



University of Bradford eThesis

This thesis is hosted in [Bradford Scholars](#) – The University of Bradford Open Access repository. Visit the repository for full metadata or to contact the repository team



© University of Bradford. This work is licenced for reuse under a [Creative Commons Licence](#).

PERFORMANCE MODELING OF CONGESTION CONTROL AND RESOURCE ALLOCATION UNDER HETEROGENEOUS NETWORK TRAFFIC

Modeling and Analysis of Active Queue Management Mechanism in the presence of
Poisson and Bursty Traffic Arrival Processes

Lan Wang

submitted for the degree
of Doctor of Philosophy

Department of Computing

University of Bradford

2010

Abstract

Along with playing an ever-increasing role in the integration of other communication networks and expanding in application diversities, the current Internet suffers from serious overuse and congestion bottlenecks. Efficient congestion control is fundamental to ensure the Internet reliability, satisfy the specified Quality-of-Service (QoS) constraints and achieve desirable performance in response to varying application scenarios. Active Queue Management (AQM) is a promising scheme to support end-to-end Transmission Control Protocol (TCP) congestion control because it enables the sender to react appropriately to the real network situation. Analytical performance models are powerful tools which can be adopted to investigate optimal setting of AQM parameters. Among the existing research efforts in this field, however, there is a current lack of analytical models that can be viewed as a cost-effective performance evaluation tool for AQM in the presence of heterogeneous traffic, generated by various network applications.

This thesis aims to provide a generic and extensible analytical framework for analyzing AQM congestion control for various traffic types, such as non-bursty Poisson and bursty Markov-Modulated Poisson Process (MMPP) traffic. Specifically, the Markov analytical models are developed for AQM congestion control scheme coupled with queue thresholds and then are adopted to derive expressions for important QoS metrics. The main contributions of this thesis are listed as follows:

- ♦ Study the queueing systems for modeling AQM scheme subject to single-class and multiple-classes Poisson traffic, respectively. Analyze the effects of the varying threshold, mean traffic arrival rate, service rate and buffer capacity on the key performance metrics.
- ♦ Propose an analytical model for AQM scheme with single class bursty traffic and investigate how burstiness and correlations affect the performance metrics. The analytical results reveal that high burstiness and correlation can result in significant degradation of AQM performance, such as increased queueing delay and packet loss probability, and reduced throughput and utilization.
- ♦ Develop an analytical model for a single server queueing system with AQM in the presence of heterogeneous traffic and evaluate the aggregate and marginal performance subject to different threshold values, burstiness degree and correlation.
- ♦ Conduct stochastic analysis of a single-server system with single-queue and multiple-queues, respectively, for AQM scheme in the presence of multiple priority traffic classes scheduled by the Priority Resume (PR) policy.
- ♦ Carry out the performance comparison of AQM with PR and First-In First-Out (FIFO) scheme and compare the performance of AQM with single PR priority queue and multiple priority queues, respectively.

Acknowledgements

First of all, I would like to express my deepest gratitude to my research supervisors, Prof. Irfan Awan and Dr. Geyong Min for their extremely valuable guidance, encouragement and discussions throughout the whole period of the research work. I do really appreciate their vast knowledge and skills in many areas, and patient advice in writing up this thesis. I could not achieve this objective without their help and support.

I would also like to express my very special thanks to my parents. It would not have been possible for me to spend this long journey full of joy without their love and constant inspirations.

Finally, I want to express my gratitude to all the friends and colleagues at the University of Bradford who constantly provided emotional support and took care of me in many aspects. I also appreciate many colleagues' valuable comments and discussions.

Publications

Chapter in Book:

- 1 L. Wang, G. Min, and I. Awan, “Effects of Bursty and Correlated Traffic on the Performance of Active Queue Management Schemes”, *Recent Advances in Modeling and Simulation Tools for Communication Networks and Services*, Springer US, September 20, pp. 349-366, 2007.

Journal Publications:

- 1 L. Wang, G. Min, and I. Awan, “Modeling and Analysis of Active Queue Management Schemes under Bursty Traffic,” *INTERNATIONAL JOURNAL OF WIRELESS INFORMATION NETWORKS (IJWIN)*, vol. 13, no. 2, pp. 161-171. 2006.
- 2 L. Wang, G. Min, and I. Awan, “Modelling and Evaluation of Congestion Control Mechanism for Different Classes of Traffic,” *Concurrency and Computation: Practice and Experience (CCPE)*, vol. 19, no. 8, pp. 1141-1156, 2007.
- 3 L. Wang, G. Min, and I. Awan, “Stochastic Modeling and Analysis of an Active Congestion Control Protocol under Differentiated Bursty Traffic”, *Journal of Interconnection Networks (JOIN)*, vol. 8, no. 4, pp. 369-385, 2007.

Conference Contributions:

- 1 L. Wang, G. Min, and I. Awan, "Analytical Modeling and Comparison of AQM-Based Congestion Control Mechanisms," *Proc. Int. Conference on High Performance Computing and Communications (HPCC'2005)*, pp. 67-76, Sorrento, Italy, September 21-23, 2005.
- 2 L. Wang, G. Min, and I. Awan, "Analysis of Active Queue Management under Two Classes of Traffic," *Proc. 21st UK Performance Engineering Workshop (UKPEW'2005)*, pp. 101-109, Newcastle, UK, July 14-15, 2005.
- 3 L. Wang, I. Awan, and G. Min, "A Performance Model for Bursty Traffic: Markov Modulated Poisson Process," *Proc. Postgraduate Research Conference in Electronics, Photonics, Communications & Networks, and Computing Science (PREP'2005)*, pp. 204-206, Lancaster, UK, March 30-April 1, 2005.
- 4 L. Wang, G. Min, and I. Awan, "Analytical Modeling and Comparison of AQM-Based Congestion Control Mechanisms," *Proc. 6th Annual PostGraduate Symposium on the Convergence of Telecommunications, Networking and Broadcasting (PGNet'2005)*, pp. 545-549, Liverpool, UK, June 27-28, 2005.
- 5 L. Wang, I. Awan, and G. Min, "A Performance Model for Bursty Traffic: Markov Modulated Poisson Process," *Proc. Sixth Informatics Workshop For Research Students*, pp. 189-191, Bradford, UK, March 23, 2005.
- 6 L. Wang, G. Min, and I. Awan, "Modelling Active Queue Management with Different Traffic Classes," *Proc. Int. Workshop on Performance Analysis and*

Enhancement of Wireless Networks (PAEWN'06 in conjunction with IEEE AINA'2006), Vienna, Austria, April 18-20, 2006.

- 7 L. Wang, G. Min, and I. Awan, "Stochastic Modeling and Analysis of GRED-I Congestion Control for Differentiated Bursty Traffic," *Proc. 21th Int. Conference on Advanced Information Networking and Applications (AINA'2007)*, IEEE Computer Society Press, Niagara Falls, Canada, May 21-23, pp. 1022 - 1030, 2007
- 8 L. Wang, G. Min, and I. Awan, "Effects of Bursty and Correlated Traffic on the Performance of Active Queue Management Schemes," *COST 285: Modelling and Simulation Tools for Research in Emerging Multi-service Telecommunications*, Surrey, United Kingdom, March 28-29, 2007.
- 9 L. Wang, G. Min, and I. Awan, "An Analytical Model for Priority-Based AQM in the Presence of Heterogeneous Network Traffic," *Proc. 22th Int. Conference on Advanced Information Networking and Applications (AINA'2008)*, IEEE Computer Society Press, GinoWan, Okinawa, Japan, March 25-28, pp. 93-99, 2008.
- 10 L. Wang, G. Min, and I. Awan, "An Analytical Model for Priority-Based AQM in the Presence of Heterogeneous Network Traffic," *Proc. 23th Int. Conference on Advanced Information Networking and Applications (AINA'2009)*, IEEE Computer Society Press, Bradford, U.K., May 26-29, accepted, 2009.

Table of Contents

Abstract	iii
Acknowledgements.....	iv
Publications	v
Table of Contents.....	viii
List of Figures.....	xii
List of Tables.....	xviii
List of Abbreviations	xix
List of Symbols.....	xxi
Chapter 1 Introduction	1
1.1 Background and Motivation.....	1
1.2 Aims and Objectives.....	3
1.3 Original Contributions	4
1.4 Outline of the Thesis.....	6
Chapter 2 Literature Reviews.....	8
2.1 Introduction	8
2.2 TCP Congestion Control.....	8
2.3 Active Queue Management.....	11
2.3.1 Queue-based AQM.....	13
2.3.2 Rate-based AQM.....	16
2.3.3 Combination of buffer and rate based AQM	17

2.4 Traffic Models	17
2.4.1 Poisson Process	18
2.4.2 Markov Modulated Poisson Process	19
2.5 Existing Analytical Models for AQM	22
2.6 Little’s Law	24
2.7 Discrete-Event Simulation	25
Chapter 3 Performance Modeling and Analysis of AQM with Non-Bursty Traffic	27
3.1 Introduction	27
3.2 AQM Scheme with Single Class Traffic	28
3.2.1 Analytical Model	28
3.2.2 Performance Validation and Evaluation.....	32
3.3 AQM Scheme with Multiple Class Traffic	36
3.3.1 Analytical Model	37
3.3.2 Performance Validation.....	44
3.3.3 Performance Evaluation	46
3.4 Summary	49
Chapter 4 Performance Modeling and Analysis of AQM with Single Class Bursty Traffic	52
4.1 Introduction	52
4.2 Analytical Model	53
4.3 Model Validation and Evaluation.....	58
4.3.1 Analytical Model	58

4.3.2 Effects of Burstiness and Correlations	63
4.4 Summary	68
Chapter 5 Performance Modeling and Analysis of AQM with Heterogeneous Traffic	70
5.1 Introductio	70
5.2 Analytical Model and Performance Measures	71
5.2.1 Proposed Markov Model	71
5.2.2 Performance Measures	76
5.3 Model Validation.....	80
5.4 Performance Analysis	84
5.5 Summary	94
Chapter 6 Performance Modeling and Analysis of Priority-Based AQM with Heterogeneous Traffic.....	96
6.1 Introduction	96
6.2 Priority-based AQM with Single Buffer.....	97
6.2.1 System Description	97
6.2.2 Analytical Model	98
6.2.3 Performance Measures	102
6.2.4 Performance Comparison between AQM and Priority-based AQM	105
6.3 Priority-Based AQM with Multiple Queues	116
6.3.1 System Description	116
6.3.2 Analytical Model	117

6.3.3 Performance Measures	122
6.3.4 Performance Comparison between Priority-based AQM with Single Queue and Multiple Queues	124
6.4 Summary	129
Chapter 7 Conclusions and Future Work	131
7.1 Conclusions	131
7.2 Future Work	135
References	137

List of Figures

Figure 2.1. Two-state MMPP model	20
Figure 2.2. Procedure of simulation program.....	26
Figure 3.1. A Model of M/M/1/K/th Queueing System	29
Figure 3.2. The dropping function	29
Figure 3.3. A State Transition Rate Diagram of M/M/1/K/ th Queueing System	30
Figure 3.4. Utilization vs the threshold subject to five different scenarios	33
Figure 3.5. The mean number of packets in the system vs the threshold subject to five different scenarios	34
Figure 3.6. The mean number of packets in the queue vs the threshold subject to five different scenarios	34
Figure 3.7. The throughput vs the threshold subject to five different scenarios	35
Figure 3.8. The mean response time vs the threshold subject to five different scenarios ...	35
Figure 3.9. The mean queueing delay vs the threshold subject to five different scenarios	36
Figure 3.10. The packet loss probability vs the threshold subject to five different scenarios.....	36
Figure 3.11. A Model of [M]2/M/1/K/th1/th2 Queueing System	37
Figure 3.12. A State Transition Rate Diagram of [M]2/M/1/K/ th1/th2 Queueing System	38
Figure 3.13. The Margianl Mean Queue Length vs th2-th1.....	47
Figure 3.14. The Margianl Throughput vs th2-th1	47
Figure 3.15. The Margianl Queueing Delay vs th2-th1	48

Figure 3.16. The Marginal Packet Loss Probability vs th_2-th_1	48
Figure 3.17. The Utilization vs th_2-th_1	48
Figure 3.18. The Aggregate Mean Queue Length vs th_2-th_1	48
Figure 3.19. The Aggregate Throughput vs th_2-th_1	49
Figure 3.20. The Aggregate Mean Queueing Delay vs th_2-th_1	49
Figure 3.21. The Aggregate Packets Loss Probability vs th_2-th_1	49
Figure 4.1. A Model of MMPP-2/M/1/K/th Queueing System	53
Figure 4.2. A State Transition Rate diagram of MMPP-2/M/1/K/th Queueing System	54
Figure 4.3. Utilization vs threshold under 4 different scenarios	60
Figure 4.4. Mean number of packets in the system vs threshold under 4 different scenarios	60
Figure 4.5. Mean number of packets in the buffer vs threshold under 4 different scenarios	61
Figure 4.6. Throughput vs threshold under 4 different scenarios	61
Figure 4.7. Mean response time vs threshold under 4 different scenarios	62
Figure 4.8. Mean queueing delay vs threshold under 4 different scenarios	62
Figure 4.9. Packet loss probability vs threshold under 4 different scenarios	63
Figure 4.10. Utilization vs r_1 with different values of c^2 and th	65
Figure 4.11. Mean number of packets in the system vs r_1 with different values of c^2 and th	65

Figure 4.12. Mean number of packets in the buffer vs r_1 with different values of c^2 and th	66
Figure 4.13. Throughput vs r_1 with different values of c^2 and th	66
Figure 4.14. Mean response time vs r_1 with different values of c^2 and th	67
Figure 4.15. Mean queueing delay vs r_1 with different values of c^2 and th	67
Figure 4.16. Packet loss probability vs r_1 with different values of c^2 and th	68
Figure 5.1. A Model of [MMPP-2]&[M]/M/1/K/th1/th2 Queueing System	72
Figure 5.2. Dropping Functions for Two Classes of Traffic	72
Figure 5.3. A State Transition Rate Diagram of [MMPP-2]&[M]/M/1/K/th Queueing System.....	73
Figure 5.4. Aggregate utilization vs five different scenarios	81
Figure 5.5. Mean number of packets in the system vs five different scenarios	81
Figure 5.6. Mean number of packets in the queue vs five different scenarios	82
Figure 5.7. Throughput vs five different scenarios	82
Figure 5.8. Mean response time vs five different scenarios	83
Figure 5.9. Mean queueing delay vs five different scenarios	83
Figure 5.10. Packet loss probability vs five different scenarios	84
Figure 5.11. Aggregate utilization vs th_1 with different values of c^2 , r_1 and mean arrival rate of Class-1 traffic	87

Figure 5.12. Aggregate throughput vs th_1 with different values of c^2 , r_1 and mean arrival rate of Class-1 traffic 88

Figure 5.13. Aggregate mean queueing delay vs th_1 with different values of c^2 , r_1 and mean arrival rate of Class-1 traffic 88

Figure 5.14. Aggregate packet loss probability vs th_1 with different values of c^2 , r_1 and mean arrival rate of Class-1 traffic 88

Figure 5.15. Marginal throughput of Class-1 vs th_1 with different values of c^2 , r_1 and mean arrival rate of Class-1 traffic 90

Figure 5.16. Marginal mean queueing delay of Class-1 vs th_1 with different values of c^2 , r_1 and mean arrival rate of Class-1 traffic 90

Figure 5.17. .Marginal packet loss probability of Class-1 vs th_1 with different values of c^2 , r_1 and mean arrival rate of Class-1 traffic 91

Figure 5.18. Marginal throughput of Class-2 vs th_1 with different values of c^2 , r_1 and mean arrival rate of Class-1 traffic 92

Figure 5.19. Marginal mean queueing delay of Class-2 vs th_1 with different values of c^2 , r_1 and mean arrival rate of Class-1 traffic 93

Figure 5.20. .Marginal packet loss probability of Class-2 vs th_1 with different values of c^2 , r_1 and mean arrival rate of Class-1 traffic 93

Figure 6.1. A Model of [MMPP-2][M]/M/1/K/ th_1 / th_2 Queueing System with PR scheme 98

Figure 6.2. State transition rate diagram of the three-dimensional Markov chain for AQM scheme with PR scheduling mechanism and a single queue for two-class traffic 99

Figure 6.3. Utilization vs th_1 with various combinations of scheduling schemes and non-bursty or bursty Class-1 traffic107

Figure 6.4. Throughput vs th_1 with various combinations of scheduling schemes and non-bursty or bursty Class-1 traffic107

Figure 6.5. Mean response time vs th_1 with various combinations of scheduling schemes and non-bursty or bursty Class-1 traffic108

Figure 6.6. Mean queueing delay vs th_1 with various combinations of scheduling schemes and non-bursty or bursty Class-1 traffic109

Figure 6.7. Packet loss probability vs th_1 with various combinations of scheduling schemes and non-bursty or bursty Class-1 traffic109

Figure 6.8. Throughput for Class-1 vs th_1 with various combinations of scheduling schemes and non-bursty or bursty Class-1 traffic111

Figure 6.9. Throughput for Class-2 vs th_1 with various combinations of scheduling schemes and non-bursty or bursty Class-1 traffic111

Figure 6.10. Packet loss probability for Class-1 vs th_1 with various combinations of scheduling schemes and non-bursty or bursty Class-1 traffic112

Figure 6.11. Packet loss probability for Class-2 vs th_1 with various combinations of scheduling schemes and non-bursty or bursty Class-1 traffic112

Figure 6.12. Fairness vs th_1 with various combinations of scheduling schemes and non-bursty or bursty Class-1 traffic113

Figure 6.13. Mean response time for Class-1 vs th_1 with various combinations of scheduling schemes and non-bursty or bursty Class-1 traffic114

Figure 6.14. Mean queueing delay for Class-1 vs th_1 with various combinations of scheduling schemes and non-bursty or bursty Class-1 traffic115

Figure 6.15. Mean response time for Class-2 vs th_1 with various combinations of scheduling schemes and non-bursty or bursty Class-1 traffic115

Figure 6.16. Mean queueing delay for Class-2 vs th_1 with various combinations of scheduling schemes and non-bursty or bursty Class-1 traffic116

Figure 6.17. A Model of [MMPP-2][M]/M/1/K1/K2/th1/th2 Queueing System with PR scheme117

Figure 6.18. State transition rate diagram of the three-dimensional Markov chain for AQM scheme with PR scheduling mechanism and multiple classes-based queue for two-class traffic119

List of Tables

Table 3.1. The parameters settings corresponding to five scenarios	32
Table 3.2. The parameter settings corresponding to five scenarios	45
Table 3.3. The analysis results of all performance metrics and the corresponding simulation results subject to five different scenarios	46
Table 4.1. The parameter settings corresponding to four scenarios	59
Table 4.2. The various combinations of c^2 and r_1 as well as the parameters of the corresponding MMPP-2 traffic	64
Table 5.1. The parameter settings corresponding to five scenarios	80
Table 5.2. The various combinations of $\bar{\lambda}^1$, c^2 and r_1 as well as the parameters of the corresponding MMPP-2 traffic	85
Table 6.1. Marginal performance metrics for Class-1 with single queue and multiple queues system corresponding to two scenarios	126
Table 6.2. Marginal performance metrics for Class-2 with single queue and multiple queues system corresponding to two scenarios	128
Table 6.3. Aggregate performance metrics with single queue and multiple queues system corresponding to two scenarios	129

List of Abbreviations

ACK	ACKnowledgment
AQM	Active Queue Management
ARED	Adaptive-RED
AVQ	Adaptive Virtual Queue
DRED	Dynamic-RED
ECN	Explicit Congestion Notification
EWMA	Exponentially-Weighted Moving Average
FCFS	First-Come First-Served
FTP	File-Transfer-Protocol
GRED	Gentle-RED
GRED-I	Gentle-RED with Instantaneous queue length
IDC	index-of-dispersion
IETF	Internet Engineering Task Force
LUBA	Link Utilization Based Approach
MMPP	Markov Modulated Poisson Process
MMPP-2	Two-state Markov Modulated Poisson Process
MQL	Mean Queue Length
PR	Preemptive Resume
QoS	Quality of Service
RED	Random Early Detection

REM	Random Exponential Marking
RFC	The IETF standards documents are called RFCs
RTT	Round Trip Time
SAVQ	Stabilized AVQ
SCV	Squared Coefficient of Variation
SRD	Short Range Dependent
SVB	Stabilized Virtual Buffer
TCP	Transmission Control Protocol
TD	Tail Drop
UDP	User Datagram Protocol
VoIP	Voice-over-IP

List of Symbols

wnd	congestion window
$rwnd$	receiver's advertised window
$ssthresh$	slow start threshold
$SMSS$	size of the largest segment the sender can transmit
L	system capacity
K	buffer capacity
th	system threshold
th_c	system threshold assigned for Class- c traffic
th_b	queue threshold
th_b^c	queue threshold assigned for Class- c traffic
λ	mean arrival rate of traffic
λ_c	mean arrival rate of Class- c traffic
Q	Infinitesimal generator of MMPP-2
Λ	arrival rate matrix of MMPP-2
δ_1	intensity of transition from state S1 to S2 of MMPP-2
δ_2	intensity of transition from State 2 to State 1MMPP-2
c^2	SCV

r_1	1-step autocorrelation coefficient of the inter-arrival times for MMPP-2
μ	mean service rate
μ_c	mean service rate for Class- c traffic
d_i	dropping probability adopted in AQM
d_i^c	dropping probability adopted in AQM for Class- c
p_i	aggregate steady state probability of system
p_{ij}, p_{ijs}	join steady state probability in system
p_i^c	marginal steady state probability for Class- c in system
p_{bi}	aggregate steady state probability in the queue
p_{bij}, p_{bajs}	join steady state probability in the queue
p_{bi}^c	marginal steady state probability for Class- c in the queue
ρ	utilization
\bar{L}	mean number of packets in the system
\bar{L}^c	mean number of Class- c packets in the system
\bar{L}_b	mean number of packets in the queue
\bar{L}_b^c	mean number of Class- c packets in the queue

\bar{T}	aggregate throughput
$\overline{T^c}$	marginal throughput of Class- c
\bar{R}	aggregate mean response time
$\overline{R^c}$	marginal mean response time of Class- c
\bar{D}	aggregate mean queueing delay
$\overline{D^c}$	marginal mean queueing delay of Class- c
\overline{PLP}	aggregate packet loss probability
$\overline{PLP^c}$	marginal packet loss probability of Class- c
F	fairness

Chapter 1

Introduction

1.1 Background and Motivation

As the Internet has been playing an ever-increasing role in the integration of other networks and multi-service applications, such as WWW, File-Transfer-Protocol (FTP), Email, video and audio streaming, Voice-over-IP (VoIP), Multi-player games and e-commerce traffic, which require differentiated Quality-of-Service (QoS) guarantees, congestion control of network traffic and provisioning of differentiated services are now of paramount importance. An effective congestion control scheme not only keeps the volume of network traffic at an acceptable value but also enables different types of Internet applications to be satisfied with the specific QoS requirements.

Transmission Control Protocol (TCP) [1] has been standardized by IETF to use four intertwined congestion control algorithms: slow start, congestion avoidance, fast retransmit and fast recovery reported in Ref. [2-4]. TCP congestion control [5-6] aims to avoid congestion collapse, provide fairness and achieve better performance (i.e., minimizing packet loss and delay, maximizing throughput and utilization). Although necessary and powerful, end-to-end TCP congestion control mechanisms are not sufficient to provide good QoS as well as effectively prevent congestion collapse in some circumstances because the limited control can be accomplished at the network edges. Therefore, some intelligent schemes are demanded in the intermediate nodes (e.g., route and switch) to complement

end-to-end congestion control. Such auxiliary schemes involve the buffer management of per-flow queue which measures and notifies the stages of congestion in router/switch and the scheduling policy which determines the packet sequence at an output port in router/switch. Traffic is controlled through the interaction between end-to-end congestion control and buffer management. Buffer management decides when to start drop packet and which packets to be dropped at congested router output port. End-to-end TCP congestion control adapts the volume of transmitted data to the current load situation by varying the congestion window as a function of packet loss rate.

The traditional approach to buffer management is named Tail Drop (TD). When the output queue is full and TD is in effect, packets are dropped until the congestion is eliminated and the queue is no longer full. Such a method can potentially cause three problems: high queueing delay, “Lock-Out” and “Full Queues” [7]. Active Queue Management (AQM), which starts dropping packets before the queue becomes full in order to notify incipient stages of congestion, has been recommended in the IETF publications [7] for overcoming the drawbacks of Tail Drop. AQM is a promising and widely applied mechanism to control the congestion occurred at a router. Two critical problems for buffer and queue management are when and how to drop packets arriving at a queueing system. In general, the former is mainly based on the queue length and the given threshold. The latter is based on dropping function used to drop packets. Both have significant impact on the average delay, system throughput and probability of packet loss.

Most existing studies [8-14] on the performance of AQM have relied on software simulation and focused on the analysis of Random Early Detection (RED), an AQM

scheme initially proposed and described in [12]. Analytical models for RED have been widely reported, but most existing models are based on the assumptions that the traffic follows the non-bursty Poisson arrival process. In real-world networks, however, traffic exhibits heterogeneous properties and differentiated service mechanisms are adopted to support various QoS requirements. In particular, the properties of burstiness and correlation have attracted many research interests. With the aim of developing cost-effective analytical tools for investigating the performance of congestion control mechanisms in the presence of heterogeneous traffic, this thesis has been dedicated to performance modeling and analysis of AQM in the presence of non-bursty Poisson and bursty Markov-Modulated Poisson Process (MMPP) [15].

1.2 Aims and Objectives

The research work in this thesis is mainly aimed at developing cost effective analytical models for performance evaluation of AQM-based congestion control. All intermediate objectives of this thesis and the steps to achieve the research aims are outlined below:

- ◆ To develop a stochastic queueing system for AQM scheme with single class non-bursty traffic.
- ◆ To further study the performance of AQM scheme in the presence of multiple classes of non-bursty traffic.
- ◆ To investigate how AQM performance is affected by the burstiness and correlation properties of network traffic.

- ♦ To develop an analytical model for performance analysis of AQM scheme under heterogeneous traffic.
- ♦ To develop an analytical Markov model of AQM coupled with Preemptive Resume (PR) priority scheduling scheme subject to multiple priority traffic classes.
- ♦ To compare performance of AQM and priority-based AQM with a single-queue buffer.
- ♦ To compare the performance of AQM coupled with PR priority scheduling scheme based on a single-queue buffer and multiple class-based queues.

1.3 Original Contributions

The original contributions of this thesis are outlined as below:

- ♦ A new performance model is developed for the AQM scheme with a buffer threshold in the presence of a single Poisson arrival process. The closed-form expressions of the steady-state probability are derived and various key performance metrics are obtained. The typical experiments are carried out to demonstrate the credibility of the model against simulation results and to investigate the effects of the traffic loads, service rate, buffer capacity and threshold on the performance metrics. The analytical results give insight into the setting of threshold value in order to satisfy a certain performance trade-off under various combinations of traffic load and buffer capacity. More details are given in Section 3.2.
- ♦ We extend the above analytical model to evaluate the performance of AQM scheme using two thresholds assigned, respectively, for two classes of traffic modeled by two independent Poisson processes. Analytical expressions for various key aggregate and

marginal performance metrics are derived and typical experiments are included to illustrate the accuracy of the proposed model. The effects of a varying threshold on all performance metrics are subsequently investigated, which enables the best threshold setting to be chosen to enable different QoS requirements of each traffic class. More details are given in Section 3.3 and [16-18].

- ♦ A two-dimensional Markov chain is introduced to model the queueing system of AQM scheme under bursty and correlated traffic modeled by a two-state MMPP. Subsequently, we derive closed-form expressions for the steady state probability and all key performance metrics. The validated model is used to discuss how the threshold value, burstiness and correlation properties of traffic affect on the desirable performance metrics and how the effects of one of these three parameters are influenced by the variation of the others. More details are given in Chapter 4 and [19-20].

- ♦ To study the performance of AQM scheme under heterogeneous traffic, we introduce the system theoretical framework in the presence of two classes of traffic that follow a bursty MMPP-2 and non-bursty Poisson process, respectively. The novel two-dimensional Markov chain is proposed to derive analytical closed form expressions for the steady state probability as well as aggregate and marginal performance metrics and to investigate the effects of input parameters of bursty traffic including the load, corresponding threshold, burstiness and correlation on the aggregate and marginal utilization, utilization, mean queueing delay and packet loss probability. More details are given in Chapter 5 and [21-23].

- ♦ A new Markov model is proposed for AQM scheme with a PR priority queue subject to two priority classes traffic and expressions of the important aggregate and marginal performance metrics are derived. The performance of AQM coupled by PR priority scheduling scheme is evaluated and compared with that of AQM with FIFO. More details are given in Section 6.2.
- ♦ A three-dimensional Markov chain is introduced for multiple class-based queues with individual thresholds for each traffic class. Based on the derived and validated expressions for the key aggregate and marginal performance metrics, we evaluate the priority-based AQM performance with multiple class-based queues and point out the advantages and disadvantages of a single PR priority queue and multiple queues, respectively, by comparing this model with the previous. More details are given in Section 6.3 and [24].

1.4 Outline of the Thesis

The rest of the thesis is organized as follows:

Chapter 2 presents a detailed literature review of congestion control technologies and AQM scheme in the Internet including their categories, advantages and disadvantages. Moreover, this chapter gives an overview of traffic models with the emphasis on properties that we will use in this thesis. The review of existing analytical models is then presented.

Chapter 3 studies the AQM schemes with Poisson process to model non-bursty traffic. Firstly, a continuous-time analytical performance model is presented for AQM scheme in the presence of single class traffic and is adopted to evaluate the performance variation due to different mean arrival rate, service rate and buffer capacity. Following this, an extended

Markovian model is proposed for AQM scheme under two classes of traffic. This chapter derives the important aggregate and marginal performance metrics and investigates the effects of varying thresholds on the performance of AQM system.

Chapter 4 proposes a Markovian model for performance evaluation of AQM mechanism using a two-state MMPP to capture the bursty traffic. It also evaluates the effects of the threshold, burstiness and correlation of the MMPP-2 traffic on the derived performance metrics of AQM system.

Chapter 5 introduces a new analytical model for AQM under two independent classes of traffic that follow bursty MMPP and non-bursty Poisson process, respectively. More specifically, this chapter investigates the effects of the average arrival rate, burstiness, correlation and the threshold of bursty traffic on the aggregate and marginal utilization, throughput, mean queueing delay and packet loss probability.

Chapter 6 presents stochastic analysis models of AQM coupled with Pre-emptive Resume (PR) priority scheduling scheme subject to heterogeneous traffic, respectively, for single-queue and multi-queue buffers. It also presents the performance comparison of FIFO and PR priority scheduling schemes, as well as single queueing and multiple queueing systems based on multiple class.

Chapter 7 concludes the thesis and highlights the future work.

Chapter 2

Literature Reviews

2.1 Introduction

This chapter provides an overview of some background knowledge relative to our project. Firstly, in Section 2.2, we review TCP's four intertwined congestion control algorithms in order to be aware of the cooperation of AQM and TCP. Secondly, we introduce in Section 2.3 the AQM scheme which is a promising congestion control scheme to support end-to-end TCP congestion control. The following aspects are covered: i) the basic idea behind AQM; ii) classification of AQM scheme; iii) some representative AQM algorithms. Thirdly, the core features of two traffic models: Poisson and MMPP, to be used in this thesis are highlighted in Section 2.4. Section 2.5 surveys existing research efforts on analytical modelling of AQM. Finally, Little's law and discrete-event simulation are briefly covered in Section 2.6 and 2.7, respectively.

2.2 TCP Congestion Control

In order to achieve high performance and avoid congestion collapse, TCP uses a number of mechanisms to control the rate of data entering the network, keeping the data flow below a rate that would trigger collapse. Modern implementation of TCP contains four intertwined algorithms: Slow Start, Congestion Avoidance, Fast Retransmit and Fast Recovery, which have been introduced in [6].

Due to the lack of considering intermediate network resources (i.e., capacity of routers and slower links), old TCPs could result in a drastic degradation in TCP throughput when two hosts (sender and receiver) are on different LANs [3, 5]. The Slow Start algorithm was first devised in [3] in order to avoid this. Instead of starting a connection with the sender injecting multiple segments into the network up to the receiver's advertised window ($rwnd$), the Slow Start algorithm controls the transmission rate based on the rate of acknowledgement returned by the receiver. Specifically, Slow Start initializes the congestion window ($cwnd$) added to the sender's TCP to one segment when establishing a new connection. Then for each transmission, the sender transmits the minimum of $cwnd$ and $rwnd$, called the transmission window, and increases $cwnd$ by one segment for each returned acknowledgement (ACK). So the window size the sender will often follows approximately an exponential increase.

With the exponential increment of $cwnd$ in Slow Start, there may be a point that one or more packets are dropped due to congestion. Under this situation, the other algorithm called Congestion Avoidance introduced below is used to slow the transmission rate. Slow start and congestion avoidance are independent algorithms, but they are implemented together. Two variables: a slow start threshold $ssthresh$ and $cwnd$, are required to determine the point at which congestion occurs. Such as, Slow Start is used if $cwnd \leq ssthresh$ and Congestion Avoidance is performed otherwise.

Both a retransmission timer expiring and the reception of duplicate ACKs are two indications of packet loss. When congestion occurs, the sender resets $ssthresh$ to be one-half of the current transmission window, which can be presented by $\min(cwnd, rwnd)$.

Specifically, $cwnd$ is reset to one segment if the congestion is indicated by a timeout. Consequently, the sender goes into Slow Start mode again in this case. In Congestion Avoidance, $cwnd$ is incremented by 1 segment per Round Trip Time (RTT), which is implemented by increasing $cwnd$ by $SMSS \times SMSS / cwnd$ each time an ACK is received, where $SMSS$ is the size of the largest segment the sender can transmit.

The basic idea of the Fast Retransmit algorithm is to use three duplicate acknowledgements (ACKs) received, instead of the retransmission timer, as an indication of segment loss to speed up the retransmission process. A lost or delayed segment enables the receiver to send a duplicate ACK as the receiver acknowledges the last continuous byte received prior to the lost or delayed segment. As for a delayed segment resulting in an out of order environment, the receiver can re-order segments before sending the latest ACK. Typically no more than two duplicate ACKs can be generated only under out of order conditions. Therefore, the appearance of three duplicate ACKs strongly indicates that at least one segment has been lost. On receiving three duplicate ACKs, the sender retransmits the lost segment (indicated by the position of the duplicate ACK in the byte stream) without waiting for the transmission timer to expire.

Following Fast Transmit, the Fast Recovery algorithm is applied to govern the transmission of new segments until a non-duplicate ACK arrives by entering Congestion Avoidance mode. The reason for not performing Slow Start algorithm is that duplicate ACKs generated after receiving a segment indicate not only segment loss but also unserious network congestion. So the sender resumes a larger transmission window instead of

reducing the flow of data abruptly (say, Slow Start). The Fast Recovery algorithm is capable of providing a higher throughput under moderate congestion conditions [5].

Specifically, the implementation of the Fast Retransmit and Fast Recovery algorithms are introduced in [6] as follows: 1) the sender is to set *ssthresh* to no more than the value $\max(\min(cwnd, rwnd)/2, 2 \times SMSS)$ when receiving the third duplicate ACK; 2) to retransmit the lost segment and set $cwnd = ssthresh + 3 \times SMSS$; 3) to increment *cwnd* by *SMSS* for each additional duplicate ACK received and to transmit a segment with a window size $\min(cwnd, rwnd)$; 4) Finally, once receiving an ACK for new data, the sender will set $cwnd = ssthresh$.

2.3 Active Queue Management

Various queue management mechanisms have been proposed to control traffic congestion and support differentiated QoS requirements. The traditional approach to queue management is to set a maximum limit on the amount of data that can be buffered. For example, in the Tail Drop (TD) approach [7], the activity to packet dropping will not start until the queue space is exhausted. When the queue becomes full, TD is in effect immediately and all forthcoming packets will be dropped until the congestion is eliminated and some space becomes available in the queue. TD is still a useful mechanism in IP routers because of its robustness and simple implementation. Unfortunately, TD often causes high packet delays, bursty packet drops, degrading system stability and bandwidth fairness in the presence of persistent congestion [8, 12, 25]. “Lock-Out” and “Full Queues” [7] are the main drawbacks of TD due to dropping packets only when the congestion has

occurred. The “Full Queues” phenomenon causes global synchronization of flows and consequently high packet delays, low link utilization as well as low overall throughput because the TD queues are always full or close to full for long periods of time so that an arriving burst will cause multiple packets to be dropped [7, 26]. Therefore, TD is not suitable for interactive network applications due to their low end-to-end delay and jitter requirements. On the other hand, as the result of synchronization, “Lock-out” phenomenon allows a single connection or a few flows to monopolize the queue space in some situations. The other two alternative queue disciplines, “Random drop on full” and “Drop front on full” can solve the “Lock-Out” problem but cannot overcome the “Full Queues” problem [7].

To deal with both problems and to provide low end-to-end delay along with high throughput, a widespread deployment of Active Queue Management (AQM) in routers has been recommended in the IETF publications [7]. To avoid the case that the buffer maintains a full status for a long time, AQM scheme starts dropping or marking packets before the queue is full in order to notify incipient stages of congestion. On receiving the congestion signal (i.e., Explicit Congestion Notification (ECN) marking or dropped packets) generated by AQM scheme, traffic sources reduce the amount of traffic injected into networks. By keeping the mean queue length small, AQM decreases the mean queueing delay and reduces the number of dropped packets, thus resulting in increased link utilisation because of absorbing more packet bursts and avoiding global synchronization. Two key issues in an AQM mechanism are when and how to drop/mark packets, which determines the incipient stage of congestion and reflects the degrees of congestion, respectively. Both have

significant effects on the key performance metrics including the utilization, throughput, mean queueing delay and packet loss probability. Some representative AQM mechanisms are overviewed below, including Random Early Detection (RED) [12], Gentle-RED (GRED) [11, 13], BLUE [27] and Random Exponential Marking (REM) [28].

2.3.1 Queue-based AQM

Random Early Detection (RED) initially proposed in [12] is the most well-known AQM algorithm. The deployment of RED was recommended by Internet Engineering Task Force (IETF) in [7], and indeed RED is widely implemented in Internet routers nowadays. RED monitors the Exponentially-Weighted Moving Average (EWMA) of the queue length and randomly drops the forthcoming packets with a dropping probability that linearly increases with the average queue length between two thresholds. In RED the exponentially weighted moving average $avg = (1 - \omega) \times avg + \omega \times q$ is used to compare with two thresholds: \min_{th} and \max_{th} . when $avg < \min_{th}$, no packets are dropped. When $\min_{th} \leq avg < \max_{th}$, the dropping probability p_b increases linearly from 0 to \max_p . Once $avg \geq \max_{th}$, p_b reaches the maximum dropping probability \max_p and the router drops all arriving packets. As the first AQM mechanism, RED was designed to address the problems exhibited in the TD policy. Uniformly random early packet marks/drops spread over time, which reduces global synchronization of traffic sources. Moreover, with RED, the average queue length is reduced to accommodate traffic bursts. Consequently, RED mechanism outperforms TD policy in terms of throughput, queueing delay and fairness [29]. The study [30] on a typical simulation scenario where four FTP sources were considered has also shown that RED

outperforms TD. Furthermore, the combination of RED and Explicit Congestion Notification (ECN) [31] can improve TCP performance over wireless networks [32-33].

Along with the in-depth studies and more scenarios considered, however, some serious drawbacks of RED have been discovered. Floyd has discussed how to set parameters for RED in 1993 [12] and in 1997 [10], respectively, and reported the complexity of the parameter settings. Authors in [34-37] have demonstrated that a specific parameter setting of RED is applicable only under a narrow range of network conditions. Some researchers [38-40] have evaluated how to adjust one of the key RED parameters (i.e., maximum drop probability) to provide good performance under dynamic network. An iteration algorithm has been proposed in [41] to optimize two RED thresholds, but the algorithm is infeasible due to its long run-time before achieving optimal values. It is hard to choose a set of parameter values to balance the trade-off between various performance measures with different scenarios. Moreover, in [42] Low *et al.* have performed a control-theoretic analysis of TCP/RED and discovered that this flow control mechanism eventually becomes unstable as RTT delay increases, or when the network capacity increases. Mikkel *et al* [9] studied the effects of RED on the performance of Web traffic using HTTP response time, a user-centric measure of performance, and found that RED cannot provide a fast response for end-users. Three key problems associated with AQM scheme: parameter setting, the insensitivity to the input traffic load variation and the mismatch between macroscopic and microscopic behaviour of queue length dynamics were surveyed in [14]. Based on the analysis of extensive experiments of aggregated traffic containing various categories of flows, Martin *et al* [11, 13] demonstrated the harm of RED due to the use of the average

queue length, especially when the average value is far away from the instantaneous queue length. The interaction between the average queue length and the sharp edge in the dropping function results in some pathology such as the increase in the drop probability of the UDP flows and the number of consecutive losses.

To overcome this drawback, Folyd and Fall [11] have suggested to use a smoothly dropping function even when the average queue length exceeds the maximum threshold but not the sharp edge and named this improved algorithm as Gentle-RED (GRED). In the GRED, the packet dropping probability p_b varies from \max_p to 1 as the avg increases from \max_{th} to $(2 \times \max_p)$. Brandauer, et al. [8] have further enhanced GRED by considering the instantaneous queue length (namely GRED-I) instead of the average queue length and varying the dropping probability smoothly from 0 to 1 between the minimum and maximum thresholds. May et al. have pointed out in [13] that it is unnecessary to choose \max_{th} smaller than the buffer capacity and $\max_p < 1$ if the instantaneous queue length is adopted. That is, it can also be viewed as GRED-I if \max_{th} equals to the buffer capacity and $\max_p = 1$. The surprising results reported in [8, 11, 13] have shown that GRED-I performs better than RED and GRED in terms of aggregate throughput, UDP loss probability, response time and the number of consecutive losses. Compared to RED, GRED appears less advantageous than GRED-I because, for RED, the averaging strategy causes more negative effects than the the sharp edge in the dropping function.

RED, GRED and GRED-I belong to queue-based AQM which aim at stabilizing the queue length. Such mechanisms measure the buffer occupancy in the buffer to determine

the severity of congestion. Extensive queue-based AQM schemes [43-56] have been proposed in the literature, such as Stabilized RED (SRED) [45], Adaptive RED (ARED) [48], Exponential-RED (E-RED) [56] etc.

2.3.2 Rate-based AQM

The second classification of AQM, named rate-based, deploys packet arrival rate to manage congestion. The goals of Rate-based AQM schemes (i.e., BLUE [27], Adaptive Virtual Queue (AVQ) [57], Stabilized AVQ (SAVQ) [58] and Link Utilization Based Approach (LUBA) [59]) are to alleviate rate mismatch between enqueue and dequeue, and thus achieve low packet loss, delay and high link utilization. As the typical representative of rate-based AQM schemes, BLUE proposed in [27] was the first attempt to couple packet loss as well as link utilization, rather than queue length, and congestion notification. BLUE can effectively send back congestion notification at the correct rate in that a probability used to mark/drop packets when they are enqueued is adjusted according to packet loss and link idle events. That is, BLUE increments the probability if arriving packets are continually dropped due to buffer overflow. Conversely, the probability is decreased if the queue becomes empty or the link is idle. Wu-chang Feng et. al. [27] have shown that BLUE would cause better performance than RED in terms of packet loss rates and buffer size requirements. On the other hand, some research efforts [60-61] have also reported that BLUE performs poorly when the buffer size is small due to no control in the average queue length and significant loss can occur when bursty traffic arrives.

2.3.3 Combination of buffer and rate based AQM

In order to achieve a tradeoff between queues stability and responsiveness, the metrics used to in some AQM schemes (i.e., Random Exponential Marking (REM) [28], Stabilized Virtual Buffer (SVB) [62] and RaQ [63]) is based on the combination of buffer and rate. For instance, REM [28] maintains a variable called “price”, which is updated based on rate mismatch (i.e., difference between input rate and link capacity) and queue mismatch (i.e., difference between queue length and target), as congestion measure. The “price” is incremented if the weighted sum of these mismatches is positive, and decremented otherwise. The weighted sum is positive when either the input rate exceeds the link capacity or there is excess backlog to be cleared, and negative otherwise. The marking probability varies exponentially with the variable “price”.

2.4 Traffic Models

Performance modeling plays a key role in analyzing AQM performance and deciding the type of AQM scheme to be implemented. Performance models in turn, require very accurate traffic models that have the ability to capture the statistical characteristics of the actual traffic on the network. Innumerable traffic models have been proposed for understanding and analyzing various statistical characteristics of traffic in real networks. This section presents the continuous-time arrival processes which have been used to model the traffic source in this thesis.

2.4.1 Poisson Process

One of the most widely used and oldest traffic models is the Poisson model which has been predominantly used for analyzing traffic in traditional telephony networks [64]. A process is referred to as Poisson process if the inter-arrival times T_n are exponentially distributed with a mean arrival rate λ and distribution function:

$$F(x) = \begin{cases} 1 - \exp(-\lambda x) & x \geq 0 \\ 0 & x < 0 \end{cases}$$

The Poisson model is attractive and popular in queueing theory due to its interesting and tractable analytical properties [65]. Primarily, the independent increment property renders a Poisson process memoryless which means that the subsequent arrival is completely independent from the previous arrivals. In addition, superposition of multiple Poisson processes with λ_i generates a new Poisson process with $\lambda = \sum \lambda_i$. Finally, a relatively simple equation $p_n = (\lambda/\mu)^n p_0$, where μ is the mean service rate and p_0 is the probability of an empty system, for steady state probability can be derived in queueing models with Poisson arrivals and exponential service.

With the in-depth study on the Internet traffic characteristics, many research efforts [66-68] have elaborated that the dominant characteristic exhibited by Internet traffic is multifaceted bursty and correlated structure. However, it has been verified that Poisson process fails to capture such Internet traffic properties [69].

2.4.2 Markov Modulated Poisson Process

The introduction of MMPP allowed the modeling of time-varying sources while keeping the analytical solution of related queueing performance tractable [15]. It has been identified in [70] that the MMPP model can be used for analyzing a mixture of voice and data traffic.

An MMPP is identified as a special case of a Markov Modulated Process (MMP) using the Poisson process as the auxiliary Markov process in which the current state of the Markov process controls the probability distribution of the traffic arrivals. In other words, MMPP is a doubly stochastic Poisson process where the arrival process is determined by an irreducible continuous-time underlying Markov chain consisting of m different states. The MMPP is generally parameterized by the infinitesimal generator \mathbf{Q} of the m -state continuous-time Markov chain and the arrival rates matrix $\mathbf{\Lambda}$.

$$\mathbf{Q} = \begin{bmatrix} -\delta_1 & \delta_{12} & \cdots & \delta_{1m} \\ \delta_{21} & -\delta_2 & \cdots & \delta_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ \delta_{m1} & \delta_{m2} & \cdots & -\delta_m \end{bmatrix} \quad (2.1)$$

$$\delta_i = \sum_{j=1, j \neq i}^m \delta_{ij} \quad (2.2)$$

$$\mathbf{\Lambda} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_m) \quad (2.3)$$

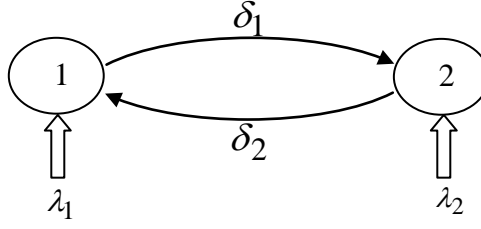


Figure 2.1. Two-state MMPP model

MMPPs are classified by the number of states, m , present in the underlying Markov Chain. The two-state MMPP (MMPP-2) as shown in Figure 2.1 has been widely used in numerous studies to model video sources and the superposition of voice sources due to its tractability [71-72]. Two states i ($i=1,2$) of MMPP-2 correspond to two different traffic arrival processes with mean rate λ_1 and λ_2 , respectively. Both δ_1 and δ_2 are the intensities of transition between State 1 to State 2. They are independent of the arrival process. That is, the duration of state i ($i=0,1$) is in accordance with an exponential distribution with mean $1/\delta_i$. The mean arrival rate $\overline{\lambda_{mmpp-2}}$ of an MMPP-2 is $(\lambda_1\delta_2 + \lambda_2\delta_1)/(\delta_1 + \delta_2)$. The infinitesimal generator \mathbf{Q} and the arrival rate matrix $\mathbf{\Lambda}$ of MMPP-2 are given as follows.

$$\mathbf{Q} = \begin{bmatrix} -\delta_1 & \delta_1 \\ \delta_2 & -\delta_2 \end{bmatrix} \quad (2.4)$$

$$\mathbf{\Lambda} = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix} \quad (2.5)$$

The steady state probability, p_i ($i=1,2$), that the MMPP-2 is in State i can be obtained from the balance equations for the underlying two-state Markov chain and is given by

$$p_i = \frac{\delta_{(3-i)}}{(\delta_1 + \delta_2)} \quad (2.6)$$

With the aim of evaluating the variation of AQM performance under different levels of burstiness and correlation of traffic modelled by MMPP-2, we present how the traffic burstiness and correlations can be expressed by the parameters of MMPP-2.

♦ ***Burstiness***

The simple but crude measure of burstiness is the ratio of the peak rate to mean rate which has the shortcoming of dependence on the interval length utilized for rate measurement. Therefore, the existing study [71] has used the index-of-dispersion for counts (IDC) [73] which is defined as the variance-to-mean ratio of the number of arrivals in an interval $[0, \tau]$ for burstiness measurement. However, IDC measure includes the autocorrelations of inter-arrival times since the number of arrivals (i.e., the numerator of the IDC) is related to the sum of inter-arrival intervals. Additionally, Squared Coefficient of Variation (SCV), c^2 , of the inter-arrival times is a more elaborate and popular measure, which can be easily obtained if the first two moments of the inter-arrival time are known. Fischer and Meier-Hellstern [15] have presented the moments of the time between arrivals in an MMPP using matrix expressions. In particular, an explicit expression of SCV for MMPP-2 has been presented in [74] as:

$$c^2 = 1 + \frac{2\delta_1\delta_2(\lambda_1 - \lambda_2)^2}{(\delta_1 + \delta_2)^2(\lambda_1\lambda_2 + \lambda_2\delta_1 + \lambda_1\delta_2)} \quad (2.7)$$

♦ ***Correlation***

In contrast, autocorrelation coefficient of the inter-arrival times is the most important parameter to measure the traffic correlations. According to the autocorrelation of a random process defined in [75], a k -step autocorrelation matrix of the inter-arrival times for an MMPP has been shown in [15]. Based on this, Cui and Nilsson [74] have provided a simplified expression (see Eq. 2.8) of the l -step autocorrelation coefficient of the inter-arrival times for MMPP-2 model.

$$r_1 = \frac{\delta_1 \delta_2 \lambda_1 \lambda_2 (\lambda_1 - \lambda_2)^2}{c^2 (\delta_1 + \delta_2)^2 (\lambda_1 \lambda_2 + \lambda_2 \delta_1 + \lambda_1 \delta_2)^2} \quad (2.8)$$

2.5 Existing Analytical Models for AQM

Though various AQM schemes have been proposed, queue-based AQM is still the most important and practical mechanism to support end-to-end congestion control in the Internet due to the consideration with the trade-off between implementation feasibility and performance effectiveness. Most existing studies on the evaluation of queue-based AQM performance are based on simulation experiments. However, analytical models are a more attractive alternative to simulation for evaluating system performance under different design spaces because of their cost-effective and time-saving properties. For instance, Kuusela and Virtamo [76] have analyzed a queuing system with RED subject to two classes of non-bursty Poisson traffic and have derived the mean queue lengths of two traffic classes. Kuusela, et al. [77] have analyzed the dynamic behavior of a single RED controlled queue interacting with a large population of idealized TCP sources. Bonald, May, and Bolot [78, 30] have developed an analytical model of RED under two classes of traffic modeled

by Poisson process and batch Poisson process, respectively. The model is used to quantify the impact of RED on the packet loss, the number of consecutively lost packets, delay and delay jitter. They have found that RED does indeed eliminate the bias against bursty traffic observed with TD and have demonstrated that RED achieves the aggregate mean queueing delay at cost of delay jitter of non bursty traffic modelled by Poisson process. Alazemi et al. [79-80] have presented a discrete time stochastic model for evaluating the performance of the RED algorithm which incorporates a two-dimensional second-order discrete-time Markov chain that captures the feedback effect of packets dropping/marking on the incoming traffic and they have derived mean system occupancy, packet drop probability, and system throughput for both marking and dropping policies. An analytical framework has been developed in [81] to evaluate the performance of an AQM router with traffic modelled by a Markov arrival process. Barbera et al. [81] have proposed a new fluid-flow-based methodology and validate the model with the consideration of a RED router subject to a constant bit-rate (CBR) traffic. A discrete-time queueing analytic model has been introduced in [82] for a new proposed dynamic random early drop (DRED) congestion control mechanism and has been used to evaluate the performance metrics including packet loss probability, average queue length, throughput and average queueing delay. Then authors have proposed a new discrete-time analytical model for two-queue nodes queueing network based on DRED algorithm in [83].

Most current existing analytical models are mainly devised for RED mechanism. As aforementioned, the combination of average queue length and sharp dropping edge is origin of drawbacks of RED. So, it is very important and necessary to use an analytical approach

to evaluate the performance of AQM schemes, such as GRED, without such weaknesses and systematically address the fundamental aspects of queue-based AQM scheme. These issues motivate our research project.

2.6 Little's Law

Little's Law provides a fundamental relationship among the average number of items in a system, the average time spent in that system (for an item) and the average arrival rate of items to the system. It states that, under steady state conditions, the average number of items equals the average arrival rate multiplied by the average time that an item spends in the system. The proof of Little's Law, first derived formally by J. D. C. Little [84], is free from assumptions of the distributions for the arrival and service process, number of servers in the system and the queueing discipline.

It is important to note that the boundary of the system in Little's Law is undefined. Therefore, the system can be an entire queueing system composed of queue and server, or it can contain only the queue. For these two systems, we can find the following equations, respectively.

$$\bar{L} = \bar{T} \bar{R} \tag{2.1}$$

$$\bar{L}_q = \bar{T} \bar{D} \tag{2.2}$$

where \bar{L} and \bar{L}_q represent the average number of packets in the system (including queue and server) and queue, respectively; \bar{R} and \bar{D} represent the average response time and queueing delay; and \bar{T} indicates the throughput of the system. Little's Law is extremely

useful to calculate an unknown parameter, which is difficult to be measured directly, based on the knowledge of the other two.

2.7 Discrete-Event Simulation

Discrete-event simulation [85] is a powerful computing technique for understanding the behavior of systems. In discrete-event simulation, the operation of a system is represented as a chronological sequence of events. The simulations developed in this thesis are programmed in JAVA and aim to validate corresponding analytical models.

Figure 2.2 shows the procedure to assist in the basic understanding of discrete-event simulation. First of all, the INITIAL function initializes all state variables which indicate the state of the system after the last simulated event occurrence. The WHILE loop enables the simulation program to run repeatedly until it becomes stable. Many particular conditions can be used for the generic “simulation not over”. Specifically, simulation clock records the time of the last event occurrence simulated is greater than 1000000 seconds. TIMING and UPDATING functions renew next events, next event time, simulation clock and the variables related to performance measures, respectively. According to the next event variable, simulation schedules two events: a packet arrival and departure. Finally, performance results are to be reported before ending simulation.

```
Begin
  INITIAL ()
  WHILE "simulation not over"
    Begin
      TIMING ()
      UPDATE ()
      IF "next event is arrival"
        ARRIVAL ()
      IF "next event is departure"
        DEPART ()
    End
  REPORT ()
End
```

Figure 2.2. Procedure of simulation program.

Chapter 3

Performance Modeling and Analysis of AQM with Non-Bursty Traffic

3.1 Introduction

As an effective and promising scheme to support traffic congestion control, AQM starts dropping packets according to a certain dropping probability once the threshold of queue length is reached in order to notify incipient stages of congestion. Therefore, both the threshold value and dropping function have great impact on the performance of AQM, which has motivated the study on stochastic analysis of queueing systems for the performance evaluation of AQM schemes. This research starts with performance modelling and analysis of AQM congestion control scheme in the presence of non-bursty Poisson arrival process.

This chapter presents two Markov models developed for AQM schemes in the presence of single class and multiple classes of non-bursty traffic in Sections 3.2 and 3.3, respectively. The accuracy of the models is verified by comparing the analytical results against those obtained from discrete-event simulator developed in JAVA programming. We derive the expressions of the important performance metrics including utilization, throughput, mean number of packets in the system, number of packets in the buffer, response time, queueing delay and packet loss probability. The first model is used to briefly evaluate the performance variations due to different mean arrival rate, service rate and

buffer capacity, whilst the second one is used to investigate the effects of varying thresholds on all aggregate and marginal performance metrics.

3.2 AQM Scheme with Single Class Traffic

In this section, we present a general analytical model for queueing system which incorporates a threshold in order to drop the arriving packets prior to periods of high congestion. Specifically, an arriving packet may be dropped randomly according to a dropping probability when the number of packets in the system reaches the threshold. Extensive experiments are carried out to validate the accuracy of the model in calculating the aforementioned various performance metrics and evaluate the performance of AQM scheme.

3.2.1 Analytical Model

Figure 3.1 depicts the queueing system of buffer in an AQM router subject to single class non-bursty traffic. Different from a single-class single-server finite queueing system $M/M/1/K$, this studied queueing system, denoted by $M/M/1/K/th$, assigns a threshold th_b in the buffer to drop an arriving packet and consequently control the injection rate of packets. The scheduling discipline is First-In First-Out (FIFO). An arriving packet can be absorbed definitely if the queue length is smaller than the threshold value. Otherwise, the system is capable of rejecting an arriving packet with a current dropping probability which is calculated using on a linear dropping function as shown in Figure 3.2.

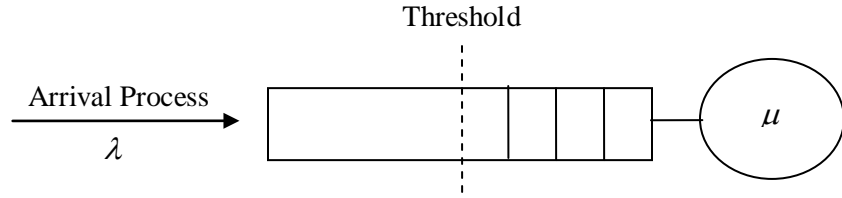


Figure 3.1. A Model of M/M/1/K/th Queueing System

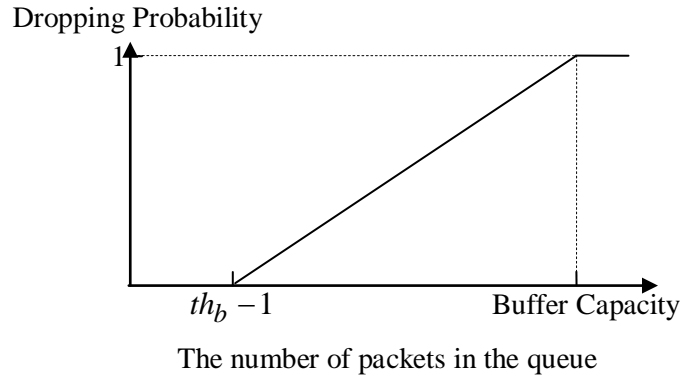


Figure 3.2. The dropping function

The proposed analytical model for the M/M/1/K/th queueing system is shown in Figure 3.3. The arrival and departure processes are modeled using two Poisson processes with mean rates λ and μ , respectively. The system capacity is L including a server and a buffer with capacity K , i.e., $L = K + 1$. The state i ($0 \leq i \leq L$) indicates the number of packets in the system. It is obvious $th = th_b + 1$. The reduction of the packet arrival rate when the system is at state i ($th \leq i \leq L$) can be viewed as a result of packets being dropped according to a certain probability. So the relationship between the reduction rate r_i and the dropping probability d_i ($0 \leq i \leq L$) is $r_i = 1 - d_i$. The calculation of d_i ($0 \leq i \leq L$) is given as follows:

$$d_i = \begin{cases} 0 & 0 \leq i < th \\ \frac{i - th + 1}{L - th + 1} & th \leq i \leq L \end{cases} \quad (3.1)$$

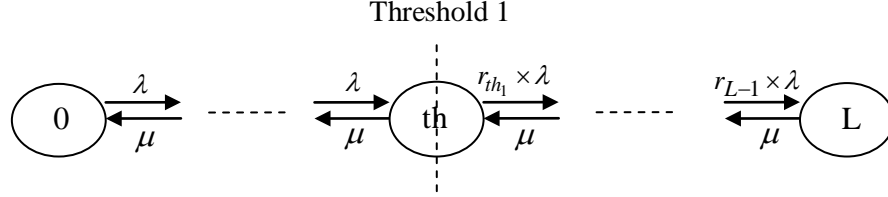


Figure 3.3. A State Transition Rate Diagram of M/M/1/K/ th Queueing System

As the fundamental of deriving the performance metrics, the probability, p_i ($0 \leq i \leq L$), that the system is at state i should be solved firstly according to the transition equilibrium between in-coming and out-going streams of each state and the normalising equation, as shown below.

$$\left. \begin{aligned} r_0 \lambda p_0 &= \mu p_1 \\ (r_i \lambda + \mu) p_i &= r_{i-1} \lambda p_{i-1} + \mu p_{i+1} \quad 1 \leq i < L \\ \mu p_L &= r_{L-1} \lambda p_{L-1} \end{aligned} \right\} \quad (3.2)$$

$$\sum_{i=0}^L p_i = 1 \quad (3.3)$$

Solving these equations, the probability p_i can be expressed as:

$$p_i = \begin{cases} \frac{1}{1 + \sum_{j=1}^L \left(\prod_{k=0}^{j-1} \frac{r_k \lambda}{\mu} \right)} & i = 0 \\ \left(\prod_{k=0}^{i-1} \frac{r_k \lambda}{\mu} \right) \times p_0 & 1 \leq i \leq L \end{cases} \quad (3.4)$$

In what follows, we will derive the performance metrics including utilization (ρ), the mean number of packets in the system (\bar{L}), mean number of packets in the buffer (\bar{L}_b),

system throughput (\bar{T}), mean response time (\bar{R}), mean queueing delay (\bar{D}) and packet loss probability (\overline{PLP}).

The server is engaged as long as the number of packets in the system is not zero. Therefore, the system utilization can be written as

$$\rho = 1 - p_0 \quad (3.5)$$

The mean number of packets in the system and queue can be calculated, respectively, as follows

$$\bar{L} = \sum_{i=0}^L (p_i \times i) \quad (3.6)$$

$$\bar{L}_b = \sum_{i=0}^{L-1} (p_{i+1} \times i) \quad (3.7)$$

Throughput is commonly defined as the average rate at which packets go through the system in the steady state. Therefore, the throughput is equal to the system service rate multiplied by utilization.

$$\bar{T} = \rho \times \mu \quad (3.8)$$

Then the mean response time and delay in the queue can be solved by Little's Law [1].

$$\bar{R} = \frac{\bar{L}}{\bar{T}} \quad (3.9)$$

$$\bar{D} = \frac{\bar{L}_b}{\bar{T}} \quad (3.10)$$

The packet loss probability consists of the probability of packet loss after the queue becomes full and that of packet dropping due to AQM scheme before the queue is full.

$$PLP = \sum_{i=0}^L p_i \times d_i \quad (3.11)$$

3.2.2 Performance Validation and Evaluation

This section presents the accuracy of the proposed model in calculating the performance metrics derived in subsection 3.2.1 and briefly analyses the effects of the variation of related parameters (i.e., threshold, buffer capacity, mean arrival rate and mean service rate) on the aforementioned performance metrics. Performance results depicted in Figures 3.4-3.10 are presented for five different scenarios with the system parameters described in Table 3.1. Figures 3.4-3.10 demonstrate the system utilization, mean number of packets in the system, mean number of packets in the buffer, throughput, mean response time, mean queueing delay and packet loss probability, respectively. The excellent match between the analysis results and the corresponding simulation results obtained from a simulator programmed in JAVA indicates the credibility of the analytical model in evaluating the performance of AQM under consideration.

	S-3.2.I	S-3.2.II	S-3.2.III	S-3.2.IV	S-3.2.V
λ	3	6	6	6	8
μ	10	10	10	7	9
K	5	5	12	12	20

Table 3.1. The parameters settings corresponding to five scenarios.

In what follows, we analyze how the variation of each parameter affects all the performance metrics. It can be viewed from Table 3.1 that the groups of the scenarios S-3.2.I and S-3.2.II, S-3.2.II and S-3.2.III, S-3.2.III and S-3.2.IV reflect the variation of the mean arrival rate, buffer capacity and mean service rate, respectively. The following figures

reveal remarkable increases in all performance metrics as the mean arrival rate rises. Secondly, the growth of the buffer capacity decreases the packet loss probability but slightly increases all the others. Thirdly, a small mean service rate (i.e., the scenario S-3.2.IV) results in a low throughput but high utilization, mean numbers of packets in the system and buffer, mean response time, mean queueing delay and packet loss probability. Finally, the decrease of the threshold value reduces the mean numbers of packets in the system and buffer, mean response time and mean queueing delay at the cost of the system utilization, throughput and packet loss probability.

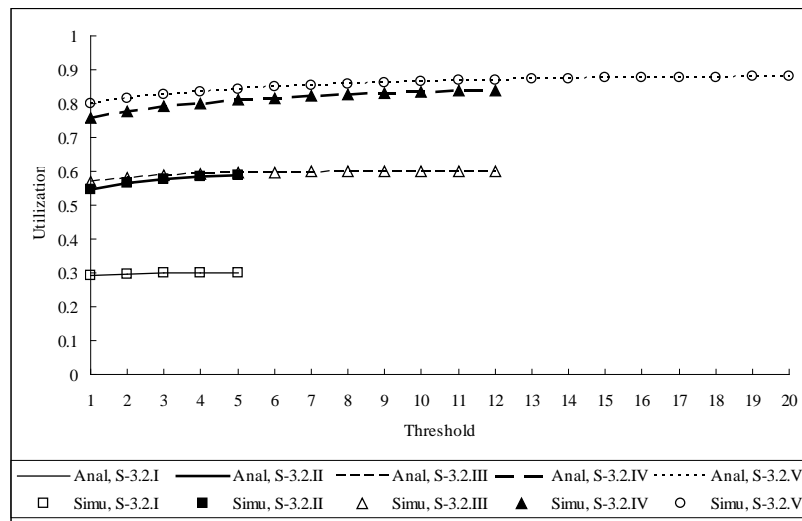


Figure 3.4. Utilization vs the threshold subject to five different scenarios.

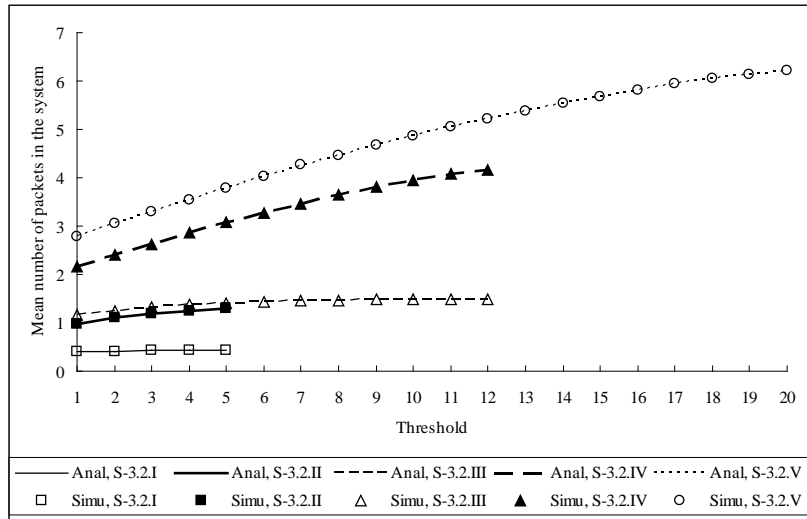


Figure 3.5. The mean number of packets in the system vs the threshold subject to five different scenarios.

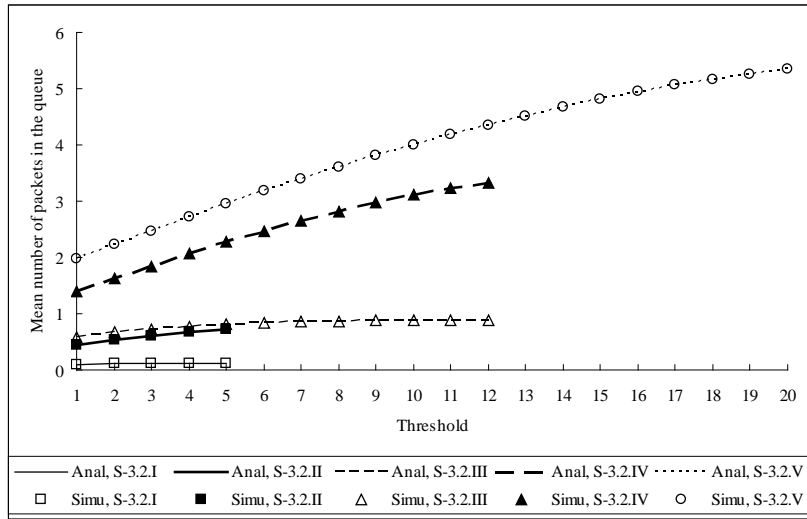


Figure 3.6. The mean number of packets in the queue vs the threshold subject to five different scenarios.

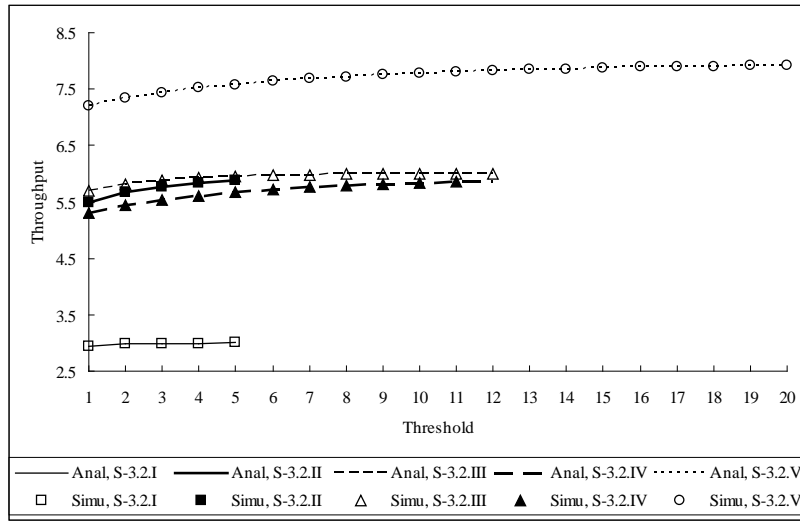


Figure 3.7. The throughput vs the threshold subject to five different scenarios.

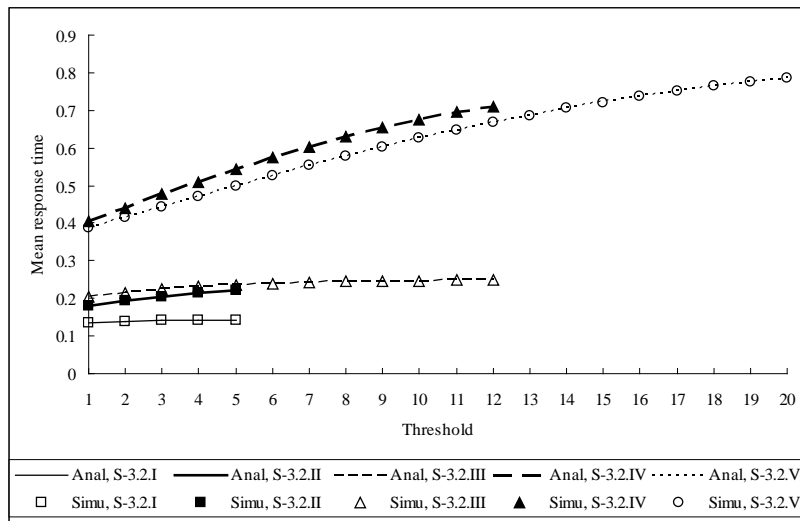


Figure 3.8. The mean response time vs the threshold subject to five different scenarios.

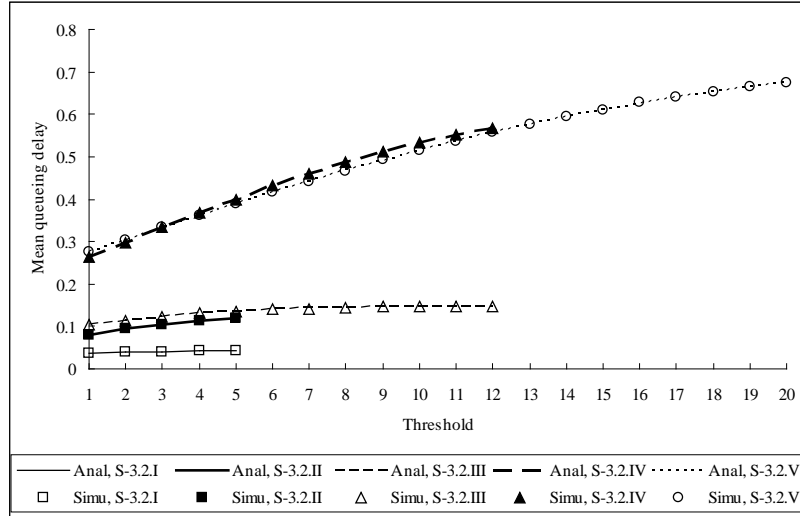


Figure 3.9. The mean queuing delay vs the threshold subject to five different scenarios.

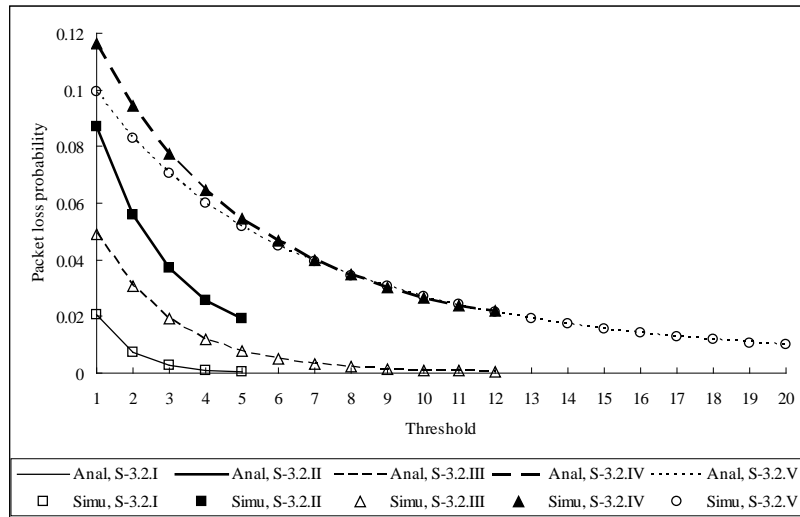


Figure 3.10. The packet loss probability vs the threshold subject to five different scenarios.

3.3 AQM Scheme with Multiple Class Traffic

Due to a variety of application requirements in communication networks, this section focuses on performance investigation of AQM congestion control scheme subject to two classes of traffic. The individual threshold is assigned to the buffer for each traffic class. A

Markov model is proposed to assist the derivation of the analytical expressions for the key aggregate and marginal performance metrics including the utilization, mean numbers of packets in the system and buffer, throughput, mean response time, mean queueing delay and packet loss probability. The investigation is concerned with the effects of a varying threshold on the aforementioned performance metrics.

3.3.1 Analytical Model

This section is to study $[M]^2/M/1/K/th_1/th_2$ queueing system where packets from two classes of traffic compete for the buffer space in an AQM router. Two thresholds th_b^1 and th_b^2 are assigned in the buffer for two class of traffic to notify the incipient stage of congestion, respectively.

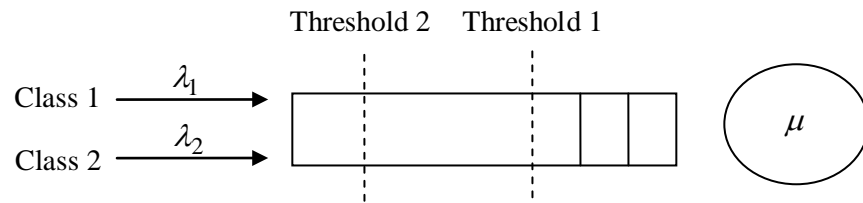


Figure 3.11. A Model of $[M]^2/M/1/K/th_1/th_2$ Queueing System

The arrival of each class c ($c = 1,2$) follows an independent Poisson process with an average arrival rate λ_c . The service time of both classes is exponentially distributed with mean $1/\mu$ and the total system capacity is $L = K + 1$, where K denotes the buffer capacity.

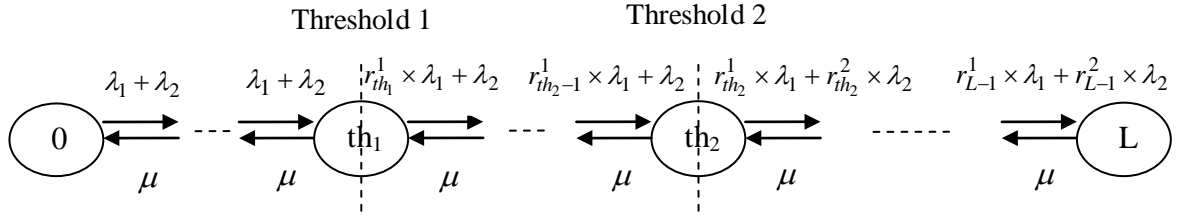


Figure 3.12. A State Transition Rate Diagram of [M]2/M/1/K/ th1/th2 Queuing System

Figure 3.12 shows the state transition rate diagram of the $[M]^2/M/1/K/th_1/th_2$ queuing system. Each state i ($0 \leq i \leq L$) represents that there are i packets in the system. Therefore, the thresholds $th_1 = th_b^1 + 1$ and $th_2 = th_b^2 + 1$, respectively. Transition from state i , ($0 \leq i \leq L-1$), to $i+1$ implies that a packet from Class-1 or Class-2 traffic enters in the system. It is well known that the superposition of independent Poisson Processes is still a Poisson process because of the memoryless property of the exponential distribution. Consequently, the aggregate mean arrival rate at state i equals the sum of current mean arrival rate of each class. In addition, the transition rate from state i to $i-1$ is μ as all packets are treated equally. The packets of class c ($c = 1, 2$) will be dropped randomly based on the dropping probability which linearly increases from 0 to 1 when the number of packets in the system grows from threshold th_c to the buffer capacity. This process can be seen as a decrease of the arriving rate with a reduction probability, r_i^c , ($0 \leq i \leq L, c = 1, 2$), which is given by

$$r_i^c = \begin{cases} 1 & 0 \leq i < th_c \\ 1 - \left(\frac{i - th_c + 1}{L - th_c + 1} \right) & th_c \leq i \leq L \end{cases} \quad (3.12)$$

The transition equilibrium equation and normalising equation can be found in Eqs. (3.13) and (3.14), respectively.

$$\left. \begin{aligned} (r_i^1 \lambda_1 + r_i^2 \lambda_2) p_0 &= \mu p_1 \\ (r_i^1 \lambda_1 + r_i^2 \lambda_2 + \mu) p_i &= (r_{i-1}^1 \lambda_1 + r_{i-1}^2 \lambda_2) p_{i-1} + \mu p_{i+1} \quad 1 \leq i < L \\ \mu p_L &= (r_{L-1}^1 \lambda_1 + r_{L-1}^2 \lambda_2) p_{L-1} \end{aligned} \right\} \quad (3.13)$$

$$\sum_{i=0}^L p_i = 1 \quad (3.14)$$

The state probability p_i ($0 \leq i \leq L$) can be given by solving those equations above:

$$p_i = \begin{cases} \frac{1}{1 + \sum_{j=1}^L \left(\prod_{k=0}^{j-1} \frac{r_k^1 \lambda_1 + r_k^2 \lambda_2}{\mu} \right)} & i = 0 \\ \left(\prod_{k=0}^{i-1} \frac{r_k^1 \lambda_1 + r_k^2 \lambda_2}{\mu} \right) \times p_0 & 1 \leq i \leq L \end{cases} \quad (3.15)$$

The derivations of the aggregate performance metrics are extremely similar to those presented in Section 3.2.1. To avoid repetition, this section demonstrates the expressions only without the detailed explanation. Eqs. (3.16)-(3.22) represent the utilization, the mean number of packets in the system, mean number of packets in the queue, throughput, mean response time, delay and packets loss probability.

$$\rho = 1 - p_0 \quad (3.16)$$

$$\bar{L} = \sum_{i=0}^L (p_i \times i) \quad (3.17)$$

$$\overline{L_b} = \sum_{i=0}^{L-1} (p_{i+1} \times i) \quad (3.18)$$

$$\overline{T} = \rho \times \mu \quad (3.19)$$

$$\overline{R} = \frac{\overline{L}}{\overline{T}} \quad (3.20)$$

$$\overline{D} = \frac{\overline{L_b}}{\overline{T}} \quad (3.21)$$

$$PLP = \sum_{i=0}^L p_i \times \frac{(1-r_i^1) \times \lambda_1 + (1-r_i^2) \times \lambda_2}{\lambda_1 + \lambda_2} \quad (3.22)$$

Next, we describe detailed derivation of the marginal performance metrics. For a system in the steady-state, the mean arrival rate equals to its throughput. So the throughput of each class can be expressed as.

$$T^c = \sum_{i=0}^{L-1} p_i \times r_i^c \times \lambda_c \quad (3.23)$$

Because both classes of traffic are served identically, the mean response time and delay of each class can be derived using Eqs. (3.24) and (3.25) [75]. The delay of a packet from class c ($c=1,2$) can be decomposed into two parts: the mean residual life due to the other packets found in service and the delay due to packets found in the queue upon its arrival. In an M/M/1/K queueing system, the mean residual life equals to the mean service time. The average response time consists of the delay and mean service time of the packet.

$$\overline{R^c} = \frac{\sum_{i=0}^{L-1} \left(\frac{r_i^c \lambda_c}{r_i^1 \lambda_1 + r_i^2 \lambda_2} \times p_{i+1} \times \frac{i+1}{\mu} \right)}{\sum_{i=0}^{L-1} \left(\frac{r_i^c \lambda_c}{r_i^1 \lambda_1 + r_i^2 \lambda_2} \times p_{i+1} \right)} \quad (3.24)$$

$$\overline{W_q^c} = \frac{\sum_{i=0}^{L-1} \left(\frac{r_i^c \lambda_c}{r_i^1 \lambda_1 + r_i^2 \lambda_2} \times p_{i+1} \times \frac{i}{\mu} \right)}{\sum_{i=0}^{L-1} \left(\frac{r_i^c \lambda_c}{r_i^1 \lambda_1 + r_i^2 \lambda_2} \times p_{i+1} \right)} \quad (3.25)$$

The ratio of the instant packet loss rate of class c ($c = 1, 2$) to the total arrival rate $\lambda_1 + \lambda_2$ is the instant packet loss probability from class c .

$$\overline{PLP^c} = \sum_{i=0}^L p_i \times \frac{(1 - r_i^c) \times \lambda_c}{\lambda_1 + \lambda_2} \quad (3.26)$$

In order to calculate the probability distribution of the marginal queue length for each class, the probability that packets of each class stay in any position in the system should be calculated firstly. There are L positions in the system with the number being $1 \dots L$ from the server to the tail of the queue. If a packet from class c ($c = 1, 2$) is allocated in the position i ($1 \leq i \leq L$) when it arrives in the system, it will experience all the positions j before i , $j \leq i$. In other words, the probability that there is a packet from class c in state j should be the sum of all the probabilities that the packet arrives in the system and is allocated at position i , $1 \leq j \leq i \leq L$. From the transition diagram above, the later

probability can be calculated intuitively as $p_{i-1} \times \frac{r_{i-1}^c \times \lambda_c}{r_{i-1}^1 \times \lambda_1 + r_{i-1}^2 \times \lambda_2}$, ($1 \leq j \leq i \leq L$). So

the probability that a packet from class c is in position i , noted as m_i^c , can be derived as:

$$m_i^c = \frac{\sum_{j=i-1}^{L-1} (p_{j+1} \times \frac{r_j^c \times \lambda_c}{r_j^1 \times \lambda_1 + r_j^2 \times \lambda_2})}{\sum_{j=i-1}^{L-1} p_{j+1}} \quad 1 \leq i \leq L \quad c = 1, 2 \quad (3.27)$$

If the number of packets from class c is q , $0 \leq q \leq L$, then the number of aggregate packets in the system should be not smaller than q . When the length of the system is l $l \geq q$, there are C_l^q combinations of two classes to make the length of class c be q . Furthermore, the probability of each combinations is different and can be calculated using m_i^c . So the marginal probability of queue length for each class c , p_q^c , can be derived as follows:

$$p_q^1 = \begin{cases} p_0 + \sum_{l=1}^L \left\{ p_l \times \left[\sum_{i=0}^{C_l^q - 1} \left(\prod_{j=0}^{l-1} m_{j+1}^q[i, j] \right) \right] \right\} & q = 0 \\ \sum_{l=q}^L \left\{ p_l \times \left[\sum_{i=0}^{C_l^q - 1} \left(\prod_{j=0}^{l-1} m_{j+1}^q[i, j] \right) \right] \right\} & 1 \leq q \leq L \end{cases} \quad (3.28)$$

$$p_q^2 = \begin{cases} p_0 + \sum_{l=1}^L \left\{ p_l \times \left[\sum_{i=0}^{C_l^q - 1} \left(\prod_{j=0}^{l-1} m_{j+1} \mathbf{B}_1^q[i, j] \right) \right] \right\} & q = 0 \\ \sum_{l=q}^L \left\{ p_l \times \left[\sum_{i=0}^{C_l^q - 1} \left(\prod_{j=0}^{l-1} m_{j+1} \mathbf{B}_1^q[i, j] \right) \right] \right\} & 1 \leq q \leq L \end{cases} \quad (3.29)$$

To represent p_q^c , two matrices \mathbf{A}_i^j and \mathbf{B}_i^j are used to describe the possible combinations and defined as:

$$\mathbf{A}_i^j = \begin{cases} (2 \ 2 \ \cdots \ 2)_{C_i^j \times i} & i = 1, 2, \dots; j = 0 \\ (1 \ 1 \ \cdots \ 1)_{C_i^j \times i} & i = 1, 2, \dots; j = i \\ \begin{pmatrix} \mathbf{a}_{(i-1) \times 1} & \mathbf{A}_{(i-1)}^{(j-1)} \\ \mathbf{\beta}_{(i-1) \times 1} & \mathbf{A}_{(i-1)}^j \end{pmatrix}_{C_i^j \times i} & i = 2, \dots, L; j = 1, 2, \dots, i-1 \end{cases} \quad (3.30)$$

$$\mathbf{B}_i^j = \begin{cases} (1 \ 1 \ \cdots \ 1)_{C_i^j \times i} & i = 1, 2, \dots; j = 0 \\ (2 \ 2 \ \cdots \ 2)_{C_i^j \times i} & i = 1, 2, \dots; j = i \\ \begin{pmatrix} \mathbf{\beta}_{(i-1) \times 1} & \mathbf{B}_{(i-1)}^{(j-1)} \\ \mathbf{a}_{(i-1) \times 1} & \mathbf{B}_{(i-1)}^j \end{pmatrix}_{C_i^j \times i} & i = 2, \dots, L; j = 1, 2, \dots, i-1 \end{cases} \quad (3.31)$$

Both matrices are of size $C_i^j \times i$. Each row of \mathbf{A}_i^j is a possible combination of Class-1 and Class-2 when the aggregate queue length is i and the queue length for class 1 is j . \mathbf{B}_i^j can be defined for class two similarly. The two basic matrices $\mathbf{a}_{i \times 1}$ and $\mathbf{\beta}_{i \times 1}$ are defined as:

$$\mathbf{a}_{i \times 1} = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}_{i \times 1} \quad i = 1, 2, \dots \quad (3.32)$$

$$\mathbf{b}_{i \times 1} = \begin{pmatrix} 2 \\ 2 \\ \vdots \\ 2 \end{pmatrix}_{i \times 1} \quad i = 1, 2, \dots \quad (3.33)$$

The relationship between \mathbf{A}_i^j and \mathbf{B}_i^j is shown as follow:

$$\gamma_{C_i^j \times i} = \begin{pmatrix} 3 \dots 3 \\ \vdots \vdots \\ 3 \dots 3 \end{pmatrix} \quad (3.34)$$

$$\mathbf{B}_i^j = \gamma_{C_i^j \times i} - \mathbf{A}_i^j \quad (3.35)$$

Similar to Eq. (3.17), the mean number $\overline{L^c}$ of packets from class c ($c = 1, 2$) in the system can be calculated by

$$\overline{L^c} = \sum_{i=1}^L (p_i \times \sum_{j=0}^i m_j^c) \quad (3.36)$$

3.3.2 Performance Validation

In this section, we validate the accuracy of the proposed models by comparing the analytical results with those obtained from the corresponding simulators programmed in JAVA. The credibility of the analytical models against simulation results are examined on

the base of different combinations of the mean arrival rate of each traffic class, mean service rate, buffer capacity and the threshold of each traffic class. Specifically, we address 5 different scenarios and the corresponding parameter settings are listed in Table 3.2 for the validation of the Markovian model. Table 3.3 presents the simulation and corresponding analytical results of all performance metrics derived in Section 3.3.1 for the five different scenarios. It can be observed that the analytical results well match the corresponding simulation results. This observation demonstrates that the developed models are very accurate in calculating various performance metrics and examining the performance of the AQM in presence of two Poisson traffic classes.

	λ_1	λ_2	S	L	Th_1	Th_2
S-3.3.I	0.2	0.2	0.5	6	2	2
S-3.3.II	0.4	0.7	1.4	10	3	5
S-3.3.III	0.6	1.2	2.0	17	8	4
S-3.3.IV	2.3	1.9	6	11	4	7
S-3.3.V	3	0.8	3.9	25	16	13

Table 3.2. The parameter settings corresponding to five scenarios.

Performance Metrics	A/S	S-3.3.I	S-3.3.II	S-3.3.III	S-3.3.IV	S-3.3.V
U	Anal	0.713906	0.747782	0.845406	0.68524	0.936557
	Simul	0.713839	0.747667	0.845527	0.685334	0.936458
\bar{L}	Anal	1.737025	2.322831	3.751313	1.933476	8.690221
	Simul	1.736909	2.321747	3.753078	1.933619	8.685808
\bar{L}_q	Anal	1.02312	1.575049	2.905907	1.248236	7.753664
	Simul	1.023071	1.57408	2.90755	1.248286	7.74935
T	Anal	0.356953	1.046895	1.690812	4.11144	3.652573
	Simul	0.356916	1.046791	1.691045	4.11197	3.652612
\bar{R}	Anal	4.86626	2.218782	2.218646	0.470267	2.379205
	Simul	4.866415	2.218004	2.219381	0.47024	2.378002
\bar{D}	Anal	2.86626	1.504497	1.718646	0.303601	2.122795
	Simul	2.866415	1.503718	1.719381	0.303574	2.121591
PLP	Anal	0.107618	0.048278	0.06066	0.021086	0.038796
	Simul	0.107638	0.048214	0.060717	0.021058	0.038788
\bar{L}^1	Anal	0.868513	0.789117	1.362643	1.022623	6.97511
	Simul	0.868296	0.78891	1.363379	1.022489	6.971768

\overline{L}_q^1	Anal	0.51156	0.525958	1.067942	0.652225	6.229305
	Simul	0.51146	0.525768	1.068665	0.652094	6.226076
T^1	Anal	0.178476	0.368423	0.589403	2.222392	2.908639
	Simul	0.178444	0.368436	0.589541	2.222441	2.908676
\overline{R}^1	Anal	4.86626	2.141876	2.311903	0.460145	2.398066
	Simul	4.866221	2.141309	2.312707	0.46008	2.396929
\overline{D}^1	Anal	2.86626	1.42759	1.811903	0.293479	2.141656
	Simul	2.866221	1.427024	1.812707	0.293413	2.140519
PLP^1	Anal	0.107618	0.078942	0.017661	0.033742	0.030454
	Simul	0.107632	0.078846	0.017663	0.033707	0.030449
\overline{L}^2	Anal	0.868513	1.533714	2.388669	0.910853	1.715111
	Simul	0.868613	1.532836	2.389699	0.91113	1.71404
\overline{L}_q^2	Anal	0.51156	1.049092	1.837965	0.596011	1.524359
	Simul	0.51161	1.048312	1.838886	0.596192	1.523274
T^2	Anal	0.178476	0.678471	1.101409	1.889048	0.743934
	Simul	0.178472	0.678355	1.101504	1.889529	0.743935
\overline{R}^2	Anal	4.86626	2.260544	2.168741	0.482176	2.305461
	Simul	4.866609	2.25966	2.169432	0.482191	2.303999
\overline{D}^2	Anal	2.86626	1.546258	1.668741	0.315509	2.049051
	Simul	2.866609	1.545374	1.669432	0.315524	2.047589
PLP^2	Anal	0.107618	0.030755	0.08216	0.005764	0.070082
	Simul	0.107645	0.030707	0.082245	0.005751	0.070059

Table 3.3. The analysis results of all performance metrics and the corresponding simulation results subject to five different scenarios.

3.3.3 Performance Evaluation

This section evaluates the effects of varying thresholds on the marginal and aggregate performance metrics including the system utilization, mean number of packets in the system (also named as mean queue length), throughput, mean queueing delay and packet loss probability. In Scenario S-3.3.VI, threshold th_1 is fixed at position 7 and threshold th_2 increases from 7 to 16 when the total system capacity is 20. In Scenario S-3.3.VII, threshold th_1 increases from 7 to 16 whereas threshold th_2 remains fixed at 16. The X-axis in all the following figures represents the difference between the threshold values, i.e., $th_2 - th_1$. In both scenarios, Class-1 and Class-2 traffic are generated, respectively, by a

Poisson process with the mean arrival rate $\lambda_1 = 0.4$ and $\lambda_2 = 0.3$. The mean service rate μ is set to be 0.8 in order to make certain that the queueing system is stable.

The marginal mean queue length, throughput, mean queueing delay and packet loss probability have been shown in Figures 3.13-3.16, respectively. It has been clearly shown in these figures that the variation of a threshold significantly affects all performance measures. In particular, as the value of a threshold increases, the number of packets of the corresponding class controlled by the threshold in the system also increases. As a consequence, its mean queue length, throughput, and mean queueing delay tend to increase whereas its packets loss probability tends to decrease. However, for the class not controlled by the fixed threshold, the throughput tends to decrease and the mean queue length, packets loss probability and mean queueing delay increase. Furthermore, it is also shown that when the value of $th_2 - th_1$ is same, the smaller values of th_1 and th_2 can reduce the marginal mean queue length and mean queueing delay of each class, but keep the marginal throughput very similar.

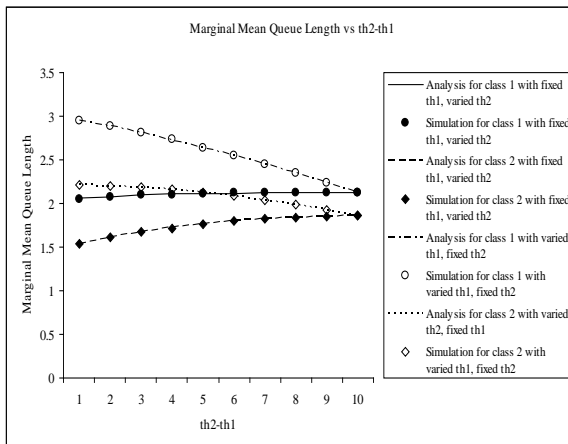


Figure 3.13. The Marginal Mean Queue Length vs $th_2 - th_1$

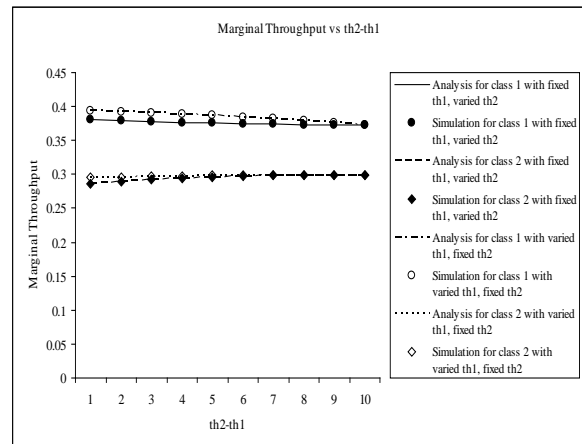


Figure 3.14. The Marginal Throughput vs $th_2 - th_1$

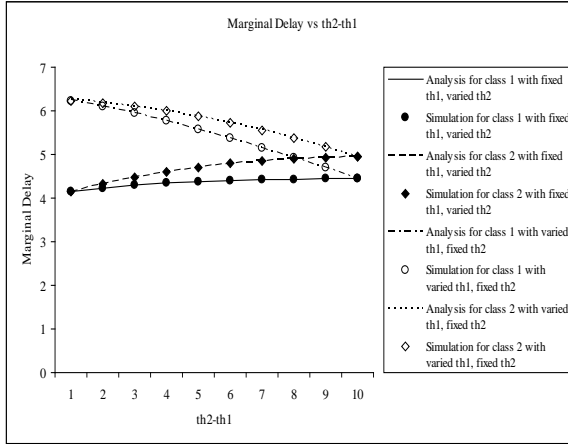


Figure 3.15. The Marginal Queueing Delay vs $th_2 - th_1$

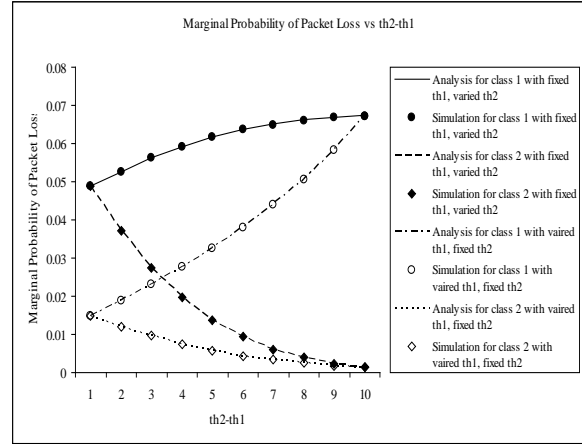


Figure 3.16. The Marginal Packet Loss Probability vs $th_2 - th_1$

The relative aggregate performance class measures have been shown in Figures 3.17-3.21. It is clear that increasing the value of a threshold enables more packets to enter into the system. So the utilization, mean queue length, throughput and mean queueing delay increase by increasing the value of a threshold but the packet loss probability reduces.

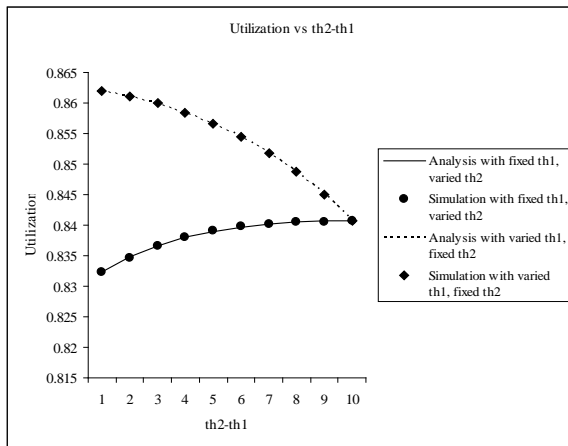


Figure 3.17. The Utilization vs $th_2 - th_1$

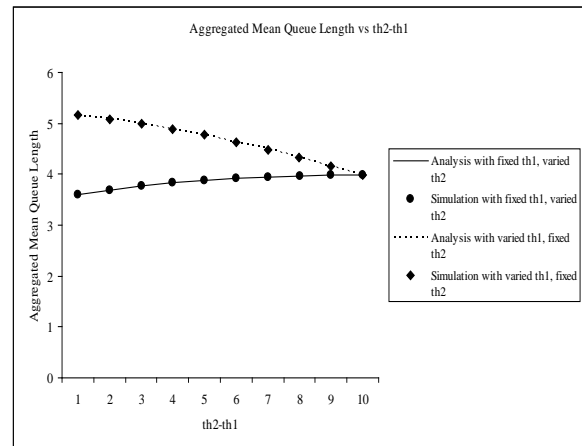


Figure 3.18. The Aggregate Mean Queue Length vs $th_2 - th_1$

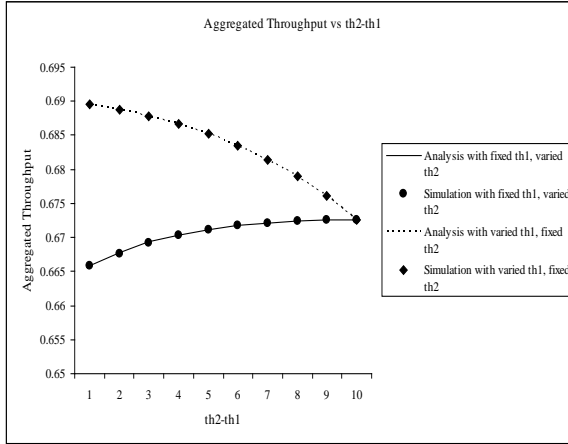


Figure 3.19. The Aggregate Throughput vs th2-th1

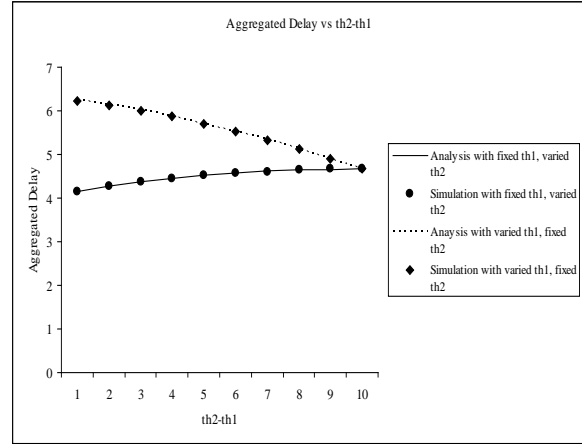


Figure 3.20. The Aggregate Mean Queueing Delay vs th2-th1

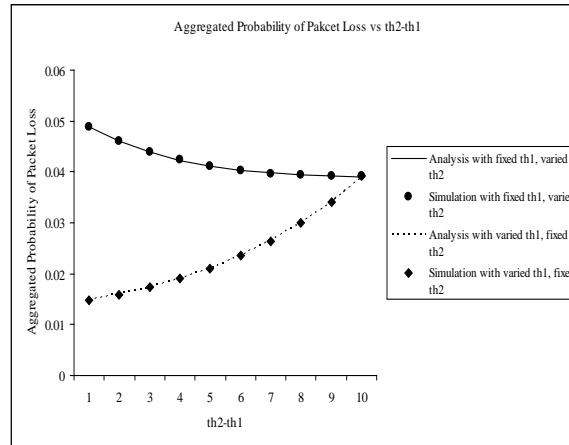


Figure 3.21. The Aggregate Packets Loss Probability vs th2-th1

3.4 Summary

This chapter has presented the stochastic analysis of single-server finite queueing systems for the performance evaluation of AQM scheme under the non-bursty Poisson arrival process. Two one-dimensional continuous time Markov models have been proposed, respectively, for AQM scheme subject to single class and multiple classes traffic. Two general analytical models are easily extended for some other AQM algorithms with

different dropping functions by re-calculating the corresponding dropping probability. Extensive experiments have demonstrated the accuracy of two models by comparing analytical results with those obtained from discrete event simulation in JAVA programming. The first model for a single class traffic is used to evaluate the system performance with different mean arrival rates, mean service rates and buffer capacity. The second model for two traffic classes is adopted to investigate the effects of varying thresholds on the aggregate and marginal performance. The main contributions of this chapter are concluded as follows:

- (i). Expressions of the marginal performance metrics, including the mean numbers of packets in the system and queue, throughput, mean response time, mean queueing delay, and packet loss probability, for a finite queueing system with thresholds under two classes of non-bursty traffic have been derived.
- (ii). A high mean arrival rate results in a high utilization and throughput, large number of packets in the system and buffer as well as packet loss probability, long response time and queueing delay.
- (iii). A small service rate decreases the throughput whilst increases all other performance metrics.
- (iv). A large buffer capacity reduces the packet loss probability but increases all other performance metrics.
- (v). The effects of the varying threshold on the traffic controlled by this threshold: as the threshold increases, the mean numbers of packets in the system and buffer,

throughput, response time and queueing delay increase, but the packet loss probability decreases.

- (vi). The effects of the varying threshold on the traffic not controlled by this threshold: as the threshold increases, the mean number of packets in the system and buffer, packet loss probability, response time and queueing delay increase, but the throughput decreases.
- (vii). The effects of the varying threshold on the aggregate system performance: as the threshold increases, the aggregate number of packets in the system and buffer, throughput, response time and queueing delay increases, but the aggregate packet loss probability decreases.

Chapter 4

Performance Modeling and Analysis of AQM with Single Class Bursty Traffic

4.1 Introduction

With the convincing evidence of traffic burstiness and correlations exhibited by key services such as compressed video, voice, etc, over modern high-speed networks, several stochastic models [86-87] have been presented to capture such traffic properties. Specifically, the well-known Markov-Modulated Poisson Process (MMPP) has been widely used for this purpose owing to its ability to model the time-varying arrival rate and to adequately capture the important correlation between inter-arrival times while still maintaining analytical tractability [15]. This chapter aims to investigate the significant effects of burstiness and correlations of bursty traffic on the performance of AQM.

A stochastic queueing model for the performance evaluation of AQM mechanism using a two-state MMPP (MMPP-2) to model the bursty traffic source is presented in Section 4.2. Then we derive the expressions of the performance metrics including the utilization, mean numbers of packets in the system and buffer, throughput, mean response time, mean queueing delay and packet loss probability. Through extensive comparisons between analytical results and those obtained from simulation experiments, we demonstrate the accuracy of the proposed model in Section 4.3. Finally, the model is adopted to investigate how the threshold value, burstiness and correlation of the MMPP-2 traffic affect

the aforementioned performance metrics of AQM system and how the effects of one of these three parameters are influenced by the variation of the others.

4.2 Analytical Model

The section introduces the system theoretical framework based on the AQM mechanism using a queue threshold in the presence of busty traffic and presents the formulation of analytical model for the corresponding queueing system under the bursty MMPP-2 arrival process. Different from Section 3.2.1, the arrivals from a bursty traffic source feed into a finite buffer in an AQM-enabled router (c.f., Figure 4.1). The system employs a threshold to inform the source of the incipient state of congestion by dropping the arriving packets according to a dropping probability when the instantaneous queue length exceeds the threshold value. The dropping probability linearly increases from 0 to 1 as the instantaneous queue length varies from the threshold value to the buffer capacity.

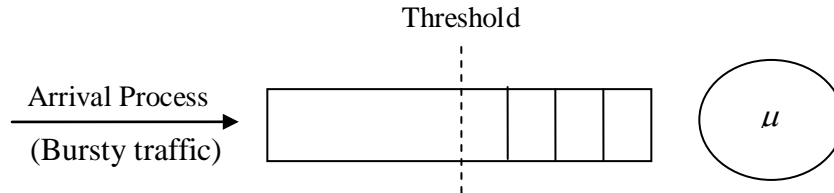


Figure 4.1. A Model of MMPP-2/M/1/K/th Queueing System

A state transition diagram of the analytical queueing model under the MMPP-2 arrival process is shown in Fig. 4.2 where State (i, j) with $(1 \leq i \leq L, 1 \leq j \leq 2)$ represents the current instantaneous queue length is i and the traffic arrival process MMPP-2 is at state j . The system capacity is given as L including a finite buffer of size K and a single server (i.e., $L = k + 1$). δ_j ($j = 1, 2$) is the transition rate out of the underlying Markov state

(i, j) to $(i, 3-j)$, where $(0 \leq i \leq L)$. The transition rate from state (i, j) to $(i-1, j)$, where $(1 \leq i \leq L, 1 \leq j \leq 2)$, is μ denoting the mean service rate. The transition from state (i, j) to $(i+1, j)$, where $(0 \leq i < L, 1 \leq j \leq 2)$, represents that a packet is injected into the system. It is noticeable that the transition rate from state (i, j) to $(i+1, j)$, where $(th \leq i < L, 1 \leq j \leq 2)$, is reduced from λ_j to $(1-d_i)\lambda_j$ since the packet dropping process can be seen as a decrease of the arriving rate with a probability $(1-d_i)$. The calculation of the packet dropping probability d_i ($0 \leq i \leq L$) is given as follows:

$$d_i = \begin{cases} 0 & 0 \leq i < th \\ \frac{i-th+1}{L-th+1} & th \leq i \leq L \end{cases} \quad (4.1)$$

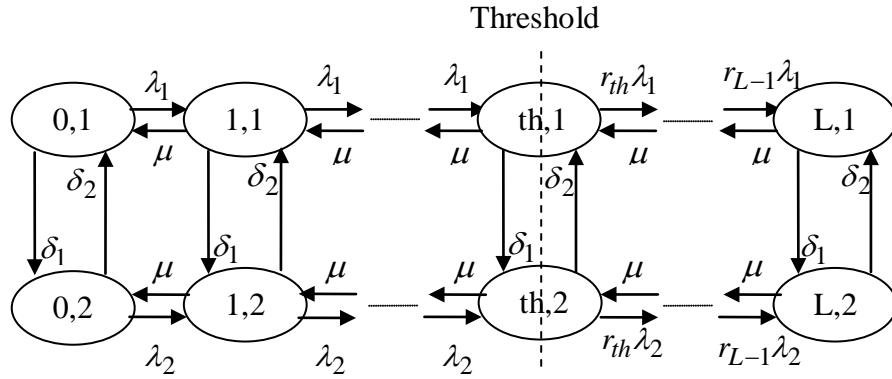


Figure 4.2. A State Transition Rate diagram of MMPP-2/M/1/K/th Queueing System.

Consider the system in the equilibrium state, the following set of equilibrium equations can be obtained directly from the state transition rate diagram. In this equilibrium case, it is clear that the rate of the input flow must be equal to that of the output flow of any given state.

$$\left. \begin{aligned}
((1-d_0)\lambda_1 + \delta_1)p_{01} &= \delta_2 p_{02} + \mu p_{11} \\
((1-d_0)\lambda_2 + \delta_2)p_{02} &= \delta_1 p_{01} + \mu p_{12} \\
((1-d_i)\lambda_1 + \delta_1 + \mu)p_{i1} &= (1-d_{i-1})\lambda_1 p_{i-1,1} + \delta_2 p_{i2} + \mu p_{i+1,1} \quad 1 \leq i \leq L \\
((1-d_i)\lambda_2 + \delta_2 + \mu)p_{i2} &= (1-d_{i-1})\lambda_2 p_{i-1,2} + \delta_1 p_{i1} + \mu p_{i+1,2} \quad 1 \leq i \leq L \\
(\delta_1 + \mu)p_{L1} &= (1-d_{L-1})\lambda_1 p_{L-1,1} + \delta_2 p_{L2} \\
(\delta_2 + \mu)p_{L2} &= (1-d_{L-1})\lambda_2 p_{L-1,2} + \delta_1 p_{L1}
\end{aligned} \right\} \quad (4.2)$$

$$\sum_{i=0}^L \sum_{j=1}^2 p_{ij} = 1 \quad (4.3)$$

where p_{ij} is the probability that the system is in State (i, j) . Solving these equations, we can find probability, p_{ij} , of each state in the Markov model as follows:

$$p_{ij} = a_{ij}m + b_{ij}n \quad (4.4)$$

with

$$a_{ij} = \begin{cases} 1 & i=0, j=1 \\ 0 & i=0, j=2 \\ \frac{(1-d_0)\lambda_1 + \delta_1}{\mu} & i=1, j=1 \\ -\frac{\delta_1}{\mu} & i=1, j=2 \\ \frac{((1-d_{i-1})\lambda_1 + \mu + \delta_1)a_{i-1,1} - (1-d_{i-2})\lambda_1 a_{i-2,1} - \delta_2 a_{i-1,2}}{\mu} & 2 \leq i \leq L, j=1 \\ \frac{((1-d_{i-1})\lambda_2 + \mu + \delta_2)a_{i-1,2} - (1-d_{i-2})\lambda_2 a_{i-2,2} - \delta_1 a_{i-1,1}}{\mu} & 2 \leq i \leq L, j=2 \end{cases} \quad (4.5)$$

$$b_{ij} = \begin{cases} 0 & i = 0, j = 1 \\ 1 & i = 0, j = 2 \\ \frac{\delta_1}{\mu} & i = 1, j = 1 \\ \frac{(1-d_0)\lambda_1 + \delta_1}{\mu} & i = 1, j = 2 \\ \frac{((1-d_{i-1})\lambda_1 + \mu + \delta_1)b_{i-1,1} - (1-d_{i-2})\lambda_1 b_{i-2,1} - \delta_2 b_{i-1,2}}{\mu} & 2 \leq i \leq L, j = 1 \\ \frac{((1-d_{i-1})\lambda_2 + \mu + \delta_2)b_{i-1,2} - (1-d_{i-2})\lambda_2 b_{i-2,2} - \delta_1 b_{i-1,1}}{\mu} & 2 \leq i \leq L, j = 2 \end{cases} \quad (4.6)$$

$$k = \lambda_1(1-d_{L-1,0})b_{L-1,1} + \delta_2 b_{L2} - (\mu + \delta_1)b_{L1} \quad (4.7)$$

$$h = (\mu + \delta_1)a_{L1} - \lambda_1(1-d_{L-1,0})a_{L-1,1} - \delta_2 a_{L2} \quad (4.8)$$

$$n = \frac{h}{k \sum_{i=0}^L \sum_{j=1}^2 a_{ij} + h \sum_{i=0}^L \sum_{j=1}^2 b_{ij}} \quad (4.9)$$

$$m = \frac{k}{h} n \quad (4.10)$$

Based on the derivation of p_{ij} , we can obtain the system performance metrics including utilization (ρ), the mean number of packets in the system (\bar{L}), numbers of packets in the buffer (\bar{L}_b), system throughput (T), response time (\bar{R}), queueing delay (\bar{W}_q) and packets loss probability (\overline{PLP}).

The server is engaged as long as the number of packets in the system is not zero. The probability that the server is idle is $(p_{01} + p_{02})$. Therefore, the system utilization can be written as

$$\rho = 1 - p_{01} - p_{02} \quad (4.11)$$

The probability that there are i ($0 \leq i \leq L$) packets in the system is $(p_{i1} + p_{i2})$, while The probability that there are i ($1 \leq i \leq K$) packets in the buffer is $(p_{i+1,1} + p_{i+1,2})$ as the system consists of the buffer and one server. Consequently, the mean number of packets in the system and buffer can be calculated, respectively, as follows

$$\bar{L} = \sum_{i=0}^L \sum_{j=1}^2 (p_{ij} \times i) \quad (4.12)$$

$$\bar{L}_b = \sum_{i=1}^L \sum_{j=1}^2 (p_{i+1,j} \times i) \quad (4.13)$$

The transition rate at which packets go through the system is μ if the server is busy and becomes 0 otherwise. Therefore, the throughput is equal to the system service rate multiplied by utilization.

$$\bar{T} = \rho \times \mu \quad (4.14)$$

Again, the mean response time and queueing delay can be solved using Little's Law [75].

$$\bar{R} = \frac{\bar{L}}{\bar{T}} \quad (4.15)$$

$$\bar{D} = \frac{\bar{L}_b}{\bar{T}} \quad (4.16)$$

The packet loss probability consists of the probability of packet loss after the buffer is full and that of packets being dropped due to AQM scheme before the buffer is full. The

difference between the average arrival rate $\overline{\lambda_{mmpp-2}}$ and the throughput can approximate the average packet loss rate. So the packet loss probability is given below:

$$PLP = \frac{\overline{\lambda_{mmpp-2}} - \overline{T}}{\overline{\lambda_{mmpp-2}}} \quad (4.17)$$

4.3 Model Validation and Evaluation

4.3.1 Model Validation

We develop a discrete-event simulator programmed in JAVA in order to validate the above analytical model. Numerous validation experiments have been carried out for different combinations of the system capacity, mean service rate and MMPP-2 input traffic. Specifically, this section presents the results of all the derived performance metrics for the following four different scenarios described in Table 4.1. The value of the threshold assigned in the buffer varies from 1 to the buffer capacity for each scenario. Figures 4.3-4.9 demonstrate the analytical and the corresponding simulation results of the utilization, mean number of packets in the system and the buffer, throughput, mean response time, mean queueing delay and the packet loss probability in the AQM system subject to MMPP-2 traffic, respectively. It can be observed from all figures that the analytical results well match the corresponding simulation results. This observation illustrates that the proposed models are very accurate in calculating various performance metrics and examining the performance of AQM under bursty MMPP traffic.

Furthermore, it can be found from all figures that, for all scenario, the packet loss probability decreases while the utilization, mean number of packets in the system and the buffer, throughput, mean response time and mean queueing delay increase as the threshold value rises. This is because such variation of the threshold value enables more packets to be injected into the system. Moreover, the figures show differential degrees of performance variation for each scenario, caused by the distinct relationships among the arrival process, service rate and system capacity. For instance, the utilization, mean response time, mean queueing delay and packet loss probability generated in Scenario S-4.4.I vary more sharply than the corresponding performance metrics produced in other scenarios as shown in Figures 4.3, 4.7-4.9, respectively.

	λ_1	λ_2	δ_1	δ_2	μ	L
S-4.4.I	1	3	0.34	0.7	2.5	4
S-4.4.II	4	18	0.9	0.12	20	10
S-4.4.III	6	4	0.5	0.8	8	19
S-4.4.IV	20	7	0.64	0.26	17	25

Table 4.1. The parameter settings corresponding to four scenarios.

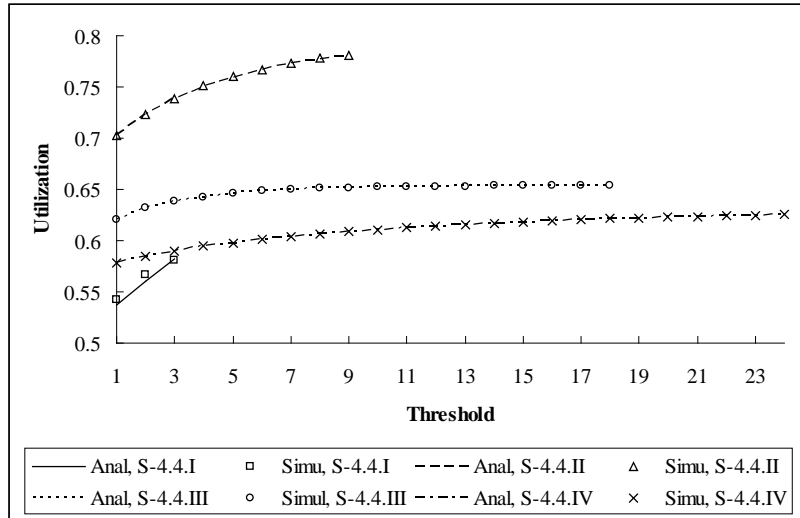


Figure 4.3. Utilization vs threshold under 4 different scenarios.

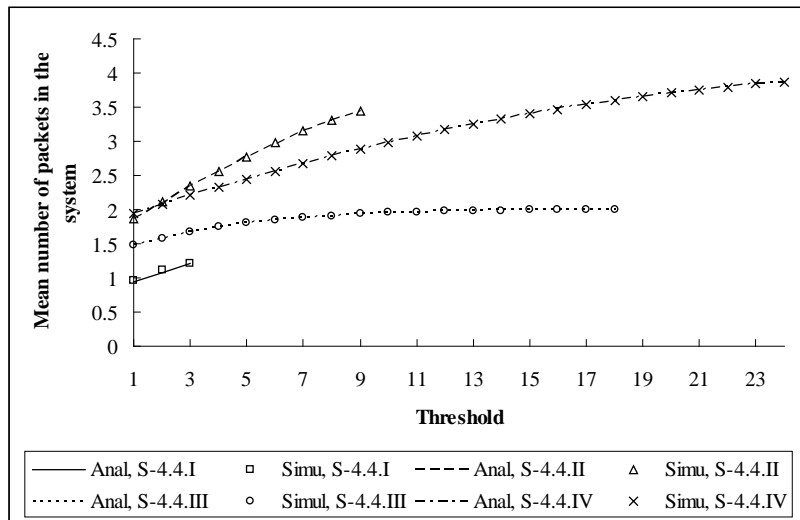


Figure 4.4. Mean number of packets in the system vs threshold under 4 different scenarios.

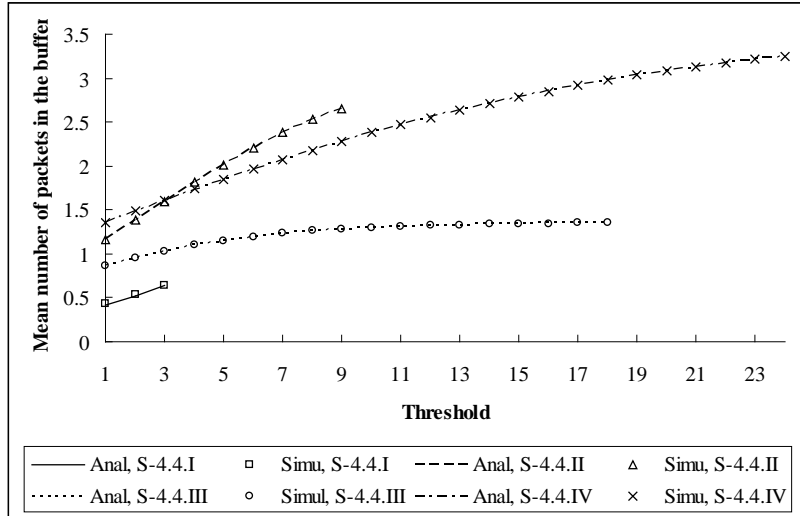


Figure 4.5. Mean number of packets in the buffer vs threshold under 4 different scenarios.

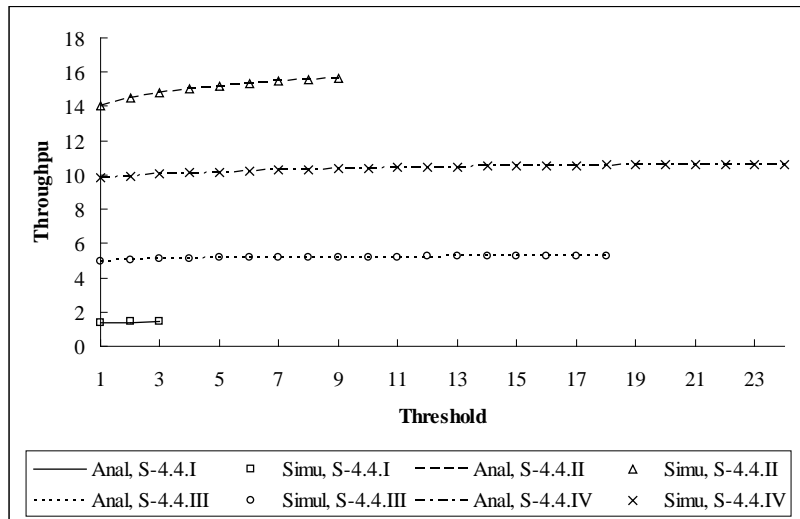


Figure 4.6. Throughput vs threshold under 4 different scenarios.

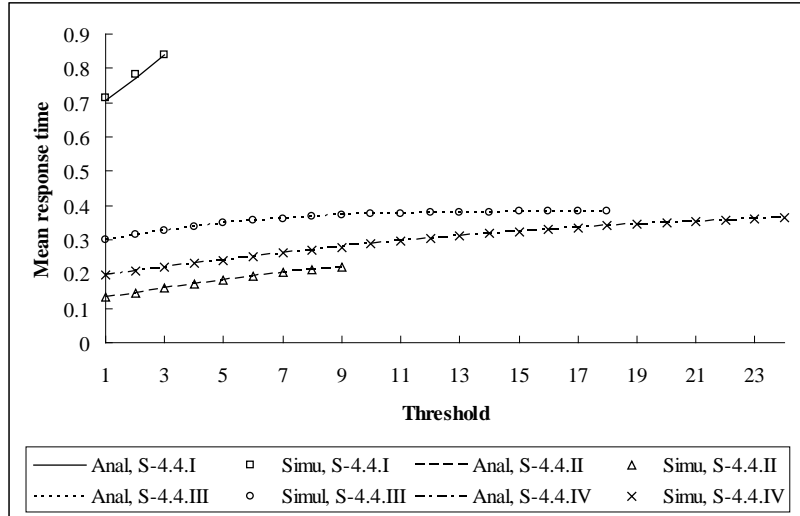


Figure 4.7. Mean response time vs threshold under 4 different scenarios.

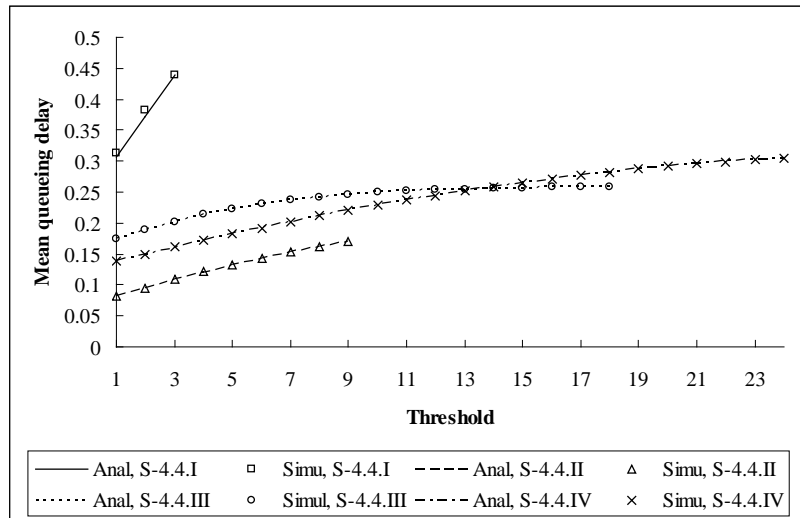


Figure 4.8. Mean queueing delay vs threshold under 4 different scenarios.

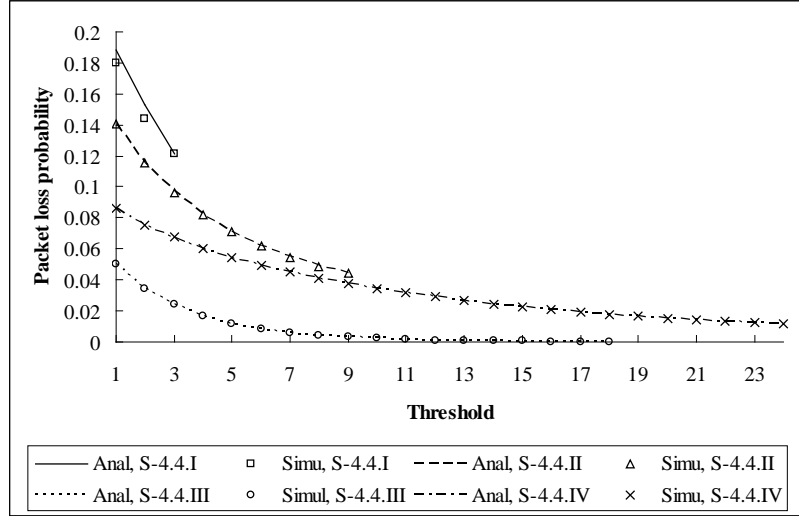


Figure 4.9. Packet loss probability vs threshold under 4 different scenarios.

4.3.2 Effects of Burstiness and Correlations

This section uses the validated analytical model to evaluate the impact of burstiness and correlations of the MMPP-2 traffic on the performance of AQM. In the following analysis, different values of the SCV and 1-step autocorrelation coefficient of MMPP-2 traffic based on Eqs. (2.7-2.8) in Section 4.2 are set to be ($c^2 = 5, 10$) and ($r_1 = 0.1, 0.2, 0.3, 0.4$), respectively. We can obtain the parameters of the MMPP-2 traffic by keeping the mean arrival rate of the MMPP-2 traffic, $\overline{\lambda_{mmpp-2}} = (\lambda_1\delta_2 + \lambda_2\delta_1)/(\delta_1 + \delta_2)$, constant at 14 and assuming $\delta_1 = \delta_2$. Table 4.2 below presents the various combinations of c^2 and r_1 as well as the parameters of the corresponding MMPP-2 traffic. The system capacities L and the mean service rate μ are 20 and 18, respectively. In addition, 3 specific values are assigned to the threshold, ($th = 8, 14, 20$), in order to compare the effects of burstiness and correlations of MMPP-2 traffic under different threshold values.

(c^2, r_1)	λ_1	λ_2	$\delta_1 = \delta_2$
(10, 0.1)	0.333381244	27.66661876	1.152941073
(10, 0.2)	0.644026712	27.35597329	0.786516853
(10, 0.3)	0.934408646	27.06559135	0.451612903
(10, 0.4)	1.206638353	26.79336165	0.144329859
(50, 0.1)	0.057947272	27.94205273	0.225526639
(50, 0.2)	0.115180896	27.8848191	0.166325276
(50, 0.3)	0.171715428	27.82828457	0.108086125
(50, 0.4)	0.227564943	27.77243506	0.05078596

Table 4.2. The various combinations of c^2 and r_1 as well as the parameters of the corresponding MMPP-2 traffic.

Figures 4.10-4.16 demonstrate the analytical results of the utilization, mean numbers of packets in the system and the buffer, throughput, mean response time, queueing delay and packet loss probability against r_1 with the different values of, c^2 and th . It can be observed from all figures below that high burstiness or correlation results in low utilization and throughput but large numbers of packets in the system and buffer, long mean response time and queueing delay as well as high packet loss probability when the other parameter settings remain unchanged, respectively. Moreover, these figures also depict that the effects of the burstiness and correlation on the performance metrics are insensitive to different threshold value. The reason is that more packets can be injected into the system before the buffer becomes full as a result of a high threshold value. Consequently, the effects of high burstiness or correlation on AQM performance are more significant with a high threshold value than those with a low one.

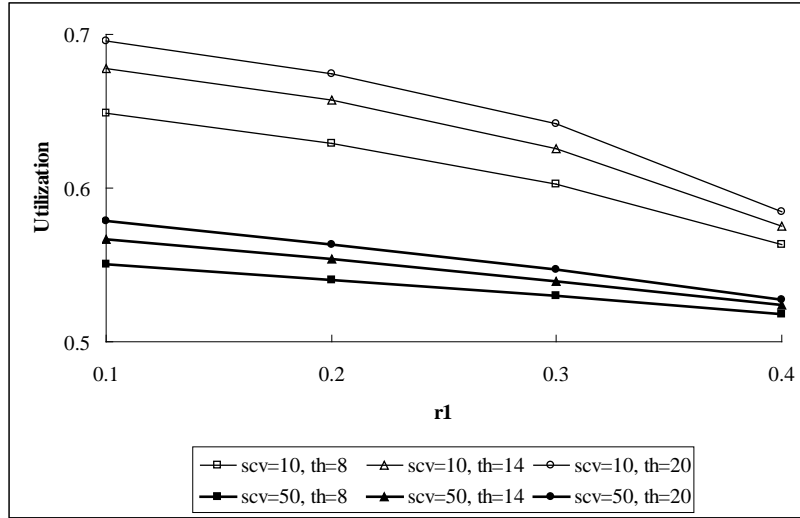


Figure 4.10. Utilization vs r_1 with different values of c^2 and th .

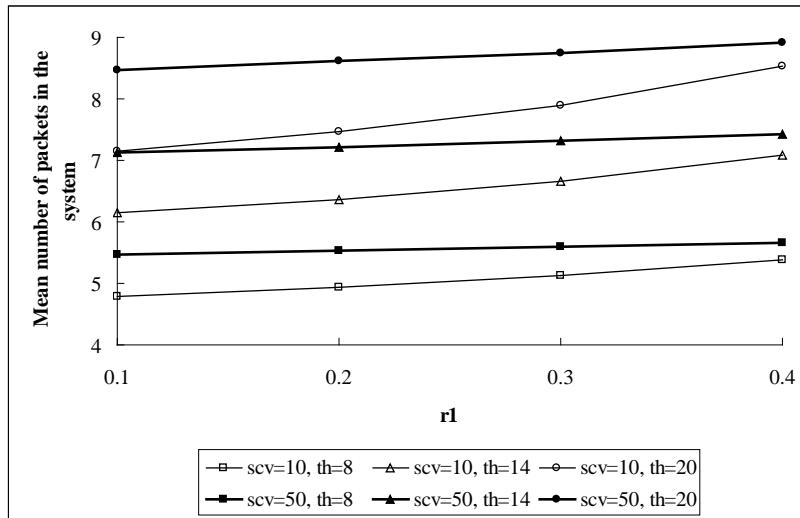


Figure 4.11. Mean number of packets in the system vs r_1 with different values of c^2 and th .

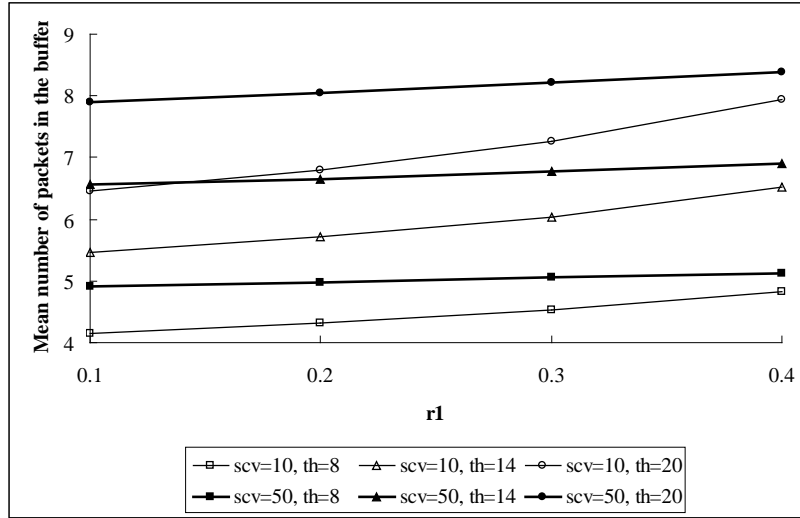


Figure 4.12. Mean number of packets in the buffer vs r_1 with different values of c^2 and th .

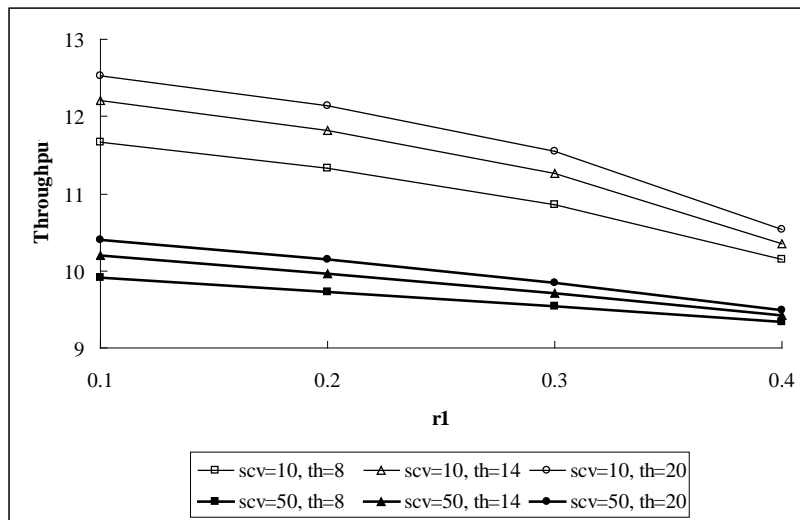


Figure 4.13. Throughput vs r_1 with different values of c^2 and th .

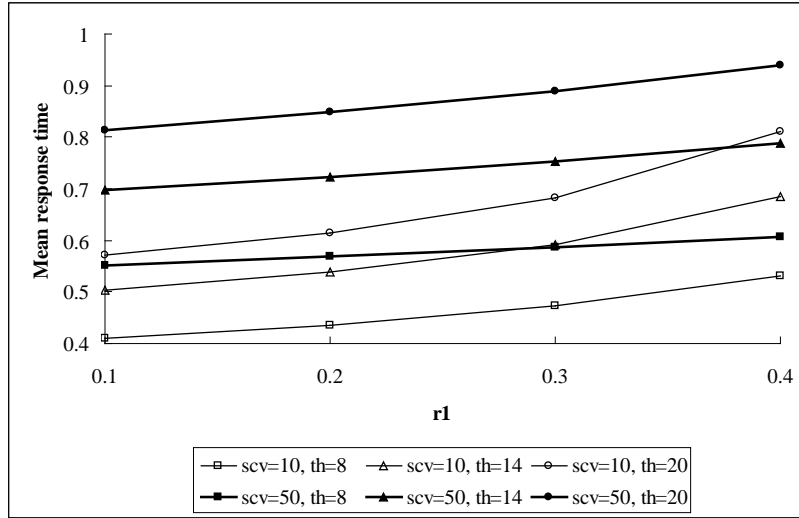


Figure 4.14. Mean response time vs r_1 with different values of c^2 and th .

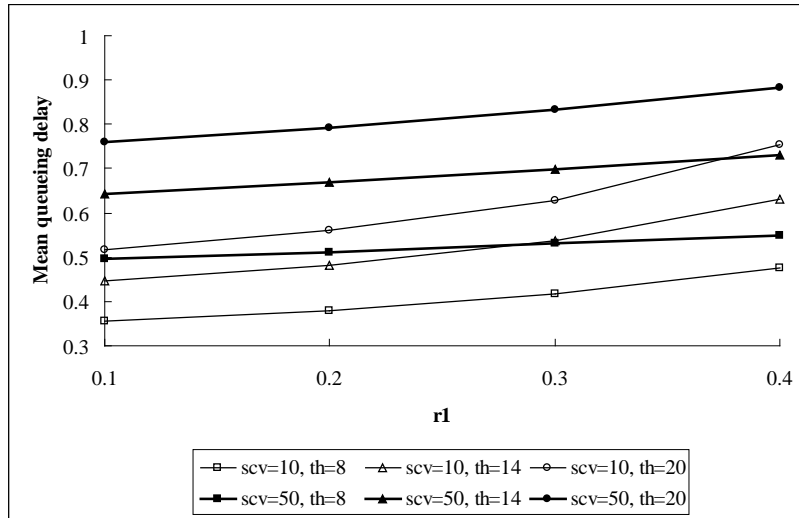


Figure 4.15. Mean queueing delay vs r_1 with different values of c^2 and th .

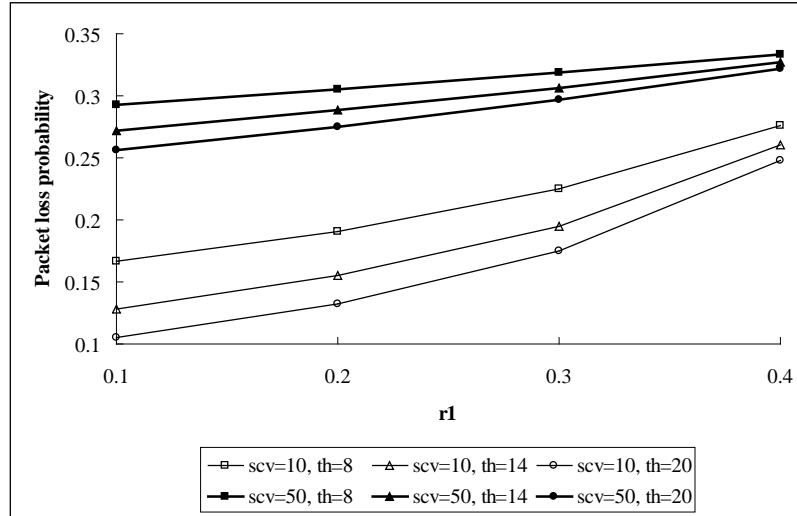


Figure 4.16. Packet loss probability vs r_1 with different values of c^2 and th .

4.4 Summary

This chapter has presented an analytical performance model of AQM system subject to bursty traffic captured by an MMPP-2. To this end, a two-dimensional Markov model has been developed to assist the derivation of several important performance metrics, including the utilization, mean numbers of packets in the system and buffer, throughput, mean response time, mean queueing delay and packet loss probability. The effectiveness and accuracy of the developed model has been demonstrated by comparing analytical results with those obtained from simulators developed in JAVA programming language. To demonstrate its application, the analytical model has been employed to investigate the effects of the burstiness and correlation of the MMPP-2 traffic. Numerical results have demonstrated that high burstiness and correlation can significantly degrade the AQM performance in terms of increasing the mean numbers of packets in the system and buffer, mean response time, mean queueing delay as well as packet loss probability and decreasing

utilization and throughput. Additionally, the effects of burstiness and correlation are also sensitive to the threshold value. More specifically, a low threshold is capable of degrading the negative effects of high burstiness and correlation on the AQM performance.

Chapter 5

Performance Modeling and Analysis of AQM with Heterogeneous Traffic

5.1 Introduction

With the development of the Internet, network applications have ranged from text-based utilities such as electronic mail and news from the early days of the Internet to the deployment of videoconferencing, multimedia streaming, the World-Wide Web, and electronic commerce. Recurrent theme relating to network traffic is the traffic burstiness and correlation exhibited by key services such as compressed video, file transfer, etc. The use of a wide variety of network applications requiring different levels of Quality of Service (QoS) motivates the effort on the study of AQM congestion control mechanism in the presence of heterogeneous traffic.

This chapter proposes a new analytical queueing model for AQM with two classes of traffic that follow respectively two different arrival processes: bursty MMPP and non-bursty Poisson process. Then we present the derivation of expressions for aggregate and marginal performance metrics including utilization, throughput, mean number of packets in the system and buffer, mean response time, queueing delay and packet loss probability. The credibility of the model is demonstrated by comparing the prediction of performance metrics from the developed model with the experimental results obtained from a simulator programmed in JAVA. This chapter highlights the effects of input parameters of bursty

traffic including the average arrival rate, burstiness, correlation and its threshold on the aggregate and marginal utilization, throughput, mean queueing delay and packet loss probability.

5.2 Analytical Model and Performance Measures

This section introduces the system theoretical framework based on AQM scheme and presents the formulation of the analytical model under heterogeneous traffic with individual thresholds.

5.2.1 Proposed Markov Model

We consider two heterogeneous classes of traffic in a single-server finite queueing system for the buffer in an AQM-enabled router with individual buffer thresholds th_b^c for each class c , ($c = 1, 2$). Different from Section 3.3.1, the arrivals of Class-1 traffic generated by a bursty traffic source are fed into a finite buffer in an AQM-enabled router (c.f., Figure 5.1). When the instantaneous queue length exceeds the threshold value, th_b^c , the arriving packets from Class- c , ($c = 1, 2$), may be rejected probabilistically according to a dropping probability aiming to inform the corresponding traffic source of the incipient state of congestion. As shown in Figure 5.2, the dropping probability for each class linearly increases from 0 to 1 as the instantaneous queue length varies from the threshold value to the buffer capacity.

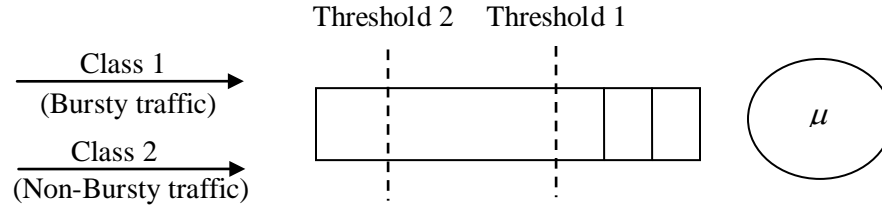


Figure 5.1. A Model of [MMPP-2]&[M]/M/1/K/th1/th2 Queueing System

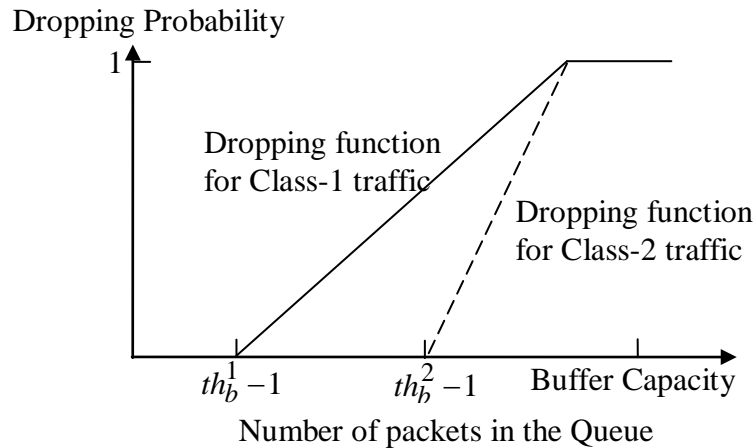


Figure 5.2. Dropping Functions for Two Classes of Traffic

In the proposed model, the first traffic class follows an MMPP-2 process for modeling burstiness (i.e., generated by voice applications) and the second class follows a Poisson process with an average arrival rate λ for modelling non-bursty traffic (i.e., text data). The

MMPP is parameterised by the infinitesimal generator $\mathbf{Q} = \begin{bmatrix} -\delta_1 & \delta_1 \\ \delta_2 & -\delta_2 \end{bmatrix}$ and the rate

matrix $\mathbf{\Lambda} = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}$. The average arrival rate of the MMPP can be calculated by

$\bar{\lambda}^1 = (\lambda_1 \delta_2 + \lambda_2 \delta_1) / (\delta_1 + \delta_2)$. The service time of both classes is exponentially distributed

with mean $1/\mu$. The buffer capacity is denoted as K , consequently the total system capacity is $L = K + 1$.

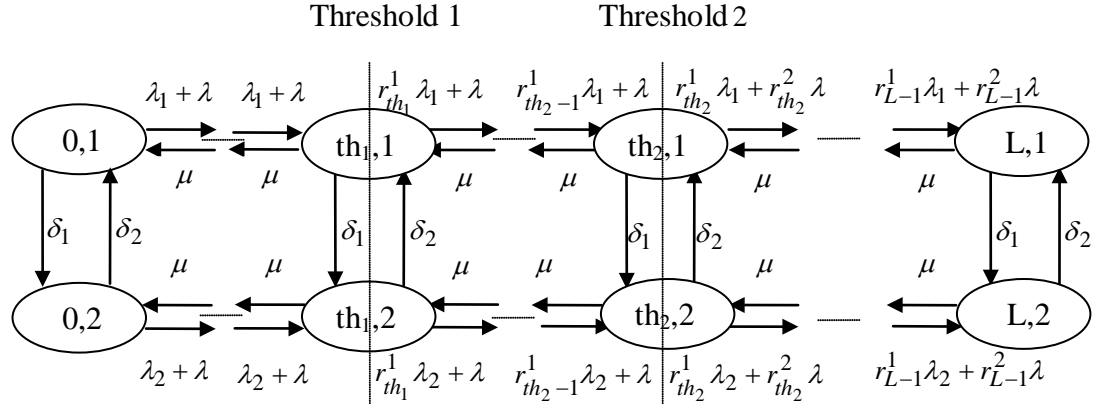


Figure 5.3. A State Transition Rate Diagram of [MMPP-2]&[M]/M/1/K/th Queueing System.

Figure 5.3 shows a state transition diagram of the analytical queueing model of the aforementioned queueing system [MMPP-2][M]/M/1/K/th1/th2. The state (i, j) with $(0 \leq i \leq L, j = 1, 2)$ represents the current number of packets in the system is i and the traffic arrival process MMPP-2 is at State j . The transition rate from State (i, j) to $(i-1, j)$, where $(1 \leq i \leq L, j = 1, 2)$, is μ denoting the mean service rate. δ_j , $(j = 1, 2)$, is the transition rate out of State (i, j) to $(i, 3-j)$, where $(0 \leq i \leq L)$. The transition from State (i, j) to $(i+1, j)$, where $(0 \leq i < L, j = 1, 2)$, represents that a packet is injected into the system. The transition rate out of State (i, j) to $(i+1, j)$, where $(0 \leq i < th_1, j = 1, 2)$, is $\lambda_j + \lambda$ because no packets are dropped. It is noticeable, however, that when the number of Class- c packets in the system exceeds the threshold th_c , i.e., when the system is at state

(i, j) , $(th_c \leq i \leq L, j = 1, 2)$, the probability that the arrivals of Class- c traffic are allowed to enter into the system is $r_i^c = (1 - d_i^c)$ since the packet dropping process can be seen as a decrease of the arriving rate. As a result, the actual arrival rate of Class- c traffic is reduced to $r_i^c \lambda_j$ from λ_j at state j of the underlying MMPP-2 Markov chain. The calculation of the packet dropping probability d_i^c ($0 \leq i \leq L, c = 1, 2$) is given as follows:

$$d_i^c = \begin{cases} 0 & 0 \leq i < th_c \\ \frac{i - th_c + 1}{L - th_c + 1} & th_c \leq i \leq L \end{cases} \quad (5.1)$$

Let p_{ij} , $(0 \leq i \leq L, j = 1, 2)$, represent the steady state probability in the state transition rate diagram. According to the transition equilibrium between in-coming and out-going streams of each state, the following group of equations can be found.

$$\left. \begin{aligned} (d_0^1 \lambda_1 + d_0^2 \lambda + \delta_1) p_{01} &= \delta_2 p_{02} + \mu p_{11} \\ (d_0^1 \lambda_2 + d_0^2 \lambda + \delta_2) p_{02} &= \delta_1 p_{01} + \mu p_{12} \\ (d_i^1 \lambda_1 + d_i^2 \lambda + \delta_1 + \mu) p_{i1} &= (d_{i-1}^1 \lambda_1 + d_{i-1}^2 \lambda) p_{i-1,1} + \delta_2 p_{i2} + \mu p_{i+1,1} \\ (d_i^1 \lambda_2 + d_i^2 \lambda + \delta_2 + \mu) p_{i2} &= (d_{i-1}^1 \lambda_2 + d_{i-1}^2 \lambda) p_{i-1,2} + \delta_1 p_{i1} + \mu p_{i+1,2} \\ (\delta_1 + \mu) p_{L1} &= (d_{L-1}^1 \lambda_1 + d_{L-1}^2 \lambda) p_{L-1,1} + \delta_2 p_{L2} \\ (\delta_2 + \mu) p_{L2} &= (d_{L-1}^1 \lambda_2 + d_{L-1}^2 \lambda) p_{L-1,2} + \delta_1 p_{L1} \end{aligned} \right\} \begin{array}{l} 1 \leq i < L \\ 1 \leq i < L \end{array} \quad (5.2)$$

$$\sum_{i=0}^L \sum_{j=1}^2 p_{ij} = 1 \quad (5.3)$$

Solving these equations, we can find probability p_{ij} of each state in the Markovian model as follows:

$$p_{ij} = a_{ij}m + b_{ij}n \quad (5.4)$$

where

$$a_{ij} = \begin{cases} 1 & i = 0, j = 1 \\ 0 & i = 0, j = 2 \\ \frac{d_0^1 \lambda_1 + d_0^2 \lambda + \delta_1}{\mu} & i = 1, j = 1 \\ -\frac{\delta_1}{\mu} & i = 1, j = 2 \\ \frac{(d_{i-1}^1 \lambda_1 + d_{i-1}^2 \lambda + \mu + \delta_1)a_{i-1,1} - (d_{i-2}^1 \lambda_1 + d_{i-2}^2 \lambda)a_{i-2,1} - \delta_2 a_{i-1,2}}{\mu} & 2 \leq i \leq L, j = 1 \\ \frac{(d_{i-1}^1 \lambda_2 + d_{i-1}^2 \lambda + \mu + \delta_2)a_{i-1,2} - (d_{i-2}^1 \lambda_2 + d_{i-2}^2 \lambda)a_{i-2,2} - \delta_1 a_{i-1,1}}{\mu} & 2 \leq i \leq L, j = 2 \end{cases} \quad (5.5)$$

and

$$b_{ij} = \begin{cases} 0 & i = 0, j = 1 \\ 1 & i = 0, j = 2 \\ -\frac{\delta_1}{\mu} & i = 1, j = 1 \\ \frac{d_0^1 \lambda_1 + d_0^2 \lambda + \delta_1}{\mu} & i = 1, j = 2 \\ \frac{(d_{i-1}^1 \lambda_1 + d_{i-1}^2 \lambda + \mu + \delta_1)b_{i-1,1} - (d_{i-2}^1 \lambda_1 + d_{i-2}^2 \lambda)b_{i-2,1} - \delta_2 b_{i-1,2}}{\mu} & 2 \leq i \leq L, j = 1 \\ \frac{(d_{i-1}^1 \lambda_2 + d_{i-1}^2 \lambda + \mu + \delta_2)b_{i-1,2} - (d_{i-2}^1 \lambda_2 + d_{i-2}^2 \lambda)b_{i-2,2} - \delta_1 b_{i-1,1}}{\mu} & 2 \leq i \leq L, j = 2 \end{cases} \quad (5.6)$$

$$k = (d_{L-1}^1 \lambda_1 + d_{L-1}^2 \lambda)b_{L-1,1} + \delta_2 b_{L2} - (\mu + \delta_1)b_{L1} \quad (5.7)$$

$$h = (\mu + \delta_1)a_{L1} - (d_{L-1}^1 \lambda_1 + d_{L-1}^2 \lambda)a_{L-1,1} - \delta_2 a_{L2} \quad (5.8)$$

$$n = \frac{h}{k \sum_{i=0}^{L-1} \sum_{j=1}^2 a_{ij} + h \sum_{i=0}^{L-1} \sum_{j=1}^2 b_{ij}} \quad (5.9)$$

$$m = \frac{k}{h} n \quad (5.10)$$

5.2.2 Performance Measures

In what follows, we will derive the aggregate system performance metrics including utilization (U), mean number of packets in the system (\bar{L}), mean number of packets in the buffer (\bar{L}_b), system throughput (T), mean response time (\bar{R}), mean queueing delay (\bar{D}), probability of packet loss (PLP) and relevant marginal performance metrics for each class of traffic.

- ♦ *The aggregate performance measures*

Following the derivation of expressions of the aggregate utilization, mean numbers of packets in the system as well as in the buffer, throughput, mean response time and mean queueing delay presented in Section 4.3, with the derived probability, p_{ij} , of each State (i, j) in the Markov model, the aggregate performance metrics are given by

$$\rho = 1 - p_{01} - p_{02} \quad (5.11)$$

$$\bar{L} = \sum_{i=0}^{L-1} \sum_{j=1}^2 (p_{ij} \times i) \quad (5.12)$$

$$\bar{L}_b = \sum_{i=1}^{L-1} \sum_{j=1}^2 (p_{i+1,j} \times i) \quad (5.13)$$

$$\bar{T} = \rho \times \mu \quad (5.14)$$

$$\bar{R} = \frac{\bar{L}}{T} \quad (5.15)$$

$$\bar{D} = \frac{\bar{L}_b}{T} \quad (5.16)$$

The aggregate packet loss rate is equal to the average aggregate arrival rate $\bar{\lambda}^1 + \bar{\lambda}^2$, (where $\bar{\lambda}^1 = (\lambda_1\delta_2 + \lambda_2\delta_1)/(\delta_2 + \delta_1)$, $\bar{\lambda}^2 = \lambda$), minus the throughput. So the packet loss probability is given below:

$$PLP = \frac{(\bar{\lambda}^1 + \bar{\lambda}^2 - T)}{\bar{\lambda}^1 + \bar{\lambda}^2} \quad (5.17)$$

♦ *The marginal performance measures*

In order to derive the marginal mean numbers of packets in the system and buffer, we calculate the marginal steady-state probability for each class c , ($c=1, 2$). Similar to Section 3.3.1, L positions in the system are numbered with $1 \dots L$ from server to the tail of the buffer. When the current instantaneous queue length including any in the server is i , ($0 \leq i < L$), the probabilities that a new arriving packet from Class-1 and Class-2 enters into the system and, of course, is allocated in the position $i+1$ are $p_i \times \frac{r_i^1 \times (\lambda_1 + \lambda_2)}{\mu}$ and, $p_i \times \frac{r_i^1 \times \lambda}{\mu}$, respectively. In addition, if a packet from class- c ($c=1,2$) is allocated in the position j ($i < j \leq L$) upon arrival in the system, it will experience the position $i+1$. In

other words, the probability that there is a packet from Class- c in the state $i+1$ should be the sum of all the probabilities that the packet arrives in the system and is allocated at position j , $1 \leq i \leq j \leq L$. So, the probabilities that position i , ($1 \leq i \leq L$), in the system is occupied by a packet from Class-1 and Class-2, noted as m_i^1 and m_i^2 , can be derived as follows, respectively:

$$m_i^1 = \frac{\sum_{j=i-1}^{L-1} (p_j \times \frac{r_j^1 \times (\lambda_1 + \lambda_2)}{\mu})}{\sum_{j=i-1}^{L-1} p_j} \quad 1 \leq i \leq L \quad (5.18)$$

$$m_i^2 = \frac{\sum_{j=i-1}^{L-1} (p_j \times \frac{r_j^2 \times \lambda}{\mu})}{\sum_{j=i-1}^{L-1} p_j} \quad 1 \leq i \leq L \quad (5.19)$$

Based on m_i^c ($c=1, 2$), the probabilities, p_i^c and $p_{b_i}^c$, that there are i packets from Class- c in the system and in the buffer can be derived, respectively, according to the method used in Section 3.3.1. So the mean number of packets from Class- c ($c=1, 2$) in the system and in the buffer can be calculated and simplified as:

$$\overline{L^c} = \sum_{i=0}^L (i \times p_i^c) = \begin{cases} \sum_{i=0}^{L-1} (i+1) \times \sum_{j=1}^2 p_{ij} \times r_i^1 \lambda_j & c=1 \\ \sum_{i=0}^{L-1} (i+1) \times \sum_{j=1}^2 p_{ij} \times r_i^1 \lambda & c=2 \end{cases} \quad (5.20)$$

$$\overline{L_b^c} = \sum_{i=0}^L (i \times p_{b_i^c}) = \begin{cases} \sum_{i=0}^{L-1} i \times \sum_{j=1}^2 p_{ij} \times r_i^1 \lambda_j & c = 1 \\ \sum_{i=0}^{L-1} (i+1) \times \sum_{j=1}^2 p_{ij} \times r_i^1 \lambda_j & c = 2 \end{cases} \quad (5.21)$$

Let us now move onto the derivation of other marginal performance metrics. Throughput is commonly defined as the average transition rate at which packets go through the system in the steady state. For a system in the steady state, the average arrival rate equals to its throughput. So the marginal throughput of each class can be expressed as:

$$T^c = \begin{cases} \sum_{i=0}^{L-1} d_i^1 (p_{i1} \lambda_1 + p_{i2} \lambda_2) & c = 1 \\ \sum_{i=0}^{L-1} d_i^2 (p_{i1} + p_{i2}) \lambda & c = 2 \end{cases} \quad (5.22)$$

Again, expressions for the marginal mean response time and mean queueing delay can be derived using Little's Law [1], respectively.

$$\overline{R^c} = \frac{\overline{L^c}}{T^c} \quad (5.23)$$

$$\overline{D^c} = \frac{\overline{L_b^c}}{T^c} \quad (5.24)$$

The ratio of the current average packet loss rate of Class- c to its original average arrival rate is the instant probability of packet loss for Class- c . The average packet loss rate can be approximated using the difference between the average arrival rate $\overline{\lambda^c}$ and the marginal throughput for Class- c . So the marginal packet loss probability can be solved.

$$PLP = \frac{\overline{\lambda^c - T^c}}{\lambda^c} \quad (5.25)$$

5.3 Model Validation

With the aim of validating the accuracy of the analytical model, a discrete-event simulator has been developed using JAVA programming. Numerous validation experiments have been carried out for different combinations of the system capacity, mean service rate, mean arrival rate of Class-2 non-bursty traffic and MMPP-2 input traffic. Specifically, this section presents the results of the derived aggregate and marginal performance metrics for the following five different scenarios described in Table 5.1. The performance results obtained from simulation experiments are illustrated and compared to analytical results in Figures 5.4-5.10. The comparison shows close consistence between simulation results with those obtained from the analytical model. This observation illustrates that the proposed models are very accurate in calculating various performance metrics and examining the performance of AQM under heterogeneous traffic.

	λ_1	λ_2	δ_1	δ_2	λ	μ	L	th_1	th_2
S-5.3.I	1	1	0.3	0.4	1	3	6	2	3
S-5.3.II	5	12	0.6	0.19	20	35	13	5	9
S-5.3.III	9	3	0.2	0.9	13	28	22	8	10
S-5.3.IV	6	28	0.27	0.7	6	23	37	11	24
S-5.3.V	17	4	0.84	0.56	2	16	54	13	42

Table 5.1. The parameter settings corresponding to five scenarios.

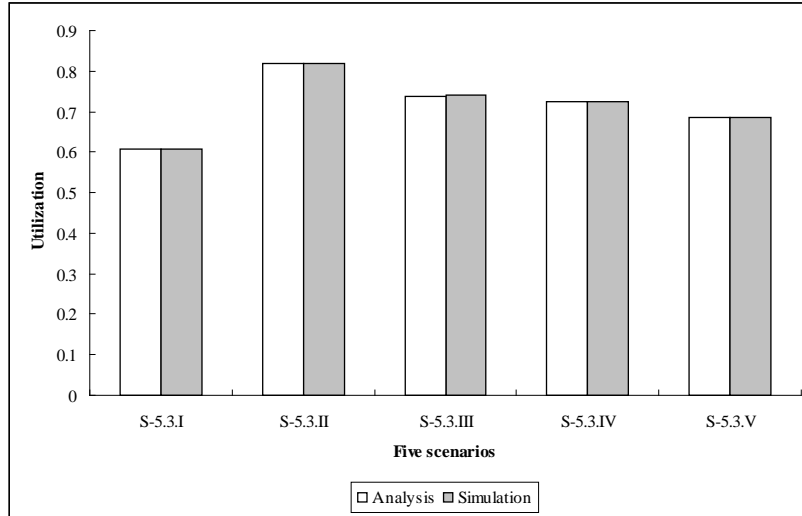


Figure 5.4. Aggregate utilization vs five different scenarios

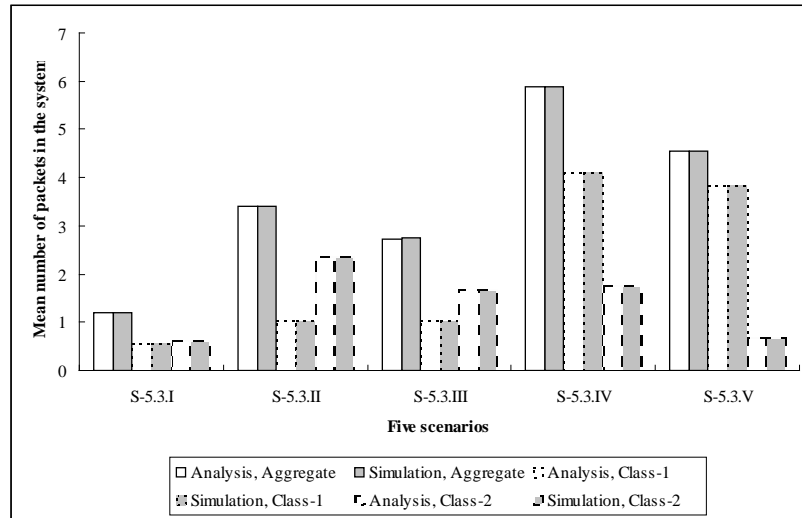


Figure 5.5. Mean number of packets in the system vs five different scenarios

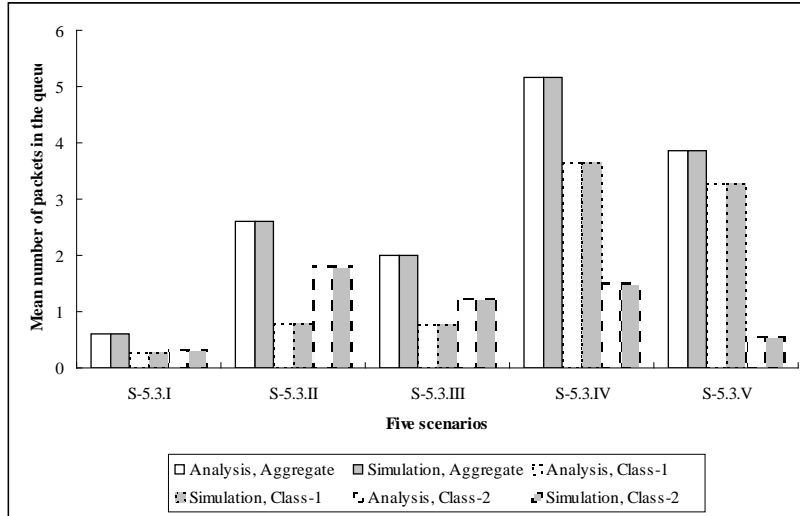


Figure 5.6. Mean number of packets in the queue vs five different scenarios

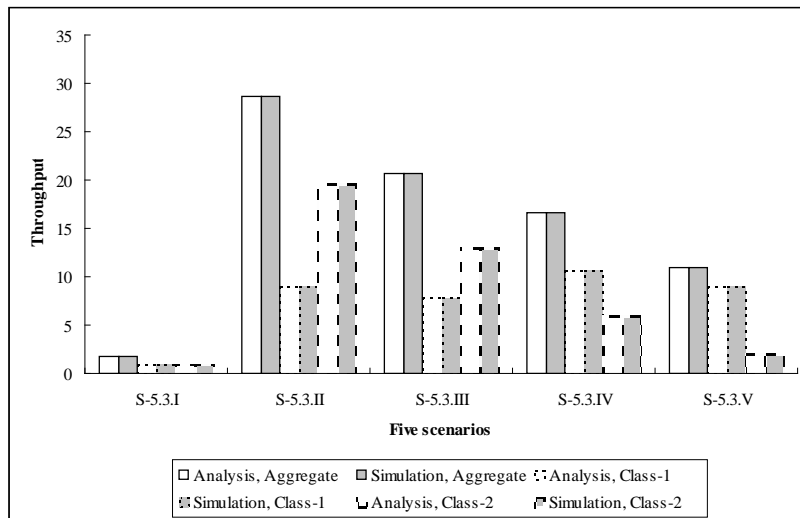


Figure 5.7. Throughput vs five different scenarios

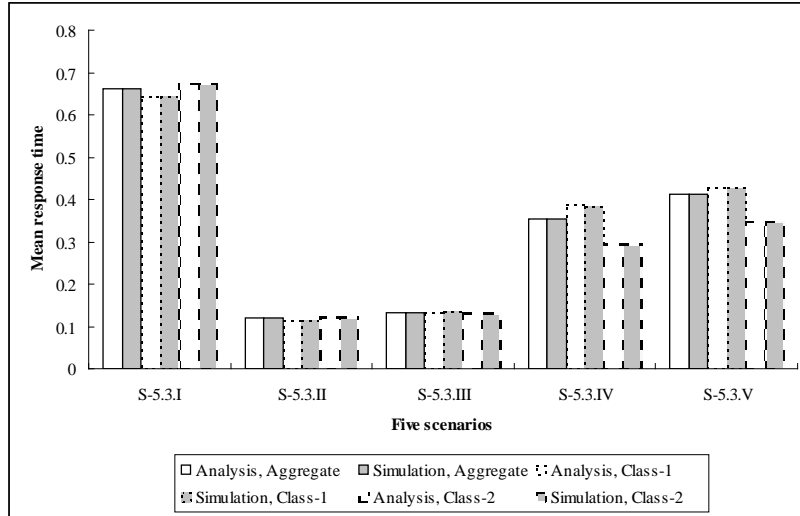


Figure 5.8. Mean response time vs five different scenarios

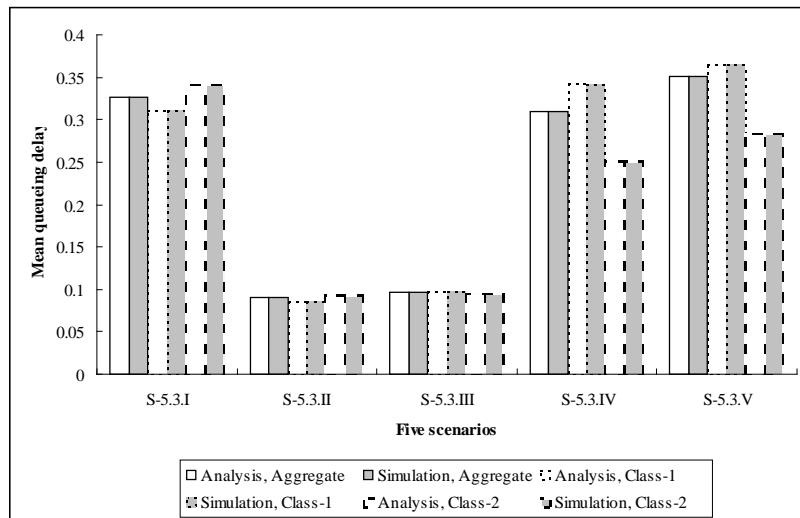


Figure 5.9. Mean queueing delay vs five different scenarios

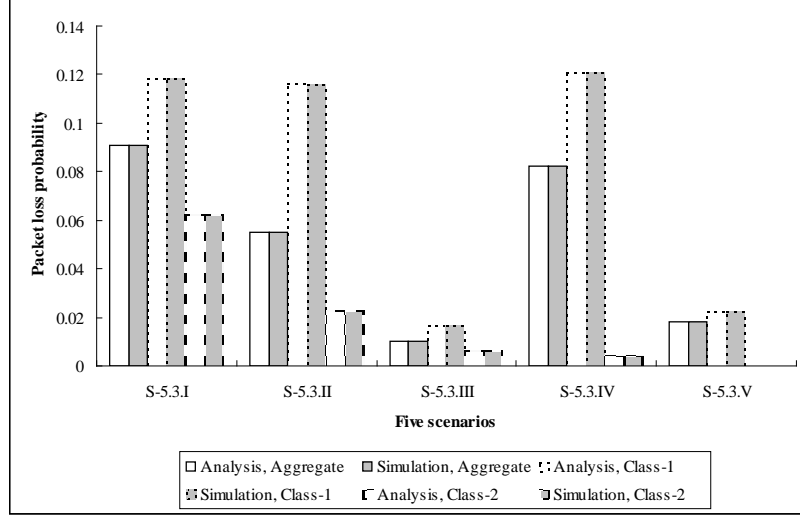


Figure 5.10. Packet loss probability vs five different scenarios

5.4 Performance Analysis

This section will use the above derived analytical model to investigate the aggregate and marginal performance of AQM, including utilization, throughput, mean queueing delay and packet loss probability, under varying threshold value, as well as working under different degrees of burstiness, correlation and rate of Class-1 traffic. In the following analysis, the value of the system capacity L is set to be 20. To focus on the effects of one threshold (i.e., th_1), the other threshold (i.e., th_2) is set to the buffer capacity in this experiment. Meanwhile the value of the first threshold th_1 varies from 1 to 19. Additionally, different values of the SCV and 1-step autocorrelation coefficient of MMPP-2 traffic based on Eqs. 4.7-4.8 in Section 4.2 are set to be $c^2 = (5, 100)$ and $r_1 = (0.01, 0.35)$, respectively. We can obtain the parameters of the MMPP-2 traffic by assuming $\delta_1 = \delta_2$ and keeping the average arrival rate of the MMPP-2 traffic, $\overline{\lambda^1}$, constant at 5 and 10, respectively. Table 5.2 below

presents the various combinations of $\bar{\lambda}^1$, c^2 and r_1 as well as the parameters of the corresponding MMPP-2 traffic. The arrivals of Class-2 traffic follow a Poisson process with a mean arrival rate $\lambda = 7.5$. In order to make sure that the queueing system works under a stable state, the mean service rate μ is set to 20 (i.e., to ensure that μ is larger than the aggregate average traffic arrival rate $\bar{\lambda}^1 + \lambda$).

$\bar{\lambda}^1$	(c^2, r_1)	λ_1	λ_2	$\delta_1 = \delta_2$
5	(5, 0.01)	0.03096005	9.96903995	1.203703704
5	(5, 0.35)	0.829711719	9.170288281	0.108695652
5	(100, 0.01)	0.0619201	19.9380799	2.407407407
5	(100, 0.35)	1.659423438	18.34057656	0.217391304
10	(5, 0.01)	0.001019992	9.998980008	0.049464559
10	(5, 0.35)	0.035332565	9.964667435	0.014586058
10	(100, 0.01)	0.002039984	19.99796002	0.098929118
10	(100, 0.35)	0.070665129	19.92933487	0.029172115

Table 5.2. The various combinations of $\bar{\lambda}^1$, c^2 and r_1 as well as the parameters of the corresponding MMPP-2 traffic.

Firstly, let us analyze the effects of varying th_1 , $\bar{\lambda}^1$, c^2 and r_1 on the aggregate performance metrics aforementioned, respectively. Figures 5.11-5.14 demonstrate the utilization, throughput, mean queueing delay and packet loss probability versus threshold th_1 with different combinations of values of c^2 , r_1 and $\bar{\lambda}^1$, respectively. It is clear that as the rate of Class-1 traffic increases, more Class-1 packets are generated for competing for the system resources (i.e., the server and the buffer). As a result, the utilization, throughput, mean queueing delay and loss probability are larger when $\bar{\lambda}^1 = 10$ than those when $\bar{\lambda}^1 = 5$. Moreover, these figures depict that high burstiness and correlation of Class-1 traffic results in the low utilization and throughput but long mean queueing delay and large packet loss

probability when the other parameter settings remain unchanged. Such trends become more remarkable when the average arrival rate of Class-1 traffic is higher.

In addition, the effects of threshold th_1 on these performance metrics are less intuitive as the variation of the threshold affects not only the average arrival rate but also the burstiness and correlation of Class-1 traffic. Because a large threshold th_1 gives rise to a high average arrival rate, burstiness and correlation of Class-1 traffic, it can be easily observed, based on above analysis, that the aggregate mean queueing delay increases with the growth of threshold th_1 as shown in Figure 5.13. Although both a large average arrival rate and a high burstiness/correlation of Class-1 traffic result in more packets to be lost due to buffer overflow, the increase of threshold th_1 reduces the aggregate packet loss probability. The reason is that the number of dropped packets decreases greatly as the threshold increases. However, because the increases of average arrival rate and burstiness/correlation of Class-1 traffic have the opposite impacts on the aggregate utilization and throughput, the effects of threshold th_1 depend on which could play a significant role. Specifically, Figures 5.11 and 5.12 illustrates that the aggregate utilization and throughput increase as threshold th_1 rises, which states that the impact of the corresponding variation of the average arrival rate of Class-1 traffic plays a dominate role. Apart from these fundamental issues, Figures 5.11-5.14 also exhibit different impact strengthes of threshold th_1 combined with various values of mean arrival rate, burstiness and correlation of Class-1 traffic. For instance, when the burstiness of Class-1 traffic is low (i.e., $c^2 = 5$), the aggregate utilization, throughput and packet loss probability with a higher

average arrival rate of Class-1 traffic (i.e., $\bar{\lambda}^1 = 10$) increase more sharply as the threshold rises than those with a low traffic rate (i.e., $\bar{\lambda}^1 = 5$), meanwhile, the effect of the threshold on the mean queueing delay when $\bar{\lambda}^1 = 10$ is more remarkable than those when $\bar{\lambda}^1 = 5$ regardless of burstiness.

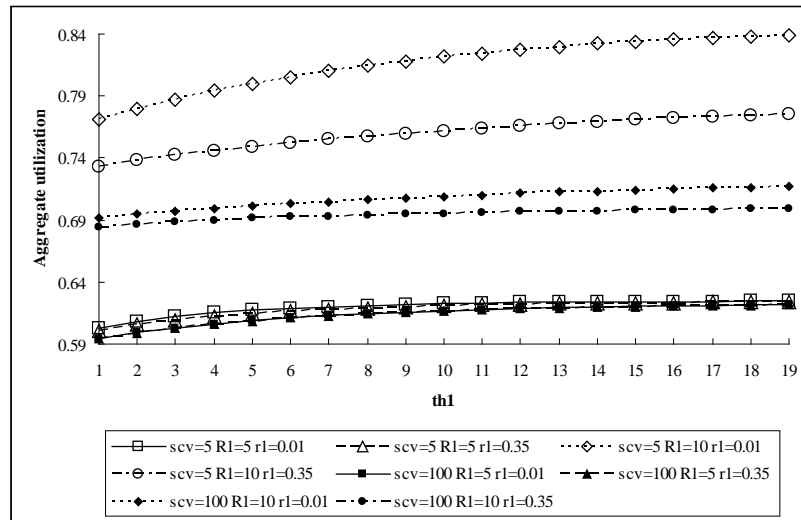


Figure 5.11. Aggregate utilization vs th_1 with different values of c^2 , r_1 and mean arrival rate of Class-1 traffic.

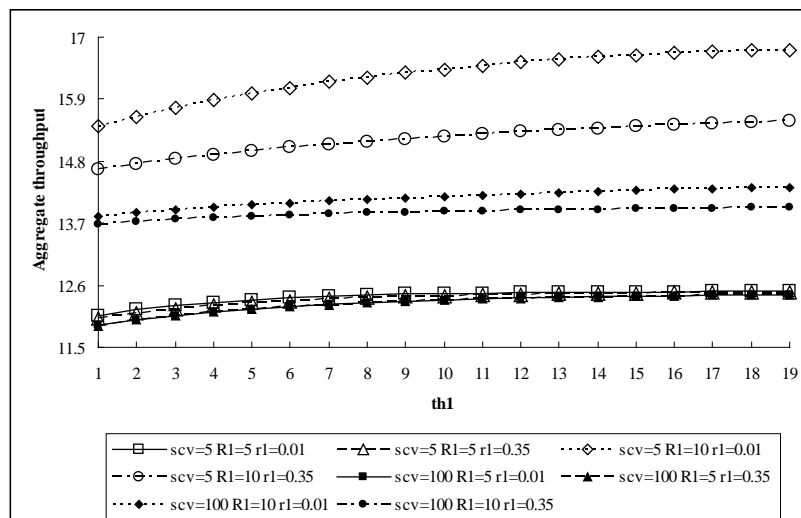


Figure 5.12. Aggregate throughput vs th_1 with different values of c^2 , r_1 and mean arrival rate of Class-1 traffic.

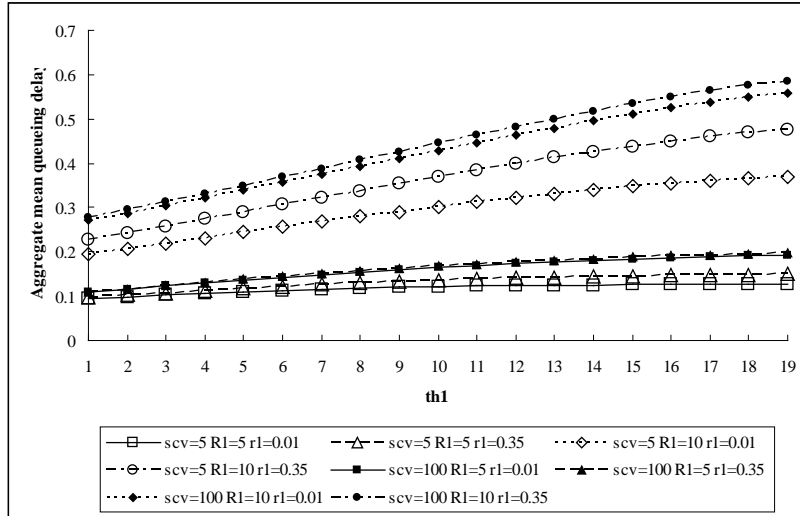


Figure 5.13. Aggregate mean queueing delay vs th_1 with different values of c^2 , r_1 and mean arrival rate of Class-1 traffic.

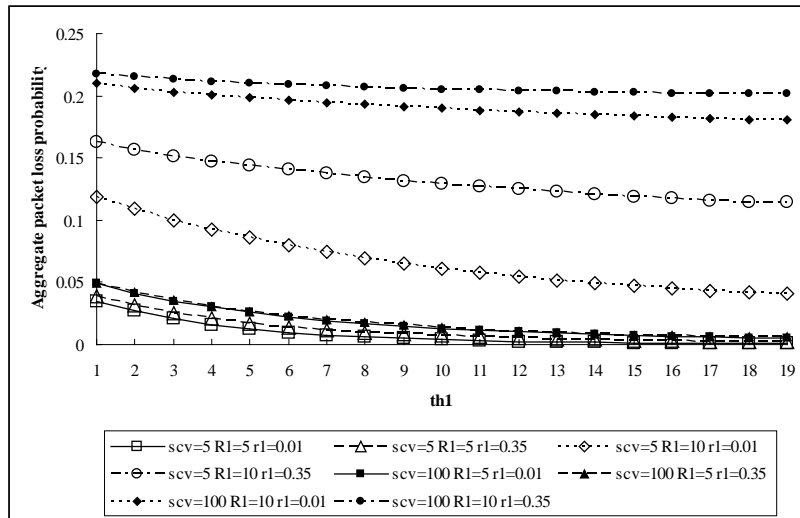


Figure 5.14. Aggregate packet loss probability vs th_1 with different values of c^2 , r_1 and mean arrival rate of Class-1 traffic.

Secondly, the variation of the marginal performance metrics for Class-1 traffic with different values of th_1 , $\overline{\lambda^1}$, c^2 and r_1 are evaluated as follows. The marginal throughput, mean queueing delay and packet loss probability for Class-1 traffic have been shown in Figures 5.15-5.17, respectively. We can find from these figures that a large average arrival rate of Class-1 traffic results in high marginal throughput, packet loss probability and long mean queueing delay for Class-1. Also, these figures illustrate that increasing the degrees of burstiness and correlation could reduce the marginal throughput for Class-1 traffic, lengthen the marginal mean queueing delay and raise the packet loss probability for Class-1 traffic. Moreover, the effects of threshold th_1 on these marginal performance metrics for Class-1 traffic are also quite similar to those on the aggregate performance metrics. The marginal throughput and mean queueing delay for Class-1 increases but the marginal packet loss probability for Class-1 decreases as a result of the growth of the threshold. In particular, if the average arrival rate of Class-1 traffic is high, the marginal throughput and packet loss probability vary more sharply when threshold th_1 is greater than 16. Besides, the effect of the threshold on the marginal mean queueing delay becomes more significant when Class-1 traffic is generated at a higher average arrival rate or with a stronger burstiness/correlation. On the other hand, Figures 5.15 and 5.17 reveal respectively that the variation of the marginal throughput and packet loss probability for Class-1 as the threshold changes are more remarkable when the average arrival rate is high (i.e., $\overline{\lambda^1} = 10$). However, if the average arrival rate of Class-1 traffic is low, the variation of these two marginal performance metrics resulted from varying threshold th_1 when the traffic burstiness is high

(i.e., $c^2 = 100$) are more remarkable than those when the traffic burstiness is low (i.e., $c^2 = 5$). Otherwise, these variations when c^2 is low are more remarkable than those when c^2 is high.

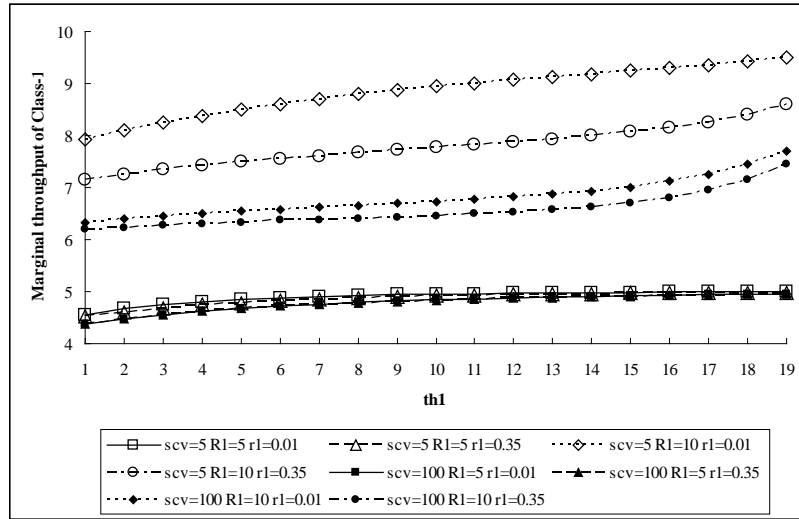


Figure 5.15. Marginal throughput of Class-1 vs th_1 with different values of c^2 , r_1 and mean arrival rate of Class-1 traffic.

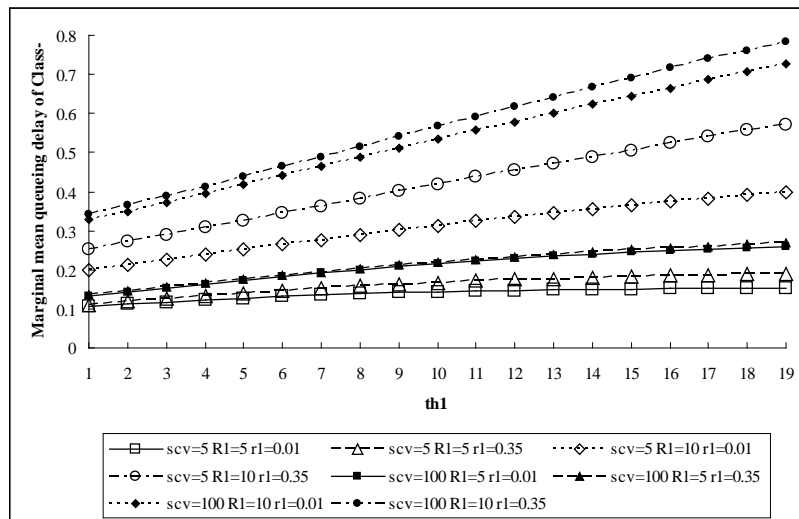


Figure 5.16. Marginal mean queuing delay of Class-1 vs th_1 with different values of c^2 , r_1 and mean arrival rate of Class-1 traffic.

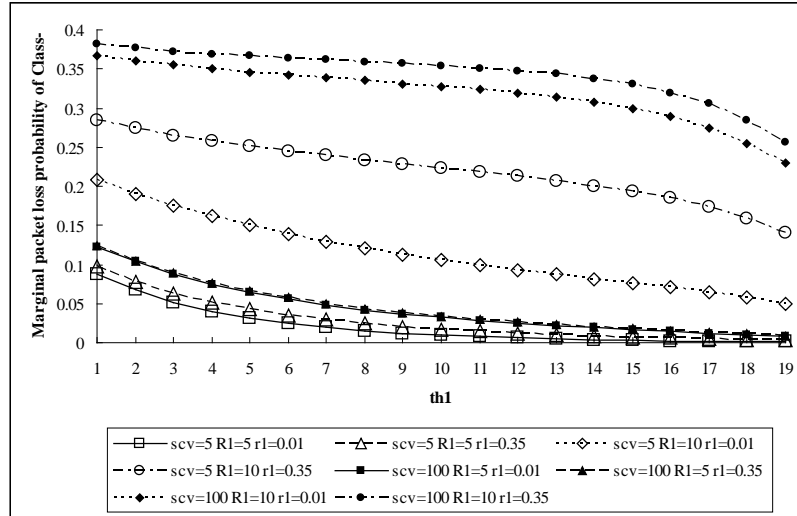


Figure 5.17. Marginal packet loss probability of Class-1 vs th_1 with different values of c^2 , r_1 and mean arrival rate of Class-1 traffic.

Finally, the last three figures present the marginal throughput, mean queueing delay and packet loss probability against varying threshold th_1 combined with different burstiness, correlation and arrival rate of Class-1 traffic. As the average arrival rate of Class-1 traffic increases, more Class-1 packets compete with Class-2 packets for the system resources. Consequently, all marginal performance metrics for Class-2 are seriously degraded. Moreover, it can be intuitively observed from Figures 5.18-5.20 that the high traffic burstiness and correlation give rise to a low marginal throughput, long mean queueing delay and large packet loss probability of Class-2. Again, high burstiness affects these marginal performance metrics of Class-2 more remarkably when the traffic correlation is low, and vice versa. A high traffic rate $\bar{\lambda}^1$ is capable of further strengthening the impact of traffic burstiness and correlation on these marginal performance metrics of Class-2 traffic. Additionally, Figures 5.18 and 5.20 indicate that the marginal throughput

decreases and packet loss probability increases accompanying the growth of threshold th_1 , respectively. Also, these figures illustrate the significant effects of traffic rate, burstiness and correlation on the two performance metrics when the threshold is big. Figure 5.19 depicts that the marginal mean queueing delay of Class-2 with a low average arrival rate of Class-1 $\bar{\lambda}^1$ increases as the threshold rises, while the delay with a high $\bar{\lambda}^1$ increases firstly and then decreases as the threshold grows.

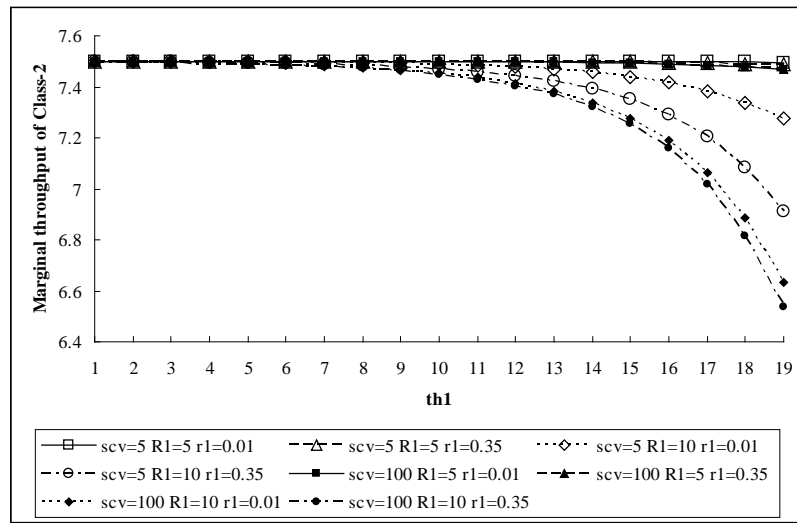


Figure 5.18. Marginal throughput of Class-2 vs th_1 with different values of c^2 , r_1 and mean arrival rate of Class-1 traffic.

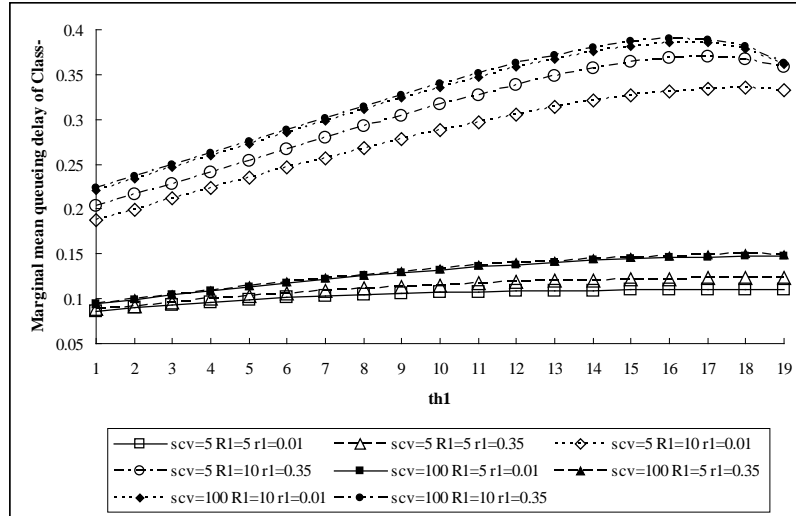


Figure 5.19. Marginal mean queuing delay of Class-2 vs th_1 with different values of c^2 , r_1 and mean arrival rate of Class-1 traffic.

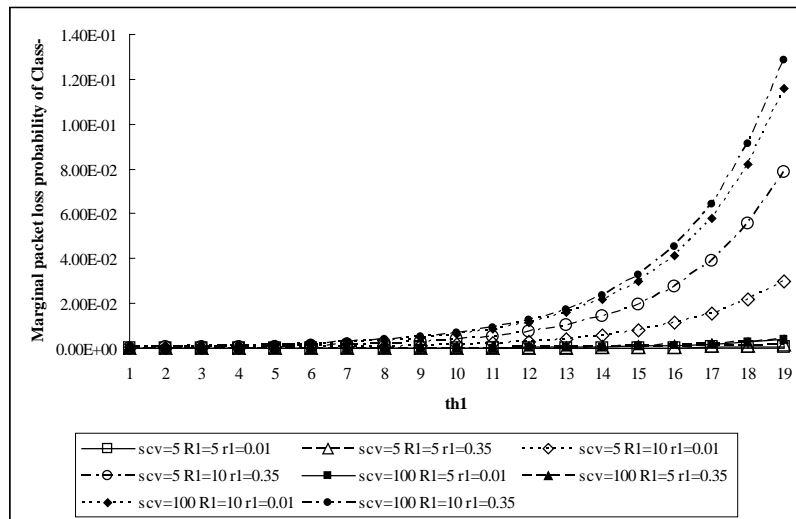


Figure 5.20. Marginal packet loss probability of Class-2 vs th_1 with different values of c^2 , r_1 and mean arrival rate of Class-1 traffic.

5.5 Summary

In this chapter, an analytical model has been developed and validated for evaluating the performance measures of a queueing system with two individual thresholds subject to heterogeneous traffic in an AQM-enabled router. This study has modeled the data and voice sources, respectively, by non-bursty Poisson Process and bursty MMPP Process. We have derived the expressions for the aggregate and marginal performance metrics including the utilization, throughput, mean numbers of packets in the system and buffer, mean response time, mean queueing delay and packet loss probability. The accuracy of the model in examining the performance of the AQM mechanism under heterogeneous traffic and calculating various performance metrics has been demonstrated by comparing analytical against simulation results. The model has been adopted to evaluate the impact of parameters related to Class-1 traffic, including the average arrival rate, burstiness, correlation and its threshold, on the aggregate and marginal utilization, throughput, mean queueing delay and packet loss probability.

Analytical results have shown that, all aggregate and marginal performance metrics change significantly as the average arrival rate of Class-1 traffic varies. For instance, as the traffic rate increases, the marginal performance metrics for Class-2 are degraded substantially. While a high traffic rate is capable of decreasing aggregate packet loss probability and the marginal one for Class-1 as well as increasing all other performance metrics. Moreover, the detrimental effects of traffic burstiness and correlation on all performance metrics have been clearly observed and reported. Finally, we analyze the uncertainty effects of the threshold assigned to Class-1 traffic on all performance metrics.

To this end, the performance evaluation of the proposed analytical model aids to find the best threshold settings and drop probability to suit a given situation; i.e., to give an appropriate trade-off between delay and packet loss probability. So settings of these parameters thus can be chosen to suit the type of service required. For example, real-time services like voice require low delay, while data services require low packet loss.

Chapter 6

Performance Modeling and Analysis of Priority-Based AQM with Heterogeneous Traffic

6.1 Introduction

AQM coupled with differentiated scheduling mechanism is an important and promising scheme for congestion control and QoS guarantee in communication networks. It is known that the non-responsive flows with which AQM is unable to cope can be regulated by effective scheduling schemes which decides on the sending order of packets in order to satisfy the QoS requirements, such as latency and fairness. The models and techniques of packet scheduling for differentiated services have been extensively studied [88-89]. Few research efforts [90-91] were devoted to investigating performance of AQM with scheduling scheme through simulation results. Due to the lack of comprehensive and analytical performance evaluation of AQM congestion control with scheduling scheme, this study develops an analytical model to investigate the performance of AQM with Pre-emptive Resume (PR) scheduling scheme in the presence of heterogeneous network traffic.

This chapter evaluates the performance for AQM with PR scheme subject to heterogeneous bursty traffic modelled by bursty MMPP and non-bursty Poisson traffic. Individual thresholds for each traffic class and PR priority scheduling mechanism are adopted to control traffic injection rate and support differential QoS. Two analytical models are proposed for the priority-based AQM system with single queue and multiple class-

based queues, respectively. This chapter presents the derivations and the expressions of the key aggregate and marginal performance metrics for these two systems. The credibility of the model is demonstrated by comparing the analytic results with those obtained through extensive simulation experiments. The first model to be developed is used to compare with that developed in Chapter 5 with the aim of investigating the effects of PR scheduling scheme. Then, the second model is adopted to analyse the AQM performance with multiple priority classes-based queues.

6.2 Priority-based AQM with Single Buffer

6.2.1 System Description

Figure 6.1 shows the model of a single-server preemptive priority queue with AQM scheme under two classes of traffic. Different from Chapter 5, two different traffic priority classes are considered here: Class-1 traffic has the higher priority than Class-2 traffic. The PR priority scheduling scheme enables an arriving Class-1 packets to pre-empt the Class-2 message which is currently occupying the server. The pre-empted Class-2 packet resumes its processing soon after the service of the Class-1 packets is complete. AQM scheme sets a single fixed threshold for each traffic class in order to inform the corresponding traffic source of the incipient state of congestion. That is, the system may reject an arriving packet from Class- c ($c = 1, 2$) according to a dropping probability if the current queue length exceeds the threshold value th_b^c . The dropping probability increases linearly from 0 to 1 as the queue length increases from $th_b^c - 1$ to the buffer capacity.

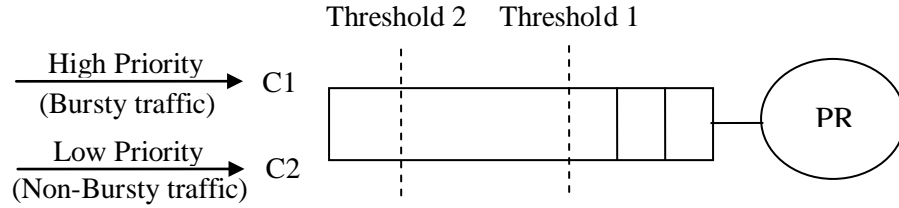


Figure 6.1. A Model of [MMPP-2][M]/M/1/K/th₁/th₂ Queueing System with PR scheme

6.2.2 Analytical Model

The arrivals of Class-1 traffic (bursty traffic) follow an MMPP-2 with the infinitesimal

generator matrix $\mathbf{Q} = \begin{bmatrix} -\delta_1 & \delta_1 \\ \delta_2 & -\delta_2 \end{bmatrix}$ and the rate matrix $\mathbf{\Lambda} = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}$. The arrivals of Class-

2 (non-bursty traffic) follow a Poisson process with the average arrival rate λ . The service time of Class- c ($c = 1, 2$) traffic is exponentially distributed with mean $1/\mu_c$, respectively.

The system capacity is $L = K + 1$, where K is the buffer size. With these assumptions, this priority queueing system is represented using the state transition rate diagram as shown in Figure 6.2. The three-dimensional Markov chain is constructed from two 2-dimensional Markov chains with one in the front layer and the other in the back (shaded) layer. The transition between the corresponding states from one layer to the other represents the transition probability between two states of MMPP-2.

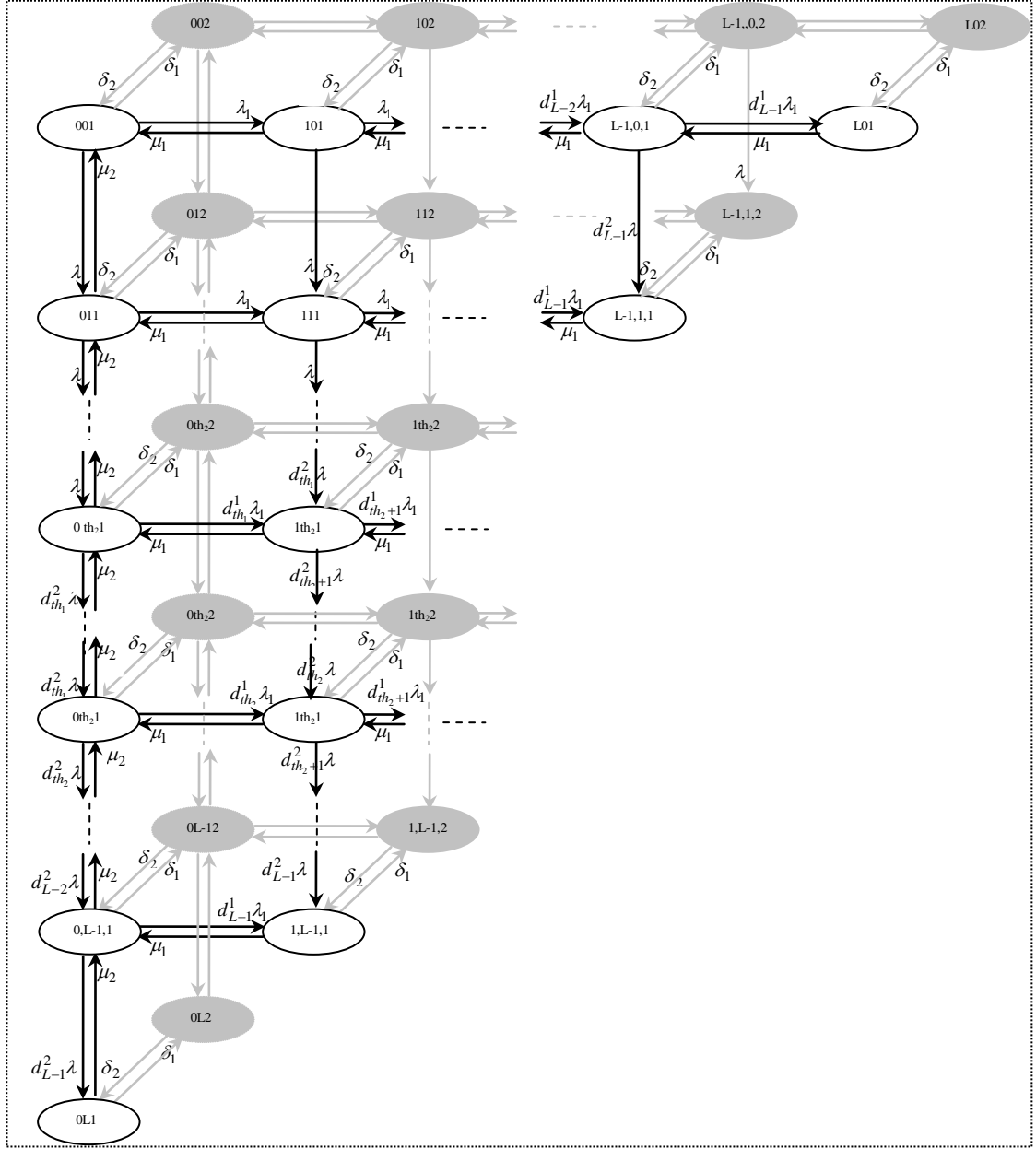


Figure 6.2. State transition rate diagram of the three-dimensional Markov chain for AQM scheme with PR scheduling mechanism and a single queue for two-class traffic.

State (i, j, s) ($0 \leq i + j \leq L, s = 1, 2$) corresponds to the situation where there are i and j packets of Class-1 and Class-2, respectively, in the system and MMPP-2 is at state s .

The transitions from State (i, j, s) to $(i+1, j, s)$, ($0 \leq i + j < L, s = 1, 2$), and from State

(i, j, s) to $(i, j+1, s)$, $(0 \leq i+j < L, s=1,2)$ imply that a packet from Class-1 and Class-2 traffic enters into the system, respectively. The transition rate out of State (i, j, s) to $(i+1, j, s)$, $(0 \leq i+j < th_1, s=1,2)$, is λ_s because no packets arriving from Class-1 are dropped before the instantaneous queue length exceeds threshold th_1 . However, when the instantaneous queue length exceeds threshold th_1 , i.e., when the system is at State (i, j, s) , $(th_1 \leq i+j < L, s=1,2)$, the probability that the arrivals of Class-1 traffic are allowed to enter into the system is denoted r_{i+j}^1 . As a result, the actual arrival rate of Class-1 traffic is reduced to $r_{i+j}^1 \lambda_s$ from λ_s at state s of the MMPP. Furthermore, when the system is at State (i, j, s) , $(th_2 \leq i+j \leq L, s=1,2)$, the actual arrival rate of Class-2 traffic is reduced to $r_{i+j}^2 \lambda$ from λ . As the packet dropping process can be viewed as a decrease of the arriving rate, the probability r_{i+j}^1 is given by:

$$r_{i+j}^c = 1 - d_{i+j}^c \quad (0 \leq i+j \leq L, c=1,2) \quad (6.1)$$

where the dropping probability d_{i+j}^c is expressed as

$$d_{i+j}^c = \begin{cases} 0 & 0 \leq i+j < th_c \\ \left(\frac{i+j-th_c+1}{L-th_c+1} \right) & th_c \leq i+j \leq L \end{cases} \quad (6.2)$$

Finally, the rate out of State (i, j, s) to $(i-1, j, s)$ $(1 \leq i \leq L, 0 \leq j \leq L-i, s=1,2)$, is equal to the service rate of Class-1 traffic, μ_1 . Class-2 packets can get service if and only if there is no Class-1 packet in the system. Therefore, the rate out of State $(0, j, s)$ to $(0, j-1, s)$

($1 \leq j \leq L, s = 1, 2$), is equal to the Class-2 traffic service rate, μ_2 . It is worth noting that the analytical model is developed in a general way by considering different service rates for two traffic classes.

The join state probability, p_{ijs} , in the three-dimensional Markov chain can be solved using the method reported in [15]. Let \mathbf{P} be the steady-state probability vector of this Markov chain, $\mathbf{P} = (p_{000}, p_{100}, p_{010}, \dots, p_{L00}, \dots, p_{0L0}, p_{001}, p_{101}, p_{011}, \dots, p_{L01}, p_{0L1})$. The infinitesimal generator matrix \mathbf{Z} of this Markov chain is of size $((L+1) \times (L+2)) \times ((L+1) \times (L+2))$. The steady-state probability vector \mathbf{P} satisfies the following equations

$$\begin{cases} \mathbf{PZ} = \mathbf{0} \\ \mathbf{Pe} = 1 \end{cases} \quad (6.3)$$

where $\mathbf{e} = (1, 1, \dots, 1)^T$ is a unit column vector of length $((L+1) \times (L+2)) \times 1$. Solving Equation (6.3) using the approach presented in [15] yields the steady-state probability vector \mathbf{P} as

$$\mathbf{P} = \boldsymbol{\alpha}(\mathbf{I} - \mathbf{X} + \mathbf{e}\boldsymbol{\alpha})^{-1} \quad (6.4)$$

where matrix $\mathbf{X} = \mathbf{I} + \mathbf{Q}/\beta$, $\beta \leq \min\{\mathbf{Q}_{ii}\}$ and $\boldsymbol{\alpha}$ is an arbitrary row vector of \mathbf{X} .

The aggregate and marginal state probabilities, p_m and p_m^c , that m packets are in the system and that m packets of Class- c are in the system can be calculated based on the solved join state probability p_{ijs} as below, respectively.

$$p_m = \sum_{i=0}^m \sum_{s=1}^2 p_{i(m-i)s} \quad 0 \leq m \leq L \quad (6.5)$$

$$p_m^c = \begin{cases} \sum_{j=0}^{L-m} \sum_{s=1}^2 p_{mjs} & c = 1, 0 \leq m \leq L \\ \sum_{i=0}^{L-m} \sum_{s=1}^2 p_{ims} & c = 2, 0 \leq m \leq L \end{cases} \quad (6.6)$$

Furthermore, the expressions of other two state probabilities p_{b_m} and $p_{b_m}^c$, that m packets are in the buffer and that m packets of Class- c are in the buffer, can be found as follows, respectively.

$$p_{b_m} = \begin{cases} \sum_{i=0}^1 p_i & m = 0 \\ \sum_{i=0}^{m+1} \sum_{s=1}^2 p_{i(m+1-i)s} & 1 \leq m \leq K \end{cases} \quad (6.7)$$

$$p_{b_m}^c = \begin{cases} \sum_{i=0}^1 \sum_{j=0}^{L-i} \sum_{s=1}^2 p_{ijs} & c = 1, m = 0 \\ \sum_{j=0}^{L-m-1} \sum_{s=1}^2 p_{(m+1)js} & c = 1, 1 \leq m \leq K \\ \sum_{i=0}^L \sum_{s=1}^2 p_{i0s} + \sum_{s=1}^2 p_{01s} & c = 2, m = 0 \\ \sum_{i=1}^{L-m} \sum_{s=1}^2 p_{ims} + \sum_{s=1}^2 p_{0(m+1)s} & c = 2, 1 \leq m \leq K \end{cases} \quad (6.8)$$

All these joint, aggregate and marginal state probabilities are useful for the following derivation of the aggregate and marginal performance metrics, respectively.

6.2.3 Performance Measures

This section derives the analytical expressions that estimate the aggregate and marginal performance metrics including utilization, throughput, mean number of packets in the

system and buffer, response time, queueing delay, packet loss probability as well as fairness.

The system utilization is equal to the probability that the server is busy. As the server is engaged as long as the number of packets in the system is not zero, the system utilization can be written as

$$\rho = 1 - p_0 \quad (6.9)$$

With the known aggregate and marginal state probabilities p_m , p_m^c , p_{b_m} and $p_{b_m}^c$, the aggregate mean number of packets in the system and buffer, \bar{L} and \bar{L}_b , as well as the marginal mean number of Class-c packets in the system and buffer, \bar{L}^c and \bar{L}_b^c can be calculated as follows, respectively.

$$\bar{L} = \sum_{i=0}^L (i \times p_i) \quad (6.10)$$

$$\bar{L}_b = \sum_{i=0}^K (i \times p_{b_i}) \quad (6.11)$$

$$\bar{L}^c = \sum_{i=0}^L (i \times p_i^c) \quad (6.12)$$

$$\bar{L}_b^c = \sum_{i=0}^K (i \times p_{b_i}^c) \quad (6.13)$$

Throughput is commonly defined as the average rate at which packets go through the system in the steady state. The aggregate throughput is equal to the addition of two marginal throughputs.

$$\overline{T}^c = \begin{cases} \mu_1 \times \sum_{i=1}^L p_i^1 & c=1 \\ \mu_2 \times \sum_{j=1}^L \sum_{s=0}^1 p_{0,js} & c=2 \end{cases} \quad (6.14)$$

$$\overline{T} = \sum_{c=1}^2 T^c \quad (6.15)$$

Little's Law [75] is adopted to calculate the aggregate and marginal mean response time (\overline{R} , \overline{R}^c) and queueing delay (\overline{D} , \overline{D}^c).

$$\overline{R} = \frac{\overline{L}}{\overline{T}} \quad (6.16)$$

$$\overline{R}^c = \frac{\overline{L}^c}{\overline{T}^c} \quad (6.17)$$

$$\overline{D} = \frac{\overline{L}_b}{\overline{T}} \quad (6.18)$$

$$\overline{D}^c = \frac{\overline{L}_b^c}{\overline{T}^c} \quad (6.19)$$

As for a stable queueing system, the packet loss probability is the ratio of the average packet loss rate, calculated as the difference between the average arrival rate and the throughput, to the corresponding original average arrival rate. The average arrival rate, $\overline{\lambda}^c$ ($c = 1, 2$), of Class- c traffic can be given by

$$\overline{\lambda}^c = \begin{cases} \frac{\lambda_1 \delta_2 + \lambda_2 \delta_1}{\delta_1 + \delta_2} & c=1 \\ \lambda & c=2 \end{cases} \quad (6.20)$$

The aggregate average arrival rate equals to $\overline{\lambda^1} + \overline{\lambda^2}$. Therefore, the expressions of the aggregate packet loss probability and the marginal one for class-c are provided as below, respectively.

$$PLP = \frac{\overline{\lambda^1} + \overline{\lambda^2} - \overline{T}}{\overline{\lambda^1} + \overline{\lambda^2}} \quad (6.21)$$

$$PLP^k = \frac{\overline{\lambda^k} - T^k}{\overline{\lambda^k}} \quad (6.22)$$

Finally, we adopt Jain's fairness index [92] to calculate the fairness of the two classes of traffic in the system as follows

$$F = \frac{(\sum_{i=1}^2 T^i)^2}{2 \times \sum_{i=1}^2 (T^i)^2} \quad (6.23)$$

6.2.4 Performance Comparison between AQM and Priority-based AQM

This section briefly investigates the effects of threshold, bursty traffic on the aforementioned aggregate and marginal system performance and compares AQM performance combined with the PR scheduling and that with FIFO scheduling scheme (c.f. Chapter 5). Meanwhile, the accuracy of the proposed analytical model is demonstrated. The performance results obtained from simulation experiments are illustrated and compared to analytical results in all the following figures. A perfect match is shown between the analytical and simulation results. This observation illustrates that the proposed model is

very accurate in calculating various performance metrics and examining the performance of the AQM mechanism combined with PR scheduling scheme under heterogeneous traffic.

This section presents the aggregate and marginal performance results shown in all figures below for the following scenario: The buffer capacity is assumed to be 10. Threshold th_2 is set to be 10, while threshold th_1 varies from 1 to 10. To investigate the priority-based AQM performance under non-bursty and bursty traffic, bursty Class-1 traffic is generated by an MMPP-2 model with the infinitesimal generated $\mathbf{Q} = \begin{bmatrix} -0.1 & 0.1 \\ 0.1 & -0.1 \end{bmatrix}$ and rate matrix $\mathbf{\Lambda} = \begin{bmatrix} 1 & 0 \\ 0 & 19 \end{bmatrix}$; meanwhile, non-bursty Class-1 traffic is generated by an MMPP-2 model with $\mathbf{Q} = \begin{bmatrix} -0.1 & 0.1 \\ 0.1 & -0.1 \end{bmatrix}$ and $\mathbf{\Lambda} = \begin{bmatrix} 10 & 0 \\ 0 & 10 \end{bmatrix}$. It is noticeable that the MMPP-2 to model Class-1 traffic degrades to non-bursty Poisson process as its arrival rates at different states are identical. Class-2 traffic is generated by a Poisson process with the average arrival rate $\lambda = 10$. Two classes of traffic are served with the mean rate $\mu = 23$ in order to make certain that queueing system is stable.

♦ *Aggregate performance metrics*

Figures 6.3-6.7 demonstrate the aggregate utilization, throughput, mean response time, mean queueing delay and packet loss probability against different values of thresholds th_1 under FIFO or PR scheduling subject to non-bursty or bursty traffic, respectively. Due to the identical average arrival rates of two classes (i.e., $\overline{\lambda^c} = 10$) as well as the identical mean service rate (i.e., $\mu_1 = \mu_2 = \mu$), it is easily understandable that the aggregate AQM

performance are unaffected by different scheduling schemes: FIFO and PR. The rest focuses on analysis and comparison of the effects of threshold and bursty traffic on the these aggregate performance metrics with FIFO and PR scheduling schemes.

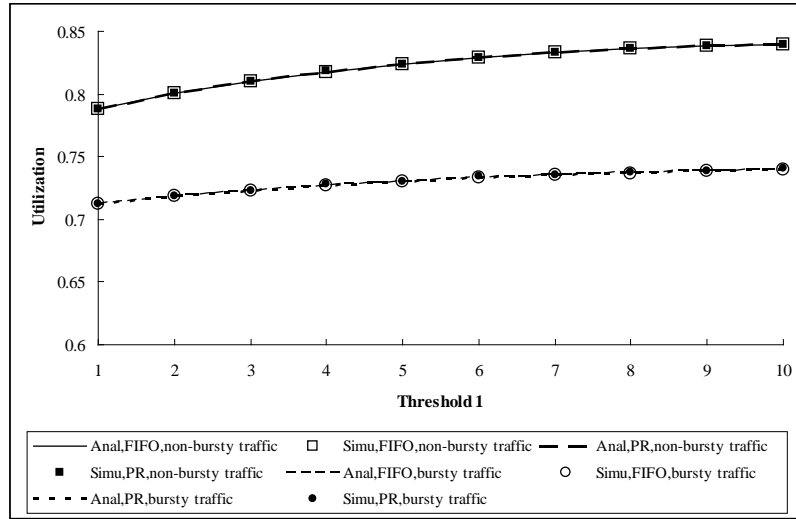


Figure 6.3. Utilization vs th_1 with various combinations of scheduling schemes and non-bursty or bursty Class-1 traffic

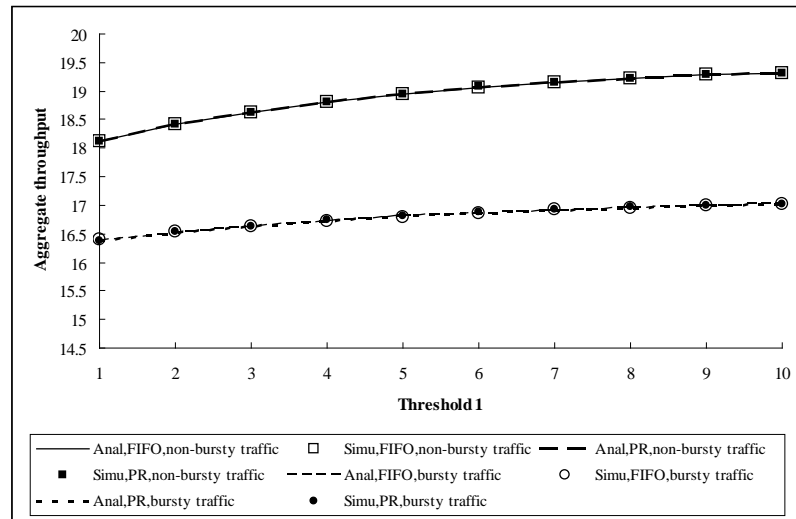


Figure 6.4. Throughput vs th_1 with various combinations of scheduling schemes and non-bursty or bursty Class-1 traffic

As shown in the above two figures, bursty traffic greatly reduces the utilization and system throughput and more remarkable reduction is demonstrated when th_1 increases. Furthermore, utilization and throughput increase with the growth of the threshold value as more Class-1 traffic can be injected into the system. The figures also demonstrate that a bit more sharp increase in utilization and throughput with the rise of threshold value when non-burst traffic is taken into account, respectively. Figures 6.5-6.6 depict that a large threshold value results in a high mean response time and long queuing delay, respectively. It can also be found that bursty traffic causes an increase in the response time and delay. Specifically, the response time and queuing delay increase more sharply with bursty traffic as the threshold value rises.

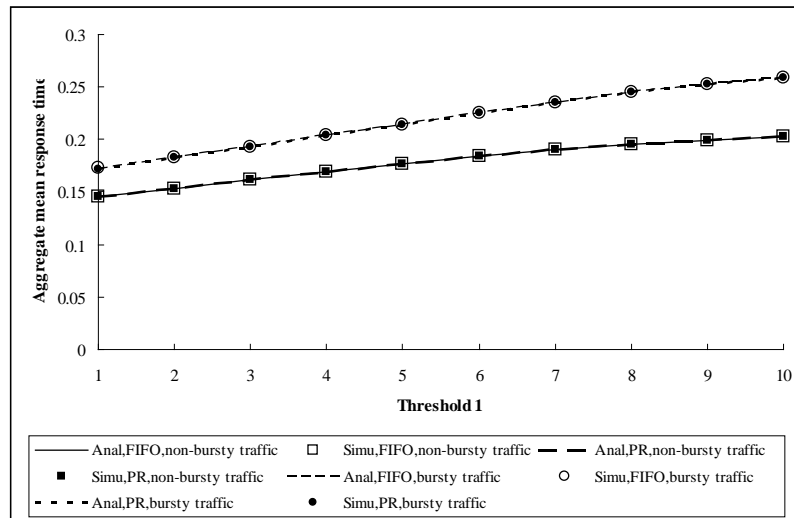


Figure 6.5. Mean response time vs th_1 with various combinations of scheduling schemes and non-bursty or bursty Class-1 traffic

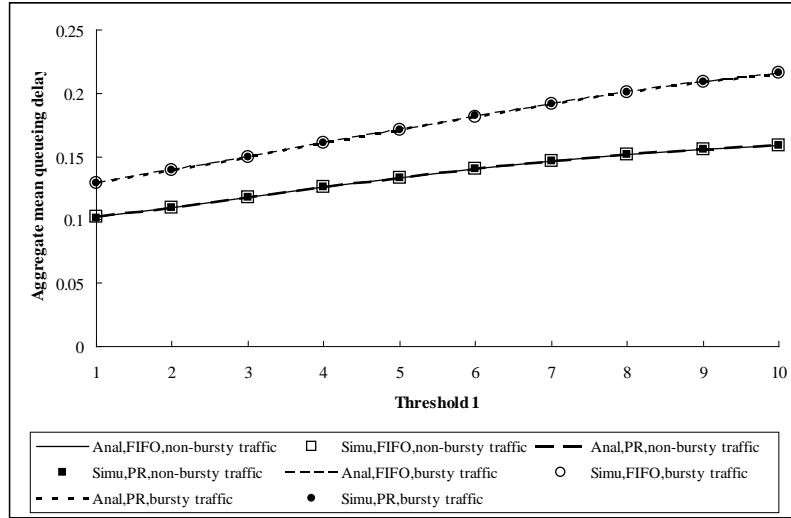


Figure 6.6. Mean queuing delay vs th_1 with various combinations of scheduling schemes and non-bursty or bursty Class-1 traffic

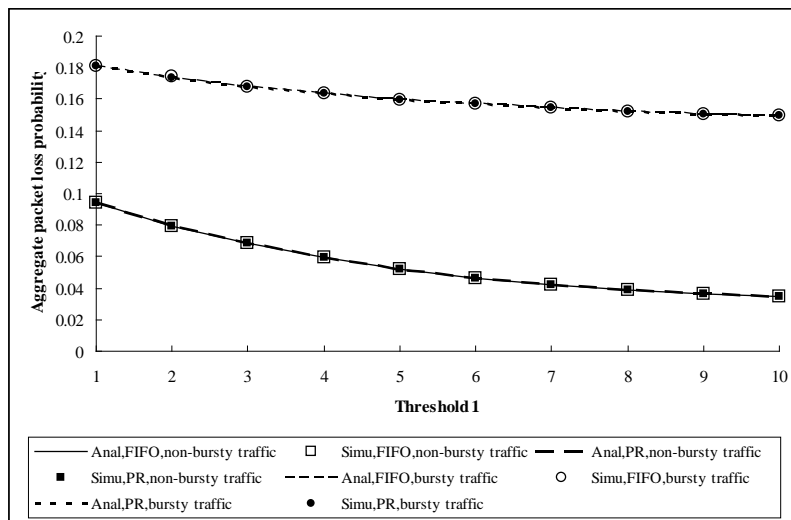


Figure 6.7. Packet loss probability vs th_1 with various combinations of scheduling schemes and non-bursty or bursty Class-1 traffic

Figure 6.7 illustrates that both bursty traffic and a smaller threshold value can increase the packet loss probability. The difference between the packet loss probability with bursty traffic and that with non-bursty traffic becomes increasingly obvious as the threshold varies from 1 to 10. It is because that the growth of threshold value enables the system to absorb

more arriving packets and, consequently, increases the burstiness degree of traffic injection into the system.

♦ *Marginal performance metrics*

We proceed to evaluate the marginal performance, including the throughput, mean response time, queueing delay, packet loss probability and fairness, of AQM coupled with PR scheduling scheme. As the average arrival rates and mean service rates of two classes of traffic are assumed to be identical respectively, adjusting the departure order of packets is unable to change the rate that traffic flow goes through the system. Therefore, as shown in Figures 6.8-6.11, AQM coupled with PR scheme results in the same marginal throughput and packet loss probability as AQM with FIFO does, respectively. It is clear that bursty traffic decreases marginal throughput but increases marginal packet loss probability. On the other hand, the throughput for Class-1 increases but that for Class-2 decreases as the threshold th_1 rises. While the packet loss probability for Class-1 decreases and that for Class-2 increases with the growth of th_1 . Figures 6.8 and 6.10 represent that the differences between the throughputs as well as the packet loss probabilities for Class-1 with non-bursty and bursty traffic are quite even when the threshold th_1 varies, respectively. However, an increase in threshold th_1 is capable of enlarging the differences between the throughput as well as the packet loss probability for Class-2 with non-bursty and bursty traffic, respectively. Additionally, Figure 6.12 demonstrates that fairness between two traffic classes increases as th_1 rises and it increases more sharply especially with bursty traffic. It is also clear that a better fairness is achieved under non-bursty traffic.

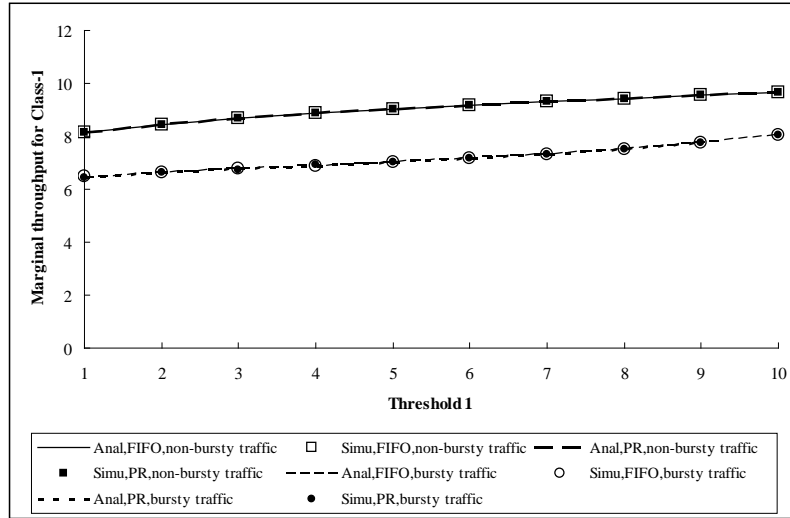


Figure 6.8. Throughput for Class-1 vs th_1 with various combinations of scheduling schemes and non-bursty or bursty Class-1 traffic

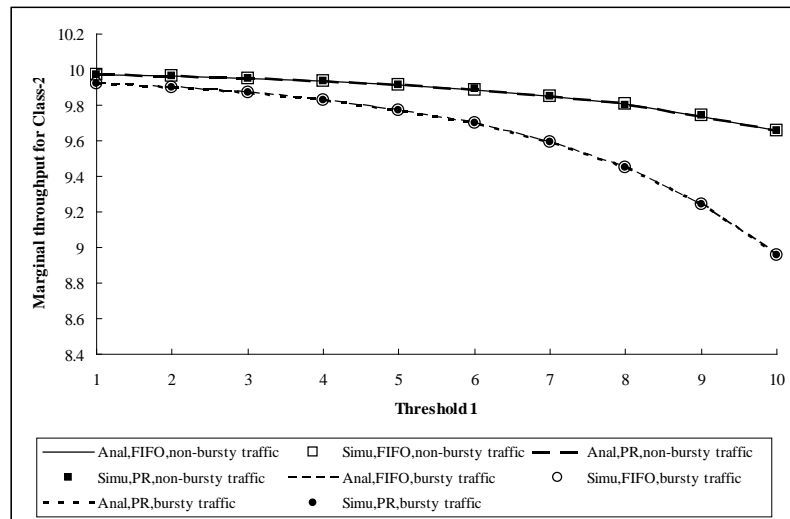


Figure 6.9. Throughput for Class-2 vs th_1 with various combinations of scheduling schemes and non-bursty or bursty Class-1 traffic

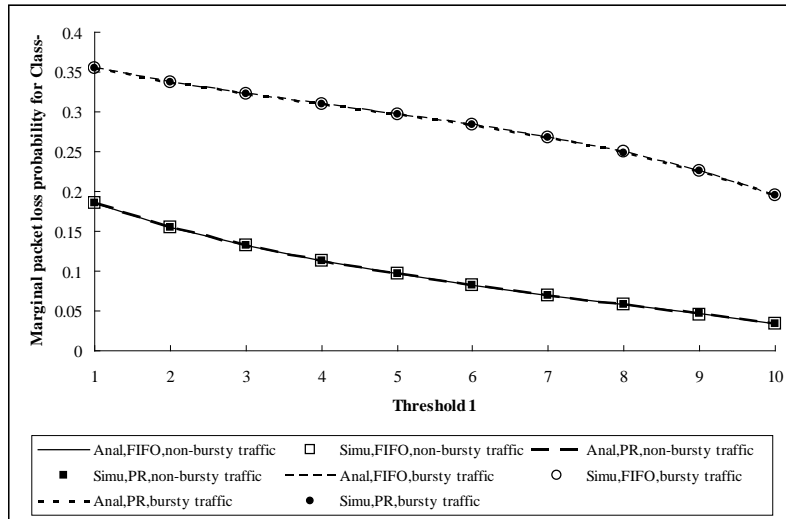


Figure 6.10. Packet loss probability for Class-1 vs th_1 with various combinations of scheduling schemes and non-bursty or bursty Class-1 traffic

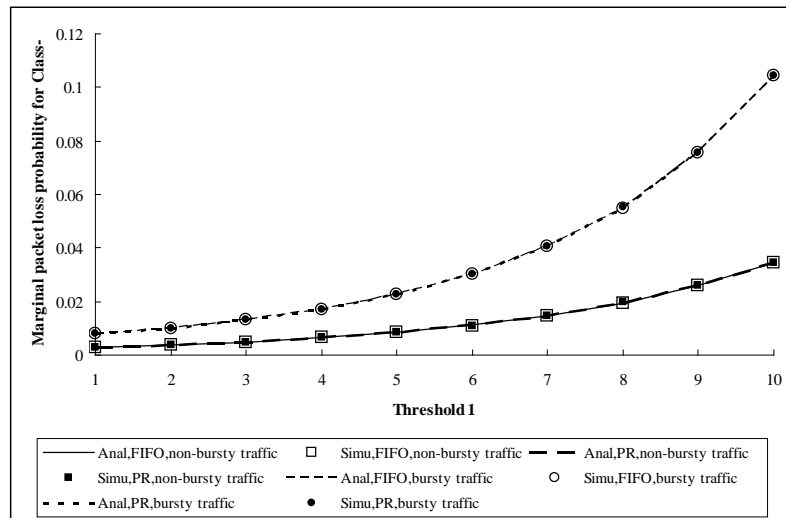


Figure 6.11. Packet loss probability for Class-2 vs th_1 with various combinations of scheduling schemes and non-bursty or bursty Class-1 traffic

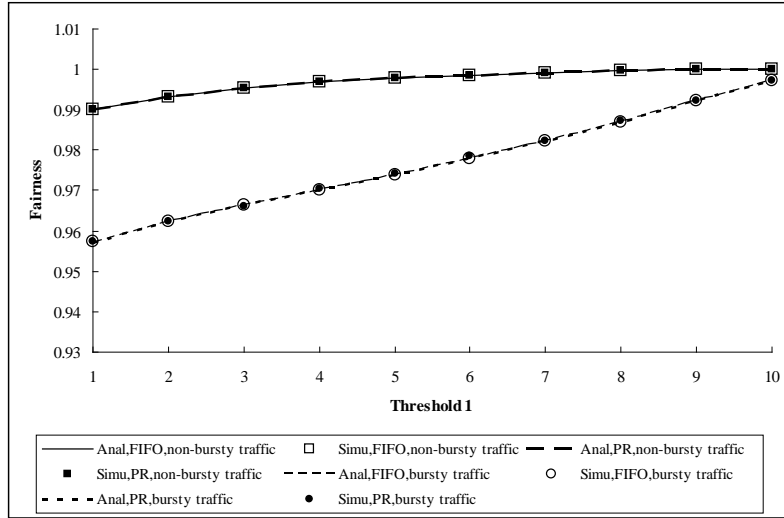


Figure 6.12. Fairness vs th_1 with various combinations of scheduling schemes and non-bursty or bursty Class-1 traffic

Figures 6.13-6.14 show that PR scheduling scheme controls the mean response time and queueing delay for Class-1 traffic better than FIFO. The reason is that PR guarantees strict priority to Class-1 traffic and consequently maintains the mean number of Class-1 packets in the system or buffer quite small. Moreover, the mean response time and queueing delay for Class-1 increase with a rising threshold value and a more apparent increasing trend can be found when the FIFO scheme is adopted. In addition, bursty traffic enables response time and queueing delay for Class-1 to increase. It can also be seen from these two figures, respectively, that effects of bursty traffic are greater with the FIFO scheme than PR. On the other hand, Figures 6.15-6.16 demonstrate the variation of mean response time and queueing delay for Class-2 against threshold th_1 under different considerations of scheduling scheme and bursty traffic. Adoption of PR scheme increases the mean response time and queueing delay for Class-2 as the cost of reducing the corresponding performance metrics for Class-1. It is easily understandable that bursty

traffic and an increasing threshold value result in a growth of response time and delay for Class-2. Different from the aforementioned analysis for Class-1, these two figures illustrate, respectively, more remarkable impact of threshold and bursty traffic on response time and queueing delay for Class-2 with PR scheme than FIFO. This is because that the PR scheme forces Class-2 packets to stay queued for a longer time, as a result, which highlights the influence of other factors, such as threshold and bursty traffic.

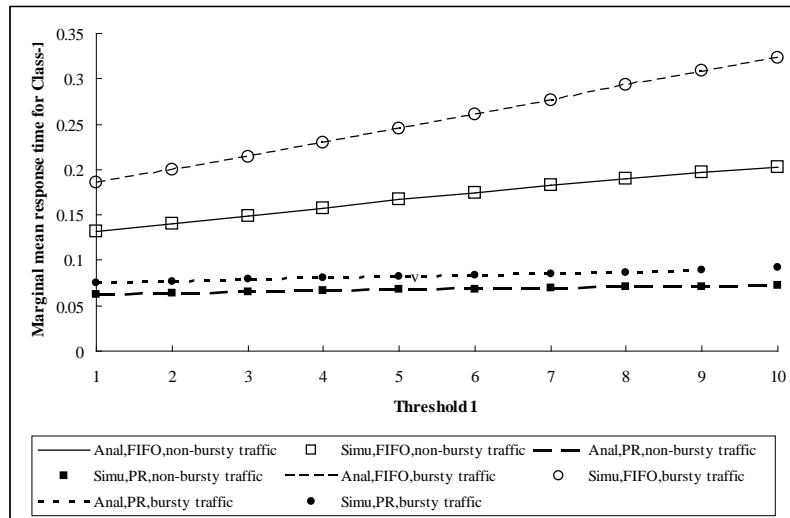


Figure 6.13. Mean response time for Class-1 vs th_1 with various combinations of scheduling schemes and non-bursty or bursty Class-1 traffic

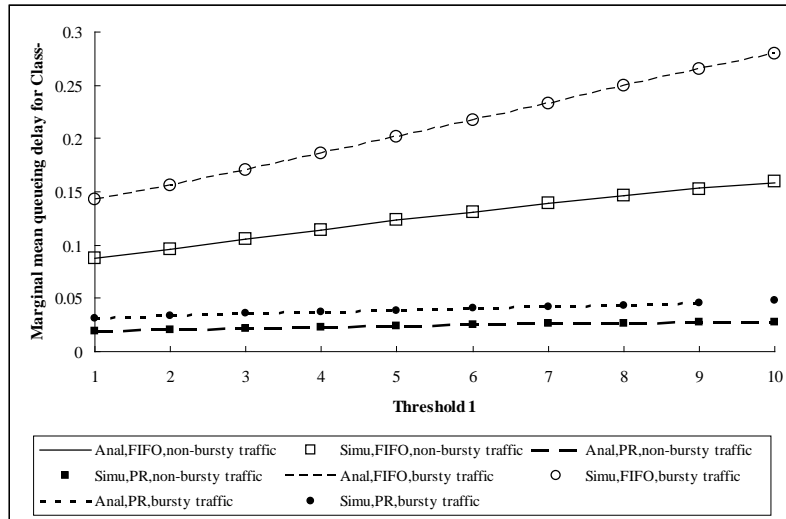


Figure 6.14. Mean queuing delay for Class-1 vs th_1 with various combinations of scheduling schemes and non-bursty or bursty Class-1 traffic

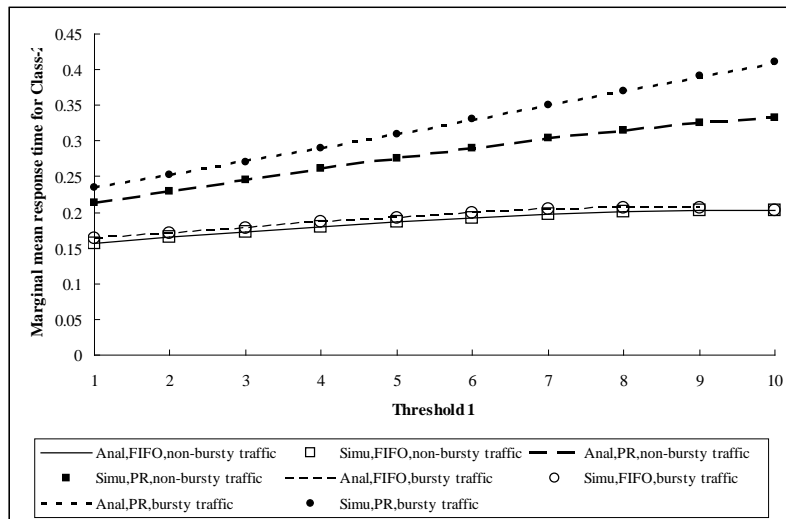


Figure 6.15. Mean response time for Class-2 vs th_1 with various combinations of scheduling schemes and non-bursty or bursty Class-1 traffic

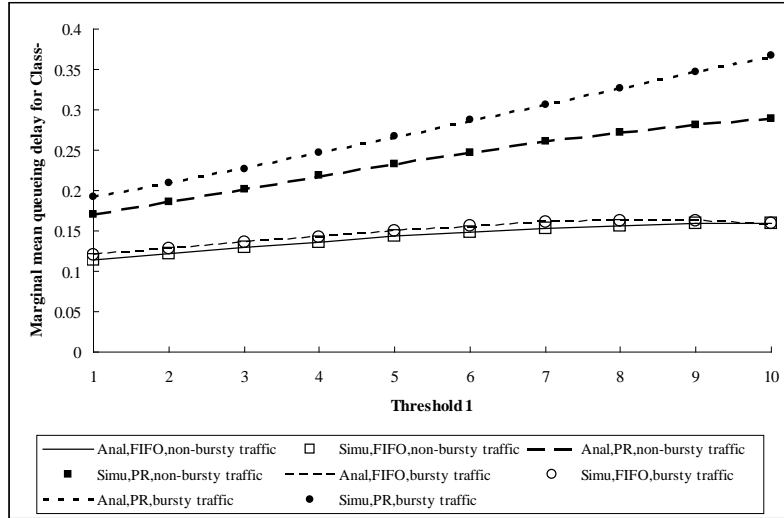


Figure 6.16. Mean queuing delay for Class-2 vs th_1 with various combinations of scheduling schemes and non-bursty or bursty Class-1 traffic

6.3 Priority-Based AQM with Multiple Queues

6.3.1 System Description

We consider a stable single server queueing system for two separate buffers in an AQM-enabled router as shown in Figure 6.17 where two priority classes of traffic wait for service at two separate finite queues, respectively. AQM scheme sets a single fixed threshold for each queue in order to control the actual rates of the corresponding traffic class injected into the system. When the current number, i , of Class- c ($c=1,2$) packets in its queue reaches the corresponding threshold, th_b^c , the forthcoming packets of this traffic class can be dropped randomly depending on the dropping probability which increases linearly from 0 to 1 as i increases from $(th_b^c - 1)$ to the queue capacity. In addition, the PR priority

scheduling mechanism is used to guarantee the high priority for bursty traffic over the non-bursty traffic.

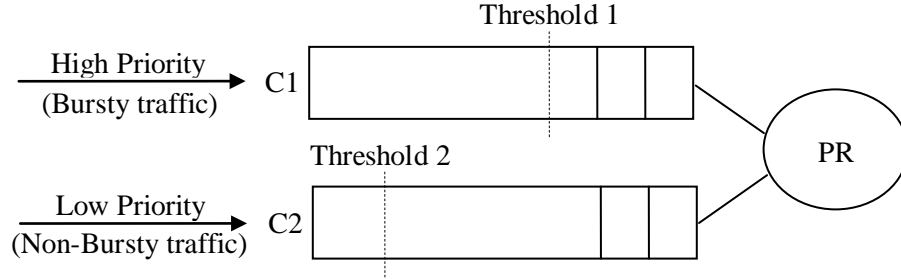


Figure 6.17. A Model of [MMPP-2][M]/M/1/K1/K2/th1/th2 Queueing System with PR scheme

6.3.2 Analytical Model

Similar to the assumptions in Section 6.2.2, an MMPP-2 model characterized by

$$\mathbf{Q} = \begin{bmatrix} -\delta_1 & \delta_1 \\ \delta_2 & -\delta_2 \end{bmatrix} \text{ as well as } \mathbf{\Lambda} = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix} \text{ and a Poisson model with the average arrival}$$

rate λ are adopted to capture the arrivals of Class-1 (bursty traffic) and Class-2 (non-bursty traffic), respectively. The service time of Class- c ($c = 1, 2$) traffic is exponentially distributed with mean $1/\mu_c$. K_c denotes the buffer size accommodating Class- c traffic.

Figure 6.18 shows the state transition diagram for this multiple-queue single-server system. Again, a three-dimensional Markov chain is constructed from two 2-dimensional Markov chains with one in the front layer and the other in the back (shaded) layer. State (i, j, s) ($0 \leq i \leq L_1, 0 \leq j \leq L_2, s = 1, 2$) corresponds to the situation where there are i and j packets of Class-1 and Class-2, respectively, in the system and MMPP-2 is at state s . All states, $(0, L_2, s)$ and (i, K_2, s) , ($1 \leq i \leq L_1, s = 1, 2$), on the bottom edge of the Markov

chain reflect the fact that the server can be occupied by Class-2 traffic if and only if there is no Class-1 packets in the system. The transitions from State (i, j, s) to $(i+1, j, s)$, $(0 \leq i < L_1, 0 \leq j \leq K_2, s = 1, 2)$, and from State (i, j, s) to $(i, j+1, s)$, $(0 \leq i \leq L_1, 0 \leq j < L_2, s = 1, 2)$ imply that a packet from Class-1 and Class-2 traffic enters into the system, respectively. Particularly, when the system is at State $(0, L_2, s)$ ($s = 1, 2$), the arrival of Class-1 packets still enables the system to immediately pre-empt the service to a Class-2 packet and put the pre-empted packet at the head of the Queue 2. Naturally, one packet (i.e., the one at the tail-end of the Queue 2) has to be discarded due to the space limitation.

The transition rate out of State (i, j, s) to $(i+1, j, s)$, $(0 \leq i < th_1, 0 \leq j \leq L_2, s = 1, 2)$, is λ_s because no packets arriving from Class-1 are dropped before the number of its packets exceeds th_1 . However, when the number of Class-1 packets in the system exceeds threshold th_1 , i.e., when the system is at State (i, j, s) , $(th_1 \leq i \leq L_1, 0 \leq j \leq K_2, s = 1, 2)$, the probability that the arrivals of Class-1 traffic are allowed to enter into the system is r_i^1 . As a result, the actual arrival rate of Class-1 traffic is reduced to $r_i^1 \lambda_s$ from λ_s at state s of the MMPP. Furthermore, when the system is at State (i, j, s) , $(1 \leq i \leq L_1, th_2 \leq j < K_2, s = 1, 2)$, the actual arrival rate of Class-2 traffic is reduced to $r_j^2 \lambda$ from λ .

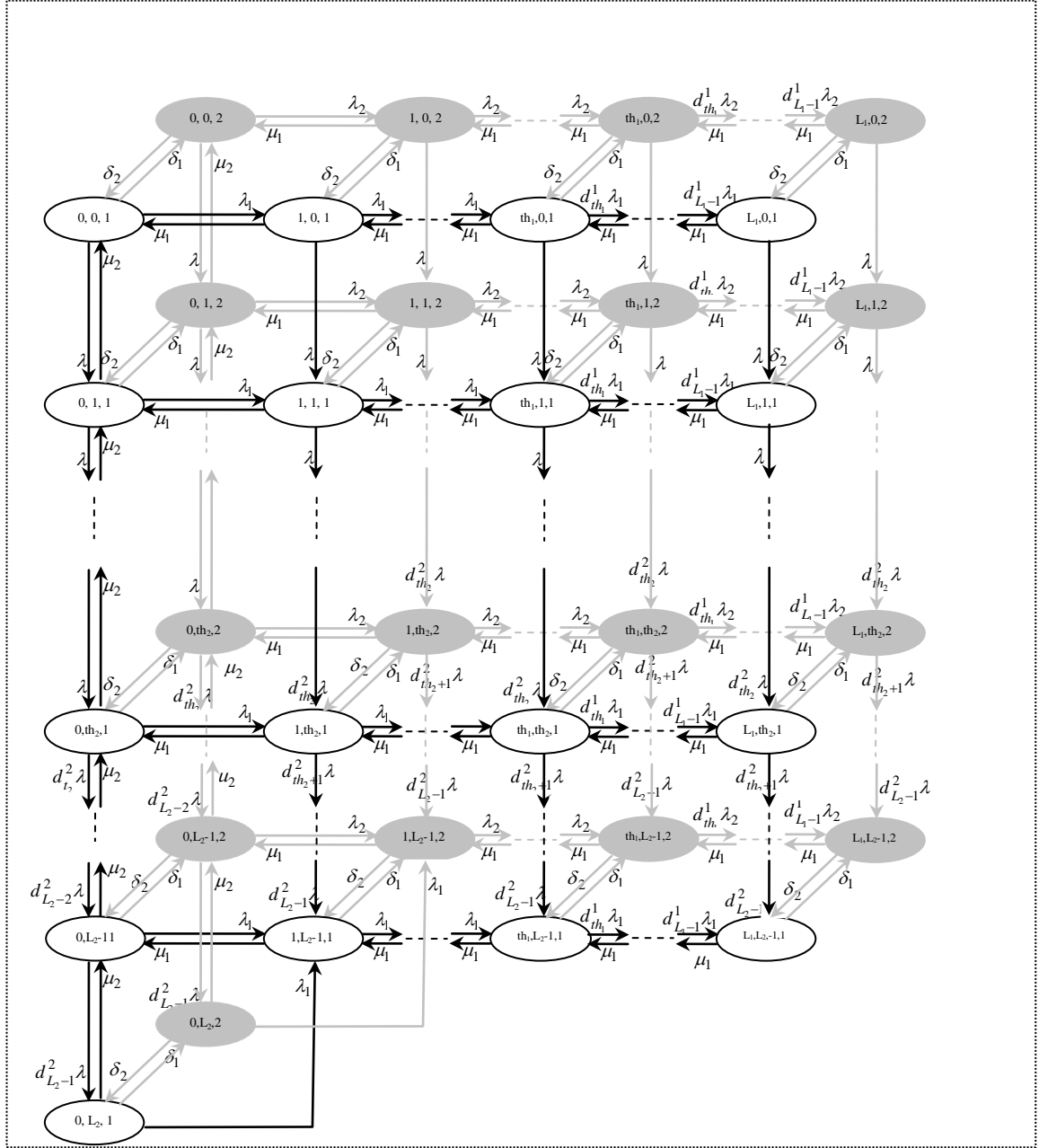


Figure 6.18. State transition rate diagram of the three-dimensional Markov chain for AQM scheme with PR scheduling mechanism and multiple classes-based queue for two-class traffic.

Meanwhile, when the system is at State $(0, j, s)$ ($th_2 + 1 \leq j < L_2, s = 1, 2$), the actual arrival rate of Class-2 traffic is reduced to $r_{j-1}^2 \lambda$ from λ in that the server is occupied by Class-2

traffic. The rate reduced probability r_i^c decreases from 1 to 0 along with the growth of the number of Class- c packets in the system.

$$r_i^c = \begin{cases} 1 & 0 \leq i < th_c \\ 1 - \left(\frac{i-th_c+1}{L_k-th_c+1}\right) & th_c \leq i \leq L_c \end{cases} \quad (6.24)$$

Finally, the rate out of State (i, j, s) to $(i-1, j, s)$ ($1 \leq i \leq L_1+1, 0 \leq j \leq K_2, s=1,2$), is equal to the service rate of Class-1 traffic, μ_1 . Class-2 packets can get service if and only if the queue of Class-1 traffic becomes empty. Therefore, the rate out of State $(0, j, s)$ to $(0, j-1, s)$ ($1 \leq j \leq L_2, s=1,2$), is equal to the Class-2 traffic service rate, μ_2 . The transition between the corresponding states from one layer to the other represents the transition probability between MMPP-2.

The join state probability, p_{ijs} , in the three-dimensional Markov chain can be solved using the method reported in [15]. Let \mathbf{P} be the steady-state probability vector of this Markov chain, $\mathbf{P} = (p_{001}, p_{011}, \dots, p_{0L_21}, p_{101}, \dots, p_{L_1K_21}, p_{002}, p_{012}, \dots, p_{L_1K_22})$. The infinitesimal generator matrix \mathbf{Z} of this Markov chain is of size $(2 \times ((L_1+1) \times L_2+1)) \times (2 \times ((L_1+1) \times L_2+1))$. The steady-state probability vector, \mathbf{P} , satisfies the following equations

$$\begin{cases} \mathbf{PZ} = \mathbf{0} \\ \mathbf{Pe} = 1 \end{cases} \quad (6.25)$$

where $\mathbf{e} = (1, 1, \dots, 1)^T$ is a unit column vector of length $(2 \times ((L_1 + 1) \times L_2 + 1)) \times 1$. Solving Equation (6.25) using the approach presented in [15] yields the steady-state probability vector \mathbf{P} as

$$\mathbf{P} = \boldsymbol{\alpha}(\mathbf{I} - \mathbf{X} + \mathbf{e}\boldsymbol{\alpha})^{-1} \quad (6.26)$$

where matrix $\mathbf{X} = \mathbf{I} + \mathbf{Q}/\beta$, $\beta \leq \min\{\mathbf{Q}_{ii}\}$ and $\boldsymbol{\alpha}$ is an arbitrary row vector of \mathbf{X} .

The aggregate state probability, p_m ($0 \leq m \leq L_1 + K_2$), that m packets are in the system is calculated as the sum of all join state probabilities p_{ijs} satisfying $i + j = m$. It is obvious that the relationship between p_m and the aggregate state probability, p_{b_m} , ($0 \leq m \leq K_1 + K_2$), that m packets are in the two buffers is

$$p_{b_m} = \begin{cases} p_0 + p_1 & m = 0 \\ p_{(m+1)} & 1 \leq m \leq (K_1 + K_2) \end{cases} \quad (6.27)$$

Moreover, the marginal state probabilities, p_m^c and $p_{b_m}^c$, that i Class- c packets are in the system and in the buffer can also be written as the following two equations, respectively.

$$p_m^c = \begin{cases} \sum_{j=0}^{K_2} \sum_{s=1}^2 p_{0js} & c = 1, m = 0 \\ \sum_{j=0}^{L_2} \sum_{s=1}^2 p_{mjs} & c = 1, 1 \leq m \leq L_1 \\ \sum_{i=0}^{L_1} \sum_{s=1}^2 p_{ims} & c = 2, 0 \leq m \leq K_2 \\ \sum_{s=1}^2 p_{0L_2s} & c = 2, m = L_2 \end{cases} \quad (6.28)$$

$$Pb_m^c = \begin{cases} p_0^1 + p_1^1 & c = 1, m = 0 \\ p_{(m+1)}^1 & c = 1, 1 \leq m \leq K_1 \\ \sum_{i=0}^{L_1} \sum_{s=1}^2 P_{ims} + P_{0(m+1)s} & c = 2, m = 0 \\ \sum_{i=1}^{L_1} \sum_{s=1}^2 P_{ims} + P_{0(m+1)s} & c = 2, 1 \leq m \leq K_2 \end{cases} \quad (6.29)$$

6.3.3 Performance Measures

Based on those joint, aggregate and marginal state probabilities, the analytical expressions of most aggregate and marginal performance metrics can be derived in a quite similar way as in Section 6.2.3. To avoid redundancies, expressions for these performance metrics including utilization (ρ), mean number of packets in the system and buffer (\bar{L} , \bar{L}^c , \bar{L}_b and \bar{L}_b^c), throughput (\bar{T} and \bar{T}^c), packet loss probability (PLP and PLP^c) and fairness (F) are listed directly as below without the detailed explanation:

$$\rho = 1 - p_0 \quad (6.30)$$

$$\bar{L} = \sum_{i=0}^{L_1+K_2} (i \times p_i) \quad (6.31)$$

$$\bar{L}^c = \sum_{i=0}^{L_c} (i \times p_i^c) \quad (6.32)$$

$$\bar{T}^c = \begin{cases} \mu_1 \times (1 - \sum_{s=0}^1 \sum_{j=0}^{L_2} p_{ojs}) & c = 1 \\ \mu_2 \times \sum_{s=0}^1 \sum_{j=1}^{L_2} p_{ojs} & c = 2 \end{cases} \quad (6.33)$$

$$\bar{T} = \sum_{c=1}^2 \bar{T}^c \quad (6.34)$$

$$PLP = \frac{\bar{\lambda}^1 + \bar{\lambda}^2 - \bar{T}}{\bar{\lambda}^1 + \bar{\lambda}^2} \quad (6.35)$$

$$PLP^c = \frac{\bar{\lambda}^c - \bar{T}^c}{\bar{\lambda}^c} \quad (6.36)$$

where, the average arrival rate, λ^c ($c = 1, 2$), of Class- c traffic can be given by

$$\lambda^c = \begin{cases} \frac{\lambda_1 \delta_2 + \lambda_2 \delta_1}{\delta_1 + \delta_2} & c = 1 \\ \lambda & c = 2 \end{cases} \quad (6.37)$$

$$F = \frac{(\sum_{i=1}^2 \bar{T}^i)^2}{2 \times \sum_{i=1}^2 (\bar{T}^i)^2} \quad (6.38)$$

Finally, let us move onto the derivation of expressions for the aggregate and marginal mean response time (\bar{R} and \bar{R}^c) and queueing delay (\bar{D} and \bar{D}^c). It is noticeable that the transition from State $(0, L_2, s)$ to $(1, K_2, s)$, ($s = 1, 2$), implicitly indicates a packet loss. That is, except departures from the server (i.e., throughput), there is the other output in this queueing system (i.e., discarding packets from Buffer 2). However, the previous derivation of aggregate and marginal mean number of packets in the system and Buffer 2 takes the whole system and the buffer into account, respectively. Therefore, the mean number of packets not to be discarded in the system or Buffer 2 should be the difference between the correspondingly known mean number of packets (i.e., \bar{L} , \bar{L}^2 , \bar{L}_b and \bar{L}_b^2) and the mean

number of discard packets. Then, Little's Law [75] may be adopted to calculate \bar{R} , \bar{R}^2 , \bar{D} and \bar{D}^2 . Moreover, discarding packets in Buffer 2 has no effects on the calculation of marginal mean response time and queueing delay for Class-1.

$$\bar{R} = \frac{\bar{L} - \sum_{s=1}^2 (P_{oL_2s} \times \frac{\lambda_1}{\lambda_1 + \delta_s + \mu_2})}{\bar{T}} \quad (6.39)$$

$$R^c = \begin{cases} \frac{\bar{L}^1}{\bar{T}^1} & c = 1 \\ \frac{\bar{L}^2 - \sum_{s=1}^2 (P_{oL_2s} \times \frac{\lambda_1}{\lambda_1 + \delta_s + \mu_2})}{\bar{T}^2} & c = 2 \end{cases} \quad (6.40)$$

6.3.4 Performance Comparison between Priority-based AQM with Single Queue and Multiple Queues

This section investigates the performance of AQM coupled with PR priority scheduling scheme with multiple class-based queues and compare the priority-based AQM performance with single queue and that with multiple queues. A good match between the analytical and simulation results shown in all tables below indicates the accuracy of the proposed model. To compare multiple queues and single queue system, the total buffer capacity is assumed to be 10, as well as the capacity (K_1) of Queue 1 in the multiple queues system varies from 1 to 9 and consequently the capacity ($K_2 = K - K_1$) of Queue 2 changes from 9 to 1 correspondingly. Both thresholds are allocated at the end of the queue in order to focus on the performance comparison between single and multiple queues

system. Two scenarios are generated as follows with the consideration of different relationships between mean arrival rates of each class. In Scenario S-6.3.I, the mean arrival rate of Class-1 traffic is lower than that of Class-2 and reversely in Scenario S-6.3.II. For instance, Class-1 traffic in Scenario S-6.3.I is generated by an MMPP-2 model with the infinitesimal generated $\mathbf{Q} = \begin{bmatrix} -0.1 & 0.1 \\ 0.1 & -0.1 \end{bmatrix}$ and rate matrix $\mathbf{\Lambda} = \begin{bmatrix} 5 & 0 \\ 0 & 5 \end{bmatrix}$; while Class-1 traffic in Scenario S-6.3.II is generated with $\mathbf{Q} = \begin{bmatrix} -0.1 & 0.1 \\ 0.1 & -0.1 \end{bmatrix}$ and $\mathbf{\Lambda} = \begin{bmatrix} 15 & 0 \\ 0 & 15 \end{bmatrix}$; moreover, in both scenarios, Class-2 traffic is generated by a Poisson process with the average arrival rate $\lambda = 10$. Two classes of traffic are served with the mean rate $\mu = 26$ in order to make certain that queueing system is stable.

Table 6.1 lists all analytical and simulation results of Class-1 performance metrics, including throughput, packet loss probability, mean response time and queueing delay, of two different queueing systems under two scenarios. In both scenarios, improved performance aforementioned for Class-1 can be achieved by splitting into two queues and adjusting each queue capacity compared with single queue system. For example, in Scenario S-6.3.I, higher throughput and less packet loss probability are demonstrated when the capacity of Queue 1 K_1 is greater than 3 but shortened mean response time and queueing delay are generated with $1 \leq K_1 \leq 3$. Specifically, mean response time and queueing delay for Class-1 are reduced approximately by 4% and 33.3%, respectively when $K_1 = 1$. Class-1 packet loss probability decreases remarkably but Class-1 throughput increases extremely slowly with the growth of K_1 . So in this case, a good combination of

queue capacities (i.e., (3,7)) can be found to provide the smaller mean response time and queueing delay and meanwhile guarantee the packet loss probability requirement. Quite similar trend of variation in corresponding performance metrics can be viewed in Scenario S-6.3.II. The combination of two queue capacities (3,7) enables the improvement of these performance metrics.

		S-6.3.I				S-6.3.II			
A/S		\bar{T}^1	PLP^1	\bar{R}^1	\bar{D}^1	\bar{T}^1	PLP^1	\bar{R}^1	\bar{D}^1
(10)	A	4.995008	0.000998	0.047551	0.00909	14.00172	0.066552	0.075071	0.036609
	S	4.99384	0.001002	0.047543	0.009087	13.99954	0.066609	0.075047	0.036581
(1,9)	A	4.849579	0.030084	0.044665	0.006203	12.38575	0.174284	0.052533	0.014071
	S	4.850546	0.030093	0.044662	0.006201	12.38446	0.174279	0.052536	0.014069
(2,8)	A	4.971239	0.005752	0.046793	0.008331	13.62957	0.091362	0.063487	0.025025
	S	4.969332	0.005757	0.046793	0.008327	13.62956	0.091353	0.063488	0.025027
(3,7)	A	4.994475	0.001105	0.047408	0.008947	14.24895	0.05007	0.071742	0.033281
	S	4.994528	0.00112	0.04741	0.00895	14.24835	0.050102	0.071742	0.033281
(4,6)	A	4.998938	0.000212	0.047568	0.009107	14.57887	0.028075	0.077779	0.039317
	S	5.000827	2.12E-04	0.047551	0.009094	14.57767	0.028062	0.077766	0.039309
(5,5)	A	4.999796	4.09E-05	0.047607	0.009146	14.76091	0.015939	0.082074	0.043613
	S	4.999807	3.93E-05	0.047587	0.009134	14.76263	0.015935	0.082071	0.043612
(6,4)	A	4.999961	7.86E-06	0.047616	0.009155	14.86332	0.009112	0.085057	0.046596
	S	4.999796	8.82E-06	0.047606	0.00915	14.86429	0.00908	0.085005	0.046553
(7,3)	A	4.999992	1.51E-06	0.047618	0.009157	14.92156	0.005229	0.087086	0.048624
	S	4.999878	7.80E-07	0.047602	0.009145	14.91882	0.005236	0.087071	0.048612
(8,2)	A	4.999999	2.91E-07	0.047619	0.009157	14.95488	0.003008	0.088441	0.049979
	S	4.999609	3.00E-07	0.047619	0.009154	14.95641	0.003003	0.088463	0.050001
(9,1)	A	5	5.59E-08	0.047619	0.009157	14.97402	0.001732	0.089332	0.05087
	S	5.00073	1.20E-07	0.047623	0.009159	14.97605	0.001722	0.089278	0.050825

Table 6.1. Marginal performance metrics for Class-1 with single queue and multiple queues system corresponding to two scenarios.

Table 6.2 illustrates the marginal throughput, packet loss probability, mean response time and queueing delay for Class-2 with two scenarios. In both scenarios, throughput increases but packet loss probability decreases and mean response time as well as queueing

delay increases first and then decrease as K_1 rises. In Scenario S-6.3.I, multiple queues system is unable to provide more throughput and less packet loss probability than single queue system. Assigning $K_1 = 1$ can result in the most closed values of throughput and packet loss probability to the corresponding performance in single queue system. Moreover, all nine combinations of queues capacities are able to reduce the mean response time and queuing delay. Therefore, multiple queues system with $K_1 = 1$ can approximately achieve the performance generated by a single queue system. In Scenario S-6.3.II, however, multiple queues system with $K_1 = 1$ increases throughput by 2% and decreases packet loss probability, response time and queuing delay by 4%, 3% and 3%, respectively.

		S-6.3.I				S-6.3.II			
Q1	A/S	\overline{T}^2	PLP^2	\overline{R}^2	\overline{D}^2	\overline{T}^2	PLP^2	\overline{R}^2	\overline{D}^2
(10)	A	9.990016	0.000998	0.111089	0.072627	9.334478	0.066552	0.426722	0.388261
	S	9.989629	0.001	0.111062	0.072602	9.336956	0.066587	0.426892	0.388424
(1,9)	A	9.98585	0.001415	0.106839	0.068378	9.599155	0.040085	0.30052	0.262059
	S	9.98634	0.001421	0.106876	0.068412	9.599272	0.040034	0.300381	0.261923
(2,8)	A	9.96882	0.003118	0.108888	0.070426	9.099861	0.090014	0.358693	0.320231
	S	9.969551	0.003115	0.108876	0.070413	9.100175	0.090047	0.358602	0.320149
(3,7)	A	9.941749	0.005825	0.107992	0.06953	8.629399	0.13706	0.370152	0.331691
	S	9.94196	0.005823	0.108043	0.069576	8.63109	0.136997	0.370049	0.33159
(4,6)	A	9.895796	0.01042	0.105706	0.067244	8.203542	0.179646	0.355326	0.316865
	S	9.897383	0.010436	0.105767	0.067302	8.204535	0.179504	0.355163	0.316702
(5,5)	A	9.814525	0.018548	0.102087	0.063625	7.787449	0.221255	0.326076	0.287615
	S	9.811198	0.018518	0.102059	0.063601	7.786954	0.22123	0.326003	0.287539
(6,4)	A	9.667339	0.033266	0.0967	0.058238	7.331453	0.266855	0.288376	0.249914
	S	9.666987	0.033295	0.096715	0.05825	7.333489	0.266481	0.288128	0.249673
(7,3)	A	9.393969	0.060603	0.088846	0.050384	6.767711	0.323229	0.244785	0.206323
	S	9.394843	0.060614	0.088826	0.050368	6.76931	0.322973	0.244634	0.206176
(8,2)	A	8.865392	0.113461	0.07753	0.039068	5.98192	0.401808	0.195432	0.156971
	S	8.86401	0.113426	0.077522	0.03906	5.982664	0.401686	0.195405	0.156952
(9,1)	A	7.766283	0.223372	0.061025	0.022563	4.717798	0.52822	0.13636	0.097899
	S	7.766867	0.223361	0.061013	0.02256	4.718826	0.528075	0.136306	0.097845

Table 6.2. Marginal performance metrics for Class-2 with single queue and multiple queues system corresponding to two scenarios.

Finally, Table 6.3 shows the analytical and simulation results of aggregate performance metrics including utilization, throughput, mean response time, packet loss probability and fairness. For both scenarios, it can be observed the absolute advantages of multiple queues system in terms of mean response time and fairness. That is, the multiple queues system with $3 \leq K_1 \leq 9$ and with $1 \leq K_1 \leq 2$ can achieve better fairness than single queue system in two Scenarios, respectively. But no other improved performance metrics can be achieved by multiple queues system. Therefore, the best performance results in both scenarios are generated by setting $K_1 = 2$ and $K_1 = 3$, respectively.

		S-6.3.I					S-6.3.II				
		U	\bar{T}	\bar{R}	PLP	F	U	\bar{T}	\bar{R}	PLP	F
(10)	A	0.5763	14.9850	0.0899	0.0009	0.9	0.8975	23.3362	0.2157	0.0665	0.9615
	S	0.5762	14.9834	0.0898	0.0010	0.8999	0.8976	23.3364	0.2158	0.0666	0.9616
(1,9)	A	0.5705	14.8354	0.0865	0.0109	0.8929	0.8455	21.9849	0.1608	0.1206	0.9841
	S	0.5706	14.8368	0.0865	0.0109	0.8930	0.8455	21.9837	0.1607	0.1205	0.9842
(2,8)	A	0.5746	14.9400	0.0882	0.0039	0.8993	0.8742	22.7294	0.1816	0.0908	0.9618
	S	0.5746	14.9388	0.0882	0.0039	0.8992	0.8741	22.7297	0.1816	0.0908	0.9618
(3,7)	A	0.5744	14.9362	0.0877	0.0042	0.9011	0.8799	22.8783	0.1842	0.0848	0.9431
	S	0.5745	14.9364	0.0877	0.0042	0.9011	0.8799	22.8794	0.1842	0.0848	0.9431
(4,6)	A	0.5728	14.8947	0.0861	0.0070	0.9024	0.8762	22.7824	0.1777	0.0887	0.9273
	S	0.5730	14.8982	0.0862	0.0070	0.9025	0.8761	22.7822	0.1776	0.0886	0.9274
(5,5)	A	0.5697	14.8143	0.0837	0.0123	0.9044	0.8672	22.5483	0.1663	0.0980	0.9127
	S	0.5695	14.8110	0.0836	0.0123	0.9045	0.8672	22.5495	0.1663	0.0980	0.9126
(6,4)	A	0.5641	14.667	0.0799	0.0221	0.9080	0.8536	22.1947	0.1522	0.1122	0.8967
	S	0.5641	14.6667	0.0799	0.0222	0.9080	0.8535	22.1977	0.1521	0.1120	0.8967
(7,3)	A	0.5536	14.3939	0.0745	0.0404	0.9147	0.8342	21.6892	0.1362	0.1324	0.8761
	S	0.5535	14.3947	0.0745	0.0404	0.9147	0.8340	21.6881	0.1362	0.1323	0.8762
(8,2)	A	0.5332	13.8653	0.0667	0.0756	0.9278	0.8052	20.9368	0.1190	0.1625	0.8448
	S	0.5332	13.8636	0.0667	0.0756	0.9279	0.8053	20.9390	0.1190	0.1624	0.8448
(9,1)	A	0.4910	12.7662	0.0557	0.1489	0.9551	0.7573	19.6918	0.1005	0.2123	0.7866
	S	0.4910	12.7676	0.0557	0.1489	0.9551	0.7573	19.6948	0.1005	0.2122	0.7866

Table 6.3. Aggregate performance metrics with single queue and multiple queues system corresponding to two scenarios.

6.4 Summary

This chapter has proposed two three-dimensional Markovian chains for AQM congestion control scheme with PR scheduling scheme in a single queue and multiple class-based systems, respectively. Two classes of traffic are generated by a non-bursty Poisson process and a bursty two-state MMPP, respectively. We have derived and evaluated the essential aggregate and marginal performance metrics including the mean number of packets in the system and in the queue, packet loss probability, mean response time, throughput, utilization and fairness. The accuracy of two models in examining the performance of the priority-based AQM under heterogeneous traffic and calculating various performance metrics has been demonstrated by comparing analytical results against simulation results obtained from a simulator programmed in JAVA.

The first proposed model has been used to evaluate AQM performance with PR scheduling scheme and a single queue and to compare PR and FIFO scheduling scheme. Particularly, the marginal mean response time and queueing delay for the high-priority traffic are improved significantly while those for the low-priority traffic are degraded remarkably due to the PR scheduling scheme. Then the second analytical model has been adopted to compare priority-based AQM performance with a single queue and multiple class-based systems. The effects of multiple queues capacities on the aggregate and marginal performance including throughput, packet loss probability, mean response time and queueing delay as well as utilization and fairness has been evaluated. We have

explained how to seek the best way to allocate the capacity to each queue according to different aggregate and marginal performance requirements for a specific scenario.

Chapter 7

Conclusions and Future Work

This chapter draws conclusions of the thesis and provides some suggestions for the future work in the performance modelling research area of congestion control mechanisms.

7.1 Conclusions

The contributions of this thesis are concluded in this section as bellow:

- ♦ All new proposed analytical models have been used to derive the expressions of essential performance metrics, including utilization, throughput, mean number of packets in the system and buffer, mean response time, queueing delay, packet loss probability and fairness, for the corresponding systems.
- ♦ A single-server finite queuing system for the performance evaluation of AQM scheme under the non-bursty Poisson arrival process has been developed. Two continuous-time Markov models have been proposed, respectively, for AQM scheme subject to single class and two classes of traffic. Closed-form expressions for corresponding performance metrics in each system have been derived. Specifically, the marginal steady state probabilities in multi-class system have been obtained.
- ♦ The model for single class traffic has been adopted to analyze the effects of mean arrival rate, mean service rate and buffer capacity on AQM performance. It can be

concluded as follows: 1) The rise of mean arrival rate enables all performance metrics to increase; 2) A high mean service rate improves most performance metrics including throughput, mean response time, queueing delay, packet loss probability and mean number of packets in the system and buffer but at the cost of utilization; 3) All performance metrics, except packet loss probability, increase as the buffer capacity enlarges. On the other hand, the model for two classes of traffic has been used to investigate the effects of thresholds on aggregate and marginal performance. Numerous experiments results have shown that a high threshold degrades the marginal performance for traffic not controlled by this threshold. Meanwhile, the varying threshold have same effects on aggregate and marginal performance for traffic controlled by it, such as, low packet loss probability but high throughput and utilization, long response time as well as queueing delay can be achieved with increase in the threshold. Furthermore, it was pointed out that, if keeping the difference between thresholds constant, the smaller values of thresholds is capable of reducing the marginal mean number of packets in the system as well as queueing delay of each class and providing similar throughput as bigger one.

- ◆ A two-dimensional Markov model has been introduced for a single-server queueing system with AQM scheme subject to bursty traffic captured by an MMPP-2. Closed-form expressions for aforementioned performance metrics have been derived and accuracy of the developed model has been demonstrated by comparing analytical results with those obtained from simulators developed in JAVA programming language. The effects of the burstiness and correlation of the MMPP-

2 traffic on performance has been investigated to demonstrate the model's application. Numerical results have demonstrated that high burstiness and correlation can significantly degrade the AQM performance in terms of increasing the mean numbers of packets in the system and buffer, mean response time, mean queueing delay as well as packet loss probability and decreasing utilization and throughput. In particular, it has been observed that high burstiness (or correlation) more remarkably affects the AQM performance if the correlation (or burstiness) is high. Additionally, the effects of burstiness and correlation are also sensitive to the threshold value. For example, a low threshold is capable of degrading the negative effects of high burstiness and correlation on the AQM performance.

- ♦ The other two-dimensional Markov model has been further developed for AQM with two individual thresholds subject to two classes of traffic modelled by a Poisson process and MMPP-2. We have adopted this model evaluate the impacts of parameters related to Class-1 traffic, including the average arrival rate, burstiness, correlation and its threshold, on the aggregate and marginal utilization, throughput, mean queueing delay and packet loss probability. It can be found that as the traffic rate grows, the marginal performance metrics for Class-2 are degraded substantially, while all values of the aggregate performance and marginal performance for Class-1 increase. Moreover, analytical results have also clearly demonstrated the detrimental impacts of traffic burstiness and correlation on all performance metrics. Lastly, the uncertainty effects of the threshold assigned to Class-1 traffic on all performance metrics have been analyzed. The analytical model is useful for assisting to find the

best set of parameters settings to suit the type of service required, for instance, real-time services like voice require low delay, while data services require low packet loss.

- ♦ A three-dimensional Markovian chains has been developed for AQM scheme with two classes traffic and PR scheduling scheme in a single queue. Two classes of traffic are generated by a non-bursty Poisson process and a bursty two-state MMPP, respectively. Adoption of PR scheme reduces the mean response time and queueing delay for Class-2 as the cost of increasing the corresponding performance metrics for Class-1. Moreover, the marginal mean response time and queueing delay for the high-priority traffic are improved significantly while those for the low-priority traffic are degraded remarkably due to the PR scheduling scheme. The other three-dimensional Markov model has been proposed for a single-server two-queues system with AQM and PR scheme. This model has been adopted to compare priority-based AQM performance with a single queue and multiple class-based systems. The effects of multiple queues capacities on the aggregate and marginal performance including throughput, packet loss probability, mean response time and queueing delay as well as utilization and fairness has been evaluated. By taking specific scenario as an example, we have explained how to seek the best way to allocate the capacity to each queue according to different aggregate and marginal performance requirements for a specific scenario.

7.2 Future Work

Moving beyond the core of the present work, there are several interesting issues and open problems that require further investigation. These are briefly outlined below.

- ♦ The present study has focused on the analysis of AQM scheme subject to one or two classes of traffic. In practical networks, traffic generated by multimedia application is classified into multiple categories. We will extend our proposed model methods to handle the AQM queuing systems subject to multiple heterogeneous traffic classes.
- ♦ More and more measurement evidences [93-94] have shown that the traffic generated by Variable-Bit-Rate (VBR) video in modern communication networks and multimedia systems exhibit extremely bursty arrival nature over a wide range of time scales. This fractal behavior of packet arrivals can be modelled using statistically self-similar or long-range-dependent processes, which have significantly different theoretical properties from those of the conventional Markovian non-memory arrival processes [95-97]. A more challenging extension of our work would be to develop the analytical model for AQM in the presence of fractal self-similar traffic.
- ♦ In wireless communications, data transmission suffers from varied signal strengths and channel bit error rates. To ensure successful packet reception under different channel conditions, various congestion control schemes based on AQM have been proposed [32, 98-99]. Our future work will aim to develop original analytical model for these new AQM congestion control schemes in the presence of error

transmission channels and use the developed models for performance analysis and resource allocation in wireless networks.

References

- [1] J. Postel, "Transmission Control Protocol," *STD 7, RFC 793*, September 1981.
- [2] R. Braden, "Requirements for Internet Hosts -- Communication Layers," *STD 3, RFC 1122*, October 1989.
- [3] V. Jacobson, "Congestion Avoidance and Control," *Computer Communication Review*, vol. 18, no. 4, pp. 314-329, Aug. 1988.
<ftp://ftp.ee.lbl.gov/papers/congavoid.ps.Z>.
- [4] V. Jacobson, "Modified TCP Congestion Avoidance Algorithm," end2end-interest mailing list, April 30, 1990. <ftp://ftp.isi.edu/end2end/end2end-interest-1990.mail>.
- [5] W. Stevens, "TCP Slow Start, Congestion Avoidance, Fast Retransmit, and Fast Recovery Algorithms", *RFC 2001*, January 1997.
- [6] M. Allman, V. Paxson, and W. Stevens, "TCP Congestion Control," *RFC 2581*, April 1999.
- [8] C. Brandauer, G. Iannaccone, C. Diot, and T. Ziegler, "Comparison of Tail Drop and Active Queue Management Performance for Bulk-data and Web-like Internet Traffic", *Proc. ISCC*, pp. 122-129, 2001.
- [7] B. Braden, D. Clark, J. Crowcroft, B. Davie, S. Deering, D. Estrin, S. Floyd, V. Jacobson, G. Minshall, C. Partridge, L. Peterson, K. Ramakrishnan, S. Shenker, J. Wroclawski, and L. Zhang, "Recommendations on Queue Management and Congestion Avoidance in the Internet," *IETF RFC2309*. 1998.

- [9] M. Christiansen, K. Jeffay, D. Ott, and F. Donelson Smith, "Tuning RED for Web Traffic", *IEEE/ACM Trans. Network*, vol. 9, no. 3, pp. 249-264, 2001.
- [10] S. Floyd, "RED: Discussions of Setting Parameters," <http://www.icir.org/floyd/REDparameters.txt>, 1997.
- [11] S. Floyd, and K. Fall, "Promoting the Use of End-to-End Congestion Control in the Internet", *IEEE/ACM Trans. Network*, vol. 7, no. 4, pp.485-472, 1999.
- [12] S. Floyd, and V. Jacobson, "Random Early Detection Gateways for Congestion Avoidance", *IEEE/ACM Trans. Network*, vol. 1, no. 4, pp. 397-413, 1993.
- [13] M. May, C. Diot, B. Lyles, and J. Bolot, "Influence of Active Queue Management Parameters on Aggregate Traffic Performance", *INRIA.RR3995*, 2000.
- [14] S. Ryu, C. Rump and C. Qiao, "Advances in Internet Congestion Control," *IEEE Communications surveys and tutorials*, vol. 5, no. 1, pp. 28-39, 2003.
- [15] W. Fischer, and K. Meier-Hellstern, "The Markov-modulated Poisson Process (MMPP) Cookbook," *Performance Evaluation*, vol. 18, no. 2, pp.149-171, 1993.
- [16] L. Wang, G. Min, and I. Awan, "Analysis of Active Queue Management under Two Classes of Traffic," *Proc. 21st UK Performance Engineering Workshop (UKPEW)*, pp. 101-109, Newcastle, UK, July 14-15, 2005.
- [17] L. Wang, G. Min, and I. Awan, "Modelling and Evaluation of Congestion Control Mechanism for Different Classes of Traffic," *Concurrency and Computation: Practice and Experience (CCPE)*, vol. 19, pp. 1141-1156, 2007.

- [18] L. Wang, G. Min, and I. Awan, "Modelling Active Queue Management with Different Traffic Classes," *Proc. Int. Conference on AINA*, pp. 442-446, Vienna, Austria, April 18-20, 2006.
- [19] L. Wang, G. Min, and I. Awan, "Analytical Modeling and Comparison of AQM-Based Congestion Control Mechanisms," *Proc. Int. Conference on High Performance Computing and Communications (HPCC)*, pp. 67-76, Sorrento, Italy, September 21-23, 2005.
- [20] L. Wang, G. Min, and I. Awan, "Modeling and Analysis of Active Queue Management Schemes under Bursty Traffic," *Journal of Wireless Information Networks*, vol. 13, no. 2, pp. 161-171, 2006.
- [21] L. Wang, G. Min, and I. Awan, "Performance Analysis of Buffer Allocation Schemes under Heterogeneous Traffic with Individual Thresholds," *Proc. 20th Int. Conference on Advanced Information Networking and Applications (AINA)*, Vienna, Austria, April 18-20, pp. 559-564, 2006.
- [22] L. Wang, G. Min, and I. Awan, "Performance Analysis of Buffer Allocation Schemes Under MMPP and Poisson Traffic with Individual Thresholds," *Cluster Computing*, vol. 10, no. 1, pp. 17-31, 2007.
- [23] L. Wang, G. Min, and I. Awan, "Effects of Bursty and Correlated Traffic on the Performance of Active Queue Management Schemes," *COST 285: Modelling and Simulation Tools for Research in Emerging Multi-service Telecommunications*, Surrey, United Kingdom, March 28-29, 2007.

- [24] L. Wang, G. Min, and I. Awan, "An Analytical Model for Priority-Based AQM in the Presence of Heterogeneous Network Traffic," *Proc. 22th Int. Conference on Advanced Information Networking and Applications (AINA)*, , GinoWan, Okinawa, Japan, March 25-28, pp. 93-99, 2008.
- [25] G. Hasegawa, and M. Murata "Analysis of Dynamic Behaviors of Many TCP Connections Sharing Tail-Drop/RED Routers," *IEEE GLOBECOM*. San Antonio, Texas, vol. 1, no. 6, pp. 1811-1815, 2001.
- [26] E. Hashem, "Analysis of Random Drop for Gateway Congestion Control," *Report LCS TR-465*, 1989.
- [27] W. C. Feng, D. D. Kandlur, D. Saha and K. G. Shin, "Blue: A New Class of Active Queue Management Algorithms," *Univ. of Michigan, Ann Arbor, Tech. Rep. CSE-TR-387-99*, 1999.
- [28] S. Athuraliya, D.E. Lapsley, and S.H. Low "Random Exponential Marking for internet congestion control" *IEEE Transactions on Network*, June 2001.
- [29] T. Bhaskar Reddy, A. Ahammed, and Reshma Banu, "Performance Comparison of Active Queue Management Techniques," *IJCSNS International Journal of Computer Science and Network Security*, vol. 9, no. 2, February 2009.
- [30] T. Bonald, M. May, and J. C. Bolot, "Analytic Evaluation of RED Performance," *Proc. IEEE INFOCOM*, Tel-Aviv, Israel, pp. 1415-1424, 2000.
- [31] K. Ramakrishnan, S. Floyd, and D. Black, "The Additional of Explicit Congestion Notification (ECN) to IP," *RFC 3168*, 2001.

- [32] S. Yi, M. Keppes, S. Garg, X. Deng, G. Kesidis, and C. Das, "Proxy-RED: An AQM Scheme for Wireless Local Area Networks," *Proc. ICCCN*, pp. 460-465, 2004.
- [33] S. Biaz, and X. Wang, "RED for Improving TCP Over Wireless Networks," International Conference on Wireless Networks, Las Vegas pp. 628—636, 2003.
- [34] C. Hollot, V. Misra, D. Towlsey, and W. Gong, "Analysis and Design of Controllers for AQM Routers Supporting TCP Flows," *IEEE Transactions on Automatic Control*, vol. 47, no. 6, pp 945-959, 2002.
- [35] C. V. Hollot, V. Misra, D. Towlsey, and W. Gong, "A Control Theoretic Analysis of RED," in *Proceedings of IEEE INFOCOM*, Anchorage, vol. 3, pp. 1510-1519, USA, April 2001.
- [36] S. H. Low, F. Paganini, J. Wang, and J. C. Doyle, "Linear stability of TCP/RED and a Scalable Control," *Computer Networks*, vol. 43, no.5, pp. 633-647, December 2003.
- [37] W. Chen and, S. H. Yang, "The Mechanism of Adapting RED Parameters to TCP Traffic," *Computer Communications*, vol. 32, no. 13-14, 2009.
- [38] W. C. Feng, D. D. Kandlur, D. Saha, and K. G. Shin, "A Self-Configuring RED Gateway," in *Proceedings of IEEE INFOCOM*, pp. 1320-1328, New York, USA, Mar 1999.
- [39] H. Sirisena, A. Haider, and K. Pawlikowski, "Auto-Tuning RED for Accurate Queue Control," in *Proceedings of IEEE GLOBECOM*, vol. 2, pp. 2010-2015, November 2002.

- [40] L. Tan, W. Zhang, G. Peng, and G. Chen, "Stability of TCP/RED Systems in AQM Routers," *IEEE Transactions on Automatic Control*, vol. 51, no. 8, pp 1393-1398, August 2006.
- [41] H. Y. Zadeh, A. Habibi, H. Jafarkhani, and C. Bauer, "Optimal Statistical Tuning of the RED parameters," in *Proceedings of IEEE ICC*, pp. 27-32, Beijing, China, May 2008.
- [42] S. Low, F. Paganini, J. Wang, S. Adlakha, and J. C. Doyle, "Dynamics of TCP/RED and a Scalable Control," in *Proceedings of INFOCOM*, , pp. 239-248, 2002.
- [43] D. Lin., R. Morris., "Dynamics of Random early Detection," *Proceedings of ACM SIGCOMM*, pp. 127-137, Octobet 1997.
- [44] M. Parris, K. Jeffay , and F. D. Smith, "Lightweight Active Router-Queue Management for Multimedia Networking" *Multimedia Computing and Networking 1999, Proceedings, SPIE Proceedings Series, Volume 3654*, pp. 162-174, San Jose, CA, January 1999.
- [45] T. J. Ott, T. V. Lakshman, and L. Wong, "SRED: Stablised RED," in *IEEE INFOCOM* , pp. 1346-1355, March 1999.
- [46] B. Zheng ,and M. Atiquzzaman, "DSRED: An active Queue Management Scheme for Next Generation Networks," *Proceedings of 25th IEEE conference on Local Computer Networks LCN*, pp. 242-251, November 2000.
- [47] J. Koo, B. Song, K. Chung, H. Lee, and H. Kahng, "MRED: a New Approach to Random Early Detection," *Proceedings of the 15th International Conference on Information Networking*, pp. 347, 2001.

- [48] S. Floyd., R. Gummadi, and S. Shenkar, "Adaptive RED: An algorithm for Increasing the robustness of RED's active Queue Management," [online] <http://www.icir.org/floyd/red.html>.
- [49] C. Wang, B. Liu, Y. Hou, K. Sohraby, and Y. Lin, "LRED: A Robust Active Queue Management Scheme Based On Packet Loss Ration," in *Proceedings of 23rd INFOCOM*, vol. 1, pp. 1-12, March 2004.
- [50] M. Claypool, R. Kinicki, and M. Hartling, "Active Queue Management for Web Traffic," *IEEE International Conference on Performance, Computing and Communication*, pp. 531-538, 2004.
- [51] L. Hu, and A. D. Kshemkalyani, "HRED: A Simple and Efficient Active Queue Management Algorithm," *13th International Conference on Computer Communications and Networking ICCCN* , pp. 387-393, October 2004.
- [52] Y. Xu, Z. Y. Wang, and H. Wang, "ARED: A Novel Adaptive Congestion Controller," *IEEE International Conference on Machine Learning and Cybernetics*, vol. 2, pp. 708-714. August 2005.
- [53] J. H. Seol, K. Y. Lee, and Y. S. Hong, "Performance Improvement of Adaptive AQM Using the Variation of Queue Length," *IEEE Region 10 Conference TENCN*, pp. 1-4, November 2006.
- [54] T. Yamaguchi, and Y. Takahashi, "A queue Management Algorithm for Fair Bandwidth Allocation", *Computer Communications*, vol. 30, no. 9, pp. 2048-2059, April 2007.

- [55] S. Chen, Z. Zhou, and B. Bensaou, "Stochastic RED and Its Applications," *ICC*, pp. 6362-6367, 2007.
- [56] S. Liu, T. Basar, and R. Srikant, "Exponential-RED: A Stabilizing AQM Scheme for Low- and High-Speed TCP Protocols," *IEEE/ACM Tran. Networking*, vol. 13, no. 5, pp. 1068-1081, 2005.
- [57] S. Kunniyur, and R. Srikant, "An Adaptive Virtual Queue (AVQ) for Active Queue Management," *IEEE/ACM Transactions on Networking*, vol. 12, pp. 286-299, April 2004.
- [58] C. Long, B. Zhao, and X. Guan, "SAVQ: Stabilized Adaptive Virtual Queue Management Algorithm," *IEEE Communications Letters*, vol. 9, no. 1, pp. 78-80, January 2005.
- [59] M. K. Agarwal, R. Gupta, and V. Kargaonkar, "Link Utilization Based AQM and its Performance," *IEEE GLOBECOM*, vol. 2, pp. 713-718, December 2004.
- [60] H. Kong, N. Ge, F. Ruan, and C. Feng, "A Control Theoretic Analysis of the AQM algorithm GREEN," *unpublished manuscript, E&E Dept. Tsinghua Univ. Beijing 100084, P.R.China*, 2002.
- [61] B. Wydrowski, and M. Zukerman, "GREEN: An Active Queue Management Algorithm for a Self Managed Internet," *Proceedings of ICC*, New York, pp. 2631-2635, 2002.
- [62] X. Deng, S. Yi, G. Kesidis, C. R. Das, "Stabilised Virtual Buffer (SVB)-An Active Queue Management Scheme for Internet Quality of Service," *IEEE GLOBECOM*, vol. 2, pp. 1628-1632, November 2002.

- [63] J. Sun, and M. Zukerman, "RaQ: A Robust Active Queue Management Scheme Based on Rate and Queue Length," *Computer Communications*, vol. 30, no. 8, pp. 1731-1741, February 2007.
- [64] B. Mandelbrot. "Self-similar Error Clusters in Communication Systems and the Concept of Conditional Stationarity," In *IEEE Transactions on Communication Technology COM-13*, pp. 71--90, 1965.
http://ieeexplore.ieee.org/xpl/freeabs_all.jsp?arnumber=1089090.
- [65] V. Frost, and B. Melamed, "Traffic Modeling for Telecommunications Networks," *IEEE Communications Magazine*, vol. 32, no. 3, pp. 70-80, March, 1994.
http://ieeexplore.ieee.org/xpl/freeabs_all.jsp?arnumber=267444.
- [66] H. Fowler and W. Leland, "Local Area Network Traffic Characteristics, with Implications for Broadband Network Congestion Management," *IEEE JSAC*, vol. 9, no. 7, pp. 1139- 1149, September, 1991.
- [67] C. Williamson, "Internet Traffic Measurement,"
<http://pages.cpsc.ucalgary.ca/~carey/papers/2001/measurements.pdf>
- [68] H. Jiang, and C. Dovrolis, "Why is the Internet Traffic Bursty in Short Time Scales," *Performance Evaluation*, vol. 33, no. 1, pp. 241-252, June 2005.
- [69] V. Paxson, and S. Floyd, "Wide-Area Traffic: The Failure of Poisson Modeling", *IEEE/ACM TON*, vol. 3, pp. 226-244, 1995.
- [70] A. Adas, "Traffic Models in Broadband Networks," *IEEE Communications Magazine*, vol. 35, no. 7, pp. 82-89, Jul. 1997.

- [71] H. Heffes, and D. M. Lucantoni, "A Markov Modulated Characterization of Packetized Voice and Data Traffic and Related Statistical Multiplexer Performance," *IEEE Journal Selected Areas Communication*, vol. 8, no. 6, pp. 856-868, 1986.
- [72] G. Min, and M. Ould-Khaoua, "Message Latency in Hypercubic Computer Networks with Bursty Traffic Pattern," *Journal of Computers and Electrical Engineering*, vol. 30, no. 3, pp. 207-222, Elsevier Science, 2004.
- [73] K. Fendick, and W. Whitt, "Measurements and Approximations to Describe the Offered Traffic and Predict and Average Workload in a Single-Server Queue," *Proc. IEEE*, vol. 77, pp. 171-194, 1989.
- [74] Z. Cui, and A.A. Nilsson, "The Impact of Correlation on Delay Performance of High Speed Networks," *Proc. Southeastern Symposium on System Theory*, pp 371-374, 1994.
- [75] L. Kleinrock, *Queueing Systems, Vol 1 Theory*. New York, John Wiley & Sons, 1975.
- [76] P. Kuusela, and J. Virtamo, "Modeling RED with Two Traffic Classes," *Proc. NTS-15*, pp. 271-282, Lund, Sweden, 2000.
- [77] P. Kuusela, P. Lassila, J. Virtamo, and P. Key, "Modeling RED with Idealized TCP Sources", *Proc. IFIP conference on Performance Modelling and Evaluation of ATM & IP Networks*, pp. 155-166, Budapest, Hungary, 2001.

- [78] A. Mokhtar, and M. Azizoglu, "A Random Early Discard Frame Work for Congestion Control in ATM Networks," *Proceeding of IEEE ATM Workshop*, Tokyo, Japan, pp. 45-50, 1999.
<http://ieeexplore.ieee.org/iel5/35/13111/00601746.pdf?isnumber=&arnumber=601746>
- [79] H. M. Alazemi, A. Mokhtar, and M. Azizoglu, "Stochastic Modeling of Random Early Detection Gateways in TCP Networks," *IEEE Global Telecommunication Conference (GLOBECOM)*, San Francisco, CA, USA, pp. 1747-1751, 2000.
- [80] H. M. Alazemi, A. Mokhtar, and M. Azizoglu, "Stochastic Approach for Modeling Random Early Detection Gateways in TCP/IP Networks," *International Conference on Communications (ICC)*, pp. 2385-2390, 2001.
- [81] M. Barbera, A. Laudani, A. Lombardo, and G. Schembra, "A New Fluid-Based Methodology to Model AQM Techniques with Markov Arrivals," *Quality of Service in Multiservice IP Networks (QoS-IP), LNCS 2601*, pp. 358-371, 2003.
- [82] H. Abdel-Jaber, M. Woodward, F. Thabtah, and A. Abu-Ali, "A Discrete-time Queue Analytical Model based on Dynamic Random Early Drop," *Proc. Int. Conf. Information Technology*, pp. 71-76, 2007.
- [83] H. Abdel-Jaber, M. Woodward, F. Thabtah, and A. Abu-Ali, "Performance evaluation for DRED discrete-time queueing network analytical model," *Journal of Network and Computer Applications*, vol. 31, no. 4, pp. 750-770, 2008.
- [84] J. D. C. Little, "A Proof of the Queueing Formula $L = \lambda W$," *Operations Research*, vol. 9, pp. 383-387, 1961.

- [85] J. BANKS, J. S. Carson, and B. L. Nelson. “Discrete-Event System Simulation”, second edition. Prentice-Hall, Englewood Cliffs, NJ, 1996.
- [86] A. M. Viterbi, “Approximate Analysis of Time-Synchronous Packet Networks,” *IEEE Journal on Selected Areas in Comrrrnzcatons (SAC)*, vol. 4, no. 6, pp. 879-890, 1986.
- [87] I. D. Moscholios, M. D. Logothetis, and G. K. Kokkinakis, “Call–burst Blocking of ON–OFF Traffic Sources with Retrials under the Complete Sharing Policy,” *Performance Evaluation*, vol. 59, no. 4, pp. 279-312, 2005.
- [88] A. E. Kamal and S. Sankaran, “A Combined Delay and Throughput Proportional Scheduling Scheme for Differentiated Services”, *Journal Computer Communications*, vol. 29, pp. 1754-1771, 2006.
- [89] Y. Z. Cho, and C. Kwan, “Analysis of the M/G/1 Queue under a Combined Preemptive/Nonpreemptive Priority Discipline,” *IEEE/ACM Trans. Communications*, vol. 41, no. 1, pp. 132-141, 1993,.
- [90] B. Suter, T. Lakshman, D. Stiliadis, and A. Choudhury, “Design Consideration for Supporting TCP with Per-flow Queueing,” *IEEE Proc. INFOCOM*, pp.299-306, 1998.
- [91] M. Nabeshima, and K.Yata, “Performance Improvement of Active Queue Management with Per-flow Scheduling,” *IEE Proc. Communication*, pp. 797-803, 2005.

- [92] R. Jain, "The Art of Computer Systems Performance Analysis," John Wiley and Sons, New York, 1991.
- [93] M. Garrett, and W. Willinger, "Analysis, Modeling and Generation of Self-Similar VBR Video Traffic," *Proceedings of SIGCOMM*, pp. 269-280, September, 1994.
- [94] W. Leland, M. Taqqu, W. Willinger, and D. Wilson, "On the Self-Similar Nature of Ethernet Traffic (Extended Version)," *IEEE/ACM Transactions on Networking*, vol. 2, no. 1, pp. 1-15, February 1994.
- [95] I. Norros. "A Storage Model with Self-Similar Input," *Queueing Systems*, vol. 16, pp. 387--396, 1994. <http://citeseer.ist.psu.edu/norros94storage.html>
- [96] T. Yoshihara, S. Kasahara, and Y. Takahashi, "Practical Time-Scale Fitting of Self-Similar Traffic with Markova-Modulated Poisson Process," *Telecommunication Systems*, vol. 17, no. 1-2, pp. 185-211, 2001.
- [97] A. Shirazinia, S.M.Safavi, and E.N. Shariati, "On the Performance of Matching MMPP to SRD and LRD Traffic Using Algorithm LAMBDA," *Proc. 3rd Int. Conf. on ICTTA*, pp. 1-6. 2008.
- [98] Y. Xue, H. Nguyen, and K. Nahrstedt, "CA-AQM: Channel-Aware Active Queue Management for Wireless Networks," *IEEE International Conference on Communications (ICC)*, pp. 4773-4778, 2007.

- [99] Q. Y. Xia, X. Jin, and M. Hamdi, "Active Queue Management with Dual Virtual Proportional Integral Queues for TCP Uplink/Downlink Fairness in Infrastructure WLANs," *IEEE Tran. Wireless Commnications*, vol. 7, no. 6, pp.2261-2271, 2008.