IEEE standard sets of 802.11e [17] and 802.16 [21], different ranges of QoS-based improvements have been involved for wireless networks at the MAC layer, in which QoS support mechanisms should be integrated into the MAC scheduling algorithms.

2.1.2 Role of Wireless MAC Protocols

Wireless channel resources are shared by multiple users for multiple services delivery. Unlike the data link control for a point-to-point communication link in a wired network, which only needs to convert the physical layer bit pipe into a higher level frame, wireless MAC protocols are also required to regulate the user's access to the channel and to achieve appropriate sharing of the channel resources among users.

MAC protocols in wireless networks are complicated as they have to handle the requirements of multi-access communication. For example, in WLANs [17], data packets are sent through a common channel so that "collisions" may occur when different users transmit their data at the same time. To coordinate the transmissions among all users, the MAC provides mechanisms for channel access control which make communications possible. A physical address (or MAC address) is assigned a unique serial number for each network adapter over which a data packet can be delivered to a given destination within a physical network. Therefore, several stations can be connected to the

same wireless physical medium to share the link [22]. In general, WLANs use a Carrier Sense Multiple Access with Collision Avoidance (CSMA/CA) MAC protocol to fit the nature of the wireless channel.

Another key function of a MAC protocol is to achieve the efficient utilisation of the wireless link and satisfy QoS requirements for each user. However, the scarce bandwidth is time-varying and location-dependent in a wireless system, and network performance suffers from high error rates due to user mobility and power constraints in their devices [23]. As a result of the special characteristics of wireless communications and complex traffic characteristics, efficient and effective scheduling algorithms may need to be developed for wireless networks. Since MAC decisions concerning resource sharing rely on the operational results of scheduling, obtaining results by modelling and performance evaluation of scheduling algorithms for wireless networks is important and more challenging than for fixed networks.

2.1.3 IEEE Standards 802.11 and 802.16

The classification of wireless networks is normally according to the geographical coverage. The wireless wide area networks (WWANs) were developed first to cover a very large area and its typical examples are mobile cellular networks and satellite networks. Then, wireless local area networks (WLANs) were developed to provide network access to a group of users within a certain local area typically based on the IEEE 802.11 standards [17] family. Recently, wireless metropolitan area networks (WMANs) have been proposed to support broadband wireless access over a metropolitan area based on the IEEE 802.16 standards [21]. In contrast to those networks with large coverage, wireless personal area networks (WPANs) only operate within a short dis-

tance of a few meters, which are based on the IEEE 802.15 standards [24] including 802.15.1 for Bluetooth, 802.15.3 for high data rate ultra-wideband networks and 802.15.4 for low-rate WPANs such as wireless sensor networks.

Relating to the issue of QoS requirements in wireless networks with integrated traffic, we discuss two standards for wireless communication which are in wide use currently, viz: the 802.11 [17] and the 802.16 [21] standards.

The 802.11 Standard was approved by the Institute of Electrical and Electronics Engineers (IEEE) in 1997, followed by the standard known as 802.11b in 1999. This set of standards was developed by the IEEE LAN/MAN Standards Committee (IEEE 802) for WLAN computer communications and contains WLAN MAC and Physical Layer (PHY) specifications. Although it was intended for indoor WLAN, a range of WLAN-based products have been used in point-to-point and point-to-multipoint outdoor solutions. In the original version of the standards, only *best effort* (BE) services were supported; there was no priority setting for different traffic types in the network, hence, no QoS differentiation was achieved. The IEEE 802.11e standard was later developed and approved as the QoS enhanced version of the original standard. A set of QoS enhancements is defined for supporting delay-sensitive applications, such as VoIP and streaming video.

In April 2002, the IEEE Standard 802.16 [21] was approved and it focused on *broadband* wireless access in metropolitan area networks. This technology provides an alternative to traditional wireline technologies such as coaxial cable networks and digital subscriber lines (xDSL) based on PSTN access networks. As multimedia applications, including data, streaming video and VoIP, are supported, each connection may have different QoS requirements.

The IEEE 802.11e and 802.16 standards enable QoS differentiation which is desirable since communication sessions usually differ in their QoS require-

ments. For instance, streaming video and VoIP conversations require a guaranteed minimum or constant throughput to ensure image and speech quality. On the other hand, other normal data applications usually can take variable transmission rates. However, these advantages for QoS differentiation make traffic scheduling and performance evaluation of scheduling algorithms more complex.

2.2 Scheduling Algorithms for Wireless Networks

A scheduling mechanism is an instrument used to treat all applications with service differentiation, which has to be implemented at routers or switches of a network. Service disciplines to support a number of services with various QoS needs are required for the scheduling mechanism in both wired and wireless networks to achieve efficient and fair scheduling among multiple users.

However, the standards mentioned in Section 2.1.3 purposely do not specify the details of a MAC protocol containing scheduling algorithms. This allows equipment manufacturers to be given the choice of implementing their own proprietary algorithms to differentiate their products from those of their competitors.

Many scheduling algorithms have been developed for wireline networks with associated performance studies, by which a certain level of QoS can be achieved by treating different types of traffic differently. Even including QoS differentiation, the existing service disciplines [25], such as fair queueing (FQ or generalized processor sharing, i.e., GPS) scheduling [26] and variants, virtual clock [27] and the earliest due date (EDD) [28], may not be adequate for wireless networks due to their different channel characteristics. In wireline links, the bandwidth is constant, while in wireless links it is time varying due

to interference, fading and shadowing. Also, due to different physical locations, different terminals may perceive different channel quality at the same time instance. This is referred to as multiuser diversity. If a wireline scheduling algorithm is directly applied, without giving consideration to channel conditions, a packet may be scheduled for transmission on a bad wireless channel so that the packet cannot be successfully transmitted. Then, the bandwidth received by the connection is not the intended bandwidth allocated by the scheduling algorithm, and in general, the QoS received by a connection is not as good as the scheduling algorithm expects to deliver.

Here we review several recently developed wireless scheduling algorithms which consider channel conditions and attempt to schedule packets for good channels. Two classes of such algorithms are considered: 1) non-opportunistic scheduling, and 2) opportunistic scheduling. Opportunistic scheduling algorithms attempt to better predict channel quality using more detailed physical layer information than the non-opportunistic algorithms.

2.2.1 Non-opportunistic Scheduling

These algorithms assume the channel is either in a *good* or a *bad* state and use per-flow queueing. Here we only provide a brief summary of these algorithms. For a more comprehensive and critical review, see [23].

Channel state-dependent packet scheduling (CSDPS) [29] is one of the earliest schemes which takes into account location-dependent and time-dependent channel conditions. Basically, packets are scheduled by a wireline scheduling algorithm, e.g., round robin (RR). The wireless system throughput is improved by temporarily disallowing the flow having a bad channel to transmit, and instead, giving the transmission opportunity to the next flow with a good chan-

nel. This introduces a fairness issue because the flow that loses its right to transmit is never compensated. This fairness problem is alleviated by adding a class-based queueing (CBQ) [30] algorithm. The resulting scheme, called CSDPS+CBQ [31, 32], restricts a flow from receiving additional bandwidth when it has already received its fair share thus ensuring fair sharing of link bandwidth.

The weighted fair queueing (WFQ) algorithm [26, 7, 31] has been widely used for wireline networks to provide fairness at the packet level. WFQ has been modified to be also applicable to the case of multiuser diversity. This WFQ modification gave rise to the algorithm called *idealized wireless fair queueing* (IWFQ) [33]. An important concept in both WFQ and IWFQ, is the so-called *finish time*, which indicates the intended time for a packet to complete its service. Accordingly, packets with earlier finish time are served first. Packets of each flow, upon arrival, are queued in a non-decreasing order of their finish times. Although the scheduler always tries to serve the packet with the smallest finish time first, this rule may not always apply under IWFQ. If the chosen packet happens to receive a bad channel, the service opportunity will be given to the packet with the next smallest finish time that also has a good channel. Packets that have lost their transmission opportunity still maintain finish times that are the smallest in the system. As a result, once their channels become good, they are selected first by the scheduler. In this way, flows that previously lost their transmission opportunity, are compensated. In brief, FQ and wireless fluid fair queueing (WFFQ) are fluid scheduling models for wireline and wireless networks, respectively. As there are implementation problems, there are WFQ and IWFQ as the corresponding packet-based emulations of FQ and WFFQ, respectively.

Another scheduling algorithm is the so-called *channel-condition indepen-*

dent packet fair queueing (CIF-Q) [34]. It approximates real service in an error-prone system to its equivalent counterpart in an error-free reference system. Then CIF-Q schedules packets as if it is operating in a wireline network environment using the so-called start-time fair queueing (SFQ) [35]. SFQ is a wireline scheduling algorithm designed to achieve fairness even if the available bandwidth in the bottleneck link varies. Unlike IWFQ, CIF-Q is designed to achieve more specific QoS objectives, viz:

- delay bound and throughput of error-free flows are guaranteed,
- short term and long term fairness are maintained among the flows,
- a minimum service guarantee for a flow which has already received more than its entitled service.

A flow is labelled *leading*, *lagging* and *satisfied* if it receives more, less or the same amount of service as it would have in the reference system, respectively. The scheduler maintains for each flow a *virtual time* that indicates the normalised amount of service time that the flow has received. The virtual time is updated according to the SFQ in the reference system. Then, the head-of-line (HOL) packet of the flow with the smallest virtual time is selected for service first. Provided that the packet has a good channel and its flow is not a leading flow that has already received its guaranteed service, the packet is transmitted. Otherwise, the transmission opportunity is given to a lagging flow. If none of the lagging flows has a good channel to transmit, then the transmission opportunity is given to a non-lagging flow.

IWFQ and CIF-Q are based on an approximate error-free reference system and do not provide explicit compensation to a flow that lost its transmission opportunity. Therefore, only error-free flows can really be treated fairly and

be provided with a guaranteed throughput. An approach that considers compensation explicitly is the so-called *server-based fair approach* (SBFA) [36]. It dedicates a certain part of the available bandwidth to compensate delayed flows using *long-term fairness server* (LTFS). If a flow has a bad channel, LTFS records the service loss. This way, SBFA can provide guaranteed throughput.

2.2.2 Opportunistic Scheduling

We have so far discussed wireless scheduling algorithms that assume a relatively simple channel model; the channel is either good or bad. However, such simplistic channel characterization may not be sufficient. Therefore, there has been recent interest in a new class of wireless scheduling algorithms, called *opportunistic scheduling* (OS) [37]. In essence, the principle of opportunistic scheduling is the same as the above-mentioned wireless scheduling algorithms in that the scheduler avoids transmitting packets having a bad channel. The difference is that OS schemes make use of more information on channel quality, and considers indicators such as estimated instantaneous carrier-to-interference ratios, supportable data rates, received signal strength indications, or bit error rates of users' links. Based on these indicators, higher priority will be given to packets with the best channel quality. This way, OS schemes improve efficiency by achieving higher channel utilisation.

2.2.3 Discussion

Although these existing wireless scheduling schemes may be applicable to different wireless networks, there are still certain scheduling issues that require attention from the research community under certain situations for delivering the promised QoS and efficient operation. Besides providing QoS, a wire-

less scheduling algorithm should also aim to maximize utilisation of wireless channels by using a minimal number of scheduling-related control messages and also minimize unproductive transmissions. Therefore, the difficulties of using such wireless scheduling algorithms are obvious, such as the complexity of computing and the evaluation of performance using analytical models.

In the downlink scheduler, for example, at the BS, the environment matches quite well with the general scenario considered in wireless scheduling algorithms. That is, the scheduler has knowledge about the channel conditions and full information of individual queues. For this reason, wireless scheduling algorithms that can deliver the required QoS of a service type are good potential candidates. However, for the uplink scheduler, normally, only limited information about each traffic queue is available. This is because in order to reduce the communication overhead, the uplink direction just transmits minimal amount of information of traffic queues. For example, the uplink scheduler may not know the packet arrival times or the packet size at the head of each queue.

The sojourn time (also known as the response time or the delay) on the performance measures for a real-time service is a very big concern. It means the time required to deliver a “message” from the origin to the destination, i.e. the total time spent in the system. Thus, the sojourn time consists of the time spent in the queue, also referred to as the waiting time, and the time spent in service. The term “message” may refer to an application layer data-unit or a network layer packet, if it has been broken down into a number of packets for transmission. As the consideration of delay strongly influences the choice of scheduling algorithms and the performance of a network, it is a critical issue to capture the nature of the delay mechanism involving characteristics of the network traffic. For wireless scheduling algorithms, the delay is even more

difficult to be traced by analytical solutions.

2.3 Processor Sharing Scheduling and its Queueing Models

Many common scheduling algorithms, are more or less related to the well-known service discipline: processor sharing (PS) [5]. The PS scheduling strategy and its inherited scheduling algorithms have been attracting prominent attention from the research community for several decades. As a convenient paradigm, the different PS models continue to play an instrumental role in the design and operations of communication systems. In the following section, we provide a detailed overview of PS models available in the literature.

Kleinrock's PS model [38] has influenced performance modelling and evaluation of resource allocation schemes since 1964. This model was used to capture the fundamental properties of a time-sharing system. Then, he proposed the PS paradigm to model round robin scheduling algorithms in time-shared computer systems [39]. In addition to time-shared systems, the PS model has also been espoused for other resource sharing systems in modelling and performance evaluation [40, 41], such as wireless networks where bandwidth-sharing is modelled. These PS models provide valuable insights into service capacity sharing under a critical assumption: the PS server is always shared in an egalitarian manner among all competing users – even if they belong to different classes and should obtain unequal shares [42], so-called Egalitarian PS (EPS). The limitation leads to a difficulty in using the model to describe heterogeneous systems.

Later, PS-related disciplines, such as GPS [6, 7, 8] and DPS [5, 6], involv-

ing the consideration of service differentiation were proposed as multiclass extensions of the single-class EPS model. These models are widely studied as suitable abstractions for performance modelling in bandwidth-sharing systems at the flow-level [40, 43]. However, the delay in such queueing systems is far less tractable than in the basic PS model. Usually, the resource-sharing mechanism for best-effort services is modelled by an EPS discipline, while the one for QoS required services is modelled by either the GPS or the DPS discipline.

Another multi-level kind of extension of the EPS policy is known as Multilevel Processor Sharing (MLPS) and was a strategy introduced by Kleinrock in the 1970's [5]. Unlike DPS and GPS which are mainly for the purpose of service differentiation, MLPS aims to improve the performance of the whole system by exploiting the variability of service demands by giving priorities to shorter requests [42].

In the following sections, we review various key results for EPS, GPS, DPS and MLPS models, highlight their inherent desirable properties from ordinary PS and discuss their capability to deliver QoS requirements.

2.3.1 Egalitarian Processor Sharing (EPS)

The original PS, also known as EPS, means that each customer is assigned an equal capacity of the processor, which is used in multi-access computer systems with competitive demands. This simple standard PS model consists of a single server and multiple users and is applicable to situations in which different users fairly receive a share of the scarce common resource. If the capacity is time-shared among all users, the service rate allocated to each one depends on the total number of users.

As the service capacity is equally shared by all users currently present in the system, EPS has been used to evaluate the flow-level performance of cellular data systems with Proportional Fairness (PF) scheduling [44]. This idea was also introduced in the study of modelling the bandwidth sharing on the Internet [40, 45] and the performance analysis of WLANs [46]. Assuming flows are Poisson arrivals, the EPS model has simple, explicit expressions for the distribution of the number of active flows in the steady state, and first-order flow-level performance metrics. These measures are insensitive to the flow size distribution.

Assume there is a Poisson arrival rate λ and random service requirements sized x . Let $\rho = \lambda \mathbf{E}(x)$ and $\rho < 1$, where $\mathbf{E}(x)$ denotes the mean of x . Sakata *et al.* [47] obtained the well-known result for the length distribution of the stationary M/G/1 EPS queue with n jobs as:

$$\pi_n = (1 - \rho)\rho^n, \text{ for } n = 0, 1, \dots \quad (2.1)$$

Accordingly, the queue length distribution is insensitive to the service requirement distribution but only depends on the mean of the service requirement distribution, i.e., $\mathbf{E}(x)$. This is also true for the mean sojourn time according to Little's law [48]. In contrast to the simple geometric queue length distribution, the sojourn time ($T(x)$ for a job with size x) distribution does not have a closed-form expression. Determining the sojourn time distribution in EPS queues turned out to be a rather challenging problem, even for exponential service requirements.

The sojourn time of EPS has been analysed in many studies. Conditioned on the initial service requirement $x > 0$, Kleinrock [38, 39] showed the mean

sojourn time for the M/M/1 EPS queue is given by:

$$\mathbf{E}[T(x)] = \frac{x}{1 - \rho}. \quad (2.2)$$

The proportional result of (2.2) reflects a certain fairness for EPS. Sakata *et al.* [47, 49] extended the case to multiple servers with generally distributed service requirements. Kleinrock also discussed a similar case in [5]. According to these results, the proportionality between $T(x)$ and x and the insensitivity between $T(x)$ and the service requirement distribution are valid. Eventually, it is also true in Cohen's generalized PS model [50].

For an EPS queue with the assumptions of a Poisson arrival process and exponential service requirements, Coffman *et al.* [51] first derived a closed-form expression for the Laplace-Stieltjes transform (LST) of the equilibrium sojourn time of an arriving job, conditioned on the service requirement $x > 0$. From these results, the first two moment expressions for the conditional sojourn time were obtained. Morrison [52] extended this work to obtain the unconditional distribution function of sojourn time, which can be used for numerical evaluation of the sojourn time distribution when the traffic intensity $\rho < 1$. When ρ is close to 1, the sojourn time distribution under the situation of heavy-traffic behaviour can also be analytically investigated.

The exact expression for the distribution of the conditional sojourn time for an EPS queue with generally distributed service requirements has been studied for many years. Several analytic solutions for the conditional sojourn time have been obtained by Yashkov [53], Ott [54], Schassberger [55], and Van den Berg and Boxma [56] via different approaches. However, because of the complexity of these transform expressions, the results are not really insightful and cannot be simply applied for computational purposes. The expression ob-

tained by Zwart and Boxma [57] is supposed to be the most explicit one in the literature, but can only recursively calculate the moments of the sojourn time. Cheung *et al.* [58, 59] studied the moments of the sojourn time when there is a service requirement $x > 0$ and derived explicit upper and lower bounds for all moments. These bounds are insensitive to the service requirement distribution but only related to the mean, hence, provide further support for the observation of EPS as a fair resource discipline. Recently, Borst *et al.* [60] gave an overview of several methods that had been developed to obtain the asymptotic equivalence and outlined the differences and similarities between these approaches.

More related work on EPS queues may be obtained from the literature under different assumptions of arrival processes and service requirement distributions: Massey [61] and Núñez-Queija [62] for M/M/1 EPS; Egorova *et al.* [63] for M/D/1 EPS; Mandjes and Zwart [64], Brandt and Brandt [65] for GI/G/1 EPS; and mostly work for M/G/1 EPS, such as Asare and Foster [66], Yashkov [67, 68, 69], Baccelli and Towsley [70], Grishechkin [71], Guillemin *et al.* [72], Bansal [73], Kitayev [74], Hampshire *et al.* [75], Kim and Kim [76].

From a practical point of view, in a standard EPS queue, small jobs cannot be blocked by large jobs since fair sharing policies prevent large jobs from hogging the server. With such a nice appealing property, the EPS discipline is a big improvement over the FCFS discipline policy where a high variation of service requirements degrades the system performance significantly.

While the EPS discipline offers crucial insights into the performance of fair resource allocation mechanisms, it is limited in analyzing and designing differentiated scheduling algorithms. Strictly speaking, EPS is applicable only if resources (e.g., bandwidth, time-slots) are shared in a perfectly fair manner. However, it may not be the case in real systems [77]. Especially, the sym-

metry properties of the EPS discipline are not suitable for use in a system which transmits data from heterogeneous classes. Hence, the GPS and the DPS disciplines have emerged as natural generalizations for modelling the performance of such service differentiation mechanisms.

2.3.2 Generalized Processor Sharing (GPS)

Cohen [50] started using the term GPS to describe an extended PS with state-dependent service rates. But the seminal work on GPS done by Parekh and Gallager [7, 8] is more broadly adopted by convention.

In the GPS model [78], all traffic sent to the server is divided into different classes, and each class i is assigned a positive weight w_i as shown in Fig. 2.1. It is assumed a traffic class i is either an individual flow or includes several flows with similar QoS requirements. The specified weight w_i is the guaranteed minimum capacity for the class i . If any class does not fully consume its reserved capacity, other classes get the right to share the excess available capacity and the same weights are used to redistribute the excess capacity as well. The guaranteed minimum rate of class i means that class i receives service at rate w_i if all classes are backlogged. Let \mathbf{B} denote the set of backlogged classes and C denote the full service capacity. We have a backlogged class $i \in \mathbf{B}$ that receives a service rate [7] of

$$\frac{w_i}{\sum_{j \in \mathbf{B}} w_j} C \geq w_i C, \quad \text{for } \sum w_i \leq 1. \quad (2.3)$$

Due to different weights assigned to different classes, GPS is designed to achieve service differentiation among heterogeneous traffic networks. For instance, in [79], GPS is proposed as a packet-based mechanism to support real-time and best effort traffic, simultaneously, in link-sharing. It significantly

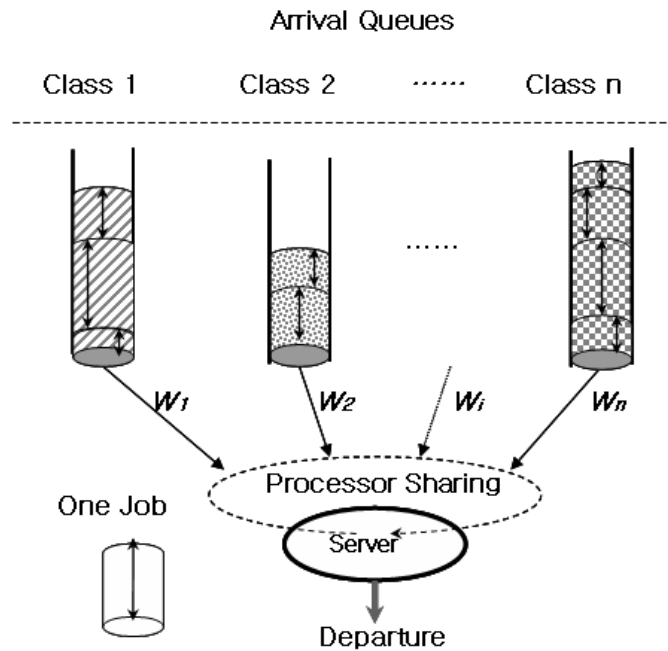


Figure 2.1: GPS model

enhances the modelling capabilities of ordinary PS. However, an underlying assumption of GPS is a requirement for infinitely divisible capacity, i.e., it is assumed that several classes can get service at the same time. It cannot be implemented with real traffic, as data traffic is composed of discrete packets which are sent sequentially. However, over a certain time scale, the fluid assumption of GPS is natural as packet sizes are very small compared to link capacity.

Approximations of GPS are required to identify implementable packet-based mechanisms. Parekh and Gallager have done well-known work to analyze packet-based GPS in [7, 8], also known as Weighted Fair Queueing (WFQ). A scheduling mechanism with WFQ calculates the service completion time of every arriving packet under GPS to determine the service order of packets. Clearly, the implementations are more complicated than with Weighted Round Robin (WRR) in which all classes are served in a fixed order.

On the other hand, results show that WFQ provides GPS within one packet transmission time and performs better than WRR.

In fact, asymmetric versions of the PS discipline including GPS and DPS have been very intractable so far, and there is no known closed-form solution for the exact distribution of the workload of an individual class, even under limited assumptions [80, 81]. Parekh and Gallager [7, 8] only provided delay guarantees for the worst-case and derived a deterministic upper bound for the delay and workload at each class by introducing a leaky bucket mechanism to control the traffic. A revisited approach in [82] uses a different type of mechanism known as the fractal leaky bucket to effectively handle Long-Range Dependent (LRD) traffic and get results for an upper bound on the delay and workload for each class. An exact analysis is only established for a special case, the two-classes GPS. Those results focussed on bounds and asymptotic approximations [83, 84, 85, 86, 87, 88, 89].

The asymptotic behaviour reveals some insight into how the system must have behaved when extremely rare events with small probabilities occur. Such events are typically used for relevant QoS measurements. Therefore, we focus on properties of stochastic majorization and stability issues but do not provide full distributional results in this section. Two asymptotic regimes are commonly of interest: (1) a large-buffer regime, which concerns the probability that the workload of an individual class in a GPS system exceeds a certain level; and (2) a many-sources regime, which concerns the probability that the buffer content exceeds a certain level due to the number of sources being increased. The investigation for these two regimes are normally approached by using specified traffic processes or multi-class traffic. Asymptotic expansions can often work as good approximations [90, 91, 92, 93].

Large-buffer asymptotics were obtained both for heavy-tailed and for light-

tailed traffic processes. Data traffic with highly variable or bursty characteristics often manifests LRD or self-similarity. Such traffic flows can behave very badly in the sense that they can grab all available capacity for a relatively long period [94, 95]. LRD or self-similar traffic can be modelled in so-called heavy-tailed distributions. In particular, traffic flows with Short-Range Dependent (SRD) or light-tailed properties will be adversely affected by heavy-tailed flows, when there is no protection mechanism. This phenomenon can be prevented in GPS by taking positive minimum weights. Yaron and Sidi [96] studied GPS stability and derived upper bounds on queue length considering GPS queues fed by exponentially bounded burstiness traffic [97]. Using the same traffic model, Zhang *et al.* [98] calculated upper bounds on the distribution of the backlog and delay for each GPS class, where the input processes of the classes can be dependent. Instead of studying the worst case behaviour, Veciana and Kesidis [99] worked on a broader class of light-tailed processes and calculated upper bounds on the logarithmic large-buffer asymptotics in a discrete time GPS system. Their results were extended by Zhang [100], where, under a similar assumption on the log-moment generating function, the exact logarithmic large-buffer asymptotics for a two-queue GPS model were derived. Then, he used these results in [101] to develop and compare several admission control schemes for both session-based and class-based service models. In [102], the large-deviations results for light-tailed traffic sources were provided. Similar work can also be found in [103] and [104]. For large-buffer asymptotics of the multiple-class GPS model, related results have been derived in [105] and [106]. Van Uiter and Borst [107] extended such work to networks of GPS queues.

Borst *et al.* [108] started the performance analysis of a multiple-class GPS system with heavy-tailed characteristics. A phenomenon referred to as

reduced-load equivalence is found when the average input rates of the classes are smaller than the GPS weights, the exact large-buffer asymptotics for a flow are shown to be equal to a flow in a system served by an isolated constant rate. Afterwards, more work on analysis of the buffer asymptotic with heavy-tailed traffic has been done by Borst *et al.* [109, 110, 111, 112], Jelenkovic and Lazar [113], Jelenkovic and Momcilovic [114], Agrawal *et al.* [115], Van Kessel *et al.* [116], Kotopoulos *et al.* [117] and Lelarge [118]. For certain assumptions on the traffic intensities, any heavy-tailed class is served at a constant rate and is immune from other classes with heavier-tailed characteristics, but only influenced by the average rates of the other classes. Beyond such regimes, however, the strong effect of heavier-tailed flows may be very obvious, known as *induced burstiness*. For a light-tailed class with a sufficient weight setting, the asymptotics combine the effect of heavy-tailed and light-tailed large-deviations behaviour, known as a *reduced-peak equivalence* [119], in an analogy with the term *reduced-load equivalence*. Similar types of qualitatively different regimes in different settings have been observed in [120, 121].

Many-sources asymptotics have been studied at the same time. Delas *et al.* [122] studied the asymptotics of the buffer occupancy distribution when buffers were accessed by a large number of stationary independent sources. The large-buffer asymptotics for a two-class GPS system in [117] has been extended to obtain many-sources asymptotics for GPS systems in [123, 124], where a GPS system with two queues and with multiple queues have been studied, respectively. Mannersalo and Norros [125] developed accurate approximations for a multiple-class GPS system with a large number of Gaussian inputs. Mandjes and Van Uitert [126, 127] further justified and refined these approximations, and established an interesting connection with tandem

queues.

So far, compared with the performance analysis of GPS with given weights, much less work has focussed on the choice of the GPS weights. Elwalid and Mitra [128] developed a framework for the design of GPS weights, based on the calculated loss probabilities of the sources. Kumaran *et al.* [129] presented an algorithm for getting on-line weight adaptation. Similar work on selecting the GPS weights can also be found in [130, 131].

For GPS, the analysis of the *actual delay* (or *sojourn time*) has been approached in only a few works. Analogously, the term *virtual delay* is used to define the time that is required for the buffer to be emptied, and is closely related to the workload. Shakkottai and Srikant [132] have shown the logarithmic many-sources asymptotics for the actual delay are equal to those for the virtual delay in a single discrete-time FCFS queue for certain assumptions on the input process, also applicable for two-queue priority models. Addie *et al.* [133] found results on the actual delay for a fluid particle in a special kind of on-off model. Borst *et al.* [134] studied the sojourn time of customers in a two-class GPS system under the assumption of a PS service discipline within the class.

For more details on GPS, readers may refer to [6, 78] and the references therein.

2.3.3 Discriminatory Processor Sharing (DPS)

Like GPS, the traffic under the DPS discipline is also divided into classes, and each class i is assigned a positive weight w_i . However, in contrast to GPS, w_i in DPS is associated with each individual job belonging to class i and not with the entire class queue, i.e., the service capacity is shared by all jobs according

to their weights. As shown in Fig. 2.2, all jobs present in the DPS system are served simultaneously at rates controlled by weights $w_i > 0, i = 1, 2, \dots$. Therefore, DPS can be reduced to an EPS if all jobs receive the same service rate, as in a multi-class extension of the EPS.

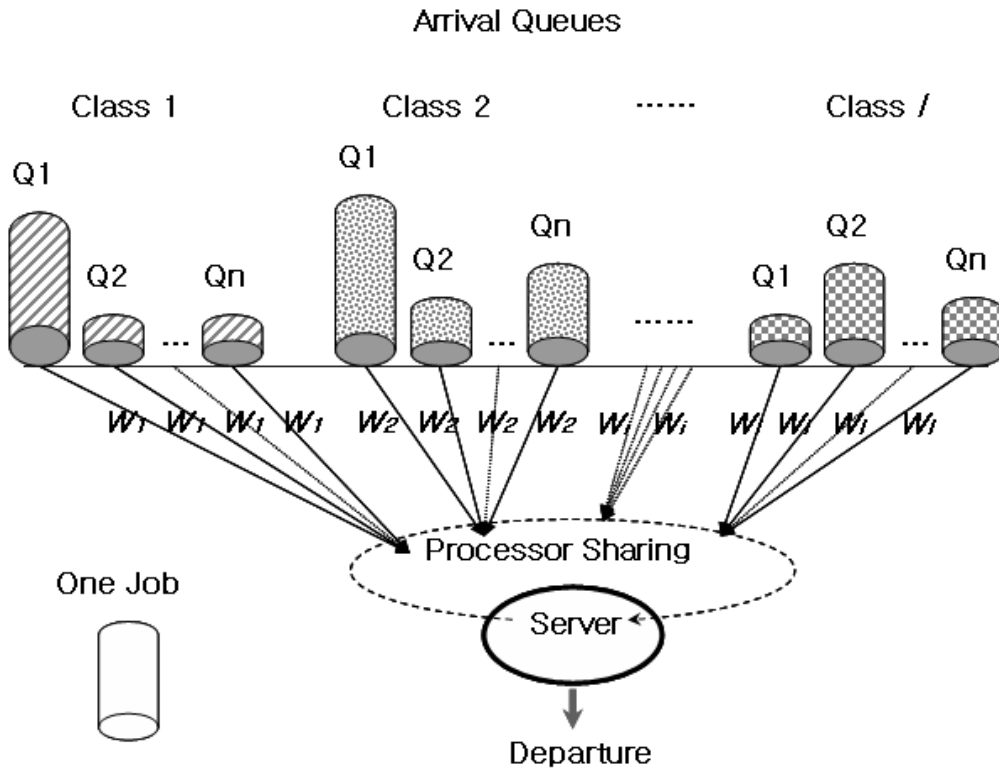


Figure 2.2: DPS model

In a DPS system, if n_i jobs are present for class $i, i = 1, 2, \dots, I$, the service rate of each class- i job R_i is [42]:

$$R_i(n_1, n_2, \dots, n_I) = \frac{w_i}{\sum_{j=1}^I w_j n_j} C. \tag{2.4}$$

As we can see, the capacity obtained by each class in the DPS is not guaranteed, and the allocated capacity to each class also depends on the number of jobs currently present in all the classes except weights w_i .

The DPS discipline provides a natural modelling approach for TCP flow-level performance with asymmetric bandwidth shares or for differentiated scheduling mechanisms. By adjusting weights, one can effectively control the instantaneous service rates of different job classes. As the weight used in DPS with regard to each active flow within the class, DPS can be an appropriate model to investigate the performance of heterogeneous systems.

Although DPS is similar to EPS with a simple model, the analysis of DPS model cannot be easily treated as an extension of the EPS model. Note that the property of insensitivity for EPS does not work in DPS, which is sensitive to specific traffic characteristics [135]. Neither the simple geometric queue length distribution, nor the tractable transform results for the sojourn time distribution, seem to exist in DPS, even for exponential service requirements. Some properties of EPS insensitivity only carry over to DPS in certain asymptotic regimes. Several asymptotic regimes will be reviewed here by reference to [42, 6, 60].

In 1967, Kleinrock proposed and studied the DPS model under another name, viz: *Priority Processor Sharing* in [39]. After his work, Fayolle, Mitrani, and Iasnogorodski [80] made an important move on DPS analysis to obtain the expected conditional sojourn times as the solution of a system of integro-differential equations for the M/G/1 DPS queue. The error contained in O'Donovan's [136] original system of equations was corrected. For the case of exponentially distributed service requirements, a closed-form expression for the conditional mean sojourn time has been derived, and the unconditional mean sojourn times can be obtained from the system of linear equations. The derivation of these equations referred to the methods used by Kleinrock *et al.* [5, 137] for a processor sharing queue with batch Poisson arrivals. It shows the asymptotic ratio (also known as the *slowdown ratio*) of the condi-

tional mean sojourn time and the service requirement in DPS is insensitive and independent of the job class. Extending this result, Avrachenkov *et al.* [138] showed that the expected unconditional sojourn time is finite and has an asymptote with slope $1/(1 - \rho)$.

The DPS queue in heavy traffic has been analyzed by Grishechkin [71] with an assumption of finite second moments for the service requirement distributions. The asymptotic regime has been discussed under heavy load conditions, i.e., $\rho = \sum_{i=1}^I \rho_i \rightarrow 1$. Rege and Sengupta [81] made a further advance in obtaining the moments of the queue length distribution as the solution to linear equations for the case of Poisson input and exponential service requirements. The joint queue length distribution has been shown to be limited. This work was extended by Van Kessel *et al.* [139] for the phase-type service requirement distributions and Kim and Kim [140] for bulk arrivals. The case when the DPS queue is in overload, i.e., $\rho > 1$, has also been studied by Altman *et al.* [141] based on techniques for EPS analysis in [142]. It shows that the queue size of any class grows asymptotically and linearly with a rate. This rate depends on the service requirement distribution in a complex manner – not just on the mean. Eventually, all jobs in the system can be finished in a finite time.

Under the time-scale decomposition regime, which was identified by Bonald and Proutiere [135], Van Kessel *et al.* [116] found that the queue length distribution turns out to be insensitive. Especially, the limiting distribution of the relatively slow dynamic traffic class is independent of the weights and also insensitive to the service time distribution, where flow sizes follow phase-type distributions. It is different with the case of generally distributed flow sizes; where the performance is affected by the distributions of flow sizes for all classes [80]. Another recent work by Boxman *et al.* [143] showed that

a time-scale decomposition approach provides a good approximation to finite capacity DPS with less computational effort even if the time-scales of classes are different.

Another asymptotic regime is for heavy-tailed service requirements. It is an extension of the asymptotic equivalence for EPS established by Zwart and Boxma [57]. Borst *et al.* [144] proved, if the service requirement distribution has finite variance, a similar tail behaviour exists between the service requirement and sojourn time for any DPS class, which is independent of the DPS weights. The additional scenarios and a broad assumption concerning service distributions have been considered in [60]. Note that these insensitivity properties do not appear for light-tailed service requirements [60]. For such cases, Egorova *et al.* [145] focused on the logarithmic estimates of sojourn-time distribution using large-deviation techniques.

More conjectures about approximate insensitivity are also considered related to statistical bandwidth sharing, and readers can refer to [146, 147, 77]. Although several decades have passed, theoretical studies relevant to DPS in the literature are still scarce and results are only under certain limiting regimes.

2.3.4 Multilevel Processor Sharing (MLPS)

Another extension of the EPS policy is the family of MLPS strategies introduced by Kleinrock [5]. The MLPS discipline is non-anticipating and work-conserving, parameterized by the specified set of level thresholds $a_0 < a_1 < \dots < a_{N+1}$ that are used for classifying jobs based on the size of their attained services, where the level $n = 0, 1, \dots, N, N + 1$ and $a_0 = 0, a_{N+1} = \infty$. If a job belongs to level n , the attained service must be at least a_{n-1} but less

than a_n . There is a strict priority applied between levels. The lower level where jobs have smaller amounts of attained service will receive the higher service priority. The EPS (simply called PS below) policy or the Foreground-Background Processor-Sharing (FBPS or FB) [148, 149] (also called as Least-Attained-Service) or FCFS may be applied for serving at each level n . Note that there are different studies for FB as an independent discipline and as a mechanism for MLPS queues. We review the literature for both cases related to the complete discussion of MLPS.

In a larger family of size-based scheduling disciplines, FB performs the smaller mean sojourn time when the service time distribution has a decreasing hazard rate (DHR) [67]. Righter [150] also proved that FB minimizes the queue size stochastically with the DHR service distribution. Righter *et al.* [151] showed that FB minimizes the mean sojourn time if the service time distribution has an increasing mean residual life (IMRL), which is a weaker condition than DHR. Wierman *et al.* [152] proved that FB is better than PS with respect to the mean delay whenever the service time distribution is the type of DHR, and vice versa if the service time distribution is the type that has an increasing hazard rate (IHR). For further details about the FB discipline, readers can refer to Nuyens and Wierman's recent survey [153]. Aalto *et al.* [154] proved that the MLPS disciplines are better than PS with respect to the mean sojourn time for the two-level case whenever the hazard rate of the service time distribution is decreasing or increasing. Afterwards, they showed that these results are valid for any MLPS discipline in [155]. Therefore, a range of MLPS disciplines seems a reasonable compromise between PS and FB, having a smaller overall mean sojourn time than PS and better fairness mean delay than FB.

Alto *et al.* [6] summarised recent results for the MLPS discipline, includ-

ing the mean sojourn time, the mean slowdown ratio, and the expected conditional sojourn time asymptotics.

- The mean sojourn times for MLPS and PS are denoted by $\mathbf{E}[T_{MLPS}]$ and $\mathbf{E}[T_{PS}]$, respectively. According to [156], MLPS with the internal disciplines FB or PS, if service requirements is IMRL, $\mathbf{E}[T_{MLPS}] \leq \mathbf{E}[T_{PS}]$; if it is decreasing mean residual life (DMRL), $\mathbf{E}[T_{MLPS}] \geq \mathbf{E}[T_{PS}]$. In contrast with [151], it showed that FB does not minimize the mean sojourn time under the situation of IMRL [157].
- There is a natural partial order among the internal disciplines applied in MLPS. For DHR or IHR service requirements, the mean sojourn time is decreased or increased which is not affected by changing internal disciplines from PS to FB, or from FCFS to PS, or by any level split in FCFS internal discipline into two adjacent FCFS levels; or by splitting level 1 with PS internal discipline into two adjacent PS levels [158]. Results of [158] proved the optimality of FB on the mean sojourn for DHR service requirements [159, 154] and quantified the reduction on the mean sojourn time by adding levels.
- Considering the mean slowdown ratio and results in [158], if a specific function for all service requirements is decreasing or increasing (see [6] 5.3), the mean slowdown ratio is decreased or increased under any conditions in the last mentioned item. Feng and Misra [159] originally proved the optimality of FB with respect to the mean slowdown ratio for DHR service requirements by the same approach.
- Comparing MLPS having PS as an internal discipline at the highest level with ordinary PS, there is an asymptotic expected conditional so-

jour time with slope $1/(1 - \rho)$ and a positive finite bias [160]. Also, the asymptotic slowdown of these disciplines is exactly the same. These results for the asymptotic expected conditional sojourn time and slowdown show that the performance of very large jobs is equivalent under both MLPS having PS and ordinary PS disciplines.

On the other hand, MLPS disciplines have recently been resurrected in some papers that focus on the differentiation between short and long TCP flows in the Internet [161, 159, 162]. Flow sizes in the Internet have been modelled by Pareto and hyperexponential distributions [163, 164]. The size-based scheduling MLPS attracts more interest as it is very often that file sizes are extremely variable and have heavy-tailed characteristics [165, 166, 163, 94].

2.4 Wireless Applications: PTT and WiMAX

2.4.1 PTT over GSM/GPRS

Push to Talk [10], a new mobile service, was first introduced by US-based network operator Nextel in 1996 [167]. Unlike a normal mobile communication in a full duplex transmission, PTT operates in a half-duplex mode to provide fast-access two-way communication between two or more communicating parties. PTT voice packets are transmitted in just one direction at any given moment, so the operation is similar to a conventional “walkie talkie”. The right to talk is transferred from one end of the conversation to another through a push of a button on the mobile terminal.

PTT features shorter call setup times, presence detection and point-to-point or point-to-multipoint communication modes. Moreover, existing mobile

phone infrastructure without any major change can deliver PTT, such as Code Division Multiple Access (CDMA), or GSM.

As PTT is an additional service that does not replace normal cellular mobile communications service, it can increase revenue to the wireless network operators. Rather than the traditional mobile service being good for long and interactive communication, PTT aims for demands of quick communications among end-users. For the significant concerns on Public Safety, or serving the need for efficient communications during times of emergency or disaster, PTT is considered as an important communication tool by government and public safety officials. Except to evolve or represent some services provided by the cellular mobile networks, PTT service benefits the customers with community-of-interest, such as younger consumers, just as Short Message Service (SMS) has been popular with this demographic. At present, PTT is also available both in Europe and in Australia.

The Open Mobile Alliance (OMA) was established to ensure mobile data service inter-operability across different devices, geographies, service providers, operators and networks. Based on the build up standards of PTT over Cellular (PoC) [168], OMA's initial work on PTT services focuses on requirements to develop specifications for an open standard to enable adoption of a PoC service over mobile networks.

In Europe and Australia, PTT services are provided through the GPRS over GSM networks that utilise the GPRS "always-on" feature, which reduces access delay. GSM uses a time division multiple access (TDMA) structure with eight slots per frame to support speech and data transmission [169]. Second-generation GSM networks deliver voice and data services with high quality and security and they have full roaming capabilities. The GSM Phase 2+ standard specifies GPRS which is a packet data communication system

using GSM physical channels [170]. However, since GPRS was originally designed for data packet transmission, there is a concern that PTT/GPRS will not meet PTT QoS requirements such as delay and jitter standards for a voice service. Therefore, it is important to verify that PTT/GPRS meets PTT QoS requirements. In Chapter 4, based on a multiservice PS model, we provide an analytical model to analyse the performance of the PTT uplink delay.

2.4.2 WiMAX

Following the increasing demand for multimedia services we have now reached a point where it is commercially justifiable to deploy wireless Internet broadband access networks. In situations where either it is difficult to use a wired technology, or its cost is too high, wireless broadband access is a viable alternative. This economic environment has led to the development of the IEEE 802.16 standards [21] for broadband wireless access in metropolitan area by Working Group 16 of the IEEE 802 committee. In addition, the WiMAX Forum [11], a non-profit industry consortium, has been formed. It is chartered to promote the technology and provide compatibility and interoperability certification of 802.16-based products. Accordingly, the term WiMAX is often used to mean the IEEE 802.16 technology.

WiMAX enables efficient and reliable broadband access by thousands of subscribers either in line-of-sight or in non-line-of-sight (NLOS) conditions. However, given that the development of the line-of-sight air interface specification has been discontinued, WiMAX technology will mainly be used in NLOS conditions. There are two modes specified in the IEEE 802.16 standard for sharing wireless media, namely, point-to-multipoint (PMP) and mesh (optional). We shall focus only on the PMP mode in this thesis. As we illustrate

in Fig. 2.3, under PMP, a single base station (BS) serves multiple SSs, each of which can be shared by multiple devices. Hence, a typical SS has multiple uplink connections generating data to be transmitted to the BS. The coverage of a BS has a NLOS range varying from 6 to 10 kms depending on the type of obstacles present (trees, huge buildings etc.). Theoretically, in a single Tx/Rx channel, WiMAX can provide data rates of up to 75 Mb/s on both the uplink and downlink channels. Communications between an SS and external networks or among SSs have to be carried out via a BS.

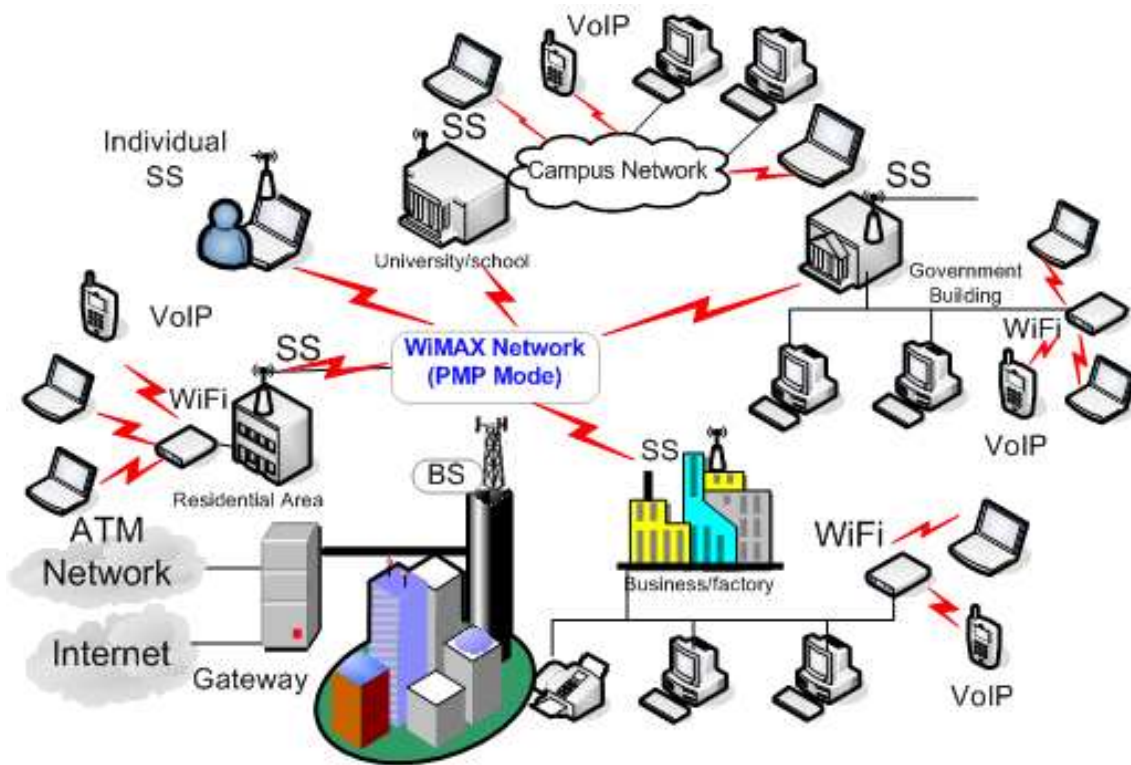


Figure 2.3: WiMAX Network (PMP Mode).

WiMAX aims to deliver efficiently various types of services including normal data, streaming video and VoIP meeting required QoS. QoS guarantee in WiMAX is provided by the scheduling mechanisms in the MAC layer of the protocol stack, where the aggregate bandwidth is granted by the BS. Based

on this, the uplink scheduler of an SS chooses a packet from the appropriate connection at each time-slot in order to deliver the required QoS. Clearly, in such a complex traffic environment, QoS delivery depends on the performance of the algorithm deployed by the uplink scheduler. However, as no specified scheduling algorithm in the standard, choosing efficient WiMAX access scheduling algorithms and having the means to evaluate their performance are important challenges for networks operators.

In Chapter 5, we propose a priority-based fair scheduling algorithm to handle both real-time and non-real-time uplink WiMAX traffic in an SS and provide an analytical model to analyse the performance of our proposed algorithm.

2.5 Conclusions

Together with their associated IEEE standards, we have described the functions of MAC protocols in wireless networks and discussed the importance of the scheduling policy and its performance model for delivery of QoS requirements at the MAC layer. We reviewed the developed wireless scheduling algorithms and then we focussed on modelling scheduling schemes. We presented an overview of work on different PS models which have been presented in the current literature. We also briefly introduced the background for two wireless applications: PTT and WiMAX, which we propose to illustrate the use of the PS scheduling policy and will be modelled by the multiservice multiqueue PS model in later chapters.

Chapter 3

Multiservice Multiqueue

Processor Sharing Models

3.1 Introduction

3.1.1 Multiservice PS Scheduling

Multiservice PS scheduling is another extension of ordinary PS scheduling policy, which is able to handle different traffic classes with different priorities. Multiple users/connections are divided into a number of groups and served in priority order. But the implementation of this priority discipline is more flexible than MLPS, as jobs do not have to be classified based on the size of their attained services. Normally, multiservice PS scheduling has been considered for use in time-sharing systems. In this traditional approach, a user from any priority holding the service token can only receive a fixed portion of service time, called the *service quantum*. Like most PS-related disciplines, multiservice PS scheduling is work-conserving.

The multiservice PS scheduling policy can be implemented in multimedia

networks without the need for complicated calculations. Thus, it is of particular benefit to wireless systems where network schedulers need a simple scheduling algorithm to make a fast decision. As the discipline provides service in a prioritized manner, real-time traffic can be protected by giving it a higher priority. Moreover, the delay performance can be modelled and analysed with a closed-form solution under certain practical assumptions.

However, for multimedia applications with heterogeneous traffic characteristics, quantum-based multiservice PS scheduling is inefficient for large size requests as the quantum is normally set according to the smallest size request. We propose a new approach which allows a priority-based service quantum to a certain group of users rather than a fixed service quantum for every group.

3.1.2 Multiservice PS Scheduling Models

Potter and Zukerman [9] proposed a multiqueue processor sharing (MPS) model for multiservice PS scheduling with the fixed service quantum. The inspiration for the MPS model came from a classical model, known as the round robin processor sharing (RRPS) model [38, 39], and the distributed queue dual bus (DQDB) protocol [172, 173, 174, 175, 176, 177].

The RRPS model introduces a round robin time-shared service system to serve a single queue, where arriving jobs with different size of service demands join the end of the single queue immediately when they arrive. When a job moves to the head of the queue, it receives a certain quantum of service (also known as the “time-slot” or “segment-time” below), and if more service is required, the incomplete part of the job goes to the tail of the queue. From this basic model, the feedback queueing model appeared as a generalization of

RRPS [5]. It replaces the single queue by a system of queues and the classical round robin model by a scheduling algorithm which determines the service order and the service quantum given to a job. As for the original RRPS, if a job needs more service after receiving its service quantum, the remaining part stays in the system until the whole job is finished. But the unfinished job re-queues at the system of queues in the feedback queueing model rather than the single queue in the RRPS model .

The basic DQDB model introduces the idea of local queues (LQs) for buffering arrivals which are waiting for transmission and a distributed queue (DQ) with a server working in priority order. The concept of a fixed service quantum is also used in this model. Moreover, at each priority, only one segment from each LQ is allowed to queue in the DQ at a time. The service manner in the DQ closely approximates RRPS under certain conditions (refer to Section IV of [9]). However, DQDB is difficult to analyse by modelling the DQ protocol. Notice that although the PS model is normally work-conserving, in [178], DQDB with bandwidth balancing [179] is considered as a non-work-conserving PS model that applies to IEEE 802.6 [172].

The MPS model falls under the category of Kleinrock's feedback queueing model and inherits some features of DQDB construction. Like DQDB, the MPS model involves several LQs and a central server with a PS queue. The server performs prioritized round robin processor sharing among these LQs. As a variation of the classical RRPS, MPS adds an infinite buffer at each LQ and uses a modified PS server handling the multi-priority case. For each priority, only one representative from a LQ can be present in the PS queue. After receiving its service quantum from the shared processor, a customer whose service is incomplete will be recycled to the tail of its own priority group within the PS queue. The MPS model extended the ordinary PS model to a

multiservice multiqueue PS model and the result relies on the conventional PS assumption, viz: an egalitarian service sharing manner among all users. An exact analysis of MPS with a closed-form solution for the mean message delay was presented in [9].

Modelling our new proposed multiservice PS scheduling, we generalize the MPS model of [9] to the case where the service quanta are different for different services. We call this new model the Priority-based Service Quanta (PBSQ) model. In the later sections, we derive an accurate approximation for the mean message delay and demonstrate the performance effect of the choice of service quantum based on implementation issues under the PBSQ model.

3.2 Multiqueue PS (MPS) Model

We review the MPS model first. In this discrete time model, time is divided into consecutive equal-length time units called *time-slots* which are related to the service time of *slots*, and the time points at the beginning of each time-slot are designated by $1, 2, 3, \dots$, so the k th time-slot is the time interval $[k, k + 1)$. Messages arriving at a priority p LQ within any time-slot are independent and identically distributed (i.i.d.) and are also independent of arrivals to other LQs. All arrivals within any time-slot are assumed to arrive at the boundary of the time-slot. Message lengths for each priority are discrete i.i.d.. Each message is assumed to consist of an integral number of units called *segments*. A segment is equal to one slot in the MPS model, and corresponds to a time-slot representing a fixed service quantum/time received by the message from the PS server at each service time. The message length distribution may be different for different priorities.

The closed-form result for the mean delay (measured in time-slots) of pri-

riority p messages was obtained as a function of the number of segments it contains, where the delay of a priority p message was analysed in two components: the waiting time in the LQ, L_p and the time spent in the PS queue until the whole message is completely transmitted, $S_p(n)$.

Based on previous studies of RRPS, an expression for $S_p(n)$ was obtained using a similar approach to that of [5] (p.168) by considering a very large test message. For the derivation of L_p , they considered an equivalence with the MPS model in terms of the average segment delay, which is the discrete-time M/G/1 with the preemptive resume priority queueing model [180]. Both queueing systems follow a strict priority discipline and are work-conserving at the segment level. Therefore, for any steady-state, statistically, if the total segment arrival processes into both systems are identical, the distribution of the total number of segments of each priority at any time-slot will be identical. By Little's formula, knowing two systems have equal average queue size and equal average arrival rate, the mean delay of segments for each priority will be equal in these systems. Notice that the segment delay distributions are different for the two models – even with equal means, as the distribution of delay usually depends on the service order. So L_p can be obtained by equating the mean of the priority p segment delay for the MPS model and the mean for the discrete-time M/G/1 with the preemptive resume priority model.

3.3 Notation

Table 3.1 provides a detailed summary of the main mathematical notation that has been used throughout this thesis for the various proposed models. Also, in Table 3.2, we provide the notation for references of [9], that are equivalent to our notation.

Symbol	Explanation
$E(X)$	The mean of the random variable X .
$\text{Var}(X)$	The variance of the random variable X .
$P(X)$	The probability of the random event X .
p	A priority in the system, $p = 1, 2, \dots, P$, where a smaller number indicates a higher priority.
a_p	A random variable representing the number of priority p messages arriving at a priority p LQ during a time-slot.
\bar{a}_p	The mean of a_p , $\bar{a}_p = E(a_p)$ [messages/time-slot].
b_p	A discrete random variable representing the priority p message size.
\bar{b}_p	The mean of b_p , $\bar{b}_p = E(b_p)$ [segments].
N_p	The length of a priority p segment in slots, $1 \leq N_p \leq N$.
$C_{a,p}^2$	The squared coefficient of variation of b_p , $C_{b,p}^2 = \text{Var}(b_p)/\bar{b}_p^2$.
$C_{b,p}^2$	The squared coefficient of variation of a_p , $C_{a,p}^2 = \text{Var}(a_p)/\bar{a}_p^2$.
M_p	The number of the LQs at priority p .
$Q_p(0)$	The mean number of priority p messages in the PS queue which have not got any service.
$F_{b,p}(n)$	The probability of $b_p \leq n$, $n = 1, 2, \dots$
$R_p(k)$	The probability that a randomly selected priority p segment is the k th in its own message, $k = 1, 2, \dots, b_p$.
L_p	The mean waiting time (or delay, as the propagation delay from an LQ to the PS queue is assumed to be zero) of priority p message in the LQ.
$S_p(n)$	The mean time of a priority p message in the PS queue till completing at least n segments service, where n is a random variable.
$D_p(n)$	The mean delay of priority p message with n segments, where n is a random variable.
λ_p	The mean total arrival rate for all priority p LQs.
ρ_p	The traffic load of priority p messages.
ε_p	$\varepsilon_p = \sum_{i=1}^p \rho_i$.
Δ_p	Compensation term for a priority p message to value the difference between the delay of PBSQ model and the delay predicted by the MPS model.
h_p	The size of MAC header [segments], $h_p < 1$.

Table 3.1: Mathematical notation used throughout this thesis

Symbol in thesis	Symbol in [9]
$p = 1, 2, \dots, P$, where 1 is the highest priority.	$p = 1, 2, \dots, H$, where H is the highest priority.
N_p denotes the length of a priority p segment in slots.	$N_p(0)$ denotes the mean number of priority p messages in the PS queue whose segments have not got any service.
$Q_p(0)$ denotes the mean number of priority p messages in the PS queue whose segments have not got any service.	
$\varepsilon_p = \sum_{i=1}^p \rho_i$	$\sigma_p = \sum_{q=p}^H \rho_q$
$\nu_i = \rho_i b_i (C_{b,i}^2 + \lambda_p C_{a,p}^2 / M_p)$ for the MPS model.	$\nu_i = b_i (C_{b,i}^2 + \lambda_p C_{a,p}^2 / M_p)$ for the MPS model.

Table 3.2: Mathematical notation used in other references

Note that L_p is not related to the message size n and $D_p(n)$ is the mean delay of priority p message with n segments, where n is a random variable. We can obtain the mean delay of all priority p messages by $\sum_{n=1}^{Max(b_p)} D_p(n)P(n)$. In [5] and [9], it has proved that $S_p(n)$ and $D_p(n)$ are linearly increasing with n in the RR model or the MPS model. We will prove that this is also true for our PBSQ model in the later part of this chapter. Therefore, if the mean message size for priority p is \bar{b}_p , the mean delay of all priority p messages can be simply represented by $D_p(\bar{b}_p)$. We will use this result directly in the following chapters of this thesis.

3.4 Correction of a Subtle Incongruity in the LQ and the PS Mean Delays of the MPS Model

In the MPS model, the mean delay of a priority p message of length n segments is given by

$$D_p(n) = S_p(n) + L_p.$$

We have reviewed the method to obtain $S_p(n)$ and L_p respectively in Section 3.2. Note that the total delay $D_p(n)$ in the MPS model has been proved correct [9]. In this section, we demonstrate just the apportionments to LQ (L_p) and to PS queue ($S_p(n)$) are slightly off and provide a correction.

3.4.1 The Subtle Incongruity for L_p and $S_p(n)$

To demonstrate the incongruity, we consider L_p obtained by the equation (13) of [9] for a simple case with only one priority that has four LQs. A constant size of 1 segment is used for all arriving messages and the arrivals follow a Poisson process with a mean arrival rate of 0.1 messages per time-slot at each LQ. Based on (13) of [9], we obtain a negative-valued result for L_p , which is -0.095. This contradicts the physical interpretation of L_p that must be non-negative.

Let us analyse the reason for this incongruity. Based on Kleinrock [5] and Appendix I of [9], equation (2) of [9] was obtained as

$$S_p(n) = n \left(\frac{Q_p(0)}{\lambda_p} \right), \quad (3.1)$$

where $Q_p(0)$ denotes the mean number of priority p messages in the PS queue which have not received any service. According to the classical method de-

scribed in Kleinrock [5] and Appendix I of [9], it has proved that $S_p(n)$ is linearly increasing with n . In (3.1), $S_p(0)$ is considered to be zero following the notation in the original round robin (RR) queue (p.168 of [5], using the symbol “ $T(0)$ ” there, representing the starting time that a customer enters the RR system).

However, there is a distinction between $S_p(0)$ in the MPS model and $T(0)$ in the RR queue. Unlike Kleinrock’s RR with a single queue, the MPS model is assumed to consist of LQs and a PS queue, where the propagation delay from an LQ to the PS queue is assumed to be zero. A message arrives at LQ first, then enters the PS queue directly if there is no message belongs to the same LQ in the PS queue; otherwise, if another message belonging to the same LQ is existed in the PS queue, the new arrival waits at the LQ, and then moves forward to the head of line (HOL) position of the LQ. For the second case, the HOL message is still waiting at its LQ when the last packet of the previous message belonging to the same LQ in the PS queue is moving out. However, following the consideration of the traditional RR, the PS queue was looked as the single queue [9]; the start time of the LQ HOL message moving into the PS queue was calculated from the start time that the last packet of the previous message belonging to the same LQ in the PS queue is served. For such a case, the LQ HOL message was supposed to be in the PS queue already, but this is incompatible with the original assumption of the MPS model in which the HOL message should still be in the LQ until the last packet of the previous message belonging to the same LQ is completely served. Therefore, $S_p(n)$ is overestimated by (2) of [9]. This period should be taken into account in L_p rather than in $S_p(n)$. Let δ_p denote the mean of this period. We give the

corrected expression for $S_p(n)$ as

$$S_p(n) = n \left(\frac{Q_p(0)}{\lambda_p} \right) - \delta_p. \quad (3.2)$$

Note that the incongruity occurs when a message has to wait for the service in its LQ. However, the probability for this case, i.e. when a message comes, another message from the same LQ is in PS queue, is complicated to analyse for different arrival and service distributions.

To provide a correction for the incongruity, we modify the original assumption of [9] to simplify the case with the incongruity. Assume that a message staying in a LQ enters the PS queue only when the PS server gives the service token to this message, i.e. the beginning time of a message in the PS queue is when the first segment of the message starts to be served. The difference of this assumption and the original MPS model assumption is only the start time when a message enters the PS queue and does not affect the total delay. But it generalizes the case to be analysed easily that we can get a solution of the incongruity. Based on RR properties summarized by Kleinrock [5] (p.169), the ratio of wasted time to service time in RR is “ $W(x)/x = \rho/(1 - \rho)$ ”, which measures how much waiting time, on average, must be sacrificed for receiving a unit of service time. Note that the ratio only relates to the traffic load “ ρ ”. This can be used to describe the mean time of priority p HOL messages waiting for the last segment of the previous message (can be from any LQ) completing service in the PS. For the multi-priority conservation system, the traffic load “ ρ ” is only due to the load of equal or higher priorities [5] (p.124); so it should be $\varepsilon_{p-1} + \rho_p \frac{M_p-1}{M_p}$ in the MPS model, as a LQ cannot contribute any traffic load to the PS queue if there is a message belonging to this LQ in the PS queue already. Under this approach, the solution of the incongruity for L_p

is

$$\frac{\varepsilon_{p-1} + \rho_p \frac{M_p-1}{M_p}}{1 - \varepsilon_{p-1} - \rho_p \frac{M_p-1}{M_p}}.$$

In Appendix A, we prove that L_p would never be negative under this approach with the above correction.

As the total delay in the MPS model is correct, for simplicity, in the following chapters, we still keep the original assumption of [9], and only use δ_p to denote the incongruity.

3.4.2 Derivation of LQ Waiting Time

We re-derive the equations which are affected by this incongruity in the MPS model. Again, consider a priority p “test” message of length x segments as in [9], whose probability of occurrence does not affect the overall statistics. Then, the time that this message spends in the PS queue, $S_p(x)$, must approach its own service requirement x , plus the time required by the total work for all messages which arrive to the PS queue during its service and waiting time but just before its last segment (i.e. the x th segment) starts service (the duration is $S_p(x) - 1$). These arriving messages include arrivals to priority p LQs other than the one that has this test message, given by $(S_p(x) - 1)\rho_p(M_p - 1)/M_p$; and the arrivals to all local queues of priorities higher than p , given by $(S_p(x) - 1)\sum_{i=1}^{p-1}\rho_i$. The corrected expression for (3) of [9] is shown as

$$S_p(x) \rightarrow x + (S_p(x) - 1)\left[\frac{M_p - 1}{M_p}\rho_p + \sum_{i=1}^{p-1}\rho_i\right], \quad \text{as } x \rightarrow \infty. \quad (3.3)$$

Defining

$$\varepsilon_p = \sum_{i=1}^p \rho_i, \quad (3.4)$$

we obtain:

$$\lim_{x \rightarrow \infty} S_p(x) = \frac{x}{1 - \varepsilon_{p-1} - \frac{M_p-1}{M_p} \rho_p}. \quad (3.5)$$

When $x \rightarrow \infty$, the slope of mean PS delay is:

$$\frac{Q_p(0)}{\lambda_p} = \frac{1}{1 - \varepsilon_{p-1} - \frac{M_p-1}{M_p} \rho_p}, \quad (3.6)$$

same as (6) of [9]. By (3.2), we obtain

$$S_p(n) = \frac{n}{1 - \varepsilon_{p-1} - \frac{M_p-1}{M_p} \rho_p} - \delta_p, \quad (3.7)$$

and have

$$D_p(n) = L_p + \frac{n}{1 - \varepsilon_{p-1} - \frac{M_p-1}{M_p} \rho_p} - \delta_p. \quad (3.8)$$

Let $R_p(k)$ for $k = 1, 2, \dots, b_p$ be the probability that a randomly selected priority p segment is the k th in its own message, where the message size b_p replaces ∞ in (9) of [9]. According to Appendix II of [9], we have

$$R_p(k) = \frac{1 - F_{b,p}(k-1)}{b_p} \quad \text{for } k = 1, 2, 3, \dots, b_p. \quad (3.9)$$

Then, the mean priority p segment delay is obtained as

$$\begin{aligned}
E[D_{seg(p)}] &= \sum_{k=1}^{b_p} [L_p + S_p(k)] R_p(k) & (3.10) \\
&= \sum_{k=1}^{b_p} [L_p R_p(k)] + \sum_{k=1}^{b_p} [S_p(k) R_p(k)] \\
&= L_p + \sum_{k=1}^{b_p} [S_p(k) R_p(k)] \\
&= L_p + \sum_{k=1}^{b_p} \left[\left(\frac{k}{1 - \varepsilon_{p-1} - \frac{M_p-1}{M_p} \rho_p} - \delta_p \right) \times \frac{1 - F_{b,p}(k-1)}{\bar{b}_p} \right] \\
&= L_p + \frac{\sum_{k=1}^{b_p} k [1 - F_{b,p}(k-1)]}{\bar{b}_p (1 - \varepsilon_{p-1} - \frac{M_p-1}{M_p} \rho_p)} - \delta_p.
\end{aligned}$$

Then, we simplify the numerator of the second term in the above result as

$$\begin{aligned}
&\sum_{k=1}^{b_p} k [1 - F_{b,p}(k-1)] \\
&= \sum_{k=1}^{b_p} [k P(b_p \geq k)] \\
&= \sum_{k=1}^{b_p} [k \sum_{j=k}^{\infty} P(b_p = j)],
\end{aligned}$$

where the random variable $b_p = 1, 2, 3, \dots$ so $j = 1, 2, \dots, \infty$. If we expand the above summation, we get:

$$\begin{aligned}
&1P(b_p = 1) + 1P(b_p = 2) + 1P(b_p = 3) + 1P(b_p = 4) + \dots \\
&\quad 2P(b_p = 2) + 2P(b_p = 3) + 2P(b_p = 4) + \dots \\
&\quad\quad 3P(b_p = 3) + 3P(b_p = 4) + \dots \\
&\quad\quad\quad \dots + \dots
\end{aligned}$$

We sum the terms vertically, the summation becomes

$$\begin{aligned}
& \sum_{k=1}^{\infty} [\text{P}(b_p = k) \sum_{j=1}^k j] \\
&= \sum_{k=1}^{\infty} [\text{P}(b_p = k) \frac{(1+k)k}{2}] \\
&= \frac{1}{2} \left[\sum_{k=1}^{\infty} k \text{P}(b_p = k) + \sum_{k=1}^{\infty} k^2 \text{P}(b_p = k) \right] \\
&= \frac{1}{2} [\bar{b}_p + \text{E}(b_p^2)] \\
&= \frac{1}{2} (\bar{b}_p + \bar{b}_p^2 + \bar{b}_p^2 C_{b,p}^2),
\end{aligned}$$

where $C_{b,p}^2 = (\text{E}(b_p^2) - \bar{b}_p^2) / \bar{b}_p^2$.

$$\begin{aligned}
E[D_{seg(p)}] &= L_p + \frac{\bar{b}_p + \bar{b}_p^2 + \bar{b}_p^2 C_{b,p}^2}{2\bar{b}_p(1 - \varepsilon_{p-1} - \frac{M_p-1}{M_p} \rho_p)} - \delta_p \\
&= L_p + \frac{1 + \bar{b}_p + \bar{b}_p C_{b,p}^2}{2(1 - \varepsilon_{p-1} - \frac{M_p-1}{M_p} \rho_p)} - \delta_p \\
&= L_p + \frac{[\bar{b}_p(1 + C_{b,p}^2) + 1]/2}{1 - \varepsilon_{p-1} - \frac{M_p-1}{M_p} \rho_p} - \delta_p. \tag{3.11}
\end{aligned}$$

Following the notation presented in Table 3.2, we rewrite (12) of [9] as

$$E[D_{seg(p)}] = \frac{\nu_p / \rho_p + \sum_{i=1}^p \nu_i / (1 - \varepsilon_p)}{2(1 - \varepsilon_{p-1})} + \frac{1}{2}. \tag{3.12}$$

Equating it with (3.11), we obtain the following corrected expression for L_p :

$$L_p = \frac{\nu_p / \rho_p + \sum_{i=1}^p \nu_i / (1 - \varepsilon_p)}{2(1 - \varepsilon_{p-1})} - \frac{[\bar{b}_p(1 + C_{b,p}^2) + 1]/2}{1 - \varepsilon_{p-1} - \frac{M_p-1}{M_p} \rho_p} + \delta_p + \frac{1}{2}. \tag{3.13}$$

Notice that the modified solution adds a term of δ_p to (13) of [9], where the bigger δ_p would increase L_p . Without this term, in [9], the $S_p(n)$ was over-

estimated and L_p was underestimated through the equivalence of (11) and (12). However, due to L_p obtained by the equivalence, the final result of the mean message delay $D_p(n)$ shown at (14) of [9] is still correct. But the new solution for the mean waiting time in LQ by adding δ_p is really matched with the model.

The physical interpretation and explanation for δ_p can be viewed as the answer to the question: What is the length of the time period spent between the LQ and the PS queue? It is based on the separate point of time in the LQ and in the PS queue. By a slightly modified assumption, our approach enables to quantify the period.

3.5 Priority-Based Service Quanta (PBSQ) Model

3.5.1 Analytical Model

We illustrate the use of our PBSQ model through the use of a two-priority (high and low) example as depicted in Fig. 3.1. Consider a centralized processor that is shared by four LQs in a prioritized PS manner. LQs 1 and 2 are exclusively loaded by low priority messages and LQs 3 and 4 by high priority messages. Each LQ is assumed to have an associated infinite buffer.

Message lengths for each priority are discrete i.i.d.. Each message is assumed to consist of an integral number of *segments*; for example, the message in LQ 1 has two segments. A segment here corresponds to a MAC layer packet representing an uninterrupted quantum of service time received by the message from the MAC PS server. A MAC layer normally operates using units called *slots* each of which consists of a fixed number of bytes. All segments of a given priority are assumed to be composed of the same integral number

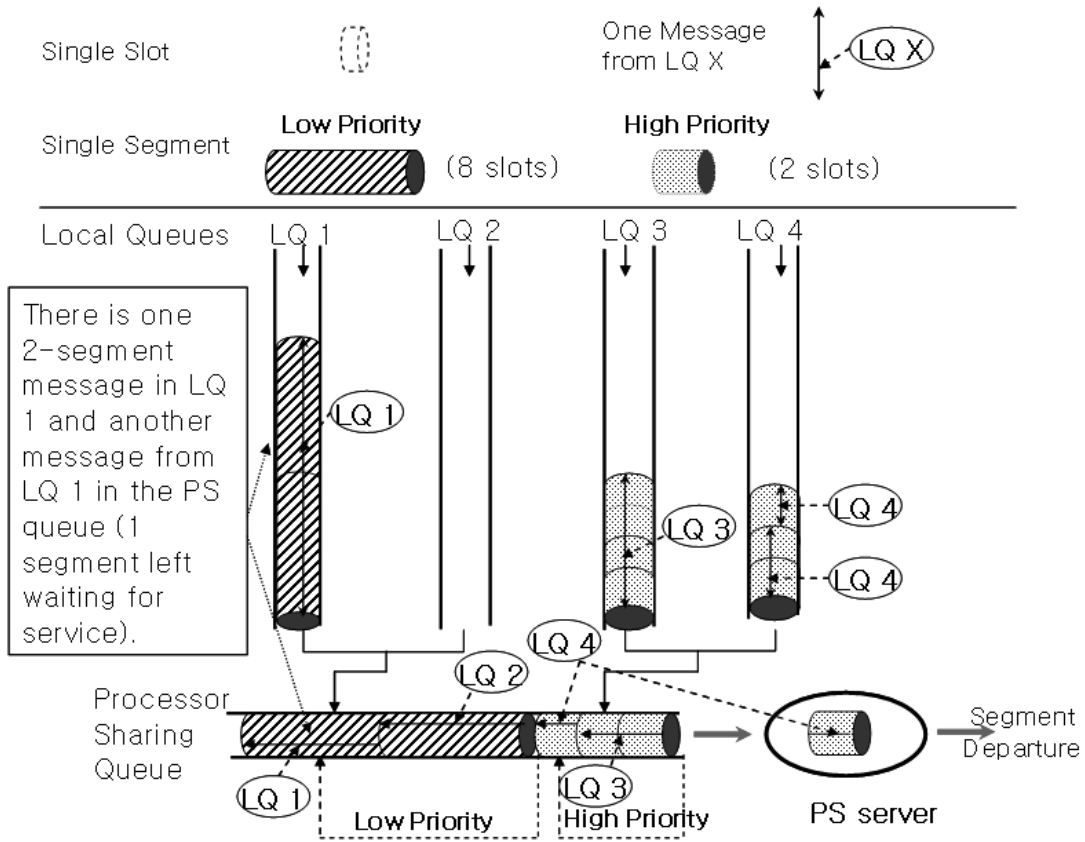


Figure 3.1: The PBSQ model for the case of two priorities.

of slots, but segments of different priorities may have different numbers of slots. For each priority p , the length of its segment is denoted by N_p [slots], $1 \leq N_p \leq N$, where N represents the maximum possible segment size. In Fig. 3.1, the high priority segment size is two slots and the low priority one is eight slots. At any time, at most one message from each LQ can be present in the PS queue, i.e. a message from an LQ cannot move into the PS queue until the previous message from the same LQ has been completely served.

As for the MPS model, time is also divided into *time-slots*. As a time-slot is related to the service time of a *slot*, the service time of a segment is known as the *segment-time*. Messages arriving at a priority p LQ within any time-slot

follow a Poisson process with a mean arrival rate of \bar{a}_p . Thus, $C_{a,p}^2$ is equal to $1/\bar{a}_p$. The mean total arrival rate of priority p messages is given by $\lambda_p = M_p \bar{a}_p$.

The PS server provides one segment of service to the message at the head of the PS queue and recycles the incomplete message to the tail of its own priority group, but ahead of all lower priority messages. In the snapshot presented in Fig. 3.1, a segment from LQ 4 is in service and the remaining segment of the same message is sent to the end of the high-priority group in the PS queue; and one segment of the incomplete LQ 3 message (2 segments) will be served next. We consider non-preemptive priority scheduling at the segment level, i.e., any new message arrival, even from a higher priority, cannot interrupt the current segment's service.

Using \bar{b}_p and $C_{b,p}^2$ defined in Table 3.1, the mean size in slots of priority p messages is given by $\bar{b}_p N_p$ and the traffic load of priority p messages is given by $\rho_p = \lambda_p \bar{b}_p N_p$, $0 < \rho_p < 1$.

3.5.2 Mean Message Delay

In [9], the delay of an arriving message is obtained by summing its waiting time in the LQ and its sojourn time in the PS queue. Notice that the method of [9] relies on the following two assumptions that we do not adopt here: (1) all messages arrive at a segment boundary; and (2) segments of different priorities are of the same size. Nevertheless, because each priority is considered separately, we can still develop an accurate approximation for the mean message delay by applying the results of [9] and then correcting the result using a compensation term. Accordingly, the mean delay of a priority p message consisting of n segments, denoted $D_p(n)$ (in units of priority p segment-time),

is given by

$$D_p(n) = L_p + S_p(n) + \Delta_p, \quad (3.14)$$

where L_p is the mean priority p message waiting time in the LQ, $S_p(n)$ is the mean time spent by priority p messages in the PS queue until the completion of at least n segments of service, and Δ_p is the compensation term defined as the difference between the delay for a priority p message and the delay predicted by the MPS model of [9].

Since, in the model of [9], each segment is equal to one slot, to use that model in our approximation it is convenient to consider the basic time unit to be the segment-time. In particular, when we evaluate the mean delay of priority p messages, we consider time to be measured in priority p segment-time units and the amount of traffic that arrives in units of priority p segments. Since other priority messages may have segments of size that are not equal to, or not an integer multiple of, the priority p segment size, and because they arrive during the priority p segment-time (at different time-slots) and not all at once (as in [9]), in our approximation we do not consider the actual arrival process, but we do fit the first two moments of our arrival process. Also, since the result of [9] is based on only two moments, our approximation can be very good. We re-describe the arrival process of priority i messages in Table 3.3.

Following the line of argument given in Section III and Appendix I of [9] and Section 3.4 of this chapter, we have the modified expression for $S_p(n)$ as

$$S_p(n) = \frac{n}{1 - \varepsilon_{p-1} - \frac{M_p-1}{M_p} \rho_p} - \delta_p, \quad (3.15)$$

where $\varepsilon_p = \sum_{i=1}^p \rho_i$.

In a similar fashion to the MPS model, we derive L_p using a discrete-

The mean length of messages in units of priority p segments	$E[b_i N_i / N_p] = b_i N_i / N_p$.
The mean message arrival rate during a priority p segment-time	$\lambda_i N_p$.
The variance of the number of messages arriving during a priority p segment-time	$N_p M_i \text{Var}(a_i) = N_p M_i a_i^2 C_{a,i}^2 = N_p M_i a_i$, as $C_{a,i}^2 = 1/a_i$ for Poisson arrivals.
The second moment of the number of priority i messages arriving during a priority p segment-time	$N_p M_i a_i^2 C_{a,i}^2 + N_p^2 \lambda^2 = N_p M_i a_i + N_p^2 \lambda^2$, as $C_{a,i}^2 = 1/a_i$ for Poisson arrivals.

Table 3.3: Priority i traffic load converted into units of priority p segments and segment-times

time non-preemptive priority queueing model [180] which is equivalent to the PBSQ model in terms of the average segment delay. Thus L_p can be obtained by equating the mean of the priority p segment delay in the PBSQ model with the mean of the discrete-time non-preemptive priority model.

According to Appendix II of [9], and in the same way as for the derivation in Section 3.4.2, we have

$$R_p(k) = \frac{1 - F_{b,p}(k-1)}{\bar{b}_p} \quad \text{for } k = 1, 2, 3, \dots, b_p,$$

$$E[D_{seg(p)}] = \sum_{k=1}^{b_p} [L_p + S_p(k)] R_p(k),$$

where $R_p(k)$ for $k = 1, 2, \dots, b_p$ is the probability that a randomly selected priority p segment is the k th in its own message. Then, the mean delay of priority p segments is obtained as

$$E[D_{seg(p)}] = L_p + \frac{[\bar{b}_p(1 + C_{b,p}^2) + 1]/2}{1 - \varepsilon_{p-1} - \frac{M_p-1}{M_p} \rho_p} - \delta_p. \quad (3.16)$$

Using the first two moments of the message arrival process in Table 3.3, we

apply (6.5) of [180] to obtain the first two moments of the arrival process for messages of priority p or higher on the priority p segment level. For priority i messages, the mean and variance of the priority p segment level is derived as

- mean = $\lambda_i N_p \bar{b}_i N_i / N_p = \lambda_i \bar{b}_i N_i = \rho_i$ [priority p segments/ priority p segment-time],
- variance = $\text{Var}(X) = E(X^2) - E(X)^2$. According to (6.5) of [180], using the results in Table 3.3, the second moment of the arrival process for messages of priority p or higher on the priority p segment level is $[(N_p M_i a_i + N_p^2 \lambda^2) - N_p \lambda] (\bar{b}_i \frac{N_i}{N_p})^2 + N_p \lambda E[(b_i \frac{N_i}{N_p})^2]$. Thus, we have

$$\begin{aligned}
& [(N_p M_i a_i + N_p^2 \lambda^2) - N_p \lambda] (\bar{b}_i \frac{N_i}{N_p})^2 + N_p \lambda E[(b_i \frac{N_i}{N_p})^2] - (N_p \lambda \bar{b}_i \frac{N_i}{N_p})^2 \\
&= [N_p M_i a_i + N_p^2 \lambda^2] (\bar{b}_i \frac{N_i}{N_p})^2 - N_p \lambda (\bar{b}_i \frac{N_i}{N_p})^2 + N_p \lambda E[(b_i \frac{N_i}{N_p})^2] - (N_p \lambda)^2 (\bar{b}_i \frac{N_i}{N_p})^2 \\
&= (\bar{b}_i \frac{N_i}{N_p})^2 [N_p M_i a_i + N_p^2 \lambda^2 - (N_p \lambda)^2] + N_p \lambda [E[(b_i \frac{N_i}{N_p})^2] - (\bar{b}_i \frac{N_i}{N_p})^2] \\
&= (\bar{b}_i \frac{N_i}{N_p})^2 N_p M_i a_i + N_p \lambda (\bar{b}_i \frac{N_i}{N_p})^2 C_{b,i}^2 \\
&= \frac{\bar{b}_i^2 N_i^2 \lambda}{N_p} + \frac{\bar{b}_i^2 N_i^2 C_{b,i}^2 \lambda}{N_p} \\
&= \frac{\rho_i \bar{b}_i N_i}{N_p} + \frac{\rho_i \bar{b}_i C_{b,i}^2 N_i}{N_p} \\
&= \rho_i \bar{b}_i \frac{N_i}{N_p} (1 + C_{b,i}^2).
\end{aligned}$$

Define $\rho_i \bar{b}_i N_i (1 + C_{b,i}^2) / N_p = \nu_i$, variance = ν_i .

Thus, the total segment arrival process for priority p messages or higher has the following first two moments:

- mean = $\sum_{i=1}^p \rho_i = \varepsilon_p$,
- variance = $\sum_{i=1}^p \nu_i$, where $\nu_i = \rho_i \bar{b}_i N_i (1 + C_{b,i}^2) / N_p$.

Parameters in (6.17) of [180]	Parameters in PBSQ model
ρ_{I_y} : the traffic load from those priorities higher than y .	ε_{p-1} : the traffic load from those priorities higher than p .
$E(\tilde{N}^{(F_y)})$ and $E([\tilde{N}^{(F_y)}]^2)$: the first two moments of the number of segments from priorities $i \leq y$.	ε_p and $\sum_{i=1}^p \nu_i + \varepsilon_p^2$: the first two moments of the number of segments from priorities p or higher arriving during a priority p segment-time, in the units of priority p segment size.
$\tilde{N}_i \bar{B}_i$: the mean arrival rate of segments from the priorities i .	ρ_i .
\bar{B}_y : the mean message/segment size depends on the equation used for calculating the mean message/segment waiting time.	1: as we calculate the mean <i>segment</i> waiting time.
\bar{N}_y and N_y^2 : the first two moments of the arriving process for priority y messages.	ρ_p and $\nu_p + \rho_p^2$: the first two moments of priority p segments arriving during a priority p segment-time.

Table 3.4: Corresponding parameter notations used in (6.17) of [180]

Then, applying the previously obtained mean and variance, we use (6.17) of [180] to calculate the mean waiting time of priority p segments. The corresponding parameters for our case using (6.17) are shown in Table 3.4. Replacement of our corresponding parameters into (6.17) of [180] and adding a one segment transmission time, we obtain the mean delay of p priority seg-

ments in units of p segment-time, giving

$$\begin{aligned}
& E[D_{seg(p)}] \\
&= \frac{1}{1 - \varepsilon_{p-1}} \left\{ \frac{\sum_{i=1}^p \nu_i + \varepsilon_p^2 - \varepsilon_p}{2(1 - \varepsilon_p)} + \sum_{i=1}^{p-1} \rho_i + \frac{\nu_p + \rho_p^2 - \rho_p}{2\rho_p} \right\} + 1 \\
&= \frac{1}{2(1 - \varepsilon_{p-1})} \left\{ \frac{\sum_{i=1}^p \nu_i}{1 - \varepsilon_p} - \frac{\varepsilon_p(1 - \varepsilon_p)}{1 - \varepsilon_p} + 2\varepsilon_{p-1} + \nu_p/\rho_p + \rho_p - 1 \right\} + 1 \\
&= \frac{1}{2(1 - \varepsilon_{p-1})} \left\{ \frac{\sum_{i=1}^p \nu_i}{1 - \varepsilon_p} - \varepsilon_p + 2\varepsilon_{p-1} + \nu_p/\rho_p + \rho_p - 1 \right\} + 1 \\
&= \frac{1}{2(1 - \varepsilon_{p-1})} \left\{ \frac{\sum_{i=1}^p \nu_i}{1 - \varepsilon_p} + \varepsilon_{p-1} - \rho_p + \nu_p/\rho_p + \rho_p - 1 \right\} + 1 \\
&= \frac{1}{2(1 - \varepsilon_{p-1})} \left\{ \frac{\sum_{i=1}^p \nu_i}{1 - \varepsilon_p} + \nu_p/\rho_p + \varepsilon_{p-1} - 1 \right\} + 1 \\
&= \frac{\sum_{i=1}^p \nu_i / (1 - \varepsilon_p) + \nu_p/\rho_p}{2(1 - \varepsilon_{p-1})} - \frac{1 - \varepsilon_{p-1}}{2(1 - \varepsilon_{p-1})} + 1 \\
&= \frac{\sum_{i=1}^p \nu_i / (1 - \varepsilon_p) + \nu_p/\rho_p}{2(1 - \varepsilon_{p-1})} + \frac{1}{2}. \tag{3.17}
\end{aligned}$$

We can obtain L_p using (3.16) and (3.17). Then, adding L_p to (3.15), we obtain

$$\begin{aligned}
L_p + S_p(n) &= \frac{\nu_p/\rho_p + \sum_{i=1}^p \nu_i / (1 - \varepsilon_p)}{2(1 - \varepsilon_{p-1})} + \\
&\quad \frac{n - [\bar{b}_p(1 + C_{b,p}^2) + 1]/2}{1 - \varepsilon_{p-1} - \frac{M_p - 1}{M_p} \rho_p} + \frac{1}{2}. \tag{3.18}
\end{aligned}$$

For the first part of the compensation, henceforth denoted $\Delta_p(1)$, when a priority p message arrives, it has to wait until the segment in progress completes its transmission – regardless of its priority – because we do not allow preemptions. At the same time, due to this delay, there are higher priority messages which arrive during the time, so the total work required by all these segments should be taken into account of $\Delta_p(1)$. Thus, for a priority p message, denoted the time waiting for the current served segment by W_p , we

have

$$\Delta_p(1) \rightarrow W_p + \Delta_p(1) \sum_{i=1}^{p-1} \rho_i.$$

We can evaluate W_p for a priority p message due to the transmission time of a priority i segment currently found in service. Since the message arrival process follows a Poisson process, this occurs with probability ρ_i (by the PASTA principle). The mean time that priority p messages have to wait until the priority i segment completes its service is estimated as half of the N_i time-slots.

We have

$$E[W_p] = \frac{\sum_{i=1}^P \rho_i N_i}{2N_p}.$$

Thus, we have

$$\Delta_p(1) = \frac{\sum_{i=1}^P \rho_i N_i}{2N_p(1 - \varepsilon_{p-1})}. \quad (3.19)$$

The part $\Delta_p(1)$ only considers compensation associated with the initial delay. It is important, however, to notice that, for the same reason that the high priority message experiences additional delay relative to the MPS model (because it has to wait until a low priority segment completes its transmission), the low priority message experiences a reduction of delay (relative to the MPS model) if a high priority message arrives during the transmission of the last segment of the low priority message. Notice that this reduction of delay, which we henceforth denote by $\Delta_p(2)$, is only relevant to the last segment. If the high priority message arrives during transmission of an earlier segment, the reduction of delay gained by the low priority message will be offset by a later delay of transmission of subsequent segments as they can only be transmitted after the transmission of the high priority message is completed and since the latter incurred initial delay, it will also be delayed in completing its transmission.

Let $J_p(x)$ for any $i < p$ be a function defined by

$$J_p(x) = \begin{cases} 1 & \text{if } x < cN_p, \\ 0 & \text{otherwise,} \end{cases}$$

where the factor c , $0 < c \leq 1$, will allow us to disregard messages of priorities higher than p when their segment size is close to that of p . Let $\lambda(p) = \sum_{i=1}^{p-1} \lambda_i J_p(N_i)$. The rate $\lambda(p)$ is the arrival rate of messages that we would like to consider in evaluating $\Delta_p(2)$. In particular, $\Delta_p(2)$ is evaluated using the time elapsed from the instant corresponding to the occurrence of an arrival (drawn from a Poisson distribution process with mean rate $\lambda(p)$) within the last segment-time of our priority p message) until the end of that segment-time.

Consider k random variables X_1, X_2, \dots, X_k that have a uniform distribution within $(0, N_p)$. The random variables represent the arrival times of new messages that come during a priority p segment-time. We have

$$P[\text{Min}(X) > x] = P(X_1 > x)P(X_2 > x)\dots P(X_k > x) = \left(1 - \frac{x}{N_p}\right)^k, \text{ for } 0 < x < N_p.$$

Hence, the density of $\text{Min}(X)$ is $k(1 - x/N_p)^{k-1}/N_p$, for $0 < x < N_p$. So we

calculate $E[\text{Min}(X_1, X_2, \dots, X_k)]$ by

$$\begin{aligned}
E[\text{Min}(X)] &= \int P(X > t) dt \\
&= \int_0^{N_p} \frac{k}{N_p} \left(1 - \frac{x}{N_p}\right)^{k-1} x dx \\
&= \int_0^{N_p} \left[1 - \left(1 - \frac{x}{N_p}\right)\right] k \left(1 - \frac{x}{N_p}\right)^{k-1} dx \\
&= k \left[\int_0^{N_p} \left(1 - \frac{x}{N_p}\right)^{k-1} dx - \int_0^{N_p} \left(1 - \frac{x}{N_p}\right)^k dx \right] \\
&= k N_p \left(\frac{1}{k} - \frac{1}{k+1} \right) \\
&= \frac{N_p}{k+1}. \tag{3.20}
\end{aligned}$$

It is known that conditioning on the number of Poisson arrivals within a segment-time, the arrival times have the same distribution as the order statistics of the same number of uniformly distributed random variables within that segment-time. Using (3.20), the mean time from the moment that the first higher priority message arrives during the last segment-time of a priority p message until the end of that segment-time, with conditioning and unconditioning of the number of high priority messages that arrive during that last segment-time, is obtained by

$$\begin{aligned}
&\sum_{k=0}^{\infty} \frac{N_p}{k+1} P(Y = k) \\
&= \frac{N_p}{N_p \lambda(p)} \sum_{k=0}^{\infty} e^{-\lambda(p)N_p} \frac{(\lambda(p)N_p)^{k+1}}{(k+1)!} \\
&= \frac{1}{\lambda(p)} (1 - e^{-\lambda(p)N_p}) \tag{3.21}
\end{aligned}$$

where Y is a Poisson random variable with parameter $\lambda(p)N_p$. Then, we obtain

the correction of the delay relative to MPS in priority p segments by

$$\Delta_p(2) \cong 1 - \frac{1}{N_p \lambda(p)} (1 - e^{-\lambda(p)N_p}), \quad 1 < p \leq P \text{ and } \lambda(p) \neq 0. \quad (3.22)$$

Since $\Delta_p(2)$ is relevant only for $i < p$, we set $\Delta_1(2) = 0$. Overall, Δ_p is estimated by

$$\Delta_p \cong \Delta_p(1) - \Delta_p(2). \quad (3.23)$$

Due to possible large variations in segment sizes for the different priorities, the compensation term Δ_p is the key to an accurate evaluation of the overall mean message delay. We can multiply the result of (3.14) by N_p to convert the delay result to the common “currency” of time-slots.

Compare the mean delay which PBSQ performs with which MPS performs, Δ_p is the only difference as PBSQ without the limitation of assumption that all messages arrive at a segment boundary. As $\Delta_1(2) \leq 1$, $\Delta_p(1)$ can have a major effect on this difference. According to (3.19), in the PBSQ model, the mean delay for a priority traffic is decreased when a bigger segment size is chosen for this priority than for other priorities. Also, under a same total loading ε_P and a same set of segment size $N_p, p = 1, 2, \dots, P$, the mean delay of priority p messages with a bigger ε_{p-1} is longer than one with a smaller ε_{p-1} .

3.5.3 Relevant Implementation Issues

The MAC layer enables the optional function of fragmentation to avoid retransmission of large frames in the presence of Radio Frequency (RF) interference. If the bit errors resulting from RF interference affect a single frame, which is quite typical of a wireless environment, it is obviously better to retransmit a smaller frame rather than a larger one. With fragmentation, a

network node can divide data messages into smaller frames according to a maximum frame length threshold set by network operators. If a message is referred to as a network layer service data unit (SDU), it needs to be turned into one or more frames known as MAC protocol data units (PDUs) for transmission over the network with added headers. MAC header formats are normally defined in the standards. For example, a typical one is the Ethernet header used predominantly for network access.

We consider the effect of a MAC header with a size h_p [segments], $h_p < 1$. When a message length b_p is divided into n segments, the mean size of MAC messages \bar{B}_p with PDU headers is given by:

$$\bar{B}_p(n) = n(1 + h_p). \quad (3.24)$$

According to (3.14) and (3.24), we can obtain an approximation for the mean message delay at the MAC layer as $D_p[\bar{B}_p(n)]$.

If a bigger segment size is chosen in the PBSQ model, the mean delay will be decreased due to the fact that adding MAC headers imposes less overhead than in the MPS model. This will be demonstrated by numerical results in the following section.

3.5.4 Model Evaluation and Numerical Results

Model Evaluation

In this section, we validate our approximation for the mean message delay using a C++ simulation program. In this simulation, there are six active stations, designated as Station 1 to Station 6 respectively, and sharing a central processor which operates as described in the PBSQ model of Section 3.5.1.

Stations 1 and 2 transmit messages at priority 3 (lowest), Stations 3 and 4 at priority 2, and Stations 5 and 6 at priority 1 (highest). The service capacity of the processor is 0.1 Gb/s. Messages of all priorities arrive in accordance with a Poisson process, with mean rates equal to 100 and 62.5 messages/s for priorities 2 and 3, respectively, and the mean rate for priority 1 is varied from 100 to 550 messages/s. Message sizes of priorities 1 and 2 are negative exponentially distributed with mean equal to 4.5 and 15 kbytes, respectively; priority 3 message lengths follow a cut-off Pareto distributed model with shape parameter = 1.1, scale parameter = 4.5 kbytes and cut-off threshold = 2 Mbytes, as in [181]. From measurements of the generated message-length deviates we obtain $C_{b,3}^2 \cong 16$. The generated message size is rounded to the nearest integral number of segments. The service quanta are 500, 5000 and 3500 bytes for priorities 1, 2 and 3, respectively. For a slot size of 500 bytes, we calculate the mean message delay for priority p traffic, $D_p[\bar{b}_p]$, using (3.18) and (3.23) with $c = 0.8$ (empirically set). Analytical and simulation results with 95% confidence intervals (which are too small to be noticed in many cases) based on a Student's t-test are presented in Fig. 3.2 and 3.3. The analytical results are in good agreement with the simulation results for all three priorities.

Comparison of the PBSQ and MPS Models

A two-priority example is considered for comparison between the mean message delays in the PBSQ and MPS models, where $M_p = 2$ for each priority. We have Poisson arrivals with a mean rate of 0.0006 per slot for low priority and a variable rate from 0.0002 to 0.001 per slot for the high priority. Message size is negative exponentially distributed with mean 142.7 slots for the high priority and cut-off Pareto distributed with a mean of 250.04 slots and $C_{b,p}^2 = 16.3$ for the low priority. We consider the message size also rounded to the nearest

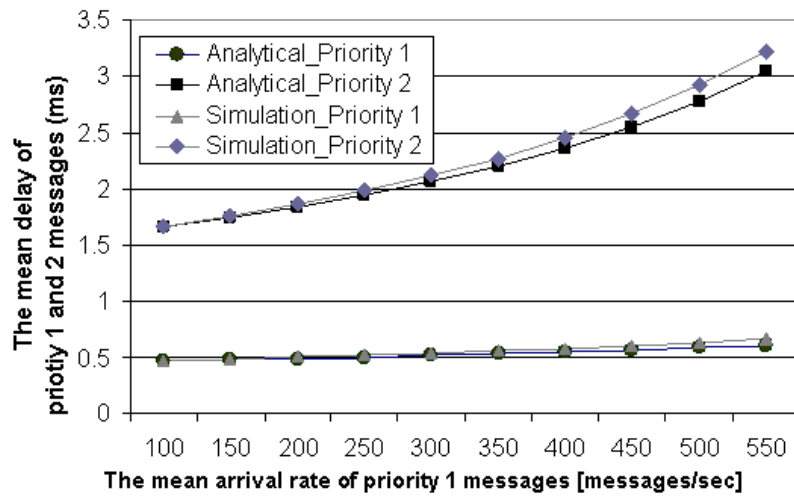


Figure 3.2: Analytical vs simulation results for the mean delay of priority 1 and 2 messages.

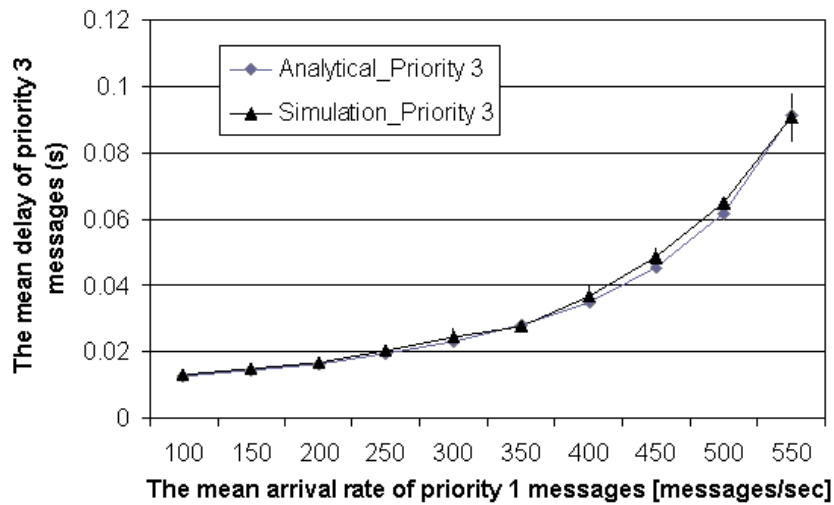


Figure 3.3: Analytical vs simulation results for the mean delay of priority 3 messages.

integral number of segments. We assume a 50-slot service quantum for the low priority and a 90-slot service quantum for the high priority traffic in the PBSQ model. The fixed service quantum in the MPS model is one slot. We consider a 0.1 slot header for each service quantum in both models. The nu-

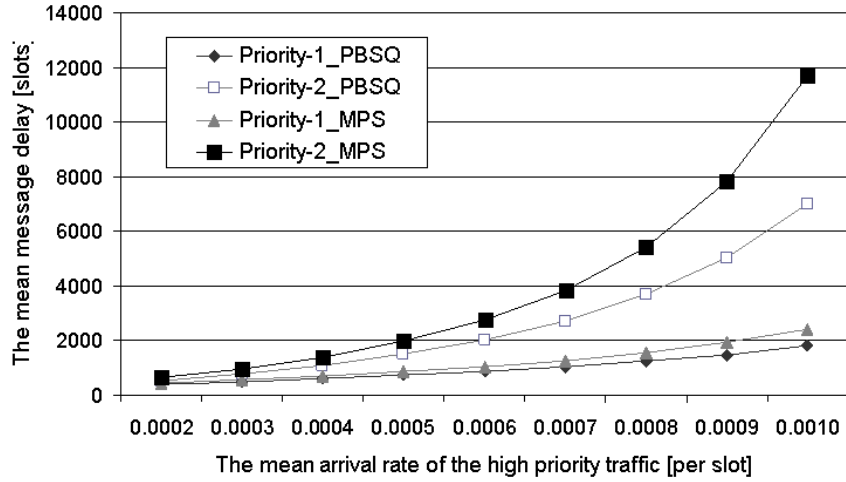


Figure 3.4: Comparison of mean message delays in the PBSQ model and in the MPS model.

merical results for the two models are plotted in Fig. 3.4. It shows the delay for the PBSQ model is significantly improved at the higher arrival rate after adding headers, although the header size is small. When the arrival rate of the high priority traffic reaches 0.001, the reduction of the mean message delay for the high priority for the PBSQ model is about 25% less than for the MPS model, and is more significant at about 40% less than the MPS model for the low priority.

3.6 Conclusions

In this chapter, we introduced a multiservice PS scheduling policy and the relevant multiservice multiqueue PS models: MPS and PBSQ. We corrected a subtle incongruity in the mean LQ waiting time given by equation (13) of [9] for the MPS model. Moreover, extending this MPS model, we described a general multiservice multiqueue PS model, that we have called the PBSQ model, in which the service quanta may be different for different services.

An accurate approximation for the mean message delay is derived and validated by a simulation study. The unique advantage of our model is that it allows variable-sized service quanta rather than the fixed service quantum scheduling mechanism assumed in other models. We also considered the implementation issue for MAC headers, we demonstrated the performance effect of the choice of service quantum and the benefits that can be achieved over a fixed service quantum scheduling mechanism through a comparison with our numerical results.

Multimedia traffic exhibits different traffic characteristics and QoS requirements. A bandwidth allocation algorithm is needed in the MAC protocol to satisfy diverse QoS requirements and to utilize bandwidth efficiently. The priority-based algorithm is a basic scheme that can be chosen as one QoS solution. Given that our closed-form solution obtained from the analytical model of such an algorithm is easily computable and captures traffic heterogeneity, so it can be incorporated as part of a practical connection admission function. The proposed model provides a simple way to evaluate performance of multimedia applications which can be used by operators in network dimensioning and management of heterogeneous traffic.

Chapter 4

Modelling and Analysis of PTT Uplink Delay in GPRS/GSM Networks

4.1 Introduction

PTT is a packetised voice service provided by mobile network operators as a value-added service. Since IP is used as a bearer, PTT service performance depends significantly on the infrastructure and technology used by carriers. We consider the case of PTT over GPRS/GSM networks, where PTT delay performance is affected by the PTT/GPRS/GSM channel sharing scheme and associated delays due to retransmissions. Moreover, the packet delay is largely due to delays experienced between the mobile station (MS) and the Base Station (BS). Therefore, an appropriate mathematical model covering the above issues is required to evaluate the performance of PTT/GPRS. There have been many publications that discuss the general performance evaluation of GPRS/GSM systems [182, 183, 184, 185, 186, 187]. However, these studies are based on

the use of GPRS for normal data services only and not for packetised speech.

By using the MPS queueing model idea, based on [170, 9, 5], we analyze the mean PTT packet delay using the partial sharing channel allocation scheme. We take into account the effect of GSM voice traffic under a quasi-stationary assumption [41]. Cases with and without provision of strict priority for PTT traffic over GPRS data are discussed. We describe the analytical model including the underlying assumptions for the MAC/RLC layer between the MS and the BS, together with numerical solutions and validate this model by simulations. Based on our analytical solution, the effect of retransmissions is evaluated. Numerical results are presented to illustrate the interactive impact of traffic loading and various design parameters on PTT packet delay. The effects of retransmission, GSM voice loading and priority assignment on mean PTT delay are also discussed.

4.2 GPRS LLC/RLC/MAC Layers

GSM uses a TDMA structure with eight time-slots (channels) per frame to support voice and data transmission [169]. Two additional support nodes: the serving GPRS support node (SGSN) and the gateway GPRS support node (GGSN), have been added into the original GSM infrastructure for GPRS packet data routing [170, 188]. To support communication between the MS and the GPRS network, the physical layer is designed to include functions for modulation/demodulation, channel coding/decoding, etc., and the data link layer is designed for establishing logical links and all of the detailed implementation of the GPRS protocols. The data link layer is split into two sub-layers, namely, logical link control (LLC) and the radio link control/medium-access control (RLC/MAC). In this section, we describe LLC/RLC/MAC func-

tions that are relevant to PTT/GPRS performance modelling.

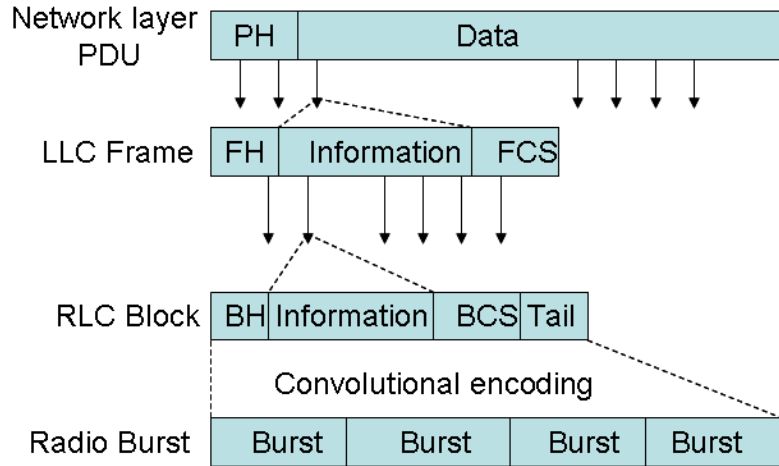


Figure 4.1: PDU segmentation into LLC frames, RLC blocks and MAC bursts in GPRS networks. (PH: PDU header; FH: frame header; BH: block header.)

- *LLC Layer* provides a logical link between the MS and the SGSN. PDUs from higher layers are segmented into variable length LLC frames (see Fig. 4.1). The LLC layer operates the control and ciphering of link-level flows either in an acknowledged or in an unacknowledged mode. In the unacknowledged mode, the LLC layer operation does not attempt to recover errors in LLC frames. In the acknowledged mode, the LLC layer enables retransmissions of erroneous LLC frames through an automatic repeat request (ARQ) mechanism using the attached frame-check sequence (FCS) within each LLC frame.
- *RLC/MAC Layer* can be considered as the RLC and MAC separately. The MAC layer handles simultaneous multiple accesses from MSs using collision detection. It employs a slotted-ALOHA-based reservation protocol [188]. The RLC lies between the LLC layer and the MAC layer. As shown in Fig. 4.1, each LLC frame is divided into several RLC data blocks and added a block-check sequence (BCS) and tail bits at

every block. Irrespective of using any of four possible channel coding schemes that are defined for GPRS, an RLC block consists of four time-slots (bursts). The RLC can also operate in either an acknowledged or unacknowledged mode. In the RLC unacknowledged mode, the RLC only reassembles RLC blocks into LLC frames without any retransmission of erroneous RLC blocks. In the acknowledged mode, the RLC uses the BCS in an RLC block to detect errors and provides an ARQ mechanism to recover them. Note that errors can be recovered by retransmission attempts both at the RLC block level and at the LLC frame level.

4.3 Analytical Model for PTT Uplink Delay

4.3.1 Partial Sharing Channel Allocation Scheme

We consider a typical GPRS/GSM network using a partial sharing scheme for the uplink. Let C be the total number of available channels for GSM and GPRS in a single cell and let g be the number of channels exclusively reserved for GPRS. The remaining $k = C - g$ channels are shared by GSM voice and GPRS packets and, in these k channels, GSM voice is assumed to have strict priority over the GPRS data traffic. Assume that voice call arrivals follow a Poisson process and their holding times are exponentially distributed. Invoking the classical Erlang M/M/k/k model, we can obtain the probability P_i of i channels ($i = 0, 1, 2, \dots, k$) being used by GSM voice calls in the steady state.

$$P_i = \frac{\frac{A^i}{i!}}{\sum_{q=0}^i \frac{A^q}{q!}}, \quad \text{for } i = 0, 1, 2, \dots, k, \quad (4.1)$$

where $A = \frac{\lambda}{\mu}$ for Poisson arrivals with the mean arrival rate λ and exponential service rate μ . Therefore, the probability Π_j of having j channels (time-slots)

available for GPRS data is given by

$$\Pi_j = P_{C-j}, \quad \text{for } j = g, g + 1, g + 2, \dots, C. \quad (4.2)$$

Since packet transmission times are much shorter than GSM voice holding times, at any state of the Markov process, we can consider that a PTT packet is serviced by a quasi-stationary [41] network.

4.3.2 Quasi-stationary Assumptions

When a network implements an architecture for differentiated services, an admission control scheme is engaged to limit the overall traffic for the stability of the network. The stability is also an important issue for implementation of the partial sharing channel allocation scheme in GPRS/GSM networks. In our case, we assume that overall service demand, including GSM voice, PTT voice and GPRS data, is less than the given capacity so that the system is ergodic. However, it can occur in some states that the number of GSM voice calls are holding more channels so that the remaining service rate can be insufficient to cope with the GPRS arrival demands. Such a situation can be referred to as *local instability* as suggested in [41]. In such a case, the GPRS traffic may suffer a long delay; consequently, authors of [41] considered the elastic traffic response time both under stable state and unstable state conditions. It has become necessary in our pseudo-stationary phase approach to include the unstable states for GPRS traffic in the model. Since the holding time of a GSM voice call is typically much greater than that of a GPRS packet, states involving local instability should not last for a significantly long period of time. However, we are able to avoid local instability by executing admission control on GSM voice calls to ensure an adequate service rate for GPRS traffic.

4.3.3 Modelling GPRS Delay

In a standard GPRS implementation, the MAC protocol enables several GPRS users to share a common transmission medium [170] and this suggests that it can be viewed as a processor sharing server. This structure was modelled in [170] by an M/G/1 PS model to estimate the PTT mean packet delay with packet service time x and link utilization ρ for the uplink of the PTT/GPRS MAC layer as follows:

$$T_x = \frac{x}{1 - \rho}. \quad (4.3)$$

Equation (4.3) assumes Poisson arrivals and considers only the service time of a test packet x and the overall utilization ρ . This can be generalized to the case of an i.i.d. arrival process and packet size distribution using the single priority version of the M/G/1 PS model of [9]. We have introduced the MPS queueing system model in Chapter 3, which consists of several distributed LQs served by a shared processor.

A BS that serves a multiplicity of users is modelled as a PS system and the GPRS packet sojourn time in the system is modelled by the time in an LQ and the time in the PS until its transmission is completed. The number of GPRS MSs sharing the channel is denoted by M . Each MS has an LQ with an infinite buffer for PTT voice or GPRS data packets with the same priority. The model is illustrated in Fig. 4.2.

As Fig. 4.1 shows, in a physical implementation of GPRS, the payload IP packets are segmented into RLC blocks. Each RLC block contains 456 bits, which is sent in a specific time-slot across four consecutive TDMA frames with a 20 ms associated transmission time [170]. As the RLC block is the smallest segment of information over the radio interface, we can measure the packet size in units corresponding to RLC blocks, as shown in Fig.4.2. Since

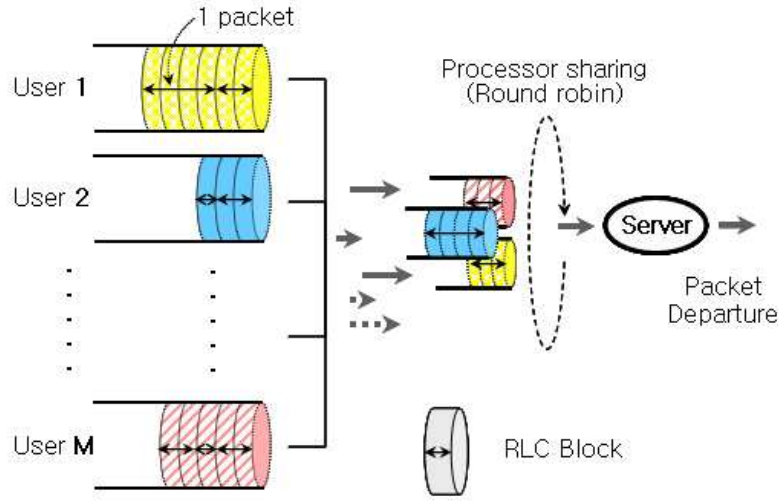


Figure 4.2: Analytical model used to analysis PTT uplink delay.

an RLC block transmission requires 20 ms, the packet size, as measured by the number of RLC blocks, also represents the packet service time.

In this single priority case, let the random variable a represent the number of packet arrivals within a RLC block associated transmission time at any MS (20 ms). We denote the mean of a as \bar{a} . Let the random variable b be the packet length measured in multiples of RLC blocks with the mean \bar{b} . Since an RLC block transmission requires 20 ms, b also represents the packet service time. At the BS, the PS server permits only one packet from each MS to join the PS queue and wait for service. If and only if there is a packet from a specific MS departing from the PS server in the last RLC block, another waiting packet from the same MS can then join the PS queue in this RLC block. For simplicity, in our model, we assume that all packets arriving within an RLC duration joining the discrete queue commence at the starting point of the next RLC block duration. Let C_a^2 and C_b^2 represent the squared coefficient of variation of a and b respectively and let $D(n)$ be the mean delay of the packet with n RLC blocks length, according to [9], $D(n)$ in the units representing the

number of RLC blocks is given by:

$$D(n) = \frac{\nu/\rho + \nu/(1-\rho)}{2} + \frac{n - [\bar{b}(1 + C_b^2) + 1]/2}{1 - \frac{M-1}{M}\rho} + \frac{1}{2}, \quad (4.4)$$

Where, $\nu = \rho\bar{b}(C_b^2 + \lambda C_a^2/M)$, $\lambda = M\bar{a}$, $\rho = \lambda\bar{b}$.

We examine the effect of GSM loading on PTT performance. As the GSM load is reduced, the PTT packet service rate increases as more time-slots are available for GPRS. Let $D(n)_j$ be the mean delay of a PTT packet with n RLC blocks given that j time-slots are available for GPRS. Equation (4.4) can be used to compute $D(n)_j$ by replacing ρ with ρ/j . Notice that the coefficient of variation of the service time is not affected by increasing the service rate. Therefore, by (4.2), with specified loads for PTT, GPRS data and GSM voice traffic, the average PTT packet delay $E[D(n)]$ is estimated by

$$E[D(n)] = \sum_{j=g}^C D(n)_j \Pi_j. \quad (4.5)$$

For the above MPS queueing model with time-varying capacity, under the quasi-stationary assumption, we consider the stability of the Markovian process that is required. Although, overall, the GPRS load is less than the given capacity, in some states of the process, the GPRS traffic experiences local instability, where $\rho/j > 1$. For such unstable cases, (4.4) cannot be used directly. With a practical approach, we replace ρ in (4.4) by $\sum_{j=g}^C (\rho/j) \Pi_j$ to obtain $E[D(n)]$. As ρ is a mean value in (4.4), this approach is justified. We can also calculate the probability of the unstable states by $\sum_{j=g}^C \Pi_j$ for any j having $\rho/j > 1$. In our approximations, we ignore the effects of these unstable states. This is consistent with a design that aims to limit the probability of being in an unstable state to a negligible value. Therefore, our approximations

are only accurate if the probability of the unstable states remains small.

4.4 Simulation Study

4.4.1 Scenario 1: Fixed Channel for GPRS Traffic

In this scenario, our analytical model is validated using the ns2 [189] simulation tool with a fixed channel for the GPRS traffic. The simulation uses CBQ objects that implement a packet-by-packet round-robin processor shared by the same priority classes. Based on the equivalent assumption for our analytical model, the CBQ buffer size in the simulation is infinite and the packet transmission time from the source node to the class-based queue is zero. We consider the following:

1. One fixed channel (equivalent to a service rate of 22.8 kbits/s) is allocated for GPRS and this is the service rate of the CBQ.
2. The same priority is given to each source node.
3. A Poisson arrival process with the same arrival rate at each source node, under the given packet arrival rate, utilisation is varied between 0.1 and 0.9. Nodes 1 and 2 generate PTT packets and the other nodes generate normal GPRS data packets.
4. Source nodes generate packets according to different packet size distributions, as described in Table 4.1. The maximum packet size is set to be 1500 bytes. As the continuous random variables (exponential and Pareto) have been assumed as models for packet sizes, the generated packet size is rounded to the nearest integral number of bytes.

5. The time when a packet is generated and when the whole packet completely departures will be recorded to calculate the packet delay. Having all packet records, the mean delay can be worked out by a statistical function of the simulation. Each set of input parameters will be run six times independently for confidence intervals.

Nodes (MS)	Packet size distribution	Mean size	Parameter
Node1	deterministic	228 bytes	-
Node2	exponential	228 bytes	-
Node3-6	Pareto	570 bytes	$\Gamma = 1.9$

Table 4.1: Packet size distributions of the source nodes

According to the same scenario, parameter settings given in Table 4.2 are used for the analytical model evaluation.

Parameters	Value
Average packet size, \bar{b} , measured by the number of RLC blocks	7.2
PTT packet size, n , measured by the number of RLC blocks	4
Squared coefficient of variation of b , C_b^2	0.5
The number of active GPRS MSs, M	6

Table 4.2: Parameters for evaluation of the model

Fig. 4.3 compares the analytical and simulation delay results. Simulation results are plotted with their 95% confidence intervals (which are too small to be noticed in many cases) based on a Student's t-test. The approximate PTT delays are in good agreement with the simulation results – although there is some divergence at high utilisation cases (e.g. 0.9).

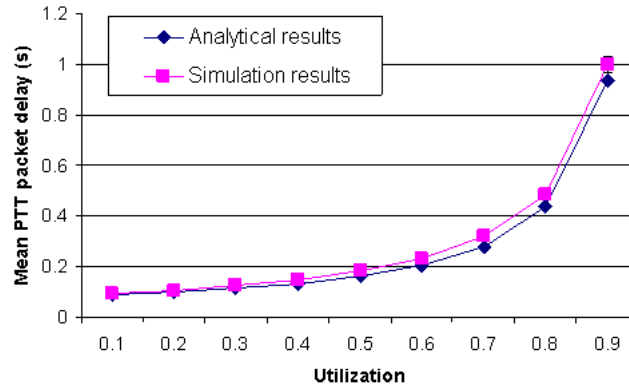


Figure 4.3: Comparison between analytical results and simulation results for PTT packet uplink delay.

4.4.2 Scenario 2: Partial Sharing Channel Scheme for GPRS and GSM Traffic

Including GSM voice traffic which shares channel resource with GPRS under using a partial sharing scheme, we present numerical results for the PTT delay based on the model described above in the case of a single cell model. We consider a cell with 24 available channels and, assuming that one of them is used for signalling and broadcasting, with an additional two being exclusively reserved for GPRS traffic, we find that this leaves 21 channels to be shared by GSM voice and GPRS traffic. Assuming a design based on a GSM voice call blocking probability of 2% and using the standard Erlang Loss formula, the traffic capacity is computed to be approximately $A = 14$ Erlangs. Using (4.2) for the standard analysis of an M/M/k/k queueing system, we calculate a set of probabilities corresponding to j available channels for GPRS, $j = 2, 3, \dots, 23$. Other parameters used in our analytical approach are listed in Table 4.3. Numerical results are plotted in Fig. 4.4.

We have developed a customised simulation tool using the C++ language

Parameters	Value
Average packet size, b , measured in units of RLC blocks	7
PTT packet size, n , in units of RLC blocks	4
Squared coefficient of variation of b , C_b^2	0.9
The number of active GPRS MSs, M	6

Table 4.3: Parameters for GPRS traffic

and employing a GUI interface. The simulation model incorporates both GSM voice traffic and GPRS traffic (which includes PTT voice and GPRS data traffic). The simulator runs as a discrete time model, the time units are optionally based around a standard GSM frame or an RLC block (4 frames). Using the same basic assumptions as for the GPRS/GSM cell described in the numerical example above, we assume two reserved channels for GPRS traffic leaving only 21 channels shared by both GSM voice and GPRS. We consider the following scenario:

1. GSM call arrivals are assumed to follow a Poisson process with a mean arrival rate 0.1167 calls per second. The channel holding time of a voice call is assumed to be exponentially distributed with a mean of 120s. When a GSM voice call arrives, it will be allocated to the next slot or RLC block (according to the option chosen) if one of 21 channels is available. Otherwise, it is blocked.
2. The same priority is given to each GPRS MS. Each GPRS MS is assumed to have an infinite buffer by setting a enough big queue size.
3. A Poisson arrival process with the same arrival rate at each GPRS MS. GPRS users generate packets according to different packet size distributions, as described in Table 4.4, then, they are placed in their individual buffers as GPRS requests. We assume one RLC block as one GPRS request. The packet transmission from a specific MS uses single-slot

FCFS, which means the packet from one user can use only one time-slot to be transmitted in the same RLC block duration.

4. Under various GPRS arrival rates, the overall system utilization of simulation process is always limited to be less than unity. We are able to do an analytical estimation of the average number of available channels for GPRS by $\sum_{j=g}^C j\Pi_j$. According to the discussion in Section 4.3.3, in the simulation, we ensure average GPRS traffic loading, i.e. $\sum_{j=g}^C \rho\Pi_j$, is always less than the average number of available channels for GPRS so that a given simulation reaches steady-state. For more details of steady-state in simulations, refer to [171].
5. For every RLC block duration, only one request from each MS is read in round robin order from users and is allocated to free channels which are not occupied by other GPRS traffic or any GSM voice traffic. If this request does not find free channels allocated, the GPRS request pointer will be recorded and used as the start point for the next round robin run in the next RLC block duration.
6. A time-slot captured by a request in a RLC block will be released in the next RLC block or TDMA frame.
7. Same as for Scenario 1, we collect delay results from simulations using a statistical function.

The other parameter settings for the simulation are the same as for the numerical example presented above, as listed in Table 4.3. Fig. 4.4 compares analytical and simulation delay results.

Simulation results are plotted with their 95% confidence intervals based on a Student's t-test. The approximate PTT delays are in good agreement

GPRS user	Packet size distribution	Mean Packet size
PTT user 1 and 2	deterministic	4 RLC blocks
Data user 3-6	exponential	8 RLC blocks

Table 4.4: Packet size distributions of GPRS users

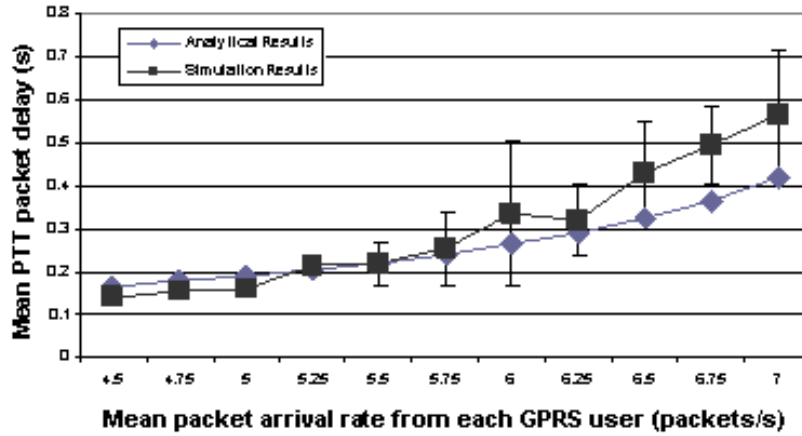


Figure 4.4: Comparison of analytical and simulation results of PTT delay.

with the simulation results although there is some divergence at high arrival rate cases (e.g. arrival rates of more than 6.5). These situations occur because of the growth of probability in the unstable states, which we have discussed before.

4.5 PTT/GPRS Retransmissions

We have introduced the retransmissions taking place at LLC level and RLC level. In this section, we quantify the effect of packet retransmissions on PTT packet delay, for when both the RLC and the LLC operate in unacknowledged mode and if channel errors occur between the source and the BS, whole GPRS data or PTT packets will be retransmitted. Thus, the average packet delay will increase. Assuming that retransmissions only occur because of errors,

and the Bit Error Rate (BER) is given, the packet error rate ϵ is estimated by

$$\epsilon = 1 - (1 - BER)^b, \quad (4.6)$$

where, b is the packet size in units of bits. The mean delay $E[D(n)]^*$ after retransmissions is based on the analytical results from the model presented above can be estimated by

$$E[D(n)]^* = \frac{E[D(n)]}{1 - \epsilon}. \quad (4.7)$$

4.6 Effects of Traffic Load for GSM and GPRS

Consider the case of two channels that are reserved for GPRS out of 23 channels where 21 channels are being shared by GSM voice and GPRS. Using the notation of (4.4) and (4.5), let the number of active GPRS users be $M = 8$, the mean GPRS packet size be $\bar{b} = 8$ RLC blocks, the PTT packet size be $n = 4$ RLC blocks, and $C_b^2 = 0.9$. Using our analytical model, we focus on the effects of GSM voice loading on PTT delay to produce the three curves given in Fig. 4.5. We consider the heaviest GPRS loading to be 2.25 Erlangs to keep the unstable state probabilities small as this is a realistic condition for practical dimensioning.

For our example, with 2.25 Erlangs of offered GPRS traffic and 18 Erlangs of GSM voice traffic, PTT packet delay is up to 400 ms, while it is less than 200 ms in low GSM voice loading cases (e.g. 14.4 or 12 Erlangs).

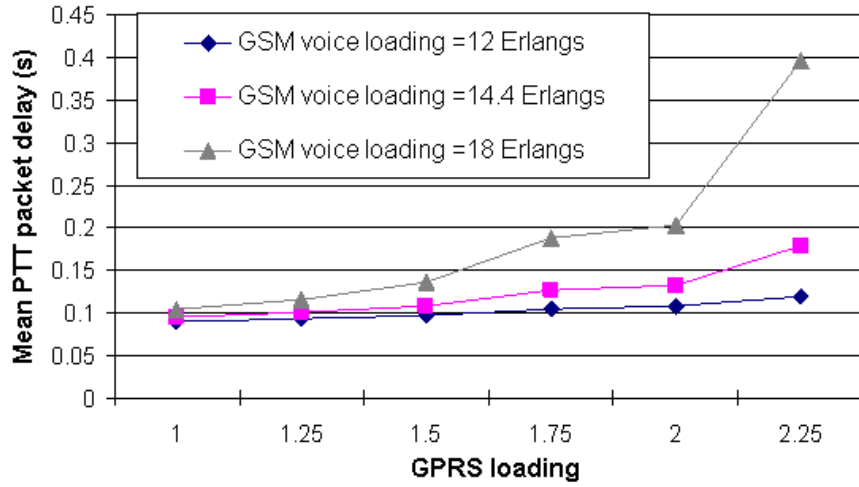


Figure 4.5: PTT packet delay under different GSM voice loads.

4.7 PTT/GPRS Priority Assignment

With strict delay requirements for PTT traffic, it may be necessary to introduce priorities. These can be implemented by combining two approaches: (1) a specified number of channels will be exclusively dedicated to GPRS, and (2) if a user requests a PTT service then they will obtain PTT-exclusive GPRS slots that will guarantee priority to PTT over other GPRS users. With these approaches, PTT packets can be protected from excessive loads from both GSM voice calls, and GPRS data packets.

Our analysis can apply to both ways of providing priority and protection for PTT traffic. In our example, two channels out of 23 are dedicated to GPRS. Since PTT packets have a higher transmission priority than other GPRS packets, where $\bar{b}_1 = \bar{b}$ is the PTT packet size, we consider $\bar{b} = n = 4$ RLC blocks, $C_b^2 = 0$ and $M = 8$. Fig. 4.6 shows the mean delay results obtained from the analytical solution for PTT with priority and without priority. Comparing with the case of PTT without priority and ($\bar{b} = 8$ RLC blocks, $C_b^2 = 0.9$), for the

mean packet size and the variation of packet size becomes smaller, PTT delay is reduced significantly. Especially, it is observed to occur under heavy traffic loading. Even under the heaviest PTT loading that we considered, PTT delay is still less than 230 ms.

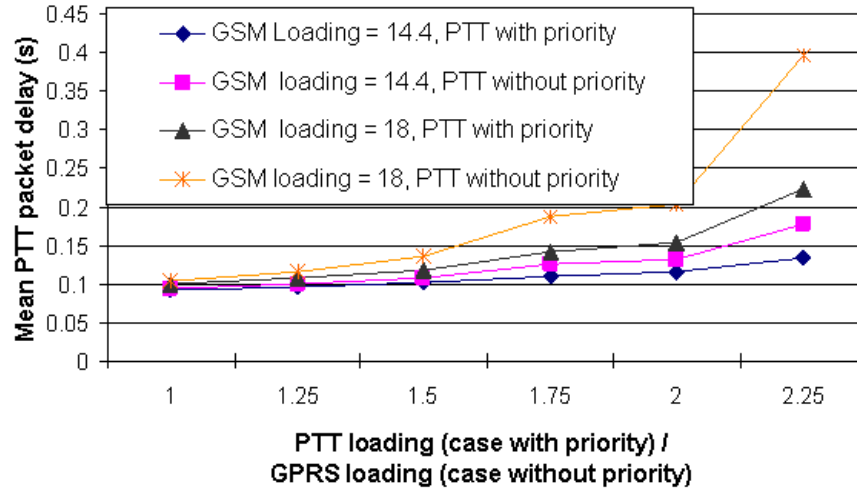


Figure 4.6: PTT packet delay with different priority assignments.

4.8 Conclusions

In this chapter, we have described an application of the multiservice multiqueue model to estimate the mean PTT packet delay in a PTT/GPRS/GSM system and validated this model by simulation. GSM voice, PTT data and normal GPRS data traffic has different traffic characteristics and QoS requirements. In the simulations, we considered a typical partial sharing channel allocation scheme implemented to distribute the bandwidth to three kinds of traffic. The approximate PTT delays obtained from our analytical model under a quasi-stationary assumption are in good agreement with the simulation results.

Although an admission control scheme is engaged to protect the stability of the network, it is still an important issue for implementation of the partial sharing channel allocation scheme in GPRS/GSM networks. In some states, the number of GSM voice calls are holding more channels so that the remaining service rate can be insufficient to cope with the GPRS arrival demands. As shown in Section 4.4.2, the simulation results under higher mean packet arrival rates have big statistic errors because of the increased the probability of instable states. For such cases, the mean delay requirements of different traffic may not always be met, as the real delay can be much bigger than the expected analytical results. Moreover, we have to consider the effect of burstiness of traffic on the probability of instable states. Even we are able to quantify the probability of instable states, some future research is required to quantify the effect of this probability on the mean delay.

Using our analytical solution, we have investigated GSM and GPRS loading, GPRS retransmissions on PTT packet delay and the improvement achieved by allowing PTT packets priority over GPRS data packets. Our results show that delay was seriously affected by high GPRS and GSM traffic loads. Moreover, the burstiness of GPRS traffic also affects PTT performance. As in packetised speech, suffering retransmissions, PTT delay also depends on channel quality. The proposed model provides a simple way for GSM operators to manage and dimension PTT traffic over GPRS/GSM networks.

Chapter 5

Modelling WiMAX Subscriber Station Uplink Delay

5.1 Introduction

The multiservice environment that WiMAX supports, possibly with multiple connections per service, is considered complex because of the various packet stream behaviours it needs to cope with. In such an environment, the packet schedulers operating at the MAC layer are very important for QoS delivery.

Four different services are specified in the IEEE 802.16 standard: unsolicited grant service (UGS), real-time polling service (rtPS), non-real-time polling service (nrtPS), and best effort (BE) [21].

- **UGS** is designed to support real-time applications with a fixed-sized packet, such as T1/E1 and VoIP without silence suppression. UGS connections are allocated fixed bandwidth at periodic intervals because of their strict delay requirements.
- **rtPS** is designed to support real-time applications with variable-sized

packets, such as with Moving Pictures Expert Group (MPEG) video and VoIP with silence suppression. Unlike UGS, it only needs a minimum reserved traffic rate and has less stringent delay requirements.

- **nrtPS** is designed for applications without any specific delay requirement but with a need for a minimum amount of bandwidth, such as File Transfer Protocol.
- **BE** is designed for applications that are delay-tolerant and does not require a minimum bandwidth.

Because of the distinct QoS characteristics of each service type, it may not be practical to use a single scheduling algorithm to handle all service types. Accordingly, a plausible solution is a two-level hierarchical approach [190]. In particular, the higher level serves these four services in the precedence order of strict priority. That is, connections of UGS are serviced first. Only after their QoS requirements have been satisfied, the connections of rtPS will be served and so on. Then, the lower level deals with scheduling of packets within each service type.

Recently, research on performance evaluation of WiMAX networks with different scheduling algorithms through the use of simulation has been conducted [190, 191, 181]. In [191, 181], deficit round robin (DRR) and WRR were chosen for the downlink and uplink schedulers, respectively. These algorithms are suitable for non-real-time data services because they focus on the throughput guarantee of data flows. Although it is possible for such rate-guaranteed algorithms to also provide a latency guarantee, these guarantees require an admission control policy based on worst-case behaviour and thus they lead to low network utilization. Therefore, they are not really suitable for real-time services which have stringent delay requirements.

We propose a priority-based fair scheduling algorithm to handle both real-time and non-real-time uplink WiMAX traffic in an SS. According to this algorithm, different service classes are assigned different priorities. Traffic is served strictly according to its priority. As traffic of non-real-time services is delay tolerant, traffic of real-time services can be protected by being assigned higher priorities. This algorithm is simple to implement. Moreover, we use a multiservice multiqueue processor sharing model to analyse the performance of our proposed algorithm. In particular, closed-form expressions for the mean *message* delay can be obtained. Using this model, the admission control module in the BS can evaluate the delay performance of various connections as part of the process of deciding whether a new connection can be admitted without sacrificing the QoS of existing connections. Note that the term “message” refers to an application layer data-unit which is broken down into a number of packets for transmission. It can represent different types of WiMAX traffic, such as a frame of streaming video, an object of a web page or a packet from a talk spurt. In other words, the delay guarantees offered to admitted connections is user-perceived delay rather than packet-level delay.

In this chapter, QoS architecture in WiMAX is reviewed first. We then focus on the scheduling component in WiMAX networks. After highlighting certain design issues for scheduling algorithms, we describe our proposed scheduling algorithm and present an analytical model for performance evaluation of our proposed scheduling algorithm for SS uplink traffic. The analytical model is verified by a simulation study. Also, the various traffic models related to the different service classes are discussed and a comprehensive set of numerical results is presented in order to illustrate the impact of traffic load and various design parameters on WiMAX message delay.

5.2 MAC Protocol of IEEE 802.16

The MAC layer of IEEE 802.16 is composed of the following three sublayers [21]:

1. *Convergence sublayer* (CS) which maps higher-layer SDUs into MAC PDUs received by the MAC common part sublayer (CPS). Currently, asynchronous transfer mode (ATM) CS and packet CS are specified in the standard. Packet CS supports all packet-based protocols such as the IP, Point-to-Point protocol (PPP) and IEEE 802.3 (Ethernet).
2. *Common part sublayer* that provides the core MAC function of medium access control, as well as the functions guaranteeing QoS.
3. *Security sublayer* which deals with security issues, which are particularly important for wireless communications. It contains two component protocols: encapsulation protocol for data encryption and privacy key management.

This section describes the functionalities of MAC CPS, which henceforth will be referred to as MAC.

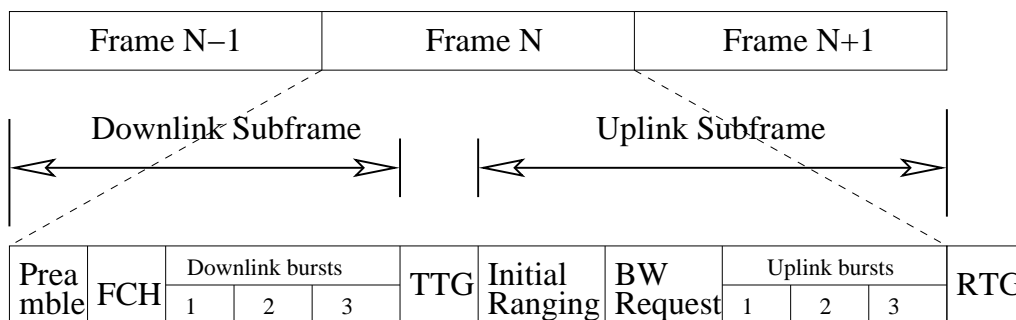


Figure 5.1: IEEE 802.16 MAC frame structure in Time Division Duplex mode; FCH: Frame Control Header, RxDS: Receiver Delay Spread clearing interval, pp. 395 and 450 [21].

The MAC protocol supports a variable length, frame-based system for receiving and transmission of data. In time division duplex mode, for example, each frame comprises a downlink and an uplink subframes separated by a Transmit/Receive transition gap (TTG) and a Receive/Transmit transition gap (RTG) as illustrated in Fig. 5.1. The downlink subframe always precedes its corresponding uplink subframe.

The MAC protocol is connection-oriented. Fig. 5.2 depicts the QoS architecture of an SS and a BS based on the 802.16 standard. Before a logical connection is established between an SS and a BS, the SS sends a connection request to the BS. The request includes information about QoS requirements including information such as bandwidth required and tolerable delay. Then, the request is analyzed by the admission control residing at the BS. The connection is admitted only if the required QoS can be satisfied without affecting the QoS of existing connections. Each accepted connection will be assigned a unique connection ID (CID) by the BS, which may be further used during the call for further bandwidth requests and QoS requirement information and will be carried by the MAC PDU headers.

After the connection has been established, traffic from each connection is generated and eventually arrives at the MAC layer where MAC PDUs are formed. The MAC PDUs are classified into different traffic queues according to their CIDs, as shown in Fig. 5.2. Such a per-connection queueing allows the scheduler to provide differentiated service to connections belonging to different service classes.

The procedures for downlink and uplink transmission are different. In the downlink, transmission is relatively simple because there is no contention as only the BS can transmit data in the downlink subframes. The downlink scheduler at the BS selects packets from the appropriate queue for the next

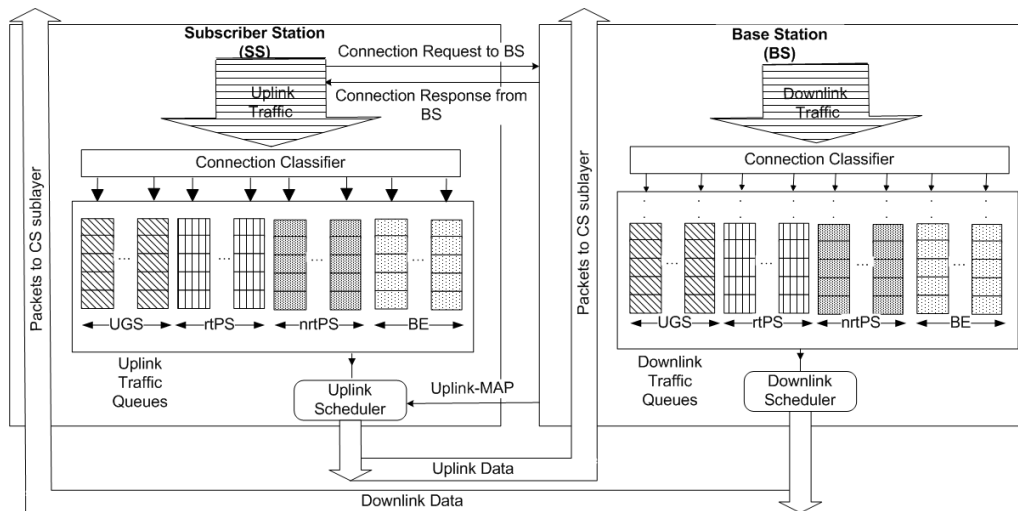


Figure 5.2: QoS architecture of IEEE 802.16

downlink frame according to the QoS parameters and the queue status of each connection. The transmission is broadcast to all SSs and each SS listens and selects only the packets destined to itself.

In regards to uplink transmission, there is a need to overcome possible contention as many SSs share the single wireless channel. In this case, only a single SS is permitted to transmit at a time. Accordingly, to avoid contention, each SS needs to be granted bandwidth before it transmits. This is facilitated by the so-called *request/grant* mechanism that allows the BS to determine the bandwidth requirements of the various SSs and grant bandwidth to them on demand. The 802.16 defines several mechanisms to solicit bandwidth requests from admitted connections:

1. *Unsolicited granting*, where a fixed amount of bandwidth on a periodic basis is specified during connection establishment. Once the connection is accepted, no subsequent bandwidth request is required. Unsolicited granting is used to support UGS connections.

2. *Unicast polling*, which allocates the bandwidth needed to transmit a bandwidth request for a polled uplink connection. Since only one uplink connection is polled at a time, it is guaranteed that the request will be received by the BS. The BS periodically grants unicast polls to rtPS connections. The polling period may optionally be specified during connection establishment.
3. *Broadcast polling*, which is issued to all uplink connections. If there are two or more uplink connections replying with their requests, collision occurs and the truncated binary exponential backoff mechanism is used to resolve the contention. Both nrtPS and BE connections request bandwidth by responding to broadcast polls from the BS.

Based on the amount of requested and granted bandwidth for each connection, the BS uplink grant scheduler estimates the residual backlog of each uplink connection, and then allocates the uplink grants to meet the negotiated QoS parameters. The resultant grant allocation for the current uplink subframe is conveyed to the SSs by the UL-MAP message carried in the frame control header (FCH) field of each downlink subframe. Although an uplink grant is allocated according to individual requests from each connection, the grants are aggregated and given to an SS to be distributed among its connections at its discretion. Therefore, upon receiving the grants, the SS uplink scheduler schedules access to each connection.

5.3 Priority-based Fair Scheduling

Based upon the above discussion, the uplink traffic arriving at an SS belongs to one of the four scheduling services. As the constant-bit-rate UGS traffic is

allocated dedicated bandwidth, it does not share the granted bandwidth with other scheduling services. Accordingly, the scheduler only needs to handle the other three scheduling services. According to their QoS requirements, rtPS, nrtPS and BE are naturally assigned to high, medium and low priorities, respectively. As shown in Fig. 5.2, each connection has its own queue. When a message arrives, it will be broken down into a number of packets. Each packet fits into one time-slot of an uplink frame. In our priority-based fair scheduling algorithm, rtPS connections are always served first. Only if all rtPS connections have no packets waiting, then will the nrtPS connections be served. Similarly, only if all rtPS and nrtPS connections have no packets waiting, will BE connections be served. Whenever there are multiple active connections of the same priority, the scheduler serves one packet from each connection in a round robin fashion.

5.4 Delay Analysis using the MPS Model

Our proposed scheduler can be analysed by the MPS model [9], which has been introduced with the corrected solution in Chapter 3. The multi-priority MPS model consists of a number of groups of distributed LQs and a central server with a PS queue. The PS server performs like an WiMAX SS scheduler and LQs like multiple connections. Group- p LQs contain M_p LQs with priority p , $p = 1, 2, \dots, P$, where P is the lowest priority. Hence, we can assign different priorities to the different WiMAX services. The central server runs prioritized round robin processor sharing among LQs by allowing no more than one message from each LQ to be present in the PS queue. Only when the service of an entire message is completed, is its LQ allowed to transfer another message into the PS queue. The model is shown in Fig. 5.3, with

$P = 3, M_1 = M_2 = M_3 = 3$ as an example.

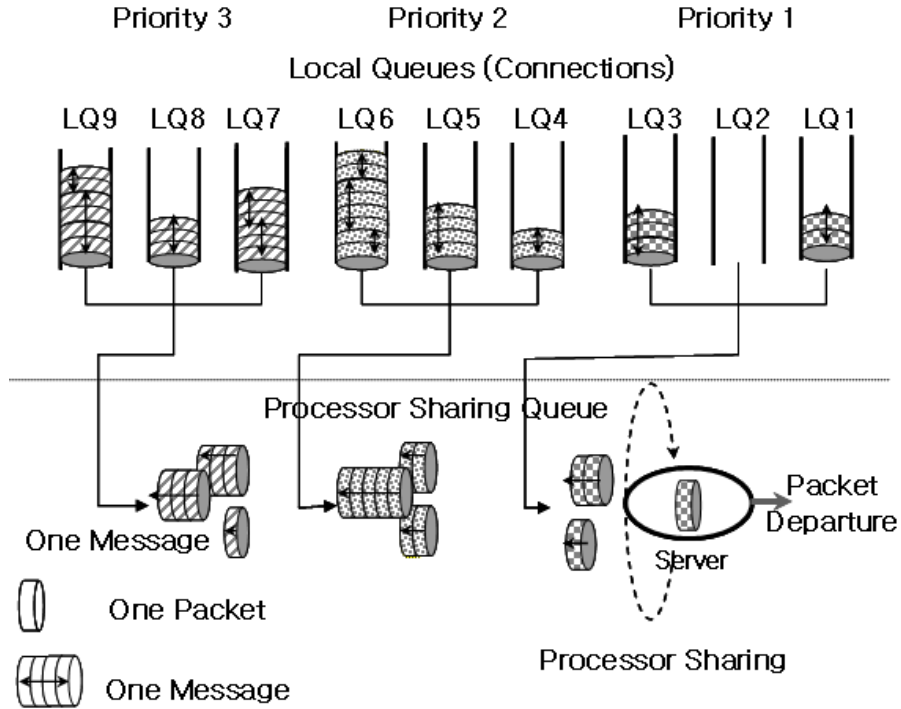


Figure 5.3: WiMAX SS delay analysis using the MPS model.

As we described in Chapter 3, MPS is a discrete-time model whereby time is divided into equal-length time-slots. It assumes that messages arriving at the LQs consist of an integral number of packets, each requiring a service time of a single time-slot. It further assumes the following:

1. For each p , the numbers of priority p messages arriving at a LQ within each time-slot are i.i.d. and are also independent of arrivals to other LQs.
2. The number of packets contained in a message (the message length) are discrete i.i.d. for each priority. The distribution of message lengths may be different for different priorities.
3. The transmission of a message can only be interrupted by messages from

higher priorities or from other connections of the same priority after the current packet is completely transmitted, i.e. until the end of this time-slot.

The mean delay $D_p(n)$ of a priority p message of length n packets is simply given by

$$D_p(n) = L_p + S_p(n), \quad (5.1)$$

where L_p is the mean time spent by a priority p message in its LQ, and $S_p(n)$ is the mean time of a priority p message, consisting of at least n packets, spends in the PS queue to complete services of n packets. Note that L_p is not related to the message length n .

Let the random variable a_p represent the number of priority p message arrivals within a time-slot to any priority p LQ. We denote the mean of a_p as \bar{a}_p . Let the random variable b_p be the priority p message length with the mean \bar{b}_p . Since a packet transmission requires a time-slot, b_p also represents the message transmission time in units of time-slot. Let $C_{a,p}^2$ and $C_{b,p}^2$ represent the squared coefficients of variation for a_p and b_p respectively. According to (3.7), we have:

$$S_p(x) = \frac{x}{1 - \varepsilon_{p-1} - \frac{M_p-1}{M_p}\rho_p} - \delta_p, \quad (5.2)$$

where, $\lambda_p = M_p\bar{a}_p$, $\rho_p = \lambda_p\bar{b}_p$ and $\varepsilon_p = \sum_{i=1}^p \rho_i$. Using the result (3.13):

$$L_p = \frac{\nu_p/\rho_p + \sum_{i=1}^p \nu_i/(1 - \varepsilon_p)}{2(1 - \varepsilon_{p-1})} - \frac{[\bar{b}_p(1 + C_{b,p}^2) + 1]/2}{1 - \varepsilon_{p-1} - \frac{M_p-1}{M_p}\rho_p} + \delta_p + \frac{1}{2}, \quad (5.3)$$

we obtain

$$D_p(n) = \frac{\nu_p/\rho_p + \sum_{i=1}^p \nu_i/(1 - \varepsilon_p)}{2(1 - \varepsilon_{p-1})} + \frac{n - [\bar{b}_p(1 + C_{b,p}^2) + 1]/2}{1 - \varepsilon_{p-1} - \frac{M_p-1}{M_p}\rho_p} + \frac{1}{2}, \quad (5.4)$$

where $\nu_p = \rho_p \bar{b}_p (C_{b,p}^2 + \lambda_p C_{a,p}^2 / M_p)$. The overall mean priority p message delay is simply given by $D_p(\bar{b}_p)$.

We assume WiMAX admission control is strictly applied in order to avoid system overload. Hence, only the admitted connections are involved in the following discussion. We also assume that an SS is granted sufficient bandwidth to serve the admitted connections. Then, this MPS model can be used directly to calculate the mean message delay under our priority-based fair scheduling algorithm. As UGS works with unsolicited granting, we do not need to consider scheduling the UGS service. Connections belonging to rtPS, nrtPS and BE are assigned with priorities 1, 2 and 3, respectively.

By (5.4), we are able to examine the impact on the message delay due to multiple WiMAX service classes. Even when the nrtPS or BE loads are changed, rtPS messages still receive the same service rate because of priority protection. In other words, the mean message delay of rtPS service is only affected by its traffic characteristics such as arrival rate, message length and their variation. However, for nrtPS and BE, their delay performance would also be affected by higher priority traffic.

5.5 Model Evaluation

In this section, the model will be validated by simulation. We first present the traffic model of a typical application of each scheduling service. We then describe the simulation environment used to validate the MPS model. Based on these traffic models, simulation and analytical results for different traffic loads are compared.

5.5.1 Traffic Models

A considerable amount of research on traffic modelling has been carried out to investigate the characteristics of different traffic sources for various communication networks [192, 193, 194, 195]. Here, we refer to their results.

VoIP with silence suppression is a typical application of an rtPS service. It is usually modelled as an exponential ON/OFF source, where the source alternates between an ON and an OFF state. For a talk spurt, the source is in the ON state during which it generates data at a constant rate. For a silence period, the source is in the OFF state implying that no data is generated. The durations of the ON and OFF states are exponentially distributed with their own mean values. For nrtPS, an example application is web access. A possible model for web access is to have the message arrivals modelled as a Poisson process and the message size following a cut-off Pareto distribution [181] with shape parameter, scale parameter and cut-off threshold. Finally, the commonly used Poisson arrival process with exponentially distributed message sizes is chosen as a model for a BE traffic source. The parameters of each traffic source used in this work are listed in Table 5.1.

5.5.2 Simulation Model

The analytical model is validated using the ns2 [189] simulation tool. In the simulation, CBQ objects implement a packet-by-packet round-robin processor sharing within the same priority class. The CBQ buffer size in the simulation is made large enough so that it corresponds to the equivalent assumptions of our analytical model. The simulation has the following settings:

- An aggregate bandwidth of 0.25 Mbits/s is assumed to be granted to and shared by all connections of an SS.

Traffic sources	Arrival process	Message size distribution	Priority level
VoIP source	Exponential ON/OFF: mean OFF period: 1.67s; mean ON period: 1.34s (one packet per 20ms during ON).	Deterministic, the size: 66 bytes.	1
Web source	Poisson, mean inter-arrival time: 5s.	cut-off Pareto: shape parameter $\alpha=1.1$; Minimum message (scale parameter): 4.5k bytes; Maximum message (cut-off threshold): 2M bytes.	2
BE traffic	Poisson.	Exponential, the mean size: 1500 bytes.	3

Table 5.1: Arrival processes and message size distributions of the traffic sources with a priority arrangement

- The above-mentioned traffic models are used to generate input traffic. The generated message size under continuous exponential and Pareto distributions is rounded to the nearest integral number of bytes.
- The operation of the scheduler follows what are discussed in Section 5.3.
- The overall utilisation is limited to be less than unity during the simulation in order to ensure that the system remains stable.
- The simulations keep the records of time when a message is generated and when the whole message is completely served, then, use them in the statistic function to calculate the mean message delay. Each set of input parameters will be run six times independently for confidence intervals.

5.5.3 Simulation and Numerical Results

We consider three scenarios according to different traffic loads. Analytical results and simulation results are compared in each scenario and presented in Fig. 5.4 to Fig. 5.7, respectively. Confidence intervals of 95% based on a Student's t-test are obtained for all of the simulation results. The range of the confidence interval of each point on the simulation curves is within 7.5% which is hardly noticeable on the figures.

Firstly, we fix the load of nrtPS and BE traffic, and investigate the effect of the load of rtPS on the mean message delay by increasing the number of ON/OFF sources. Traffic of lower priorities (nrtPS and BE) contribute about 50% of the load, and each priority has one source. We set the mean inter-arrival time of the BE traffic to be $0.133s$, the mean message sizes of nrtPS and BE traffic as 20.6 kbytes and 1500 bytes, $C_{b,2}^2 = 8$ and $C_{b,3}^2 = 1$ respectively. The number of VoIP sources increases from one to ten and each of these has the same traffic parameter values: $\bar{b}_1 = 66$ bytes and $C_{b,1}^2 = 0$. Results obtained from the simulation and analytical models are shown in Fig. 5.4 and Fig. 5.5. It can be seen that the simulation and analytical results are in good agreement, particularly for the nrtPS and BE traffic. Some divergency at high VoIP loading as the traffic of multiple VoIP connections modelled as Poisson for obtaining analytical results, in fact, has some correlation. But for nrtPS and BE message sizes are big, this effect is not significant.

In the second scenario, rtPS and BE traffic are fixed but the nrtPS traffic is changed by increasing the number of connections. We have four VoIP sources and one BE source in this scenario. Parameters of the BE traffic are set to be the same as in the first scenario: the mean inter-arrival time is $0.133s$ and the mean message size is 1500 bytes. The number of web sources increases

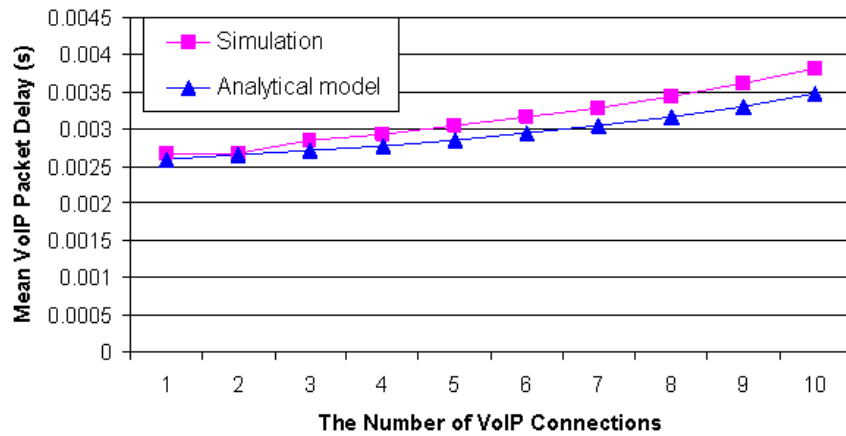


Figure 5.4: Comparison between analytical results and simulation results for VoIP packet delay.

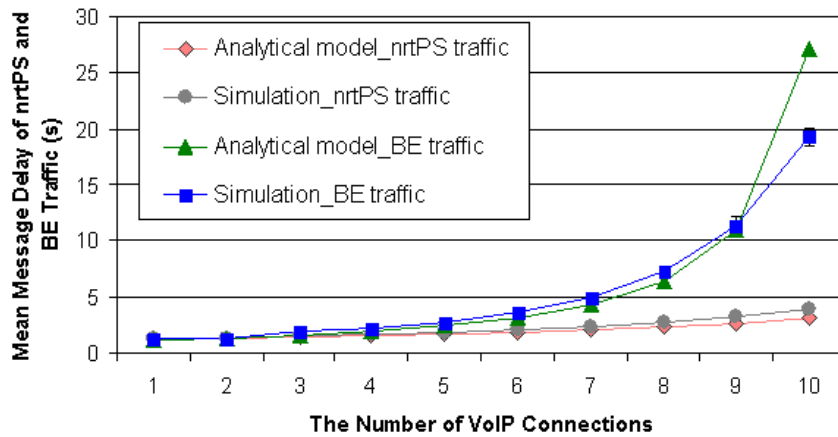


Figure 5.5: nrtPS and BE message delays under different rtPS traffic loads.

from one to four. Fig. 5.6 shows the delay for rtPS traffic is not affected much by the change of nrtPS loads. However, as the web source traffic features big message sizes and big variations in size, even increasing one connection leads to a dramatic growth in the BE traffic delay.

Finally, we keep the same load for rtPS and nrtPS traffic. Three VoIP sources and two web sources are used in this scenario. Another three sources of BE traffic change the total load with variable inter-arrival times. Fig. 5.7

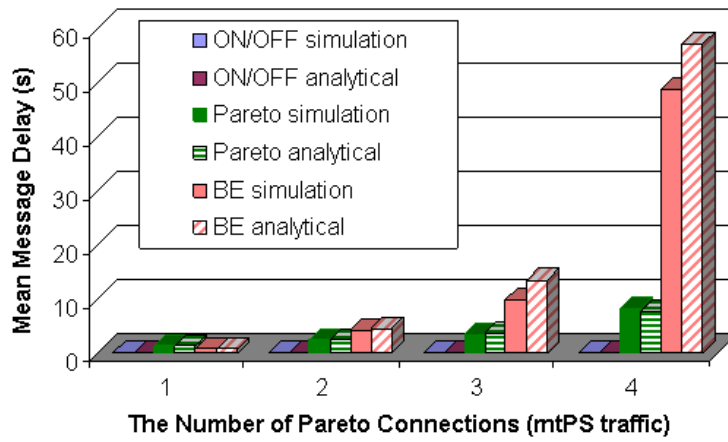


Figure 5.6: WiMAX message delays under different nrtPS traffic loads.

is plotted for this case. As expected, the loading of BE traffic does not affect the message delays of other higher priority traffic.

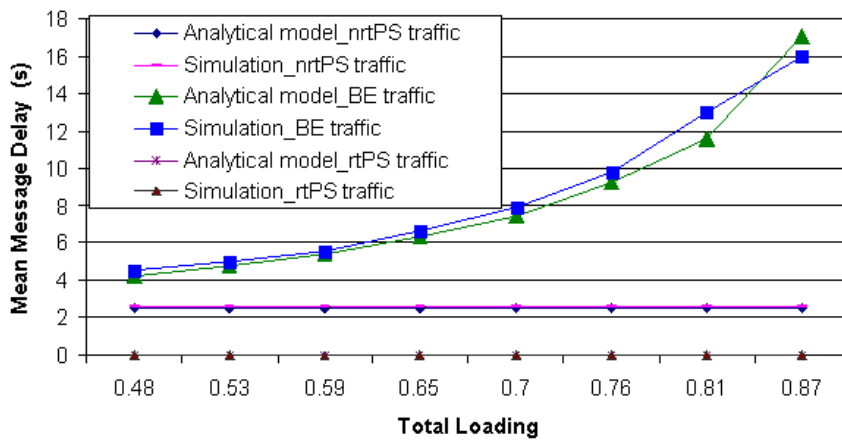


Figure 5.7: WiMAX message delays under different BE traffic loads.

5.6 Conclusions

In this chapter, we have presented a priority-based fair scheduling algorithm for WiMAX uplink traffic at SS and modelled it by a multiservice multiqueue

processor sharing model, MPS queueing system, so that a closed-form approximation is available to approximate WiMAX mean message delays for various service classes. The model has been validated by simulation.

To apply the model, it is important to choose suitable traffic models to represent the characteristics of different traffic sources. At the same time, the first two moments of arrival processes and message size distributions can be available for the analytical solution. In the simulation, we use VoIP silence suppression as for the case of rtPS service. Note that VoIP packet arrivals are not i.i.d.. However, the traffic of multiple VoIP connections can be modelled as Poisson according to the research in [192]. In the scenario one of simulations, we can see that some divergency between the simulation results and the obtaining analytical results at high VoIP loading, as the real traffic arrival process has some correlation. But the effect is not significant since the Poisson approximation is reasonable and nrtPS and BE message sizes are big.

Using the model we have studied various effects of traffic loading on the performance of various services. Our results show that, as expected, the mean message delay can be seriously affected by high traffic loads; however, the delay for real-time services can be maintained at an acceptable level if they are protected by a priority mechanism from the non-real time traffic loads. Moreover, we have also demonstrated the effect of traffic burstiness on delay performance. The proposed model provides a tractable method for operators to manage and dimension their WiMAX systems for uplink multimedia traffic.

Chapter 6

Conclusions

This dissertation has been an investigation of MAC modelling issues in the design and development of wireless networks with multimedia services. The main objective has been to study multiservice multiqueue PS models and to apply these models to analyse delay performance in wireless applications. Current scheduling algorithms for multimedia wireless networks have been considered mostly for non-realtime services and faced the problem of complexity of computation and lack of the analytical models to evaluate the performance. We proposed a multipriority PS discipline at the scheduler of wireless networks to handle multiple services, so the delay requirement of real-time service can be protected by the priority assignment. This multipriority PS scheduling policy is simple to implement at the scheduler. Moreover, the analytical models are suitable for dimensioning purposes.

We investigated multiservice multiqueue models, MPS and PBSQ, to analyse the delay performance in a system following the multipriority PS discipline. The analytical models incorporate important parameters to characterise different network traffic. Based on a comprehensive description of the MPS model, we demonstrate and correct a subtle incongruity for the delay

in a local queue of [9]. Also, we developed the PBSQ model to extend the analysis of the ordinary multiservice PS scheduling policy to the analysis of a priority-based service quanta scheduling policy, where a priority-based service quantum is given by the PS server at each time rather than a fixed quantum. A good closed-form approximation of the mean message delay was obtained in Chapter 3.

We validated the new model using simulations and carried out a numerical study to demonstrate the effects of different parameters. Results show that the priority-based service quanta scheduling policy can improve delay performance of large size requests if they are given a larger service quantum than ordinary multiservice PS scheduling. It does so, however, at the cost of an additional delay for smaller sized requests if they have relatively smaller service quanta. Also, an efficiency increase can be achieved by reducing the overhead.

The work recommends a simple, but practical, choice for wireless MAC scheduling. At the same time, multiservice multiqueue PS models are available to apply to the general performance analysis of MAC protocols with QoS specifications. When a MAC structure is capable of service differentiation, the scheduler can assign priorities for different traffic and perform as a multipriority PS server to transmit data. The above mentioned analytical results can be used to answer two fundamental questions for the CAC, i.e., how to properly allocate the bandwidth required for transmission over a multimedia wireless link, and if the CAC should accept the new connections. The decision is subject to realistic traffic conditions and meeting specified QoS requirements.

We carried out studies using the multiservice multiqueue model for wireless applications, the following areas have been considered in detail: MAC layer QoS architectures including issues in the design and performance spec-

ification; network stability for analysis; service types and the alterations to their traffic models. Using the research results presented, the thesis chapters have established analytical methods leading to the development of specific guidelines for the dimensioning of wireless multimedia networks.

In PTT over GPRS/GSM networks, the MPS model has been applied to estimate the delay of PTT packets. We considered a common scenario for GPRS/GSM networks, where the bandwidth is given to all traffic and the partial sharing channel allocation scheme is adopted to handle GSM voice calls and GPRS (including PTT) packets. Under a quasi-stationary assumption, the use of the MPS model is invoked for all GPRS/GSM traffic. The PTT voice packet delay can be estimated by our analytical model. Moreover, for GPRS traffic, the cases of *with and without priority* for PTT service have been studied using the same model. The work shows that there is good protection for the PTT delay by reserving of channels and through priority assignment. Our solution for the mean packet delay can also be used directly to quantify the effect of retransmissions.

We also applied the MPS model to WiMAX networks. Based on the QoS architecture, we proposed a priority-base fair scheduling for the SS scheduler. Through modelling of this scheduling algorithm's operation at SS under the traffic from specified WiMAX services, some guidelines for the delay analysis have been provided for WiMAX provisioning. A closed-form approximation of the mean message delay is obtained for various WiMAX service classes. The model enables the performance study of the effects due to traffic loading and burstiness. We demonstrated that our proposed scheduling algorithm provides an effective protection for the delay of real-time traffic. Although high traffic loads and burstiness seriously affect the mean message delay, the delay of real-time services can be kept at an acceptable level. Therefore, a tractable

method for use by WiMAX operators is available for system management and dimensioning.

6.1 Summary of Contributions

A final list of the contributions by this thesis is given below:

- Development of a generalised multiservice multiqueue PS model that allows different service quanta to be used for different types of services.
- Correction of an incongruity in the literature associated with the determination of the delay for the developed MPS model.
- Development of a multiservice multiqueue PS model to analyse PTT packet delay that takes into account GSM voice traffic and GPRS data traffic in GPRS/GSM networks.
- An analytical solution for PTT retransmissions in GPRS/GSM networks.
- A study of PTT packet delay in PTT over an GPRS/GSM network using ns-2 and C++ simulators.
- Development of a priority-based scheduling discipline applied at the SS in WiMAX networks.
- An analytical study of the message delay under priority-based scheduling in a WiMAX network.
- Design and implementation of a simulation model for WiMAX message delays using ns-2.

6.2 Future Research

Significant improvements to the development of multiservice multiqueue PS models have been made in this work. However, there are still some interesting issues remaining that need to be further studied and addressed. In conclusion, we comment on some of these areas for future research.

6.2.1 Enhanced Models

It might be possible to derive an exact expression for the PBSQ model, although whether such an exact expression would significantly improve on the accuracy of the model is questionable.

We have shown the efficacy of the multiservice multiqueue PS models to accurately predict delay performance through using them for wireless applications considered in this thesis, but it is not conclusive proof of their universal suitability. In the specific approach taken in this thesis, there is further work remaining in the development of more comprehensive extended models under more general assumptions, such as traffic models, that are applicable to a wider range of wireless systems.

Thus far, the PBSQ model has been developed only under the assumption of a Poisson arrival process. It is our hope that an extension can be made for more general arrival processes and service distributions. If there is a more generalised model only using the first two statistical moments of real traffic, the model can be used as a tool whenever an accurate traffic model is available. Validation of such models could be done by simulation using data traces. Such work will largely improve the applicability of the model.

In the PBSQ model, the parameter N_p is set arbitrarily. We know that large values of N_p will affect the delay results. An alternative effort that can

be made is relating N_p to some measurable statistic, such as the statistic of an aggregate data stream. If it is not feasible to estimate the parameters of the real aggregate traffic, an estimate may be derived from the requirements of individual sources contributing to the traffic.

6.2.2 More Applications

Two applications of multiservice multiqueue PS models have been addressed in Chapters 4 and 5 of this thesis. However, some other wireless applications of these models exist, such as the IEEE 802.11e protocol. The standard defines two new MAC modes to support up to eight-priority traffic classes that map directly to the RSVP protocol. Therefore, our multiservice PS scheduling policy can be implemented at the scheduler of 802.11e networks and the delay performance can be modelled.

6.2.3 Effects of Further MAC Scheduling and Modelling Issues

Although focusing on a single layer is the usual approach for protocol design and analysis, recent research has involved capturing the interactions across different layers, such as the interaction between the MAC and PHY layers. Some of these research results are being proposed for MAC scheduling as a potential way to improve the performance of wireless networks. As cross-layer interaction affects overall system performance, there are many open questions of quantitative analysis for the performance modelling of such approaches. A further study of new models should involve the effect of these interactions.

In wireless networks, since traffic demands might change rapidly and the mechanism does not always perform under stable channel conditions, it is dif-

difficult to achieve end-to-end QoS just by working on a single simple scheduling algorithm. One must consider a good combination of CAC and scheduling algorithms as well. The CAC and its co-worked scheduling policy are equally important. Hence, it is expected that the optimal solution should take into account these issues at the MAC layer, though some foundation modelling work is essential for network operators.

Bibliography

- [1] A. J. Paulraj *et al.*, “An overview of MIMO communications - a key to gigabit wireless,” *Proceedings of the IEEE*, vol. 92, pp. 198–218, 2004.
- [2] R. VanNee and R. Prasad, *OFDM for Wireless Multimedia Communications*. Norwood, MA, U.S.A. Artech House, Inc., 2000.
- [3] C. Berrou and A. Glavieux, “Near optimum error correcting coding and decoding: turbo-codes,” *IEEE Transactions on Communnications*, vol. 44, no. 10, pp. 1261–1271, October 1996.
- [4] A. S. Tanenbaum, *Computer Networks (4th ed.)*. Prentice Hall, Inc., 2003.
- [5] L. Kleinrock, *Queueing Systems Volume II: Computer Applications*. John Wiley & Sons Inc., 1976.
- [6] S. Aalto *et al.*, “Beyond processor sharing,” *ACM SIGMETRICS Performance Evaluation Review*, vol. 34, pp. 36–43, March 2007.
- [7] A. K. Parekh and R. G. Gallager, “A generalized procoessor sharing approach to flow control in integrated services networks: the single-node case,” *IEEE/ACM Transactions on Networking*, vol. 1, no. 3, pp. 344–357, June 1993.

- [8] A. K. Parekh and R. G. Gallager, "A generalized procoessor sharing approach to flow control in integrated services networks: the multiple node case," *IEEE/ACM Transactions on Networking*, vol. 2, no. 2, pp. 137–150, April 1994.
- [9] P. Potter and M. Zukerman, "Analysis of a discrete multi-priority queueing system involving a central shared processor serving many local queues," *IEEE Journal on Selected Areas in Communications*, no. 2, pp. 194–202, February 1991.
- [10] L. A. DaSilva *et al.*, "The resurgence of push-to-talk technologies," *IEEE Communications Magazine*, pp. 48–55, 2006.
- [11] "The WiMAX Forum." [Online]. Available: <http://www.wimaxforum.org/>.
- [12] S. Chan *et al.*, "Algorithms for WiMAX Scheduling," in *Proceedings of APMC 2008*, pp.1-4, Hong Kong, December 2008.
- [13] Y. Wang, M. Zukerman, and R. Harris, "PTT packet delay analysis for GPRS/GSM links," *IEEE Communications Letters*, no. 6, pp. 456–458, June 2006.
- [14] Y. Wang, M. Zukerman, and R. Harris, "Modeling PTT uplink in GPRS/GSM networks," in *Proceedings of IEEE VTC 2006-Spring*, pp. 420–424, May 2006.
- [15] Y. Wang *et al.*, "Priority-based fair scheduling for multimedia WiMAX uplink traffic," in *Proceedings of IEEE ICC 2008*, pp. 301–305, May 2008.

- [16] Y. Wang *et al.*, “A priority-based Processor Sharing Model for TDM Passive Optical Networks,” *Submitted to IEEE Journal on Selected Areas in Communications*, August 2009.
- [17] “Wireless local area networks - the working group for WLAN standards,” *IEEE STD 802.11*, 1999.
- [18] J. R. Gallardo, P. Medina, and W. Zhuang, “QoS mechanisms for the MAC protocol of IEEE 802.11 WLANs,” *Wireless Networks*, vol. 13, no. 3, pp. 335–349, 2007.
- [19] J. Jin and K. Nahrstedt, “QoS specification languages for distributed multimedia applications: a survey and taxonomy,” *IEEE MultiMedia*, vol. 11, no. 3, pp. 74–87, July-September 2004.
- [20] “QoS Parameters.” [Online]. Available: <http://www.wireless-center.net/Mobile-and-Wireless/767.html/> .
- [21] “Air interface for fixed broadband wireless access systems,” *IEEE STD 802.16 - 2004*, October 2004.
- [22] D. Bertsekas and R. Gallager, *Data networks (2nd ed.)*. Prentice-Hall, Inc., 1992.
- [23] Y. Cao and V. O. K. Li, “Scheduling algorithms in broad-band wireless networks,” *Proceedings of the IEEE*, vol. 89, no. 1, pp. 76–87, January 2001.
- [24] “IEEE 802.15 Working Group.” [Online]. Available: <http://www.ieee802.org/15/> .

- [25] H. J. Chao and X. Guo, *Quality of Service Control in High-Speed Networks*. John Wiley & Sons, Inc., 2002.
- [26] A. Demers, S. Keshav, and S. Shenker, "Analysis and simulation of a fair queueing algorithm," *SIGCOMM'89*, vol. 19, no. 4, pp. 1–12, September 1989.
- [27] L. Zhang, "Virtual clock: a new traffic control algorithm for packet switching networks," *SIGCOMM Comput. Commun. Rev.*, vol. 20, no. 4, pp. 19–29, 1990.
- [28] D. Ferrari and D. C. Vexma, "A scheme for real-time channel establishment in wide-area networks," *IEEE Journal on Selected Areas in Communications*, vol. 8, pp. 368–379, April 1990.
- [29] P. Bhagwat *et al.*, "Enhancing throughput over wireless LANs using channel state dependent packet scheduling," in *Proceedings of INFOCOM'96*, vol. 3, pp. 1133–1140, 1996.
- [30] S. Floyd and V. Jacobson, "Link-sharing and resource management models for packet networks," *IEEE/ACM Transactions on Networking*, vol. 3, pp. 365–386, August 1995.
- [31] H. Zhang, "Service disciplines for guaranteed performance service in packet-switching networks," *Proceedings of the IEEE*, vol. 83, no. 10, pp. 1374–1396, October 1995.
- [32] C. Fragouli, V. Sivaraman, and M. B. Srivastava, "Controlled multimedia wireless link sharing via enhanced class-based queueing with channel-state-dependent packet scheduling," in *Proceedings of INFOCOM'98*, pp. 572–580, March 1998.

- [33] S. Lu, V. Bharghavan, and R. Srikant, “Fair scheduling in wireless packet networks,” *IEEE/ACM Transactions on Networking*, vol. 7, pp. 473–489, August 1999.
- [34] T. S. E. Ng, I. Stoica, and H. Zhang, “Packet fair queueing algorithms for wireless networks with location-dependent errors,” in *Proceedings of IEEE INFOCOMM’98*, vol. 3, pp. 1103–1111, March 1998.
- [35] P. Goyal, H. M. Vin, and H. Cheng, “Start-time fair queueing: a scheduling algorithm for integrated services packet switching networks,” *IEEE/ACM Transactions on Networking*, vol. 5, no. 5, pp. 690–704, October 1997.
- [36] P. Ramanathan and P. Agrawal, “Adapting packet fair queueing algorithms to wireless networks,” in *Proceedings of ACM/IEEE MOBICOM*, pp. 1 – 9, 1998.
- [37] A. Gyasi-Agyei and S. Kim, “Cross-layer multiservice opportunistic scheduling for wireless networks,” *IEEE Communications Magazine*, vol. 44, no. 6, pp. 50–57, June 2006.
- [38] L. Kleinrock, “Analysis of a time-shared processor,” *Naval Research Logistics Quarterly*, vol. 11, pp. 59–73, 1964.
- [39] L. Kleinrock, “Time-shared systems: a theoretical treatment,” *Journal of the ACM*, vol. 14, no. 2, pp. 242–261, 1967.
- [40] S. Ben Fredj *et al.*, “Statistical bandwidth sharing: a study of congestion at flow level,” *SIGCOMM Comput. Commun. Rev.*, vol. 31, no. 4, pp. 319–332, 2001.

- [41] F. Delcoigne, A. Proutiere, and G. Regni, "Modeling integration of streaming and data traffic," *Performance Evaluation*, vol. 55, pp. 185–209, February 2004.
- [42] E. Altman, K. Avrachenkov, and U. Ayesta, "A survey on discriminatory processor sharing," *Queueing System*, vol. 53, no. 1-2, pp. 53–63, 2006.
- [43] J. W. Roberts, "A survey on statistical bandwidth sharing," *Computer Networks*, vol. 45, no. 3, pp. 319–332, 2004.
- [44] S. Borst, "User-level performance of channel-aware scheduling algorithms in wireless data networks," *IEEE/ACM Transactions on Networking*, vol. 13, no. 3, pp. 636–647, 2005.
- [45] T. Bonald *et al.*, "Insensitivity results in statistical bandwidth sharing," *17th ITC*, December 2001.
- [46] R. Litjens *et al.*, "Performance analysis of wireless LANs: an integrated packet/flow level approach," *18th ITC*, 2003.
- [47] M. Sakata, S. Noguchi, and J. Oizumi, "Analysis of a processor shared model for time sharing systems," in *Proceedings of 2nd Hawaii International Conference on System Sciences*, pp. 625–628, January 1969.
- [48] J. D. C. Little, "A proof of the queueing formula $L = \lambda W$." *Operations Research*, vol. 9, no. 3, pp. 383–387, May-June 1961.
- [49] M. Sakata, S. Noguchi, and J. Oizumi, "An analysis of the M/G/1 queue under round-robin scheduling," *Operations Research*, vol. 19, no. 2, pp. 371–385, March-April 1971.

- [50] J. W. Cohen, "The multiple phase service network with generalized processor sharing," *Acta Information*, vol. 12, no. 3, pp. 245–284, October 1979.
- [51] E. G. Coffman, R. R. Muntz, and H. Trotter, "Waiting time distributions for processor-sharing systems," *Journal of the ACM*, vol. 17, no. 1, pp. 123–130, 1970.
- [52] J. A. Morrison, "Response-time distribution for a processor-sharing system," *SIAM Journal on Applied Mathematics*, vol. 45, no. 1, pp. 152–167, 1985.
- [53] S. F. Yashkov, "A derivation of response time distribution for a M/G/1 processor sharing queue," *Problems Control Inform. Theory* 12, pp. 133–148, 1983.
- [54] T. J. Ott, "The sojourn-time distribution in the M/G/1 queue with processor sharing," *Journal of Applied Probability*, vol. 21, no. 2, pp. 360–378, June 1984.
- [55] R. Schassberger, "A new approach to the M/G/1 processor-sharing queue," *Advances in Applied Probability*, vol. 16, no. 1, pp. 202–213, March 1984.
- [56] J. L. Van den Berg and O. J. Boxma, "The M/G/1 queue with processor sharing and its relation to a feedback queue," *Queueing Systems*, vol. 9, pp. 365–402, October 1991.
- [57] A. P. Zwart and O. J. Boxma, "Sojourn time asymptotics in the M/G/1 processor sharing queue," *Queueing Systems*, vol. 35, no. 1–4, pp. 141–166, 2000.

- [58] S. K. Cheung, J. L. vandenBerg, and R. J. Boucherie, "Insensitive bounds for the moments of the sojourn time distribution in the M/G/1 processor-sharing queue," *Queueing Systems*, vol. 53, no. 1-2, pp. 7–18, June 2006.
- [59] S. K. Cheung, *Processor-Sharing Queues and Resource Sharing in Wireless LANs*. PhD thesis, University of Twente, Netherlands, 2007.
- [60] S. Borst, R. Nunez-Queija, and B. Zwart, "Sojourn time asymptotics in processor sharing queues," *Queueing Systems*, vol. 53, no. 1-2, pp. 31–51, June 2006.
- [61] W. A. Massey, "Asymptotic analysis of the time dependent M/M/1 queue," *Mathematics of Operations Research*, vol. 10, pp. 305–327, 1985.
- [62] R. Nunez-Queija, "Sojourn times in a processor sharing queue with service interruptions," *Queueing Systems*, vol. 34, no. 1-4, p. 351-386, 2000.
- [63] R. Egorova, A. P. Zwart, and O. Boxma, "Sojourn time tails in the M/D/1 processor sharing queue," *Probability in the Engineering and Informational Sciences*, vol. 20, no. 3, pp. 429-446, 2006.
- [64] M. Mandjes and A. P. Zwart, "Large deviations of sojourn times in processor sharing queues," *Queueing Systems*, vol. 52, no. 4, pp. 237–250, April 2006.
- [65] A. Brandt and M. Brandt, "A sample path relation for the sojourn times in G/G/1-PS systems and its applications," *Queueing Systems*, vol. 52, no. 4, pp. 281–286, April 2006.

- [66] B. K. Asare and F. G. Foster, "Conditional response times in the M/G/1 processor sharing system," *Journal of Applied Probability*, vol. 20, no. 4, pp. 910-915, December 1983.
- [67] S. F. Yashkov, "Processor-sharing queues: Some progress in analysis," *Queueing Systems*, vol. 2, no. 1, pp. 1-17, June 1987.
- [68] S. F. Yashkov, "Mathematical problems in the theory of processor-sharing queueing systems," *Journal of Mathematical Sciences*, vol. 58, no. 2, pp. 101-147, January 1992.
- [69] S. F. Yashkov, "On a heavy traffic limit theorem for the M/G/1 processor sharing queue," *Stochastic Models*, vol. 9, no. 3, pp. 467-471, 1993.
- [70] F. Baccelli and D. Towsley, "The customer response times in the processor sharing queue are associated," *Queueing System*, vol. 7, no. 3-4, pp. 269-282, November 1990.
- [71] S. Grishechkin, "On a relationship between processor-sharing queues and Crump-Mode-Jagers branching processes," *Advances in Applied Probability*, vol. 24, no. 3, pp. 653-698, 1992.
- [72] F. Guillemin and J. Boyer, "Analysis of the M/M/1 queue with processor sharing via spectral theory," *Queueing Systems*, vol. 39, no. 4, pp. 377-397, 2001.
- [73] N. Bansal, "Analysis of the M/G/1 processor-sharing queue with bulk arrivals," *Operations Research Letters*, vol. 31, no. 5, pp. 401-405, September 2003.
- [74] Y. Kitayev, "The M/G/1 processor-sharing model: transient behavior," *Queueing Systems*, vol. 14, no. 3-4, pp. 239-273, September 1993.

- [75] R. C. Hampshire, M. Harchol-Balter, and W. A. Massey, "Fluid and diffusion limits for transient sojourn times of processor sharing queues with time varying rates," *Queueing Systems*, vol. 53, no. 1-2, pp. 19–30, June 2006.
- [76] J. Kim and B. Kim, "The processor-sharing queue with bulk arrivals and phase-type services," *Performance Evaluation*, vol. 64, no. 4, pp. 277–297, May 2007.
- [77] Y. Wu, C. Williamson, and J. Luo, "On processor sharing and its applications to cellular data network provisioning," *Performance Evaluation*, vol. 64, no. 9-12, pp. 892–908, October 2007.
- [78] M. J. G. Van Uitert, *Generalized Processor Sharing Queues*. PhD thesis, Eindhoven University of Technology, Netherlands, 2003.
- [79] J. C. R. Bennett and H. Zhang, "Hierarchical packet fair queueing algorithms," *IEEE/ACM Transactions on Networking*, vol. 5, no. 5, pp. 675-689, October 1997.
- [80] G. Fayolle, I. Mitrani, and R. Iasnogorodski, "Sharing a processor among many job classes," *Journal of the ACM*, vol. 27, pp. 519-532, July 1980.
- [81] K. M. Rege and B. Sengupta, "Queue-length distribution for the discriminatory processor-sharing queue," *Operations Research*, vol. 44, no. 4, pp. 653-657, July-August 1996.
- [82] F. M. Pereira, N. L. S. Fonseca, and D. S. Arantes, "On the performance of Generalized Processor Sharing servers under long-range dependent traffic," *Computer Networks*, vol. 40, pp. 413-431, October 2002.

- [83] G. Fayolle and R. Iasnogorodski, "Two coupled processors: the reduction to a Riemann-Hilbert problem," *Probability Theory and Related Fields*, vol. 47, no. 3, pp. 325-351, January, 1979.
- [84] A. G. Konheim, I. Meilijson, and A. Melkman, "Processor-sharing of two parallel lines," *Journal of Applied Probability*, vol. 18, no. 4, pp. 952–956, December 1981.
- [85] J. W. Cohen and O. J. Boxma, *Boundary Value Problems in Queueing System Analysis*. The North Holland, 1983.
- [86] F. Guillemin *et al.*, "Analysis of the fluid weighted fair queueing system," *Journal of Applied Probability*, vol. 40, no. 1, pp. 180-199, March, 2003.
- [87] F. Guillemin and D. Pinchon, "Analysis of generalized processor-sharing systems with two classes of customers and exponential services," *Journal of Applied Probability*, vol. 41, no. 3, pp. 832–858, 2004.
- [88] C. Knessl and J. A. Morrison, "Heavy traffic analysis of two coupled processors," *Queueing Systems*, vol. 43, no. 3, pp. 173–220, March 2003.
- [89] T. Osogami, M. Harchol-Balter, and A. Scheller-Wolf, "Analysis of cycle stealing with switching times and thresholds," *Performance Evaluation*, vol. 61, no. 4, pp. 347 – 369, August 2005.
- [90] A. I. Elwalid *et al.*, "Fundamental bounds and approximations for ATM multiplexers with applications to video teleconferencing," *IEEE Journal on Selected Areas in Communications*, vol. 13, pp. 1004-1016, 1995.
- [91] G. L. Choudhury, D. M. Lucantoni, and W. Whitt, "Squeezing the most out of ATM," *IEEE Transactions on Communications*, vol. 44, no. 2, pp. 203-217, February 1996.

- [92] E. Willekens and J. L. Teugels, "Asymptotic expansions for waiting time probabilities in an M/G/1 queue with long-tailed service time," *Queueing Systems*, vol. 10, no. 4, pp. 295–311, December 1992.
- [93] J. Abate, G. L. Choudhury, and W. Whitt, "Waiting-time tail probabilities in queues with long-tail service-time distributions," *Queueing Systems*, vol. 16, no. 3-4, pp. 311–338, September 1994.
- [94] W. E. Leland *et al.*, "On the selfsimilar nature of Ethernet traffic," *IEEE/ACM Transactions on Networking*, vol. 2, pp. 1–15, February 1994.
- [95] A. Paxson and S. Floyd, "Wide area traffic: the failure of Poisson modeling," *IEEE/ACM Transactions on Networking*, vol. 3, pp. 226-244, June 1995.
- [96] O. Yaron and M. Sidi, "Generalized Processor Sharing networks with exponentially bounded burstiness arrivals," in *Proceedings of IEEE INFOCOM'94*, pp. 628–634, February 1994.
- [97] C. Kluppelberg, "Subexponential distributions and integrated tails," *Journal of Applied Probability*, vol. 25, no. 1, pp. 132-141, 1988.
- [98] Z. L. Zhang, D. Towsley, and J. Kurose, "Statistical analysis of the Generalized Processor Sharing scheduling discipline," *IEEE Journal of Selected Areas in Communications*, vol. 13, pp. 1071-1080, August 1995.
- [99] G. deVeciana and G. Kesidis, "Bandwidth allocation for multiple qualities of service using Generalized Processor Sharing," *IEEE Transactions on Information Theory*, vol. 42, pp. 268-272, January 1996.

- [100] Z. L. Zhang, "Large deviations and the Generalized Processor Sharing scheduling for a two-queue system," *Queueing Systems*, vol. 26, no. 3-4, pp. 229–245, November 1997.
- [101] Z. L. Zhang *et al.*, "Call admission control schemes under Generalized Processor Sharing scheduling," *Telecommunication Systems*, vol. 7, pp. 125-152, 1997.
- [102] Z. L. Zhang, "Large deviations and the Generalized Processor Sharing scheduling for a multiple-queue system," *Queueing Systems*, vol. 28, no. 4, pp. 349–376, 1998.
- [103] D. Bertsimas, I. C. Paschalidis, and J. N. Tsitsiklis, "Large deviations analysis of the generalized processor sharing policy," *Queueing Systems*, vol. 32, no. 4, pp. 319-349, 1999.
- [104] L. Massoulié, "Large deviations estimates for polling and weighted fair queueing service systems," *Advances in Performance Analysis*, vol. 2, pp. 103-128, 1999.
- [105] P. Dupuis and K. Ramanan, "A Skorokhod problem formulation and large deviation analysis of a processor sharing model," *Queueing Systems*, vol. 28, no. 1-3, pp. 109-124, May 1998.
- [106] T. Mikosch *et al.*, "Is network traffic approximated by stable Levy motion or fractional brownian motion?" *Annals of Applied Probability*, vol. 12, no. 1, pp. 23–68, February 2002.
- [107] M. J. G. Van Uitert and S. C. Borst, "A reduced-load equivalence for Generalised Processor Sharing networks with long-tailed input flows," *Queueing Systems*, vol. 41, no. 1-2, pp. 123-163, June 2002.

- [108] B. C. Borst, O. J. Boxma, and P. R. Jelenkovic, "Generalized processor sharing with long-tailed traffic sources," *In Proceedings of IEEE ITC-16*, pp. 345-354, 1999.
- [109] B. C. Borst, and O. J. Boxma, and P. R. Jelenkovic, "Reduced-load equivalence and induced burstiness in gps queues with long-tailed traffic flows," *Queueing Systems*, vol. 43, no. 4, pp. 273-306, 2003.
- [110] S. C. Borst, O. J. Boxma, and M. J. G. VanUitert, "The asymptotic workload behavior of two coupled queues," *Queueing Systems*, vol. 43, no. 1-2, pp. 81-102, January 2003.
- [111] S. C. Borst, S. C. Mandjes, and M. J. G. VanUitert, "Generalized processor sharing queues with light-tailed and heavy-tailed input," *IEEE/ACM Transactions on Networking*, vol. 11, no. 5, pp. 821-834, October 2003.
- [112] B. C. Borst, O. J. Boxma, and P. R. Jelenkovic, "Asymptotic behavior of generalized processor sharing with long-tailed traffic sources," *in Proceedings of INFOCOM 2000*, pp. 912-921, 2000.
- [113] P. R. Jelenkovic and A. A. Lazar, "Asymptotic results for multiplexing subexponential on-off processes," *Advances in Applied Probability*, vol. 31, no. 2, p. 394-421, 1999.
- [114] P. R. Jelenkovic and P. Momcilovic, "Network multiplexer with generalized processor sharing and heavy-tailed on-off flows," *in Proceedings of IEEE ITC-17*, pp. 719-730, 2001.

- [115] R. Agrawal, A. M. Makowski, and P. Nain, "On a reduced load equivalence for fluid queues under subexponentiality," *Queueing Systems*, vol. 33, no. 1-3, pp. 5–41, 1999.
- [116] G. Van Kessel, R. Nunez-Queija, and S. C. Borst, "Differentiated bandwidth sharing with disparate flow sizes," in *Proceedings of IEEE INFOCOM 2005*, vol. 4, pp. 2425–2435, 2005.
- [117] C. Kotopoulos, N. Likhanov, and R. R. Mazumdar, "Asymptotic analysis of the GPS system fed by heterogeneous long-tailed sources," in *Proceedings of IEEE INFOCOM 2001*, vol. 4, pp. 299-308, 2001.
- [118] M. Lelarge, "Asymptotic behavior of Generalized Processor Sharing queues under subexponential hypothesis," *Report RR-4339, INRIA Rocquencourt*, 2001.
- [119] S. C. Borst and A. P. Zwart, "A reduced-peak equivalence for queues with a mixture of light-tailed and heavy-tailed input flows," *Advances in Applied Probability*, vol. 35, no. 3, pp. 739-805, 2003.
- [120] O. J. Boxma, Q. Deng, and A. P. Zward, "Waiting-time asymptotics for the M/G/2 queue with heterogeneous servers," *Queueing Systems*, vol. 40, pp. 5–31, February 2002.
- [121] O. J. Boxma and I. A. Kurkova, "The M/G/1 queue with two service speeds," *Advances in Applied Probability*, vol. 33, no. 2, p. 520-540, 2001.
- [122] S. Delas, R. R. Mazumdar, and C. P. Rosenberg, "Tail asymptotics for HOL priority queues handling a large number of independent stationary sources," *Queueing Systems*, vol. 40, pp. 183-204, March 2002.

- [123] C. Kotopoulos and R. R. Mazumdar, "Many sources asymptotics for a 2-buffer system with Generalized Processor Sharing," *Preprint, Purdue University*, 2002.
- [124] C. Kotopoulos and R. R. Mazumdar, "Buffer occupancy and delay asymptotics in multi-buffered systems with Generalized Processor Sharing handling a large number of independent traffic streams," *Preprint, Purdue University*, 2002.
- [125] P. Mannersalo and I. Norros, "GPS schedulers and Gaussian traffic," in *Proceedings of INFOCOM 2002*, vol. 3, pp. 1660-1667, 2002.
- [126] M. Mandjes and M. J. G. Van Uitert, "Sample-path large deviations for generalized processor sharing queues with gaussian inputs," in *Proceedings of INFOCOM 2002*, vol. 61, pp. 225-256, July 2005.
- [127] M. Mandjes and M. J. G. VanUitert, "Sample-path large deviations for tandem and priority queues with gaussian inputs," *The Annals of Applied Probability*, vol. 15, pp. 1193-1226, 2005.
- [128] A. I. Elwalid and D. Mitra, "Design of Generalized Processor Sharing schedulers which statistically multiplex heterogeneous QoS classes," in *Proceedings of INFOCOM 1999*, vol. 3, pp. 1220-1230, 1999.
- [129] K. Kumaran *et al.*, "Novel techniques for the design and control of Generalized Processor Sharing schedulers for multiple QoS classes," in *Proceedings of INFOCOM 2000*, vol. 2, pp. 932-941, 2002.
- [130] P. M. D. Lieshout, M. Mandjes, and S. C. Borst, "GPS scheduling: selection of optimal weights and comparison with strict priorities," *ACM*

- SIGMETRICS Performance Evaluation Review*, vol. 34, pp. 75–86, June 2006.
- [131] A. Panagakis *et al.*, “Optimal call admission control on a single link with a GPS scheduler,” *IEEE/ACM Transactions on Networking*, vol. 12, no. 5, pp. 865–878, October 2004.
- [132] S. Shakkottai and R. Srikant, “Many-sources delay asymptotics with applications to priority queues,” *Queueing Systems*, vol. 39, no. 2-3, pp. 183–200, October 2001.
- [133] R. Addie, P. Mannersalo, and I. Norros, “Most probable paths and performance formulae for buffers with Gaussian input traffic,” *European Transactions on Telecommunications*, vol. 13, pp. 183–196, 2002.
- [134] S. Borst, R. Nunez-Queija and M. J. G. vanUitert, “User-level performance of elastic traffic in a differentiated-services environment,” *Performance Evaluation*, vol. 49, pp. 507–519, September 2002.
- [135] T. Bonald and A. Proutiere, “On stochastic bounds for monotonic processor sharing networks,” *Queueing Systems*, vol. 47, pp. 81–106, May–June 2004.
- [136] T. M. O’Donovan, “Direct solutions of M/G/1 processor-sharing models,” *SIGCOMM’89*, vol. 22, no. 6, pp. 1232–1235, November–December 1974.
- [137] L. Kleinrock, R. R. Muntz, and E. Rodemich, “The processor sharing queueing model for time-shared systems with bulk arrivals,” *Networks Journal*, vol. 1, no. 1, pp. 1–13, 1971.
- [138] K. E. Avrachenkov *et al.*, “Discriminatory processor sharing revisited,” *in Proceedings of INFOCOM 2005*, vol. 2, pp. 784–795, March 2005.

- [139] G. V. Kessel, R. Nunez-Queija, and S. C. Borst, "Asymptotic regimes and approximations for Discriminatory Processor Sharing," *ACM SIGMETRICS Performance Evaluation Review*, vol. 32, pp. 44–46, September 2004.
- [140] J. Kim and B. Kim, "Sojourn time distribution in the M/M/1 queue with discriminatory processor-sharing," *Performance Evaluation*, vol. 58, no. 4, pp. 341-365, December 2004.
- [141] E. Altman, T. Jimenez, and D. Kofman, "DPS queues with stationary ergodic service times and the performance of tcp in overload," in *Proceedings of IEEE INFOCOM 2004*, vol. 2, pp. 975– 983, March 2004.
- [142] A. Jean-Marie and P. Robert, "On the transient behavior of the processor sharing queue," *Queueing Systems*, vol. 17, no. 1-2, pp. 129-136, March 1994.
- [143] O. J. Boxma, N. Hegde, and R. Nunez-Queija, "Exact and approximate analysis of sojourn times in finite discriminatory processor sharing queues," *International Journal of Electronics and Communications*, vol. 60, no. 2, pp. 109-115, 2006.
- [144] S. C. Borst, D. V. Ooteghem, and A. P. Zwart, "Tail asymptotics for discriminatory processor sharing queues with heavy-tailed service requirements," *Performance Evaluation*, vol. 61, no. 2-3, pp. 281-298, July 2005.
- [145] R. Egorova, M. R. H. Mandjes, and A. P. Zwart, "Sojourn time asymptotics in processor sharing queues with varying service rate," *Queueing Systems*, vol. 56, no. 3-4, pp. 169–181, August 2007.

- [146] L. M. T. Bonald, "Impact of fairness on Internet performance," in *Proceedings of ACM SIGMETRICS*, pp. 82-91, June 2001.
- [147] S. Fredj *et al.*, "Statistical bandwidth sharing: A study of congestion at flow level," in *Proceedings of ACM SIGCOMM*, August 2001.
- [148] E. G. Coffman and L. Kleinrock, "Feedback queueing models for time-shared systems," *Journal of the Association for Computing Machinery*, vol. 15, no. 1, pp. 549-576, 1968.
- [149] L. E. Schrage, "The queue M/G/1 with feedback to lower priority queues," *Management Science*, vol. 13, no. 7, pp. 466-474, 1967.
- [150] R. Righter and J. G. Shanthikumar, "Scheduling multiclass single server queueing systems to stochastically maximize the number of successful departures," *Probability in the Engineering and the Informational Sciences*, vol. 3, pp. 323-333, 1989.
- [151] J. G. S. R. Righter and G. Yamazaki, "On extremal service disciplines in single-stage queueing systems," *Journal of Applied Probability*, vol. 27, no. 2, pp. 409-416, June 1990.
- [152] A. Wierman, N. Bansal, and M. Harchol-Balter, "A note on comparing response times in the M/GI/1/FB and M/GI/1/PS queues," *Operations Research Letters*, vol. 32, no. 4, pp. 73-76, 2004.
- [153] M. Nuyens and A. Wierman, "The Foreground-Background queue: a survey," *Performance Evaluation*, vol. 65, pp. 286-307, March 2008.
- [154] S. Aalto, U. Ayesta, and E. Nyberg-Oksanen, "Two-level processor-sharing scheduling disciplines: mean delay analysis," *ACM SIGMETRICS Performance Evaluation Review*, vol. 32, pp. 97-105, June 2004.

- [155] S. Aalto, U. Ayesta, and E. Nyberg-Oksanen, "M/G/1/MLPS compared to M/G/1/PS," *Operations Research Letters*, vol. 33, pp. 519–524, September 2005.
- [156] S. Aalto, "M/G/1/MLPS compared with M/G/1/PS within service time distribution class IMRL," *Mathematical Methods of Operations Research*, vol. 64, pp. 309-325, October 2006.
- [157] S. Aalto and U. Ayesta, "On the nonoptimality of the foreground-background discipline for IMRL service times," *Journal of Applied Probability*, vol. 43, pp. 523-534, 2006.
- [158] S. Aalto and U. Ayesta, "Mean delay analysis of multi-level processor-sharing disciplines," in *Proceedings of IEEE INFOCOM 2006*, pp. 1–11, April 2006.
- [159] H. Feng and V. Misra, "Mixed scheduling disciplines for network flows," *ACM SIGMETRICS Performance Evaluation Review*, vol. 31, pp. 36-39, September 2003.
- [160] K. E. Avrachenkov, U. Ayesta, and P. Brown, "Batch arrival processor-sharing with application to multi-level processor-sharing scheduling," *Queueing Systems*, vol. 50, no. 4, pp. 459–480, August 2005.
- [161] K. E. Avrachenkov *et al.*, "Differentiation between short and long TCP flows: predictability of the response time," in *Proceedings of INFOCOM 2004*, vol. 2, pp. 762– 773, March 2004.
- [162] L. Guo and I. Matta, "Differentiated control of web traffic: a numerical analysis," in *Proceedings of SPIE ITCOM 2002*, 2002.

- [163] M. E. Crovella and A. Bestavros, “Self-similarity in World Wide Web traffic: evidence and possible causes,” in *Proceedings of SIGMETRICS/RICS’96*, pp. 160–169, 1996.
- [164] A. Feldmann and W. Whitt, “Fitting mixtures of exponentials to long-tail distributions to analyze network performance models,” in *Proceedings of IEEE INFOCOM’97*, vol. 3, pp. 1096-1104, April 1997.
- [165] V. Anantharam, “Scheduling strategies and long-range dependence,” *Queueing System*, vol. 33, no. 1-3, pp. 73–89, 1999.
- [166] O. J. Boxma and B. Zwart, “Tails in scheduling,” *ACM SIGMETRICS Performance Evaluation Review*, vol. 34, no. 4, pp. 13–20, March 2007.
- [167] S. Griffin, “Push-To-Talk an overview of the past, present and future,” November 2004. [Online]. Available: <http://www.iee.org/OnComms/PN/communications/Griffin.pdf>.
- [168] “Open mobile alliance.” [Online]. Available: <http://www.openmobilealliance.org/>.
- [169] M. Mouly and M. Pautet, *The GSM System for Mobile Communications*, 1st ed. France: Palaiseau, 1992.
- [170] G. Sanders *et al.*, *GPRS Network*. John Wiley & Sons Ltd, 2003.
- [171] A. Law and D. Kelton, *Simulation Modeling and Analysis (3th ed.)*. McGraw-Hill, Inc., 2000.
- [172] “DQDB Metropolitan Area Network, doc. no. p802.6/d14,” *IEEE 802.6 Proposed Standard*, July 1990.

- [173] Z. L. Budrikis *et al.*, “QPSX: A queued packet and synchronous circuit exchange,” in *Proceedings of ICC 1986*, pp. 288–293, 1986.
- [174] J. L. Hullett and P. Evans, “New proposal extends the reach of metro area nets,” *Data Communications Magazine*, pp. 139–147, February 1988.
- [175] R. M. Newman and J. L. Hullett, “Distributed queueing: A fast and efficient packet access protocol for QPSX,” in *Proceedings of ICC’86*, pp. 294–299, 1986.
- [176] R. M. Newman, Z. L. Budrikis, and J. L. Hullett, “The QPSX man,” *IEEE Communications Magazine*, vol. 26, pp. 20–28, 1988.
- [177] M. Zukerman and P. G. Potter, “The effect of eliminating the STANDBY state on DQDB performance under overload,” *International Journal of Digital and Analog Cabled Systems*, vol. 2, no. 3, pp. 179–186, 1989.
- [178] R. C. Hung, *Mean packet delay of DQDB and two of its enhancement*. Master thesis, Monash University, Australia, 1998.
- [179] E. L. Hahne, A. K. Choudhry, and N. F. Maxemchuk, “DQDB networks with and without bandwidth balancing,” *IEEE Transaction on Communications*, vol. 40, no. 7, pp. 1192–1204, July 1992.
- [180] I. Rubin and Z. Tsai, “Message delay analysis of multiclass priority TDMA, FDMA, and discrete-time queueing systems,” *IEEE Transaction on Information Theory*, vol. 35, no. 3, pp. 36–43, May 1989.
- [181] C. Cicconetti *et al.*, “Performance evaluation of the IEEE 802.16 MAC for QoS support,” *IEEE Transactions on Mobile Computing*, vol. 6, no. 1, pp. 26–38, January 2007.

- [182] M. A. Marsan, P. Laface, and M. Meo, "Packet delay analysis in GPRS systems," in *Proceedings of IEEE INFOCOM 2003*, pp. 970–978, March 2003.
- [183] P. Venkatesh and H. Sirisena, "Modeling for dimensioning and configuring general packet radio service networks," in *Proceedings of IEEE WCNC*, pp. 1510–1515, 2004.
- [184] H. H. Liu, J. L. Wu, and W. C. Hsieh, "Delay analysis of integrated voice and data service for GPRS," *IEEE Communications Letters*, vol. 6, no. 8, pp. 319–321, August 2002.
- [185] Y. Cao, H. R. Sun, and K. S. Trivedi, "Performance analysis of reservation media-access protocol with access and serving queues under bursty traffic in GPRS/EGPRS," *IEEE Transactions on Vehicular Technology*, vol. 52, no. 6, pp. 1627–1641, November 2003.
- [186] B. Baynat and P. Eisenmann, "Toward an Erlang-like law for GPRS/EDGE network engineering," in *Proceedings of IEEE ICC*, no. 1, June 2004.
- [187] K. Premkumar and A. Chockalingam, "Performance analysis of RLC/MAC and LLC Layers in a GPRS protocol stack," *IEEE Transactions on Vehicular Technology*, vol. 53, no. 5, pp. 1531–1546, September 2004.
- [188] G. Brasche and B. Walke, "Concepts, services, and protocols of the new GSM phase 2+ general packet radio service," *IEEE Communications Magazine*, vol. 35, no. 8, pp. 94–104, August 1997.

- [189] “The NS simulator and the documentation.” [Online]. Available: <http://www.isi.edu/nsnam/ns/>.
- [190] K. Wongthavarawat and A. Ganz, “Packet scheduling for QoS support in IEEE 802.16 broadband wireless access systems,” *International Journal of Communication Systems*, vol. 16, pp. 81–96, February 2003.
- [191] C. Cicconetti, *et al.*, “Quality of service support in IEEE 802.16 networks,” *IEEE Network Magazine*, vol. 20, pp. 50–55, March-April, 2006.
- [192] K. Sriram and W. Whitt, “Characterizing superposition arrival processes in packet multiplexers for voice and data,” *IEEE Journal on Selected Areas in Communications*, no. 6, pp. 833–846, September 1986.
- [193] A. Golaup and H. Aghvami, “A multimedia traffic modeling framework for simulation-based performance evaluation studies,” *Computer Networks*, vol. 50, pp. 2071–2087, August 2006.
- [194] Z. Sun *et al.*, “Internet QoS and traffic modelling,” *IEE Proceedings Software, Special Issue on Performance Engineering*, vol. 151, no. 5, pp. 248–255, October 2004.
- [195] M. Molina, P. Castelli, and G. Foodis, “Web traffic modelling exploiting TCP connections’ temporal clustering through HTML-REDUCE,” *IEEE Network Magazine*, vol. 14, no. 3, pp. 46–55, May 2000.

Abbreviations

ATM	Asynchronous transfer mode
BE	Best Effort
BER	Bit Error Rate
BS	Base Station
CAC	Connection Admission Control
CBQ	Class-Based Queue
CDMA	Code Division Multiple Access
CID	Connection ID
CIF-Q	Channel-condition Independent Packet Fair Queueing
CPS	Common part sublayer
CS	Convergence sublayer
CSDPS	Channel state dependent packet scheduling
DHR	Decreasing hazard rate
DMRL	Decreasing mean residual life
DPS	Discriminatory Processor Sharing
DQ	Distributed queue
DQDB	Distributed Queue Dual Bus
EDD	Earliest Due Date
EPS	Egalitarian Processor Sharing
FBPS or FB	Foreground-Background Processor-Sharing
FCFS	First Come First Served
FCH	Frame control header
FQ	Fair Queueing

GGSN	Gateway GPRS support node
GPRS	General Packet Radio Service
GPS	Generalized Processor Sharing
GSM	Global System for Mobile Communications
HOL	Head of line
IHR	Increasing hazard rate
IMRL	Increasing mean residual life
IP	Internet Protocol
IWFQ	Idealized Wireless Fair Queueing
LCFS	Last Come First Served
LQ	Local queue
LRD	Long-Range Dependent
LST	Laplace-Stieltjes transform
MAC	Medium Access Control
MIMO	Multiple-input Multiple-output
MLPS	Multilevel Processor Sharing
MPS	Multiqueue Processor Sharing
MS	Mobile Station
NLOS	Non-line-of-sight
OFDM	Orthogonal Frequency-division Multiplexing
OMA	Open Mobile Alliance
OSI	Open Systems Interconnection
PBSQ	Priority-Based Service Quanta
PDU	Protocol data units
PF	Proportional Fairness
PHY	Physical Layer
PMP	Point-to-multipoint
PoC	PTT over Cellular
PPP	Point-to-Point protocol
PS	Processor Sharing

PTT	Push to Talk
QoS	Quality of Service
RF	Radio Frequency
RLC	Radio link control
RR	Round robin
RRPS	Round robin processor sharing
RTG	Receive/Transmit transition gap
SDU	Service data unit
SIRO	Service-in-random-order
SMS	Short Message Service
SRD	Short-Range Dependent
SS	Subscriber Station
SGSN	Serving GPRS support node
TCP	Transport Control Protocol
TDMA	Time division multiple access
TTG	Transmit/Receive transition gap
VoIP	Voice over IP
WFFQ	Wireless Fluid Fair Queueing
WFQ	Weighted Fair Queue
WiMAX	Worldwide Interoperability for Microwave Access
WLAN	Wireless Local Area Networks
WRR	Weighted Round Robin
xDSL	Digital Subscriber Line

List of Figures

2.1	GPS model	31
2.2	DPS model	36
2.3	WiMAX Network (PMP Mode).	45
3.1	The PBSQ model for the case of two priorities.	62
3.2	Analytical vs simulation results for the mean delay of priority 1 and 2 messages.	75
3.3	Analytical vs simulation results for the mean delay of priority 3 messages.	75
3.4	Comparison of mean message delays in the PBSQ model and in the MPS model.	76
4.1	PDU segmentation into LLC frames, RLC blocks and MAC bursts in GPRS networks. (PH: PDU header; FH: frame header; BH: block header.)	80
4.2	Analytical model used to analysis PTT uplink delay.	84
4.3	Comparison between analytical results and simulation results for PTT packet uplink delay.	88
4.4	Comparison of analytical and simulation results of PTT delay.	91
4.5	PTT packet delay under different GSM voice loads.	93

<i>LIST OF FIGURES</i>	148
4.6 PTT packet delay with different priority assignments.	94
5.1 IEEE 802.16 MAC frame structure in Time Division Duplex mode; FCH: Frame Control Header, RxDS: Receiver Delay Spread clearing interval, pp. 395 and 450 [21].	99
5.2 QoS architecture of IEEE 802.16	101
5.3 WiMAX SS delay analysis using the MPS model.	104
5.4 Comparison between analytical results and simulation results for VoIP packet delay.	110
5.5 nrtPS and BE message delays under different rtPS traffic loads.	110
5.6 WiMAX message delays under different nrtPS traffic loads. . .	111
5.7 WiMAX message delays under different BE traffic loads.	111

List of Tables

3.1	Mathematical notation used throughout this thesis	52
3.2	Mathematical notation used in other references	53
3.3	Priority i traffic load converted into units of priority p segments and segment-times	65
3.4	Corresponding parameter notations used in (6.17) of [180] . . .	67
4.1	Packet size distributions of the source nodes	87
4.2	Parameters for evaluation of the model	87
4.3	Parameters for GPRS traffic	89
4.4	Packet size distributions of GPRS users	91
5.1	Arrival processes and message size distributions of the traffic sources with a priority arrangement	108

Appendix

Two theorems in the MPS model

We prove the solution for L_p with our correction [see Section 3.4.1]

$$\frac{\varepsilon_{p-1} + \rho_p \frac{M_p-1}{M_p}}{1 - \varepsilon_{p-1} - \rho_p \frac{M_p-1}{M_p}}.$$

never be negative in the MPS model, i.e. the mean waiting time of priority p message in the LQ $L_p \geq 0$. Following the notation in Table 3.1, the priority $p = 1, 2, \dots, P$ and 1 is the highest priority. For a priority p , we have

M_p the number of LQs $M_p = 1, 2, 3, \dots$;

\bar{b}_p the mean of packet size b_p , where $b_p = 1, 2, 3, \dots$, so $\bar{b}_p \geq 1$;

\bar{a}_p the mean of arrival rate a_p , where $\bar{a}_p < 1$ because $b_p \geq 1$ and $0 < \rho_p < 1$;

$\text{Var}(a_p) = E(a_p^2) - \bar{a}_p^2$ the variance of a_p ;

$\text{Var}(b_p) = E(b_p^2) - \bar{b}_p^2$ the variance of b_p ;

$C_{b,p}^2$ the squared coefficient of variation of b_p ;

$C_{a,p}^2$ the squared coefficient of variation of a_p ;

$\lambda_p = M_p \bar{a}_p$;

$\rho_p = \lambda_p \bar{b}_p$, and $0 < \rho_p < 1$;

$\varepsilon_p = \sum_{i=1}^p \rho_i$;

$\nu_p = \rho_p \bar{b}_p (C_{b,p}^2 + \lambda_p C_{a,p}^2 / M_p)$.

Theorem 1: $L_p < L_{p+1}$

Firstly, we prove the higher priority has the shorter waiting time than the lower priority, i.e. $L_p < L_{p+1}$, when the parameters of traffic are same for two priorities. The modified L_p with the incongruity correction is:

$$L_p = \frac{\nu_p/\rho_p + \sum_{i=1}^p \nu_i/(1-\varepsilon_p)}{2(1-\varepsilon_{p-1})} - \frac{[\bar{b}_p(1+C_{b,p}^2) + 1]/2 - (\varepsilon_{p-1} + \rho_p \frac{M_p-1}{M_p})}{1 - \varepsilon_{p-1} - \rho_p \frac{M_p-1}{M_p}} + \frac{1}{2}. \quad (1)$$

As priorities p and $p+1$ with same parameters, $\rho_{p+1} = \rho_p$ and $\nu_{p+1} = \nu_p$. Moreover, $\varepsilon_{p+1} - \varepsilon_p = \rho_p$. Comparing L_{p+1} with L_p , we obtain:

$$L_{p+1} - L_p = \frac{\nu_p}{(1-\varepsilon_{p+1})(1-\varepsilon_p)(1-\varepsilon_{p-1})} - \frac{\rho_p[\bar{b}_p(1+C_{b,p}^2) - 1]/2}{(1-\varepsilon_p + \frac{\rho_p}{M_p})(1-\varepsilon_{p+1} + \frac{\rho_p}{M_p})}. \quad (2)$$

For $M_p = 1, 2, \dots$, we have

$$\begin{aligned} L_{p+1} - L_p &> \frac{\nu_p}{(1-\varepsilon_{p+1})(1-\varepsilon_p)(1-\varepsilon_{p-1})} - \frac{\rho_p[\bar{b}_p(1+C_{b,p}^2) - 1]/2}{(1-\varepsilon_p)(1-\varepsilon_{p+1})} \\ &= \frac{\nu_p/\rho_p - (1-\varepsilon_{p-1})[\bar{b}_p(1+C_{b,p}^2) - 1]/2}{(1-\varepsilon_{p+1})(1-\varepsilon_p)(1-\varepsilon_{p-1})}. \end{aligned}$$

When the system is stable, $\sum_{i=1}^P \rho_i = \varepsilon_P < 1$, we just need to prove the top of above equation never being negative. For $0 \leq \varepsilon_{p-1} < 1$, we also have

$$\begin{aligned} &\frac{\nu_p}{\rho_p} - (1-\varepsilon_{p-1})\left[\frac{\bar{b}_p(1+C_{b,p}^2)}{2} - \frac{1}{2}\right] \\ &\geq \bar{b}_p C_{b,p}^2 + \bar{b}_p \bar{a}_p C_{a,p}^2 - \frac{\bar{b}_p}{2} - \frac{\bar{b}_p C_{b,p}^2}{2} + \frac{1}{2} \\ &= \frac{1}{2} \bar{b}_p C_{b,p}^2 + \bar{b}_p (\bar{a}_p C_{a,p}^2 - \frac{1}{2} + \frac{1}{2\bar{b}_p}). \end{aligned}$$

As $\bar{b}_p C_{b,p}^2/2 \geq 0$, we only need to test if:

$$\bar{a}_p C_{a,p}^2 - \frac{1}{2} + \frac{1}{2\bar{b}_p} \geq 0.$$

We have:

$$\begin{aligned} \bar{a}_p C_{a,p}^2 - \frac{1}{2} + \frac{1}{2\bar{b}_p} &= \frac{\mathbb{E}(a_p^2) - \bar{a}_p^2 - \bar{a}_p}{\bar{a}_p} + 1 - \frac{1}{2} + \frac{1}{2\bar{b}_p} \\ &= \frac{\mathbb{E}(a_p^2) - \bar{a}_p}{\bar{a}_p} + \left(1 - \bar{a}_p - \frac{1}{2} + \frac{1}{2\bar{b}_p}\right). \end{aligned}$$

As $\bar{a}_p = \rho/(M_p \bar{b}_p)$, we have:

$$\begin{aligned} 1 - \bar{a}_p - \frac{1}{2} + \frac{1}{2\bar{b}_p} &= 1 - \frac{\rho}{M_p \bar{b}_p} - \frac{1}{2} + \frac{1}{2\bar{b}_p} \\ &= \frac{\bar{b}_p - 2\rho/M_p + 1}{2\bar{b}_p} \\ &= \frac{\bar{b}_p - \rho/M_p + 1 - \rho/M_p}{2\bar{b}_p} > 0. \end{aligned}$$

as $\bar{b}_p \geq 1$ and $\rho/M_p < 1$. Then, we only need to prove if $\mathbb{E}(a_p) - \bar{a}_p \geq 0$. As

$$\begin{aligned} \bar{a}_p &= \sum_{a_p: \mathbb{P}(a_p) > 0} a_p \mathbb{P}(a_p), \\ \mathbb{E}(a_p^2) &= \sum_{a_p: \mathbb{P}(a_p) > 0} a_p^2 \mathbb{P}(a_p), \end{aligned}$$

and a_p is an integer random number of arrivals, $a_p = 1, 2, 3, \dots$, $\mathbb{E}(a_p^2) \geq \bar{a}_p$.

Therefore, we prove

$$\begin{aligned} \bar{b}_p (\bar{a}_p C_{a,p}^2 - \frac{1}{2} + \frac{1}{2\bar{b}_p}) &> 0, \\ \text{and } L_{p+1} - L_p &> 0. \end{aligned} \tag{3}$$

Thus, we have the theorem for all priorities p in the MPS model:

$$L_p < L_{p+1}. \quad (4)$$

Theorem 2: $L_p \geq 0$ for all priority p

Secondly, we prove $L_p \geq 0$ for all priority p based on Theorem 1 obtained before. If for the highest priority 1, $L_1 \geq 0$, for all priorities p , we have $L_p \geq 0$.

Using the modified result in (1), for the highest priority, we have

$$\begin{aligned} L_1 &= \frac{\nu_1/\rho_1 + \nu_1/(1-\rho_1)}{2} - \frac{[\bar{b}_1(1+C_{b,1}^2) + 1]/2 - \frac{M_1-1}{M_1}\rho_1}{1 - \frac{M_1-1}{M_1}\rho_1} + \frac{1}{2} \\ &= \frac{\nu_1/\rho_1}{2(1-\rho_1)} - \frac{[\bar{b}_1(1+C_{b,1}^2) + 1]/2 - \frac{M_1-1}{M_1}\rho_1}{1 - \frac{M_1-1}{M_1}\rho_1} + \frac{1}{2} \\ &= \frac{\frac{\nu_1}{\rho_1}(1 - \frac{M_1-1}{M_1}\rho_1) - [\bar{b}_1(1+C_{b,1}^2) + 1](1-\rho_1) + 2\frac{M_1-1}{M_1}\rho_1(1-\rho_1) + (1-\rho_1)(1 - \frac{M_1-1}{M_1}\rho_1)}{2(1-\rho_1)(1 - \frac{M_1-1}{M_1}\rho_1)} \\ &= \frac{(\bar{b}_1C_{b,1}^2 + \bar{a}_1\bar{b}_1C_{a,1}^2)(1-\rho_1 + \frac{\rho_1}{M_1}) + (2\frac{M_1-1}{M_1}\rho_1 - \bar{b}_1 - \bar{b}_1C_{b,1}^2 - 1 + 1 - \frac{M_1-1}{M_1}\rho_1)(1-\rho_1)}{2(1-\rho_1)(1 - \frac{M_1-1}{M_1}\rho_1)} \\ &= \frac{(\bar{b}_1C_{b,1}^2 + \bar{a}_1\bar{b}_1C_{a,1}^2)(1-\rho_1) + \frac{\rho_1}{M_1}(\bar{b}_1C_{b,1}^2 + \bar{a}_1\bar{b}_1C_{a,1}^2) + (\frac{M_1-1}{M_1}\rho_1 - \bar{b}_1 - \bar{b}_1C_{b,1}^2)(1-\rho_1)}{2(1-\rho_1)(1 - \frac{M_1-1}{M_1}\rho_1)} \\ &= \frac{(1-\rho_1)(\bar{a}_1\bar{b}_1C_{a,1}^2 - \bar{b}_1 + \frac{M_1-1}{M_1}\rho_1) + \frac{\rho_1}{M_1}(\bar{b}_1C_{b,1}^2 + \bar{a}_1\bar{b}_1C_{a,1}^2)}{2(1-\rho_1)(1 - \frac{M_1-1}{M_1}\rho_1)}. \end{aligned}$$

Let us see if the top of above equation is not negative.

$$\begin{aligned}
& (1 - \rho_1)(\bar{a}_1 \bar{b}_1 C_{a,1}^2 - \bar{b}_1 + \frac{M_1 - 1}{M_1} \rho_1) + \frac{\rho_1}{M_1} (\bar{b}_1 C_{b,1}^2 + \bar{a}_1 \bar{b}_1 C_{a,1}^2) \\
= & (1 - \rho_1) \left(\bar{b}_1 \left[\frac{E(a_1^2) - \bar{a}_1^2}{\bar{a}_1} - 1 + \bar{a}_1 (M_1 - 1) \right] \right) + \frac{\rho_1}{M_1} (\bar{b}_1 C_{b,1}^2 + \bar{a}_1 \bar{b}_1 C_{a,1}^2) \\
= & (1 - \rho_1) \left(\bar{b}_1 \left[\frac{E(a_1^2) - \bar{a}_1}{\bar{a}_1} - \bar{a}_1 + \bar{a}_1 (M_1 - 1) \right] \right) + \bar{b}_1 \bar{a}_1 (\bar{b}_1 C_{b,1}^2 + \bar{a}_1 \bar{b}_1 C_{a,1}^2) \\
= & (1 - \rho_1) \left(\bar{b}_1 \left[\frac{E(a_1^2) - \bar{a}_1}{\bar{a}_1} + \bar{a}_1 (M_1 - 1) \right] \right) + \bar{b}_1 \bar{a}_1 (\bar{b}_1 C_{b,1}^2 + \bar{a}_1 \bar{b}_1 C_{a,1}^2 - 1 + \rho_1) \\
= & (1 - \rho_1) \left(\bar{b}_1 \left[\frac{E(a_1^2) - \bar{a}_1}{\bar{a}_1} + \bar{a}_1 (M_1 - 1) \right] \right) + \bar{b}_1 \bar{a}_1 (\bar{b}_1 C_{b,1}^2 + \bar{b}_1 \frac{E(a_1^2) - \bar{a}_1^2}{\bar{a}_1} - 1 + \rho_1) \\
= & (1 - \rho_1) \left(\bar{b}_1 \left[\frac{E(a_1^2) - \bar{a}_1}{\bar{a}_1} + \bar{a}_1 (M_1 - 1) \right] \right) + \bar{b}_1 \bar{a}_1 (\bar{b}_1 C_{b,1}^2 + \bar{b}_1 \frac{E(a_1^2) - \bar{a}_1}{\bar{a}_1} - \bar{a}_1 + \rho_1).
\end{aligned}$$

In above equation, $1 - \rho_1 > 0$, $M_1 \geq 1$ and $E(a_1^2) \geq \bar{a}_1$. Also, as $M_1, \bar{b}_1 \geq 1$, $\rho_1 \geq \bar{a}_1$. Thus, we have

$$(1 - \rho_1)(\bar{a}_1 \bar{b}_1 C_{a,1}^2 - \bar{b}_1 + \frac{M_1 - 1}{M_1} \rho_1) + \frac{\rho_1}{M_1} (\bar{b}_1 C_{b,1}^2 + \bar{a}_1 \bar{b}_1 C_{a,1}^2) \geq 0, \tag{5}$$

and

$$L_1 \geq 0. \tag{6}$$

According to Theorem 1, we have

$$L_p < L_{p+1}.$$

Therefore, we can prove

$$L_p \geq 0, \text{ for all } p. \tag{7}$$