

**PERFORMANCE ANALYSIS OF AN ASYNCHRONOUS  
TRANSFER MODE MULTIPLEXER WITH MARKOV  
MODULATED INPUTS**

**A DISSERTATION  
SUBMITTED TO THE DEPARTMENT OF  
ELECTRICAL AND ELECTRONICS ENGINEERING  
AND THE INSTITUTE OF ENGINEERING AND SCIENCE  
OF BILKENT UNIVERSITY  
IN PARTIAL FULFILLMENT OF THE REQUIREMENTS  
FOR THE DEGREE OF  
DOCTOR OF PHILOSOPHY**

**TK  
7872  
E5  
A43  
1993**

**By  
Nail AKAR  
August 1993**

PERFORMANCE ANALYSIS OF AN ASYNCHRONOUS  
TRANSFER MODE MULTIPLEXER WITH MARKOV  
MODULATED INPUTS

A DISSERTATION

SUBMITTED TO THE DEPARTMENT OF ELECTRICAL AND ELECTRONICS ENGINEERING  
AND THE INSTITUTE OF ENGINEERING AND SCIENCE  
OF BILKENT UNIVERSITY  
IN PARTIAL FULFILLMENT OF THE REQUIREMENTS  
FOR THE DEGREE OF  
DOCTOR OF PHILOSOPHY

By

Nail Akar

August 1993

NAİL AKAR  
İstanbul, Türkiye

TK

7872

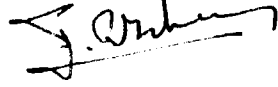
.ES

ALH2

1991

3 024375

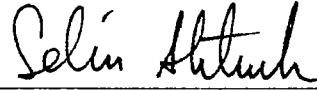
I certify that I have read this thesis and that in my opinion it is fully adequate, in scope and in quality, as a dissertation for the degree of Doctor of Philosophy.



---

Erdal Arıkan, Ph. D. (Supervisor)

I certify that I have read this thesis and that in my opinion it is fully adequate, in scope and in quality, as a dissertation for the degree of Doctor of Philosophy.



---

Selim Aktürk, Ph. D.

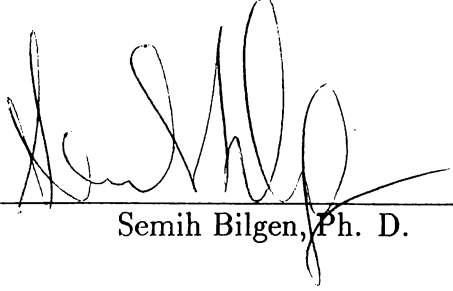
I certify that I have read this thesis and that in my opinion it is fully adequate, in scope and in quality, as a dissertation for the degree of Doctor of Philosophy.



---

Ender Ayanoglu, Ph. D.


I certify that I have read this thesis and that in my opinion it is fully adequate, in scope and in quality, as a dissertation for the degree of Doctor of Philosophy.



---

Semih Bilgen, Ph. D.

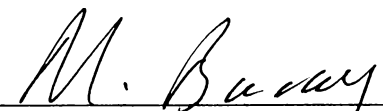
I certify that I have read this thesis and that in my opinion it is fully adequate, in scope and in quality, as a dissertation for the degree of Doctor of Philosophy.



---

Levent Onural, Ph. D.

Approved for the Institute of Engineering and Science:



---

Mehmet Baray, Ph. D.  
Director of Institute of Engineering and Science

# Abstract

## PERFORMANCE ANALYSIS OF AN ASYNCHRONOUS TRANSFER MODE MULTIPLEXER WITH MARKOV MODULATED INPUTS

Nail Akar

Ph. D. in Electrical and Electronics Engineering

Supervisor:

Assoc. Prof. Dr. Erdal Arıkan

August 1993

Asynchronous Transfer Mode (ATM) networks have inputs which consist of superpositions of correlated cell streams. Markov modulated processes are commonly used to characterize this correlation. The first step through gaining an analytical insight in the performance issues of an ATM network is the analysis of a single channel. One objective of this study is the performance analysis of an ATM multiplexer whose input is a Markov modulated periodic arrival process. Based on the transient behavior of the  $nD/D/1$  queue, we present an approximate method to compute the queue length distribution accurately. The method reduces to the solution of a linear differential equation with variable coefficients. Another general traffic model is the Markov Modulated Poisson Process (MMPP). We employ Padé approximations in transform domain for the deterministic service time distribution in an MMPP/D/1 queue so as to compute the distribution of the buffer occupancy. For both models, we also provide algorithms for analysis in the case of finite queue capacities and for computation of effective bandwidth.

**Keywords:** ATM, statistical multiplexing, fluid models, Markov modulated processes, traffic control in ATM networks, effective bandwidth.

# Özet

## MARKOV MODÜLELİ GİRDİLERLE BESLENEN BİR EŞZAMANSIZ AKTARIM MODU ÇOĞULLAYICISININ BAŞARIM ANALİZİ

Nail Akar

Elektrik ve Elektronik Mühendisliği Doktora

Tez Yöneticisi:

Doç. Dr. Erdal Arıkan

Ağustos 1993

Eşzamansız Aktarım Modu (ATM) ağlarının girdileri ilintili paket akışlarından oluşur. Bu ilintiyi tarif edebilmek için genel olarak Markov modüleli süreçler kullanılmaktadır. ATM ağlarını kavrayabilmek için öncelikle tek bir ATM çoğullayıcısının başarım analizini yapmak gerekir. Bu çalışmanın amaçlarından biri girdisi Markov modüleli periyodik varış süreci olan bir ATM çoğullayıcısının başarım analizini yapmaktır. Bu analizi yapabilmek için  $nD/D/1$  kuyruğunun geçici davranışına dayanarak kuyruk uzunluğu dağılımını bulan yaklaşık bir yöntem önerilmektedir. Bu dağılım ise doğrusal ve değişken katsayılı türevsel bir denklemin çözümüyle elde edilir. ATM ağları için genel olarak kullanılan bir başka trafik modeli ise Markov modüleli Poisson sürecidir (MMPP). MMPP/D/1 kuyruğunun dağılımını hesaplamak amacıyla sabit servis zamanı için dönüşüm uzayında Padé yaklaşımları kullanılmıştır. Bu iki model için ayrıca sonlu kuyruk kapasiteleri durumunu inceleyen ve eşdeğer bant genişliği hesaplayan yöntemler önerilmiştir.

### Anahtar

**sözcükler:** ATM, istatistiksel çoğullama, sıvı akış modelleri, Markov modüleli süreçler, ATM ağlarında trafik denetimi, eşdeğer bant genişliği.

# Acknowledgement

I would like to express my deepest gratitude to Dr. Erdal Arıkan for his supervision and invaluable advices in all the steps of the development of this work. His encouragement and his motivating approach based on deadlines contributed a lot in completing my PhD study.

I would also like to thank to Dr. Abdullah Atalar, chairman of our department, for that he encouraged me to go on completing my PhD study in critical and desperate moments and didn't even hesitate once in giving me both moral and technical support.

I am grateful to Dr. Ender Ayanođlu, for his collaboration, guidance, and invaluable advices throughout this study. From miles and miles away, he has always supported me in all aspects even in his busiest moments.

Fortunate to have friends Adil Baktır, Cem Ođuz, Ogan Ocalı, and Gzde Bozdađı for that we have been together here and we will continue to be.

Finally, my sincere thanks are due to my family for their love, patience, and continuous moral support throughout my whole graduate study.



# Contents

Abstract	i
Özet	ii
Acknowledgement	iii
Contents	iv
List of Figures	vii
List of Tables	x
<b>1 Introduction</b>	<b>1</b>
1.1 Asynchronous Transfer Mode . . . . .	1
1.2 Statistical Multiplexing . . . . .	3
1.3 Call Admission Control . . . . .	7
1.4 Traffic Modeling and ATM Multiplexer Performance Analysis . . . . .	8
1.5 Objectives and Outline of the Thesis . . . . .	16

<b>2</b>	<b>Markov Modulated Fluid Sources</b>	<b>21</b>
2.1	Problem Formulation and Analysis . . . . .	23
2.2	An Alternative Formulation . . . . .	27
<b>3</b>	<b>Markov Modulated Periodic Arrival Process</b>	<b>33</b>
3.1	$nD/D/1$ Queue . . . . .	36
3.2	An Approximation to the Transient Behavior of the $nD/D/1$ Queue . . .	37
3.3	MMPAP/D/1 Queue . . . . .	43
3.3.1	Numerical Examples . . . . .	48
3.4	Finite Buffers . . . . .	62
3.5	Effective Bandwidth . . . . .	63
<b>4</b>	<b>Padé Approximations in the Analysis of the MMPP/D/1 System</b>	<b>69</b>
4.1	Transient Analysis of the M/G/1 Queue . . . . .	72
4.2	MMPP/G/1 Queue . . . . .	74
4.3	Padé Approximations in the MMPP/D/1 Queue . . . . .	76
4.3.1	Numerical Examples . . . . .	81
4.4	Computational Aspects . . . . .	83
4.5	MMPP/D/1/K Queue . . . . .	93
4.5.1	Numerical Examples . . . . .	98
4.6	Effective Bandwidth . . . . .	99

<b>5 Conclusions and Suggestions for Future Work</b>	<b>105</b>
<b>Vita</b>	<b>114</b>

# List of Figures

1.1	Approximate ATM traffic performance requirements. . . . .	3
1.2	ATM network model. . . . .	4
1.3	Cells multiplexed on a single link. . . . .	5
1.4	$N \times N$ non-blocking ATM switch: output queueing solution. . . . .	5
1.5	Variable bit rate sources . . . . .	9
1.6	2-state Markov model for an on/off source. . . . .	10
1.7	Birth-death model for the superposition of $N$ on/off sources. . . . .	10
1.8	Statistical multiplexing of on-off sources (MMPAP/D/1 queue). . . . .	11
1.9	The MMPP/D/1 queue. . . . .	12
1.10	Statistical multiplexing of on-off sources (Poisson arrivals during on periods). . . . .	12
1.11	Statistical multiplexing of two-state fluid sources. . . . .	13
1.12	$\sum D_i/D/1$ queue. . . . .	14
1.13	$nD/D/1$ queue. . . . .	15
3.1	Comparison of approximations for the expected value of the queue length for the case $R = 10$ and $n = 8$ (underload). . . . .	39

3.2	Performance of the proposed approximation for the expected value of the queue length for the case $R = 10$ and $n = 12$ (overload). . . . .	40
3.3	Comparison of approximations for the expected value of the queue length for the case $R = 48$ and $n = 40$ (underload). . . . .	41
3.4	Performance of the proposed approximation for the expected value of the queue length for the case $R = 48$ and $n = 50$ (overload). . . . .	42
3.5	Comparison of the queue length survivor function for our proposed method with simulation results and the fluid flow approximation ( $N = 15$ , $R = 10$ , utilization = 0.52). . . . .	49
3.6	Comparison of the queue length survivor function for our proposed method with simulation results and the fluid flow approximation ( $N = 20$ , $R = 10$ , utilization = 0.70). . . . .	50
3.7	Comparison of the queue length survivor function for our proposed method with simulation results and the fluid flow approximation ( $N = 60$ , $R = 48$ , utilization = 0.44). . . . .	52
3.8	Comparison of the queue length survivor function for our proposed method with simulation results and the fluid flow approximation ( $N = 90$ , $R = 48$ , utilization = 0.66). . . . .	53
3.9	Comparison of the queue length survivor function for our proposed method with simulation results and the fluid flow approximation ( $N = 120$ , $R = 48$ , utilization = 0.88). . . . .	54
3.10	Queue length survivor function for $N = 45$ , $N = 75$ , and $N = 105$ when $L_b = 16250$ bytes. . . . .	58
3.11	Queue length survivor function obtained via $j^{th}$ -order approximations for $N = 75$ and $L_b = 16250$ bytes. . . . .	59

3.12	Queue length survivor function for $N = 45$ , $N = 75$ , and $N = 105$ when $L_b = 500$ bytes. . . . .	60
3.13	Queue length survivor function obtained via $j^{th}$ -order approximations for $N = 75$ and $L_b = 500$ bytes. . . . .	61
4.1	Performance comparison of the Padé approximations in terms of the queue length survivor function ( $N = 8$ , utilization = 0.28). . . . .	84
4.2	Performance comparison of the Padé approximations in terms of the queue length survivor function ( $N = 10$ , utilization = 0.35). . . . .	85
4.3	Performance comparison of the Padé approximations in terms of the queue length survivor function ( $N = 15$ , utilization = 0.52). . . . .	86
4.4	Performance comparison of the Padé approximations in terms of the queue length survivor function ( $N = 20$ , utilization = 0.70). . . . .	87
4.5	Cell loss rate approximations ( $N = 40$ , utilization = 0.29). . . . .	100
4.6	Cell loss rate with respect to the buffer size obtained by Padé approximation (2,2) as $N$ is varied. . . . .	101

# List of Tables

1.1	Some services and their characteristics. . . . .	2
1.2	A brief survey of teletraffic analysis of ATM multiplexers. . . . .	20
3.1	Comparison of approximations of the mean waiting time with the simulation results for the case $R = 10$ . . . . .	51
3.2	Comparison of approximations of the mean waiting time with the simulation results for the case $R = 48$ . . . . .	55
4.1	Performance comparison of the Padé approximations in terms of the mean waiting time. . . . .	82

# Chapter 1

## Introduction

### 1.1 Asynchronous Transfer Mode

The Asynchronous Transfer Mode (ATM) is considered by CCITT, International Consultative Committee for Telephone and Telegraph, (now the International Technological Union - Telecommunications Section, or ITU-TS) as the preferred transfer mode for B-ISDN (Broadband Integrated Services Digital Network) [26]. Unlike traditional networks, the B-ISDN will be required to support a wide mix of services (e.g., voice, low- and high-speed data, image and video) over a common ATM transport network. In an ATM based network, all information is conveyed using fixed size packets (called “cells”). To achieve high speed integrated transport, the ATM network adopts a simplified transport protocol based on hardware cell switching [7],[48].

A basic factor that favors ATM is its capability to handle “bursty” traffic via the use of statistical multiplexing. Bursty calls generate traffic at high rates for short periods of time and traffic at much lower rates at other times [22]. Burstiness of a call is simply described in [7] as the ratio between the maximum and the average information rates during the holding time of the call. The average bit rate and the burstiness are important measures to describe the traffic stream associated with a



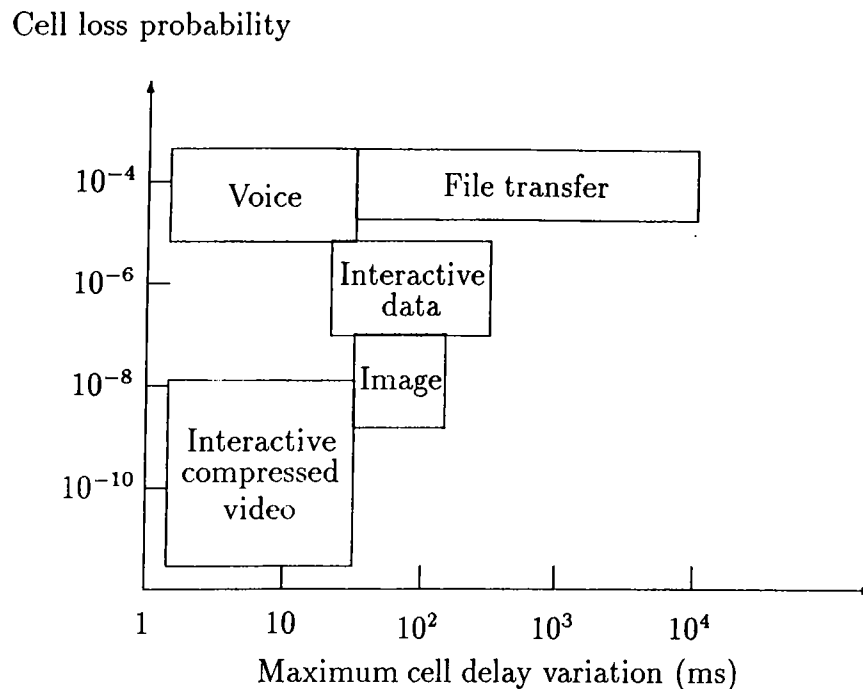
particular service. These two measures of interest actually depend on particular coding and compression techniques used to transport a service. Table 1.1 shows these traffic parameters for certain broadband services in order to demonstrate the diversification of traffic characteristics in the B-ISDN [7]. A service type is a set of services that have the same Quality of Service (QoS) requirements. ATM should satisfy the different

Service type	Mean bit rate	Burstiness
Voice	32 kbits/s	2
Interactive data	1-100 kbits/s	10
Bulk data	1-10 Mbits/s	1-10
Standard quality video	20-30 Mbits/s	2-3
High definition TV	100-150 Mbits/s	1-2
High quality video telephony	~ 2 Mbits/s	5

**Table 1.1:** Some services and their characteristics.

QoS requirements of different services. These requirements are usually measured in terms of maximum delays and cell loss rates. Figure 1.1 shows approximate delay and loss requirements for some expected services [22],[58]. Services such as voice and real time video have strict delay requirements. If cells are not delivered within their delay requirements, they are considered lost due to the real time nature of the services. Delay jitter, the standard deviation concerning delays, should also be small so that the information can be reconstructed in a continuous fashion. In many cases, a certain amount of loss is tolerable although lost cells will have some adverse effects on real-time traffic. Data traffic, such as transfer of files, can generally be characterized by a flexible delay requirement and a strict loss sensitivity.

ATM cells consist of an header and an information field. These cells are transmitted over a virtual circuit and routing is performed based on the information in the header. The cell transmission time is equal to a slot length and slots are allocated to a call on a demand basis. Since bursty traffic does not require continuous allocation of the bandwidth at its peak rate, a large number of bursty traffic sources can share the bandwidth, thus increasing resource utilization. ATM can also support continuous



**Figure 1.1:** Approximate ATM traffic performance requirements.

bit-rate services by allocating bandwidth based on their bit rates. This multiplexing could lead to more efficient use of resources, but may require new kinds of bandwidth management and traffic control. In the next section, we address the statistical multiplexing concept in more detail.

## 1.2 Statistical Multiplexing

An ATM network model is shown in Figure 1.2. The bit stream from an individual source is first segmented into cells at the edge of the network and a header is attached to each cell. The cells are then transported to the destination through the network and the bit stream is reconstructed at the receiving side by stripping the header and “playing out” the cells. In both the access nodes and the output buffers of the intermediate switches, the key factors in performance deterioration are cell losses and delays due to

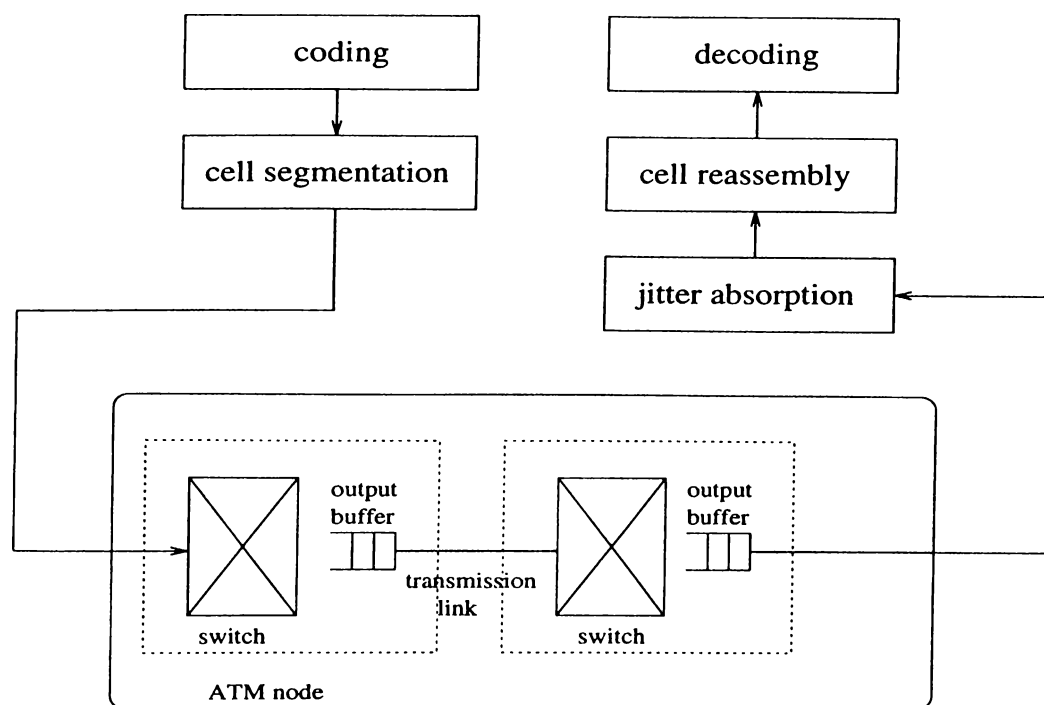
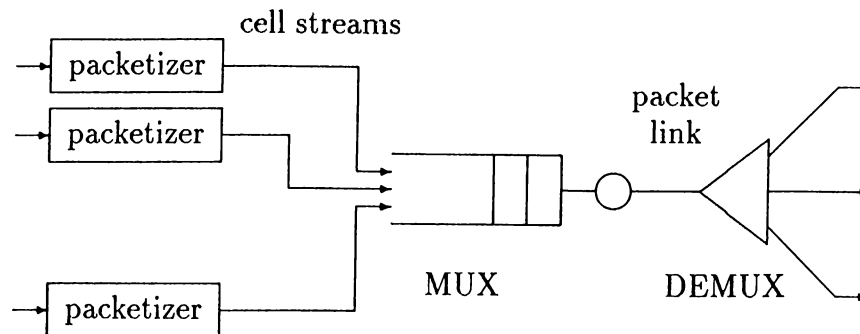


Figure 1.2: ATM network model.

queueing. There are other performance deterioration factors (e.g., cell segmentation delay, cell loss, and cell misdelivery due to header field errors due to transmission, etc.) which are independent of incoming traffic and are out of scope of this dissertation. Our objective in this dissertation is to obtain a fundamental understanding of the queueing characteristics when traffic from several bursty sources are multiplexed on network links. We study this problem for a concentrator where a single link carries multiplexed cell streams (shown in Figure 1.3).

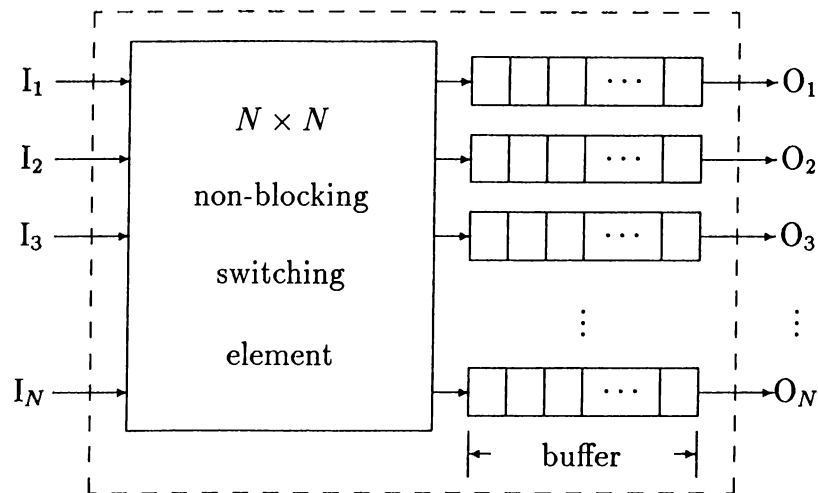
Let us further consider an ATM switch (shown in Figure 1.4) to understand how statistical multiplexing takes place inside an ATM network. We will describe the basic properties of the switch that will yield our ATM multiplexer model.

There are actually different queueing schemes proposed for an ATM switch depending on where the queues are employed (i.e., inputs or outputs). Input queueing solution has a significant throughput limitation [21]. We therefore consider the output queueing solution [21],[53] in which there is a reserved buffer for each output port and the incoming



**Figure 1.3:** Cells multiplexed on a single link.

cells are allowed to use these reserved buffers of the output ports they are destined for. This is in contrast with completely shared buffers (central queueing) [21] where the total



**Figure 1.4:**  $N \times N$  non-blocking ATM switch: output queueing solution.

memory is common to all connections. The approach in output queueing is that, cells of different inputs destined to the same output can be transferred through the switching element during one cell time. However, only a single cell may be served by an output in a cell time, causing possible output contention. This contention is solved by queues which are located at each output of the switching element. The switching device is assumed to be internally non-blocking in the sense that no cell is blocked in the switching fabric

when being transferred to the output ports, cell blocking is only due to possible buffer overflows. The control of the output queues is based on a simple FIFO (first-in-first-out) discipline to ensure that cells belonging to a certain connection will remain in the correct sequence. Priority mechanisms that will change this control are out of scope of this dissertation.

With this kind of an operation, the incoming lines from bursty traffic sources are said to be “statistically multiplexed” on the output port they are destined for. Statistical multiplexing occurs when the capacity of the output channel is less than the sum of the connection peak bandwidths, but is larger than their average total bandwidth requirement. The statistical gain is the factor by which the sum of the peak bandwidths can exceed the output channel’s capacity while satisfying the QoS requirements, or in other words, the throughput gain in using statistical multiplexing instead of deterministic multiplexing (e.g., time or frequency division multiplexing). Statistical multiplexing, therefore, relies on the input channels being bursty due to variable information transfer rates. This statistical gain directly depends on the bandwidth allocation and traffic characteristics of the input channels.

Achieving any statistical gain results in a nonzero probability of cell level overload or congestion. Congestion can be eliminated to a limited extent by using large storage capacity buffers. The buffers will absorb excess information until the sum of the input rates drops below the output rate of the multiplexer. The larger the buffer, the greater the overload that can be absorbed, but this occurs at the expense of large queueing delays which cannot be tolerated by real time applications. Therefore, this delay constraint makes it inconvenient to use very large buffer sizes which would have ensured very low probabilities of buffer overflow.

It is the risk of potential cell losses and delays in a high-speed network which necessitates new traffic control schemes. Teletraffic analysis is necessary to clarify the fundamental properties of statistical multiplexing in ATM networks and to develop effective bandwidth management and congestion control [9],[27]. The next section briefly addresses a particular congestion control strategy which is called the *call admission*

*control.*

### 1.3 Call Admission Control

The design of B-ISDN based on ATM technology depends on the definition of an effective traffic control mechanism capable of guaranteeing required quality of service for a wide variety of connection types. The term “traffic control” includes the actions of routing and resource allocation, necessary for setting up virtual connections as well as the protective measures required to maintain throughput in overload situations [43]. The high transmission speeds and the widely differing traffic characteristics and quality of service requirements require novel procedures for congestion control in ATM networks.

Many of the congestion control schemes developed for existing networks fall into the class of *reactive control*. Reactive control reacts to the congestion after it happens and tries to bring the degree of network congestion to an acceptable level. Due to high transmission speeds, reactive control is, in general, found to be ill-suited for use in ATM networks [2],[62]. Unlike reactive control where control is invoked upon the detection of congestion, *preventive control* does not wait until congestion occurs but attempts to prevent the network from reaching an unacceptable level of congestion. The most common and effective approach is to control the traffic at the entry points to the network (e.g., access nodes). This approach is especially effective due to the connection-oriented feature of ATM networks. With connection-oriented transport, a decision to admit new traffic can be made based on the knowledge of the state of the route which the traffic would follow [59].

One of the major preventive controls is call admission control which determines whether to accept or reject a new connection at the time of call set-up. When a new connection is requested, the network examines the call’s performance requirements (e.g., acceptable end-to-end delay and cell loss probability) and traffic characteristics (e.g., peak rate, mean rate, mean burst length, etc.). The network then examines the current

load and decides whether or not to admit the new call. A call admission policy therefore limits the number of calls in the system so as to give proper QoS guarantees to different services.

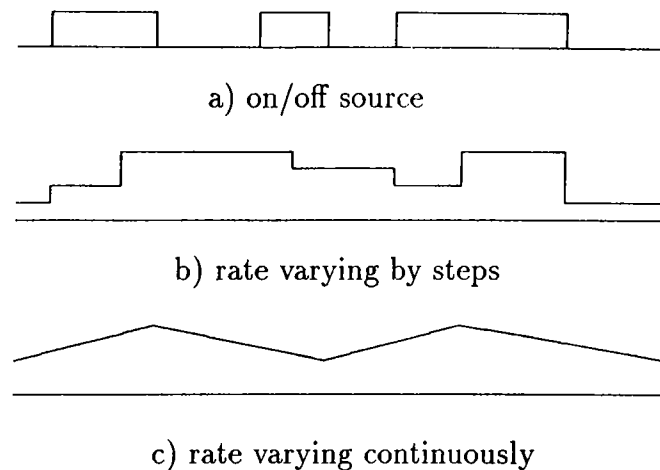
A conservative admission policy (e.g., peak rate bandwidth allocation) allows relatively low loading on the network links and minimizes the probability of cell level congestion. Such an approach, on the other hand, results in a higher level of call blocking relative to a more aggressive admission policy [17]. Therefore, efficient call admission procedures are required, especially for users with predictable traffic parameters, in order to provide an adequate use of network resources. In this dissertation, we will also consider a particular call admission policy that depends on the notion of *effective bandwidth*. For various models, it has been shown that an effective bandwidth can be associated with each source, and that the queue can deliver its performance guarantee by limiting the sources served so that their effective bandwidths sum to less than the capacity of the link.

The characterization of statistical gain mentioned in the preceding section and the definition of an effective bandwidth depend critically on how the traffic is generated. We now give a survey on traffic modeling and ATM multiplexer performance analysis with special emphasis on multiplexers fed by sources as Markov modulated rate processes.

## 1.4 Traffic Modeling and ATM Multiplexer Performance Analysis

When variable bit rate sources (VBR sources) are multiplexed in an ATM network, there arise queues fed by a particular form of correlated arrival process. Accurate traffic modeling is necessary to characterize this arrival process which is composed of a superposition of packet streams generated by these variable bit rate sources. Depending on the bit rate variability, these sources may be classified as (Figure 1.5):

- on-off sources,
- more general piecewise constant rate sources,
- continuously varying rate sources.

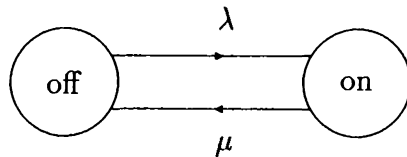


**Figure 1.5:** Variable bit rate sources

Many forms of data-, speech- and image-based communication are expected to exhibit output of the first kind while the latter two may be more typical to multi-media and VBR video communications [40]. In this dissertation, we will rather focus on on-off type source modeling.

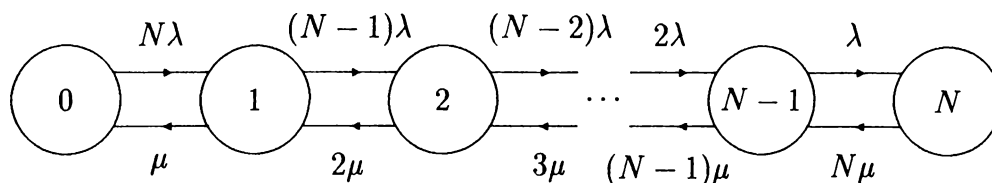
Bit rate variability manifests itself in the network by the changing frequency of cell arrivals. Sources employing constant bit rate coding schemes transmit cells periodically at a frequency determined by their bit rate. On-off sources emit cells periodically during activity periods, or “bursts”, of variable length alternating with silence times, also of variable length. The superposition of on-off sources has been studied, notably, in the context of packetized speech [6],[19] for which the silence times and the activity times are modeled to be exponentially distributed with means  $1/\lambda$  and  $1/\mu$ , respectively [5],[61]. The bit stream belonging to an on-off source is therefore characterized by a 2-state continuous-time Markov chain (Figure 1.6). This 2-state model can easily be extended





**Figure 1.6:** 2-state Markov model for an on/off source.

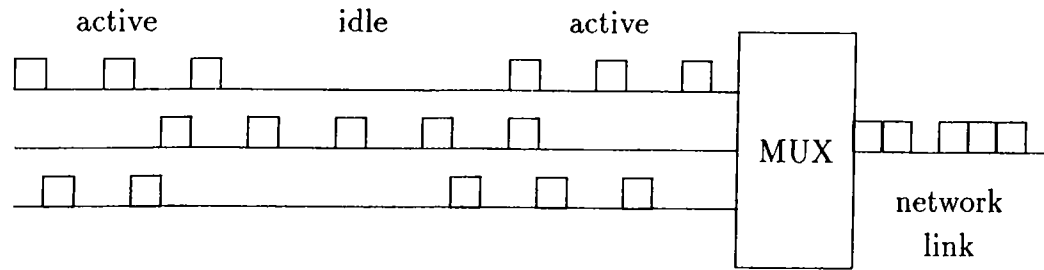
to construct an  $N$ -state Markov chain to describe the superposition process of  $N$  on-off sources of the same type (Figure 1.7). The birth-death model in Figure 1.7 might



**Figure 1.7:** Birth-death model for the superposition of  $N$  on/off sources.

also be used to characterize a single video source without scene changes [2],[46]. In case scene changes are taken into account, the above model should be extended to a multi-dimensional birth-death process [46].

Let us focus our attention to the birth-death process. In an arbitrary state, say  $n$ , of the Markov process whose state holding time is exponentially distributed with parameter  $\sigma_n = (N - n)\lambda + n\mu$ ,  $n$  sources independently transmit cells periodically with the same period. Generally, we call such an arrival process a Markov Modulated Periodic Arrival Process (MMPAP). When such a process is offered to a deterministic server, we call the resulting system the MMPAP/D/1 queue (Figure 1.8). This queueing system has turned out to be one of the most challenging problems of teletraffic theory in recent years due to its practical significance in the ATM context. It has long been known that the apparently convenient device of assuming that the superposition of a large number of independent on-off sources yields a Poisson arrival process can lead to quite inaccurate results [6],[47]. More accurate queueing models must take into account the correlated nature of the cell arrival process which possesses basically two kinds of correlation [47]:

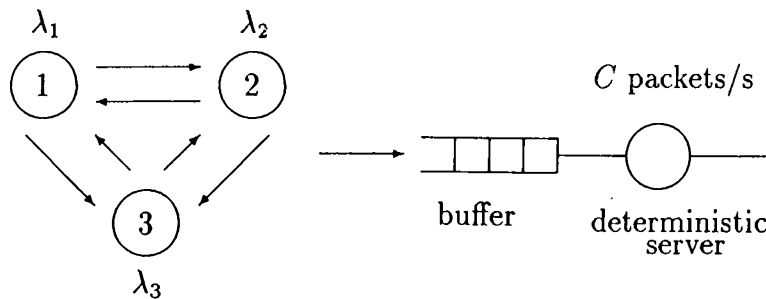


**Figure 1.8:** Statistical multiplexing of on-off sources (MMPAP/D/1 queue).

- negative correlation of cell arrivals in successive time slots due to the periodic cell emissions of active sources,
- positive correlation between the average arrival rates in successive periods of length greater than the inter-cell time of the multiplexed sources.

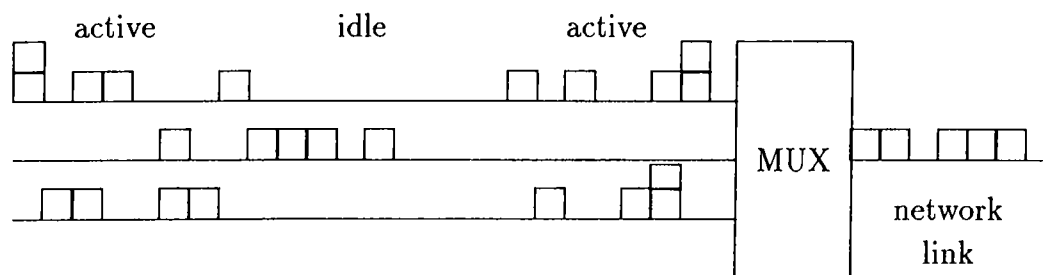
Various modeling approaches in the literature attempt to account for these correlation effects while providing computationally tractable performance analysis schemes. A promising approach is to approximate the superposition nonrenewal point process by a renewal process [47] in which positive correlations are accounted for by the choice of the second moment of the packet interarrival time distribution.

An approach which has proved more popular is to approximate the arrival process by a Markov Modulated Poisson Process (MMPP): the arrival process is governed by the evolution of a discrete-space Markov process; when in state  $n$ , cells are generated according to a Poisson process with intensity  $\lambda_n$ . The resulting queue is called the MMPP/D/1 queue since the packet lengths are fixed in the ATM environment (see Figure 1.9). The more general queueing system named the MMPP/G/1 queue for which packet service times have a general distribution is solved algorithmically in [20] using matrix geometric methods [39]. The technique suggested by Neuts [39] is iterative and has been criticized in [6] to have a slow convergence rate. A 2-state MMPP is proposed in [20] where four parameters (state transition rates and the two arrival intensities) of the MMPP are chosen to match four particular arrival process characteristics of the



**Figure 1.9:** The MMPP/D/1 queue.

superposition process. Other choices for the four fitting parameters have also been proposed in [3],[35] to yield more accurate results. In [25], the superposition of  $N$  on-off sources is modeled by an  $N$ -state MMPP where the arrival intensity is simply proportional to the number of active sources, in other words, an on-off source is assumed to generate packets with respect to a Poisson process in activity times. Figure 1.10 demonstrates the underlying multiplexing system in Ide's work [25]. MMPP models are



**Figure 1.10:** Statistical multiplexing of on-off sources (Poisson arrivals during on periods).

also employed in [45] to characterize not only the packetized voice traffic but also a superposed video arrival process.

On the computation side, when the number of states of the Markov chain increase, numerical problems occur in solving the state equations of the MMPP/D/1 queue to determine the performance measures of interest. Spectral expansion techniques [13] are

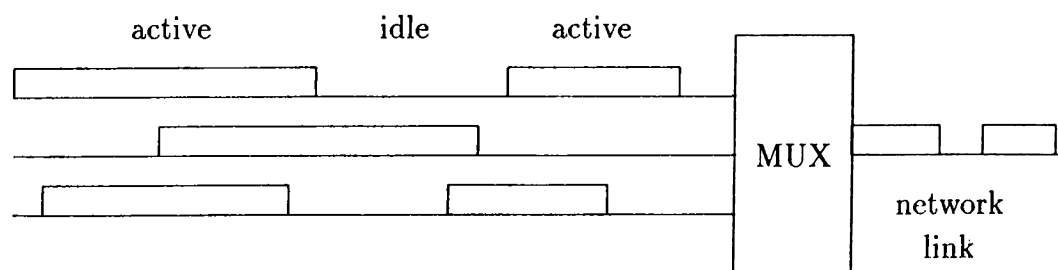
shown to reduce this complexity in the MMPP/M/1 framework for which packet lengths are assumed to be exponentially distributed. The deterministic service time is in general hard to tackle in the MMPP framework. Certain Erlang distributions are therefore used to approximate the deterministic service time distribution in [45],[56],[60].

The use of point process models, such as the MMPP, can be criticized on two counts [40]:

- they do not accurately represent short term correlation effects,
- performance evaluation remains complex.

Simpler models, which also capture the long-term correlation characteristics of the arrival process, are obtained through the so-called fluid flow approximations. In these models, the cell arrivals are approximated by uniform and continuous arrival of fluid, in other words, the concept of packetization is absent. This appears to be a reasonable approximation when the cell interarrival times are small compared to the time between arrival rate changes.

Fluid flow models have attracted the attention of many researchers in the telecommunications literature due to their simplicity. The superposition of a finite number of on/off sources is considered in [1] where the arrival rate is modulated with respect to the state of a Markov chain as in MMPP (see Figure 1.11). A computationally



**Figure 1.11:** Statistical multiplexing of two-state fluid sources.

efficient algorithm is also given, however, the model does not give accurate results for low

to moderate traffic when packet layer contention dominates over burst layer contention [36]. The model proposed in [1] is extended for the finite buffer case in [54] to solve for the information loss rate, a critical value in ATM networks. In [49], the authors give a general algebraic theory for separable Markov Modulated Fluid Sources (MMFS). This actually removes the restriction of the on-off type modeling of a single source. In addition, the work presented in [49] is capable of treating a superposition of nonidentical MMFS, thus allowing multi-state and multi-class traffic into the buffer. The common feature of continuous time fluid flow models is that the solution to the queue length distribution is given in terms of a linear differential equation with constant coefficients. Discrete time models with correlated input described in [34] are also the members of the family of fluid flow approximations. In spite of their shortcoming in accurate traffic modeling, many extensions of fluid flow models have been proposed to analyze more sophisticated queueing systems (e.g., queues with overload control [11],[63]). In [46], the authors employ fluid flow models to evaluate the performance of a statistical multiplexer fed with variable bit rate video sources.

The negative correlation between cell arrivals in successive slots is a local phenomenon occurring while the composition of active sources remains constant. When the overall arrival rate remains below multiplex capacity, the system behaves like the so-called  $\sum D_i/D/1$  queue: a superposition of independent periodic sources of possibly different periods and random phase is offered to a deterministic server (see Figure 1.12). The

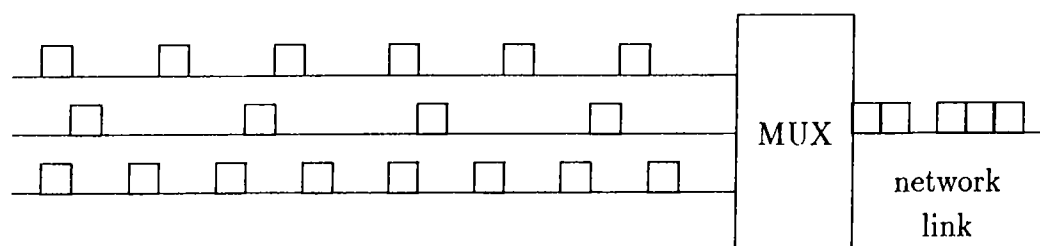


Figure 1.12:  $\sum D_i/D/1$  queue.

system is called the  $nD/D/1$  queue if all the users have an identical period (Figure

1.13). The  $nD/D/1$  queueing system is solved in [8] and revisited more recently in

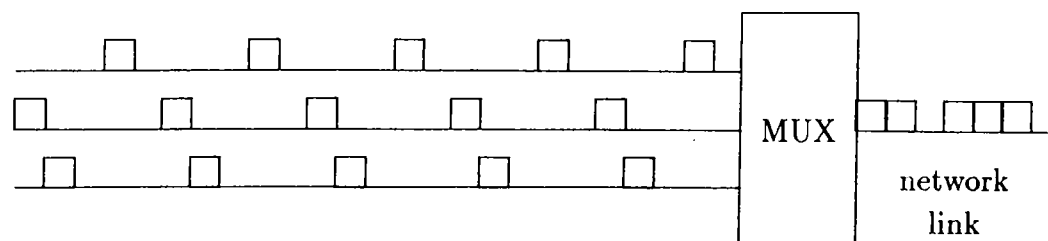


Figure 1.13:  $nD/D/1$  queue.

[4],[44],[57]. Among these approaches, the technique in [4] based on the Ballot theorems [52] seems to be the most efficient one in terms of computational complexity. The more general superposition of sources with different periods is considered in [31] and [44], where accurate approximate formulas for the queue length distribution are derived.

The concept of effective bandwidth has been used to propose admission control policies in ATM based networks. In [23] and [24], Hui has shown that for a simple model of an unbuffered resource, the probability of resource overload can be held below a desired level by requiring that the number of calls  $N_i$  accepted from sources of class  $i$ ,  $i = 1, 2, \dots, m$ , satisfies

$$\sum_i e_i N_i < C,$$

where  $C$  is interpreted as the capacity of the resource, and  $e_i$  is the effective bandwidth of each source of class  $i$ . Kelly [28], Gibbens and Hunt [16], Guerin, Ahmadi and Naghshineh [18], and Elwalid and Mitra [12] offer different approaches to effective bandwidth for buffered resources. Kelly finds effective bandwidth for GI/G/1 queues (queues with general and independent interarrival and service time distribution). In [16], effective bandwidth of on-off type fluid sources is derived for the asymptotic regime of large buffers and small buffer overflow probabilities. Guerin et al. [18] independently obtain the formulas in [16] and extend them through heuristics. In [12], the authors extend the results of [16] to multi-state sources both in the MMFS and MMPP/M/1 frameworks. They show that the effective bandwidth of a Markovian source is the

maximal real eigenvalue of a matrix derived from the source and channel characteristics, and of dimension equal to the number of states.

## 1.5 Objectives and Outline of the Thesis

In this dissertation, we have considered the queueing analysis of a statistical multiplexer which plays a fundamental role in the performance evaluation of ATM networks. The system of interest is a FIFO buffer located at one of the output ports of an ATM switch which is capable of multiplexing variable bit rate sources. What makes the problem challenging is that the interarrival times of the incoming cell streams to the multiplexer are correlated. Teletraffic modeling approaches attempt to characterize this correlation to provide computationally tractable analysis schemes. For voice and video sources, it has widely been accepted that, the arrival rate of information to the multiplexer changes with respect to the state of an underlying continuous-time, discrete-state Markov process. These type of arrivals are called Markov modulated rate processes. In this model, we also need to specify the distribution of the interarrival times of the cell arrivals whose rates are governed by a Markov process in order to have a complete characterization of the input traffic. Among the continuous-time approaches, fluid flow models, periodic arrival processes, and Poisson processes are essentially used in the literature to capture this cell generation process.

The case of a buffer offered with a Markov modulated periodic arrival process is the most accurate model for a wide variety of input traffic types, including voice, video, and interactive data. Despite the accuracy in traffic modeling, no exact solution is available for the so-called MMPAP/D/1 queue. Fluid flow approximations and MMPP-based approaches are among the most popular techniques that attempt to give a solution for the buffer occupancy or the waiting time in this system. These proposed methods in general suffer from inaccuracy since they are incapable of capturing the short term cell scale fluctuations. Fluid models have especially attracted the attention of many researchers in this field due to the ease of computation of the performance measures of

interest despite the inaccuracies encountered in low to moderate traffic regimes. MMPP based models (i.e., MMPP/D/1 queueing system) seem to be more appropriate in terms of accuracy but they suffer from numerical problems especially when the number of states in the Markov chain are large. Our main goal in this dissertation is improving the accuracy of the fluid flow approximations by better traffic modeling but preserving its ease of computation. In other words, our objective is to derive the queue length distribution in both the MMPAP and MMPP frameworks while making use of fluid flow techniques.

We now describe the contributions of the thesis and the significance of the results we have obtained.

a) *Fluid flow approximations.*

- A new derivation of the queue length distribution is provided in transform domain.
- The underlying method in this derivation is readily extendible to more sophisticated queueing systems (i.e., MMPP/D/1 queue), basically an appropriate characterization of the transient behavior of a simpler system (i.e., M/G/1 queue) is required.

b) *ATM multiplexer analysis offered with a superposition of on-off sources.*

We extend the fluid flow technique by incorporating also the short-term cell layer fluctuations, in an approximate way, within the same model. For the case when the system is momentarily underloaded and the number of active sources is fixed, a simple relation is derived. This relation shows that, over complete periods, the queue length evolves as the maximum of a fluid flow term and the queue length in equilibrium. This relation is then used to obtain the following results.

- Via a linear interpolation of the queue length for the  $nD/D/1$  queue which is exactly known at certain time epochs, a new approximation is proposed for the MMPAP/D/1 queue.



- The solution to the queue length distribution is given in terms of a linear differential equation as in the fluid queue. The difference is that, in the fluid queue, the coefficients of the differential equation are constants, in the solution presented here, they are variable.
- This approximation captures the short term cell scale fluctuations and is therefore able to approximate the queue length distribution accurately irrespective of the utilization in the system. Assessment of the approximation's performance is made via a numerical study of a packetized voice multiplexer.
- The solution procedure is quite similar to fluid flow approximations, the essential difference being the determination of a certain linear operator obtained by a number of matrix exponentiations and matrix multiplications. Methods that can decrease the computational effort in computing this linear operator are presented through numerical examples.
- The case of finite buffers is also investigated. The underlying method is based on an extension of [54] where fluid flow approximations are used to solve for a packetized voice multiplexer of finite size.
- An effective bandwidth may be assigned to an MMPAP in the asymptotic regime of large buffers and small overflow probabilities which is the same as assigned to an MMFS.

c) *MMPP/D/1 queue.*

We provide a novel proof for the transform expression of the unfinished work in an MMPP/G/1 queue based on Takács' integro-differential equation that describes the transient behavior of the M/G/1 queue. The deterministic service time distribution is then approximated by several Padé approximations of different orders in transform domain. A Padé approximation is simply a rational function for which a number of first coefficients of its Taylor series expansion match with those of the original function. In our case, the original function is the Laplace transform of the probability density function of the deterministic service time. The number of coefficients to be matched depends on the order of the particular Padé

approximation. The significance of these results lie under the fact that the algebraic theory developed for Markov modulated fluid sources [49] and the MMPP/M/1 system [13] is readily extendible to the MMPP/D/1 queue using the Padé theory. Our results are:

- Instead of Erlang distributions, Padé approximations in transform domain are employed for the deterministic service time which give more accurate results when the computational complexities of these proposed methods are forced to be the same. The underlying reason is that, use of Padé approximations allows one to exactly match the higher order moments of the deterministic service time distribution whereas Erlang distributions don't have this nice property. To give an example, the zero variance of the deterministic service time which plays a critical role in the performance of the queueing system can be captured by a simple Padé approximation. On the other hand, no matter how one can choose the degree of approximation in using Erlang distributions, the zero variance cannot be captured exactly.
- A simple relationship between the fluid flow models and the MMPP/D/1 queue in transform domain is obtained via the use of Padé theory.
- The approximations proposed for the MMPP/D/1 system follow closely the fluid flow methodology and may benefit from the results obtained in the literature for the fluid models. This benefit is shown to be possible if finding computationally efficient algorithms is of concern. In order to demonstrate the viability of this benefit, a procedure is given when the input traffic is a superposition of many 2-state MMPP's of the same type.
- The extension to finite buffers (i.e., MMPP/D/1/K queue) is also presented. The computational complexity of the proposed algorithm is independent of the buffer size and therefore, the computation is tractable even for large buffer sizes.
- An effective bandwidth assignment is shown to be possible for an MMPP in the asymptotic regime of large buffers and small overflow probabilities.

We believe that Table 1.2 will be helpful in clarifying the issues encountered in traffic modeling and performance evaluation of ATM networks. This table attempts to summarize the previous work and the methods we propose in certain perspectives to yield an easy understanding. In this table, we present the queueing models used in certain references and in this dissertation.

reference	no. of act. sources	arrival type	serv. time distr.	no. traf. classes	model app./exact	solution app./exact
Anick [1]	Markov mod.	fluid	fluid	1	app.	exact
Stern [49]	Markov mod.	fluid	fluid	>1	app.	exact
Heffes [20]	Markov mod.	Poisson	general	1	app.	exact
Elwalid [13]	Markov mod.	Poisson	exp.	>1	app.	exact
Bhargava [4]	fixed	periodic	determ.	1	exact	exact
Roberts [44]	fixed	periodic	determ.	>1	exact	app.
Chapter 3	Markov mod.	periodic	determ.	1	exact	app.
Chapter 4	Markov mod.	Poisson	determ.	>1	app.	app.

**Table 1.2:** A brief survey of teletraffic analysis of ATM multiplexers.

The organization of the material is as follows. Chapter 2 is devoted to the analysis of an ATM multiplexer with MMFS models. We then examine the MMPAP/D/1 queue in Chapter 3 and propose an approximate technique to evaluate the queue length distribution in this system. The objective of Chapter 4 is the analysis of the MMPP/G/1 queue, and in particular, the MMPP/D/1 system. Conclusions and suggestions for future work are given in Chapter 5.

## Chapter 2

# Markov Modulated Fluid Sources

In ATM networks, information arrives to the multiplexer at a rate which fluctuates randomly, often with a high degree of correlation in time as explained in the preceding chapter. Accurate capture of these statistical fluctuations is facilitated by modeling the time-varying arrival rate to be governed by a Markov process. If the information arrives uniformly on each line of the multiplexer with a rate controlled by the state of the Markov process and the server similarly removes information from the queue uniformly, then this model is generally called the Markov modulated fluid model and finds its roots in the works of [1],[15]. This model is also called the uniform arrival and service model (UAS) in the packetized voice framework [54].

The performance of the multiplexer when the traffic offered is fixed, has two distinct components corresponding to congestion phenomena, which are generally referred to as cell layer congestion and burst layer congestion [42]. Let us have in mind a superposition of homogeneous on-off sources. Cell layer congestion occurs due the simultaneous arrival of cells from independent sources when the overall cell arrival rate due to active sources is less than the multiplex capacity. Burst layer congestion occurs when the overall arrival rate exceeds the multiplex capacity; buffer content continues to grow as long as the arrival rate excess exists.

As many authors have noted, the fluid flow models are well matched to the ATM environment at the burst layer [11],[40],[37]. Several major reasons have been mentioned:

- the small and uniform cell size and the constant interarrival time of the cells in a burst (periodic packet arrivals for continuous bit oriented (CBO) sources) fit easily in the fluid framework and are difficult to handle in the queueing framework,
- the computational complexity encountered in solving the fluid models in the finite buffer case does not depend on the buffer size while this complexity increases in the queueing model.

The major disadvantage of the fluid model is that, it cannot handle the short-term queue length increases at the cell layer since it removes the concept of packetization from the real arrival process. This is actually why the fluid flow approximation techniques generally do not produce accurate results in light to moderate traffic regimes particularly when the packet layer contention dominates over the burst layer contention. One of the main goals of this dissertation is to improve the accuracy of the fluid flow approximation by refining upon the source model while taking advantage of the ease in computation encountered in fluid models.

The organization of this chapter is as follows. First, the buffer occupancy and queueing delay expressions are obtained in a general Markov modulated setting. Then, the case of a superposition of two-state on-off sources being fed into a multiplexer is discussed. Finally, a new mathematical formulation is developed in this particular case which yields an expression for the stationary queue length distribution. The formulation here can easily be generalized to buffers with more sophisticated input traffic models including Markov modulated Poisson sources and Markov modulated periodic sources. These models will be investigated in the forthcoming chapters in which the relationships and performance comparison of these models and the fluid flow models will be examined. This is one of the reasons why we include a brief presentation of Markov modulated fluid sources in this dissertation. Except for the alternative mathematical formulation that we propose, the exposition that follows is mainly based on [49].

## 2.1 Problem Formulation and Analysis

Consider a buffer with arrival rate  $\lambda(\mathbf{S}(t))$  where  $\mathbf{S}(t)$  is the state of a finite irreducible Markov process at time  $t$ . Let the service rate be  $C$ . Let  $X(t)$  (non-negative random variable) be the buffer content at time  $t$ . Within the fluid flow framework, the behavior of  $X(t)$  in the infinite buffer case is described by

$$\frac{dX}{dt} = \lambda(\mathbf{S}(t)) - C, \quad X > 0. \quad (2.1)$$

Without any loss of generality,  $\mathbf{s} \in \mathbf{S}$  is assumed to be integer-valued, that is;

$$\mathbf{S}(t) \in \{0, 1, 2, \dots, N\}.$$

In view of ATM multiplexers, the size of the Markov chain,  $N + 1$ , depends on the total number of individual sources that can be multiplexed on a common link. In the sequel, we will describe this dependence when a number of homogeneous on-off sources are statistically multiplexed.

Now let

$$P(t, s, x) = Pr\{\mathbf{S}(t) = s, X(t) \leq x\}.$$

Since the modulating Markov process is finite and irreducible, its equilibrium probabilities

$$\pi_s = \lim_{t \rightarrow \infty} Pr\{\mathbf{S}(t) = s\}$$

exist. The mean arrival rate  $\bar{\lambda}$  to the buffer is expressed as

$$\bar{\lambda} = \sum_{\mathbf{s} \in \mathbf{S}} \pi_s \lambda(\mathbf{s}),$$

by which we can find the system utilization

$$\rho = \bar{\lambda}/C.$$

A necessary and sufficient condition for the existence of equilibrium probabilities  $F(s, x)$  for the joint process  $(\mathbf{S}, X)$  in the infinite buffer case is  $\rho < 1$ . We therefore assume that this condition is fulfilled in which case

$$F(s, x) = \lim_{t \rightarrow \infty} P(t, s, x).$$

Let  $M(s, u)$  be the transition rate from state  $u$  to state  $s$  for the underlying modulating process for  $u \neq s$ , and define

$$M(s, s) = - \sum_{u \neq s} M(u, s).$$

The forward Kolmogorov differential equation defining the function  $P(t, s, x)$  for this system is [49]

$$\frac{\partial P}{\partial t} + d(s) \frac{dP}{dx} = \sum_u M(s, u) P(t, u, x), \quad (2.2)$$

where

$$d(s) = \lambda(s) - C.$$

In order to find the equilibrium probabilities of the joint process, we set  $\frac{\partial P}{\partial t} = 0$  in equation (2.2) to obtain

$$d(s) \frac{\partial}{\partial x} F(s, x) = \sum_u M(s, u) F(u, x). \quad (2.3)$$

The equation (2.3) represents a set of  $N + 1$  linear ordinary differential equations which, with suitable boundary conditions, can be solved uniquely for  $F(\cdot)$ . Without loss of generality, we assume  $\lambda(s) \neq C$  for each  $s$ , otherwise the set of equations in (2.3) become singular. In this case, one equation becomes algebraic and may be removed. Denoting now

$$\begin{aligned} F(x) &= \left[ F(0, x) \quad F(1, x) \quad \cdots \quad F(N, x) \right]^T, \\ D &= \text{diag}\{d(j)\}, \quad j = 0, 1, \dots, N, \\ M &= [M(i, j)], \quad i, j = 0, 1, \dots, N, \end{aligned}$$

equation (2.3) can be rewritten as

$$D \frac{d}{dx} F(x) = M F(x), \quad (2.4)$$

where  $M$  is the transpose of the infinitesimal generator matrix for the underlying Markov process and  $D$  is called the drift matrix. The solution to (2.4) then takes the form

$$F(x) = \sum_n a_n \exp(z_n x) \phi_n,$$

where each pair  $(z_n, \phi_n)$  satisfies the eigenvalue-eigenvector problem

$$zD\phi = M\phi.$$

Now let  $S_-$  and  $S_+$  be the set of states such that  $\lambda(s) < C$  and  $\lambda(s) > C$ , respectively. Also let  $d_-$  and  $d_+$  be the cardinality of the corresponding sets. It is well-known that [1],[49], if the Markov chain is reversible then the differential system described by (2.4) has real eigenvalues, only one at the origin,  $d_+$  negative, and  $d_- - 1$  positive.

The boundary conditions can easily be formed by observing that

- 1)  $F(\infty) = \pi \triangleq \left[ \pi_0 \ \pi_1 \ \cdots \ \pi_N \right]^T$  is the stationary distribution of the underlying Markov process.
- 2)  $a_n = 0$  for  $z_n > 0$ , otherwise the solution for the stationary queue length distribution grows without bound.
- 3) For  $s \in S_+$ , the queue is always increasing, so the queue length cannot be zero. Therefore  $F(s, 0) = 0$  for  $s \in S_+$ .

Employing these boundary conditions, one can obtain the unique solution for the differential equation (2.4). The resulting queue length cumulative distribution function (cdf) is then written by the following expression:

$$Pr\{\text{queue length} \leq x\} = \sum_{n=0}^N F(n, x). \quad (2.5)$$

The problem dealt with is in fact a standard eigenvalue problem and a solution subject to the boundary conditions is, in principle, straightforward. However, it becomes intractable because of its size since the number of equations (e.g.,  $N + 1$  in the above framework) can range from hundreds to tens of thousands in typical situations in ATM based networks. Therefore, special structure of the system equations should be taken into account in order to avoid numerical problems.

We now consider the multiplexing of several calls of on-off type onto a single link with capacity  $C$ . Let  $P$  denote the peak rate of one call in packets/sec. The link is shared



by  $N$  statistically identical and independent calls alternating between active and idle periods, which are assumed to be exponentially distributed with mean values  $\mu^{-1}$  and  $\lambda^{-1}$ , respectively (see Figure 1.11). Each call generates information at a rate  $P$  when active and at rate zero when idle. This model has indeed been used for packet voice with speech detection [6],[36],[54].

The number of active calls at time  $t$ ,  $\mathbf{S}(t)$ , is represented as a continuous-time birth-death process. When  $\mathbf{S}(t) = n$ , the mean arrival rate to the multiplexer is  $Pn$ . If  $p(n, m) \triangleq M(m, n)$  is defined to be the transition rate from state  $n$  to state  $m$ , the birth and death rates are given [54] by

$$\begin{aligned} p(n, n+1) &= (N-n)\lambda, \quad n = 0, 1, \dots, N-1, \\ p(n, n-1) &= n\mu, \quad n = 1, 2, \dots, N. \end{aligned}$$

We also define the total probability flow rate out of state  $n$ ,  $\sigma_n$ ,

$$\sigma_n = (N-n)\lambda + n\mu.$$

Within this framework, (2.4) holds with

$$D = \text{diag}\{Pn - C\}, n = 0, 1, \dots, N,$$

and

$$M = \begin{bmatrix} -\sigma_0 & p(1,0) & & & & & & \\ p(0,1) & -\sigma_1 & p(2,1) & & & & & \\ & p(1,2) & -\sigma_2 & p(3,2) & & & & \\ & & & \ddots & & & & \\ & & & & p(N-2, N-1) & -\sigma_{N-1} & p(N, N-1) & \\ & & & & & p(N-1, N) & -\sigma_N & \end{bmatrix}. \quad (2.6)$$

In [1], this particular structure of the infinitesimal generator matrix is made use of in order to evaluate explicitly the eigenvalues and the eigenvectors of the associated differential system, thus providing a computationally efficient method for the analysis of

a statistical multiplexer in case a single class of traffic is present. There the eigenvalue problem is reduced to a set of uncoupled quadratic equations for this birth-death process.

In practice, the queueing delay distribution may be of greater interest. In fact, the buffer occupancy corresponds, with a change of scale, to the virtual waiting time (delay seen by an arriving cell). Taking into account the change of scale, we have

$$Pr\{\text{delay} \leq t\} = \frac{1}{\alpha N} \sum_{n=0}^N nF(n, Ct), \quad (2.7)$$

where  $\alpha \triangleq \frac{\lambda}{\lambda + \mu}$ , is the average fraction of active calls. This can be verified by observing that an arbitrary cell arrives in state  $n$  with probability  $\frac{\pi_n n}{\alpha N}$ . Then,

$$\begin{aligned} Pr\{\text{cell delay} \leq t\} &= \sum_n Pr\{\text{cell delay} \leq t, \text{cell arrives at state } n\} \\ &= \sum_n Pr\{\text{cell delay} \leq t \mid \text{cell arrives at state } n\} \frac{\pi_n n}{\alpha N} \\ &= \sum_n Pr\{\text{queue length} \leq Ct \mid \text{chain state} = n\} \frac{\pi_n n}{\alpha N} \\ &= \frac{1}{\alpha N} \sum_n nF(n, Ct) \end{aligned}$$

The next section is devoted to our alternative formulation of the same problem using transform domain techniques. The significance of this formulation will be clear when we extend it to more general Markovian sources in the subsequent chapters.

## 2.2 An Alternative Formulation

Consider the same traffic model. Let  $X(t)$  be the buffer content and  $\mathbf{S}(t)$  be the state of the Markov chain at time  $t$ . We then define the following stationary probabilities (as  $t \rightarrow \infty, \Delta t \rightarrow 0^+$ ):

$$F_b(n, x) = Pr\{\mathbf{S}(t) = n\} \bar{F}_b(n, x), \quad (2.8)$$

where

$$\bar{F}_b(n, x) = Pr\{X(t) \leq x \mid \mathbf{S}(t + \Delta t) = n, \mathbf{S}(t) \neq \mathbf{S}(t + \Delta t)\}$$

and

$$F_e(n, x) = Pr\{\mathbf{S}(t) = n\} \bar{F}_e(n, x), \quad (2.9)$$

where

$$\bar{F}_e(n, x) = Pr\{X(t) \leq x \mid \mathbf{S}(t + \Delta t) \neq \mathbf{S}(t), \mathbf{S}(t) = n\}.$$

Note that, since  $\mathbf{S}(t)$  is the state of a continuous-time Markov chain, given  $\mathbf{S}(t)$ , the buffer content  $X(t)$  is independent of  $\mathbf{S}(t + \Delta t)$ . This fact yields

$$\bar{F}_e(n, x) = Pr\{X(t) \leq x \mid \mathbf{S}(t) = n\},$$

and we therefore write

$$F_e(n, x) = Pr\{X(t) \leq x, \mathbf{S}(t) = n\}. \quad (2.10)$$

To interpret,  $\bar{F}_b(n, x)$  is the equilibrium probability that the queue length is less than  $x$  given that a state transition to state  $n$  is about to occur. Similarly,  $\bar{F}_e(n, x)$  is the stationary probability that the queue length is less than  $x$  given that a state transition from state  $n$  is about to occur. In other words, we observe the queue length at the time epochs when state transitions occur and henceforth define the corresponding random variables. Recall that the state holding time at state  $n$  is exponentially distributed with parameter  $\sigma_n$ , which is in fact, the total flow rate out of state  $n$ . Conditioning on the state holding time and by exploiting the fluid flow model (i.e., queue length changes with a rate  $C - Pn$  at state  $n$ ), we can now write

$$F_e(n, x) = \int_0^\infty F_b(n, x + (C - Pn)t) \sigma_n \exp(-\sigma_n t) dt, \quad x \geq 0. \quad (2.11)$$

One can verify by using the equality (2.11) the following relationships:

$$F_e(n, x) = \begin{cases} F_b(n, x) * \left(\frac{-\sigma_n}{Pn - C} \exp\left(\frac{\sigma_n x}{C - Pn}\right) u(-x)\right), & Pn < C \\ F_b(n, x) * \left(\frac{\sigma_n}{Pn - C} \exp\left(\frac{\sigma_n x}{C - Pn}\right) u(x)\right), & Pn > C \end{cases} \quad (2.12)$$

where  $*$  is the convolution operator and  $u(\cdot)$  is the unit step function. In case  $Pn < C$ , the equality holds for  $x \geq 0$ , but the term on the right-hand side may be nonzero for

$x < 0$  whereas  $F_e(n, x)$  must equal zero in this interval. In other words, the expression (2.12) suggests that  $F_e(n, x)$  is the orthogonal projection of the term (when  $Pn < C$ )

$$F_b(n, x) * \left( \frac{-\sigma_n}{Pn - C} \exp\left(\frac{\sigma_n x}{C - Pn}\right) u(-x) \right)$$

onto the positive  $x$ -axis. On defining  $\hat{F}_b(n, s)$  and  $\hat{F}_e(n, s)$  as the Laplace transforms of  $F_b(n, x)$  and  $F_e(n, x)$ , respectively,  $\hat{F}_e(n, s)$  turns out to be the Toeplitz operator with symbol  $H$  operating on  $\hat{F}_b(n, s)$  [14], where

$$\begin{aligned} H(s) &= \int_{-\infty}^{\infty} \frac{-\sigma_n}{Pn - C} \exp\left(\frac{\sigma_n x}{C - Pn}\right) u(-x) dx, \\ &= \frac{\sigma_n}{Pn - C} \frac{1}{\left(s + \frac{\sigma_n}{Pn - C}\right)}. \end{aligned}$$

Then, in case  $Pn < C$  we have [14]

$$\hat{F}_e(n, s) = [H(s)\hat{F}_b(n, s)]_s \quad (2.13)$$

where  $[\cdot]_s$  denotes the stable part of transform  $[\cdot]$ . To explain, since  $F_e$  is a nonnegative random variable, the unstable part of the above transform corresponding to negative queue lengths should be removed when  $Pn < C$ . In regard of this,

$$\hat{F}_e(n, s) = \begin{cases} \frac{\sigma_n}{Pn - C} \frac{\hat{F}_b(n, s) - \hat{F}_b(n, \frac{-\sigma_n}{Pn - C})}{s + \frac{\sigma_n}{Pn - C}} & \text{if } Pn < C, \\ \frac{\sigma_n}{Pn - C} \frac{\hat{F}_b(n, s)}{s + \frac{\sigma_n}{Pn - C}} & \text{if } Pn > C. \end{cases} \quad (2.14)$$

Remark that

$$F_e(n, 0) = \frac{\sigma_n}{C - Pn} \hat{F}_b(n, \frac{\sigma_n}{C - Pn}). \quad (2.15)$$

Our objective now is to express  $F_b$ 's in terms of  $F_e$ 's. For this purpose, we rewrite  $\bar{F}_b(n, x)$  in equation (2.8) as  $t \rightarrow \infty$ ,  $\Delta t \rightarrow 0^+$ :

$$\begin{aligned} \bar{F}_b(n, x) &= \frac{\sum_{m \neq n} Pr\{X(t) \leq x, \mathbf{S}(t + \Delta t) = n, \mathbf{S}(t) = m\}}{\sum_{m \neq n} Pr\{\mathbf{S}(t + \Delta t) = n \mid \mathbf{S}(t) = m\} Pr\{\mathbf{S}(t) = m\}} \\ &= \frac{\sum_{m \neq n} Pr\{X(t) \leq x \mid \mathbf{S}(t + \Delta t) = n, \mathbf{S}(t) = m\} Pr\{\mathbf{S}(t + \Delta t) = n \mid \mathbf{S}(t) = m\} \pi_m}{\sum_{m \neq n} p(m, n) \pi_m \Delta t} \\ &= \frac{\sum_{m \neq n} \bar{F}_e(m, x) p(m, n) \pi_m}{\sum_{m \neq n} p(m, n) \pi_m} \end{aligned}$$

Multiplying the last equality by  $\pi_n$  and recalling the balance equations of the Markov process:

$$\pi_n \sigma_n = \sum_{m \neq n} p(m, n) \pi_m$$

we have

$$\sigma_n F_b(n, x) = \sum_{m \neq n} p(m, n) F_e(m, x),$$

or, in transform domain,

$$\sigma_n \hat{F}_b(n, s) = \sum_{m \neq n} p(m, n) \hat{F}_e(m, s). \quad (2.16)$$

Substituting equations (2.15) and (2.16) into (2.14) and solving for  $\hat{F}_e(n, s)$ , one finally obtains

$$(sI - D^{-1}M) \hat{F}_e(s) = \left[ F_e(0, 0) \quad F_e(1, 0) \quad \cdots \quad F_e(C_0, 0) \quad 0 \quad \cdots \quad 0 \right]^T,$$

where  $\hat{F}_e(s)$  is the Laplace transform of  $F_e(x)$  and  $C_0$  is the largest integer  $n$  such that  $Pn < C$ . This transform equation is actually the transform domain equivalent of the equation (2.4) with the imposed boundary conditions. One can now easily write down the buffer occupancy cdf;

$$\begin{aligned} Pr\{\text{queue length} \leq x\} &= \sum_{n=0}^N Pr\{\text{queue length} \leq x, \text{chain state} = n\}, \\ &= \sum_{n=0}^N F_e(n, x). \quad (\text{by definition (2.10)}) \end{aligned}$$

This kind of an alternative formulation in terms of transform domain equations provides a major advantage; it forms a basis for obtaining similar results for queues and point processes in which the traffic sources are Markov modulated Poisson processes or Markov modulated periodic sources which are the topics of the forthcoming chapters. The stationary probability definitions for  $F_b(n, x)$  and  $F_e(n, x)$  will be the same as well as the interconnecting equations (2.16) for these upgraded models. The underlying reason is that these interconnecting equations are only dependent upon the modulating Markov chain but not the type of arrivals (i.e., Poisson, periodic, etc.). What will

mainly differ is the relation between  $F_b(n, x)$  and  $F_e(n, x)$  as in equation (2.14) which will critically depend on how cell generation takes place. The main approach is to obtain a counterpart to equation (2.14) for the MMPAP/D/1 and the MMPP/D/1 queues through the transient behaviors of the  $nD/D/1$  and  $M/D/1$  systems, respectively.

The emerging high-speed networks, particularly the ATM-based broadband ISDN, are expected to integrate through statistical multiplexing large numbers of traffic sources having a broad range of burstiness characteristics. The fluid flow model is suggested to be a prime instrument for analyzing such systems since it handles the essential characteristics of the traffic process at the burst layer. With this model, besides a single class of traffic with each connection having two states, multi-state and multi-class traffic feeding finite buffers with overload control are also examined in the literature [11],[32] with computationally tractable algorithms. Despite being computationally tractable and extendible for analysis to more complicated queueing systems encountered in ATM networks, fluid flow models do not generally give accurate results for low to moderate loads. In the subsequent chapters, we attempt to overcome this drawback in accuracy by using more accurate source modeling, such as the MMPP, but using the same analytical methods used for solving the fluid models.

A typical instrument for controlling congestion is the admission control which limits the number of calls and guarantees a grade of service determined by the cell loss probability in the multiplexer. Fluid flow models have made it possible to assign an effective bandwidth to each source which is an explicitly identified, simply computed quantity, varying between the mean and peak bit rates of the source depending on its burstiness and the grade of service requirements of the call [12],[16], [18]. This quantity has been shown in the above-mentioned references to yield efficient call admission procedures in the natural asymptotic regime of small cell loss probabilities and large buffer sizes. This in turn enables us to extend the model and analysis to a network of channels using approximations such as the Erlang fixed point procedure for a standard circuit-switched network [29]. One other objective of this study is that the use of the same mathematical framework as in fluid flow models will make it possible to assign an

effective bandwidth to calls of more sophisticated type (e.g., MMPP type).

## Chapter 3

# Markov Modulated Periodic Arrival Process

In this chapter, we focus on a particular category of multiplexers whose inputs consist of periodic packet streams. Periodic packet generation is a major feature of continuous bit stream oriented (CBO) sources. With fixed-length packets as in ATM networks, each CBO source generates packets periodically with the period being the packetization time. Depending on the nature and the bit rate of the underlying application, the periods separating successive packets (or the packetization times) can widely differ. The queueing system that models the sharing of a network link by such incoming connections is actually a single server queue with periodic arrivals and deterministic service times. We call this system as an  $nD/D/1$  queue ( $n$  denotes the number of connections) when all connections have an identical period and a  $\sum D_i/D/1$  queue where multiple periods are allowed to coexist.

In our queueing model, it is assumed that a silence detection mechanism exists for CBO sources in the sense that each user alternates between active (on) and idle (off) times of variable length. Sources generate packets periodically at a constant rate during active times and they generate no data during idle times. The 2-state continuous-time Markov chain model will be used to describe the above-mentioned traffic stream (Figure



1.6). In this model, the idle times and the activity times are exponentially distributed with means  $1/\lambda$  and  $1/\mu$ , respectively. The  $N$ -state Markov chain (Figure 1.7) now describes the superposition process of  $N$  on/off sources where the state of the Markov chain is defined to be the number of active sources. In an arbitrary state, say  $n$ , of the Markov chain whose state holding time is exponentially distributed with parameter  $\sigma_n = (N - n)\lambda + n\mu$ ,  $n$  sources independently transmit cells with an identical period. In general, we call this arrival process a Markov Modulated Periodic Arrival Process (MMPAP). Even though the Markov process that governs an MMPAP is arbitrary in the above definition, throughout this chapter we rather focus on the birth-death model (Figure 1.7) due to its practical significance. This way of traffic modeling enables us to appropriately characterize a packet stream originated by a fixed-bit rate coding scheme employed on a 2-state on/off source with a silence detection mechanism.

Given the network link speed and the traffic parameters of an individual source, we are interested in the probability distribution of the buffer content as a function of the number of users. This distribution is derived for a discrete-time queueing system which operates in a slotted fashion; a slot defines the base unit for data generation and data transmission. In particular, in each slot, the link is capable of transmitting one packet and an active source generates at most one packet within that time. Such a slotted operation is an adequate representation for ATM networks where data transmission is in the form of fixed-size packets (cells). Incoming cells are then transmitted on the network link and stored in the multiplexing buffer when the aggregate input rate exceeds the capacity of the link (Figure 1.8).

The method developed here is valid for discrete-time queueing schemes where the modulating process is a continuous-time Markov chain. This choice is due to the discrete-time operation of ATM multiplexers and the continuous-time nature of the fluid flow approximations on the basis of which we make the performance comparisons. This is significant because, our approach combines the discrete-time nature of periodic arrivals in a slotted system and the continuous-time nature of the underlying Markov chain. The framework presented here can readily be reformulated to cover other models (e.g., both

the multiplexer and the chain work in continuous-time (or in discrete-time)).

Consider now the superposition of on-off sources in Figure 1.8 which is offered to the ATM multiplexer; when the instantaneous arrival rate is less than the link rate and the number of active sources is fixed, the queueing system behaves as the  $nD/D/1$  queue. The change in the number of active sources ( $n$ ) which is governed by the Markov chain in Figure 1.7 forces us to examine the transient behavior of the  $nD/D/1$  queue which turns out to have a crucial role in our analytical approach. Actually, the focus of this chapter is on the derivation of relationships between fluid sources and CBO sources, arrival rates of which are Markov modulated in the same manner, through an approximation of the transient behavior of the  $nD/D/1$  queue. This approximation is mainly based on an interpolation of the queue length of the  $nD/D/1$  queue whose distribution is exactly known at certain epochs. The case of multiple periods have not been investigated due to the lack of exact results in the literature for the steady-state distribution of the queue length of the  $\sum D_i/D/1$  queue. The solution to the buffer content distribution for the overall problem is then reduced to the solution of a linear differential equation with variable coefficients whereas in fluid flow approximations, the corresponding equation is simply linear with constant coefficients. Numerical results are given in order to demonstrate the performance of our proposed performance analysis scheme. Finally, effective bandwidth calculation of on-off sources based on this scheme is presented.

The method used in solving for the steady-state distribution of the queue length for the Markov modulated periodic arrival case is composed of two main stages. The first stage consists of an approximation to the transient behavior of the discrete-time  $nD/D/1$  queue in a continuous-time framework. In the second stage, we extend our results for the  $nD/D/1$  queue to solve for the continuous-time Markov model which characterizes the input traffic. Let us then first consider the  $nD/D/1$  queue.

### 3.1 $nD/D/1$ Queue

Throughout this section, we assume that the number of active users ( $n$ ) is fixed. In our queueing model, the time axis is slotted, where each time slot is as long as the transmission time of a single packet. The packets arriving to the queue are served on a first-come-first-serve basis and the queue has infinite size. Each one of the  $n$  active sources transmits fixed length packets with a period of  $R$  slots, independent of other sources. In an arbitrary frame of  $R$  slots, each input source's packet can be in any of these  $R$  slots with equal probability. The source rate in packets/sec is denoted by  $P$  and the service rate of the buffer is denoted by  $C$ , which actually equals to  $PR$  packets/sec. Without loss of generality, we assume that the departures take place at the beginning of slots, and arrivals during slots. We define the following random variables for  $k = 1, 2, \dots, R$ ,

$$\begin{aligned} Q_k &= \text{queue length at the end of } k^{\text{th}} \text{ slot,} \\ a_k &= \text{number of arrivals in the } k^{\text{th}} \text{ slot.} \end{aligned}$$

Note that

$$a_1 + a_2 + \dots + a_R = n.$$

The queueing discipline is the following:

$$Q_k = \begin{cases} Q_0 & \text{if } k = 0 \\ \max(Q_{k-1} - 1, 0) + a_k & \text{if } k > 0 \end{cases}$$

By iteration on  $k$ , one can check using algebraic manipulations that

$$\begin{aligned} Q_1 &= \max(a_1, Q_0 + a_1 - 1) \\ Q_2 &= \max(a_2, a_1 + a_2 - 1, Q_0 + a_1 + a_2 - 2) \\ &\vdots \\ Q_R &= \max(\tilde{Q}_n, Q_0 + n - R) \end{aligned} \tag{3.1}$$

where the random variable  $\tilde{Q}_n$  is defined via

$$\tilde{Q}_n = \max_{0 \leq j < R} \left( \sum_{l=R-j}^R a_l - j \right). \tag{3.2}$$

Note that if  $n < R$ ,  $Q_r \rightarrow \tilde{Q}_n$  as  $r \rightarrow \infty$  for fixed  $n$ . Let us first focus our attention on the probability distribution for  $\tilde{Q}_n$ . This steady-state queue length distribution when  $n < R$  is explicitly given in [4];

$$Pr(\tilde{Q}_n > q) = \sum_{x=1}^{n-q} \frac{R-n+x}{R-x} \binom{n}{q+x} \left(\frac{x}{R}\right)^{q+x} \left(1 - \frac{x}{R}\right)^{n-q-x}. \quad (3.3)$$

Note that the buffer cannot contain more than  $n$  packets in the steady-state, that is  $Pr(\tilde{Q}_n > q) = 0$ ,  $q \geq n$ . We also define the cumulative distribution function (cdf) of  $\tilde{Q}_n$

$$\tilde{Q}_n(q) \triangleq Pr(\tilde{Q}_n \leq q).$$

In [4], the change in the number of active sources,  $n$ , is assumed to happen slowly. In this case, the queue length reaches its steady-state distribution  $\tilde{Q}_n(\cdot)$  whenever  $n < R$ . The equation (3.3) is then sufficient to compute the distribution of the queue length and the queueing delay assuming that  $n$  does not exceed  $R$  so that the queueing system is stable. Considering the Markov model (Figure 1.7) which governs the number of incoming active sources, we have two significant observations: 1) there are possible overload states ( $n > R$ ) in which case there is no limiting distribution 2) even for the underload states ( $n < R$ ), the state holding time may not be long enough for the queue length to reach its steady-state distribution  $\tilde{Q}_n(\cdot)$  before the Markov chain makes a transition to another state. Therefore, an accurate capture of the transient behavior of the  $nD/D/1$  queue turns out to be the major issue for our purposes. This problem is addressed in the next section.

## 3.2 An Approximation to the Transient Behavior of the $nD/D/1$ Queue

In order to obtain the queue length evolution equations for  $n < R$ , we iterate on equation (3.1) on an  $R$ -slot basis so that by periodicity of arrivals we have

$$Q_{kR} = \max(\tilde{Q}_n, Q_0 + k(n - R)), \quad k = 1, 2, \dots \quad (3.4)$$

There is, in fact, a strong interconnection between periodic models and fluid flow models. In the latter models, information is assumed to arrive uniformly to the multiplexer and the server similarly removes information from the queue, in a continuous manner. The computational tractability and buffer size independent solvability of fluid flow approximation techniques suggest a further study of this interconnection.

If we define  $Q(t)$  as the queue length at time  $t$ , the fluid flow approximations suggest that [1]:

$$Q(t) = \max(0, Q_0 + (Pn - C)t). \quad (3.5)$$

Note that  $Q(t)$  may take noninteger values due to the absence of the concept of packetization in fluid models.

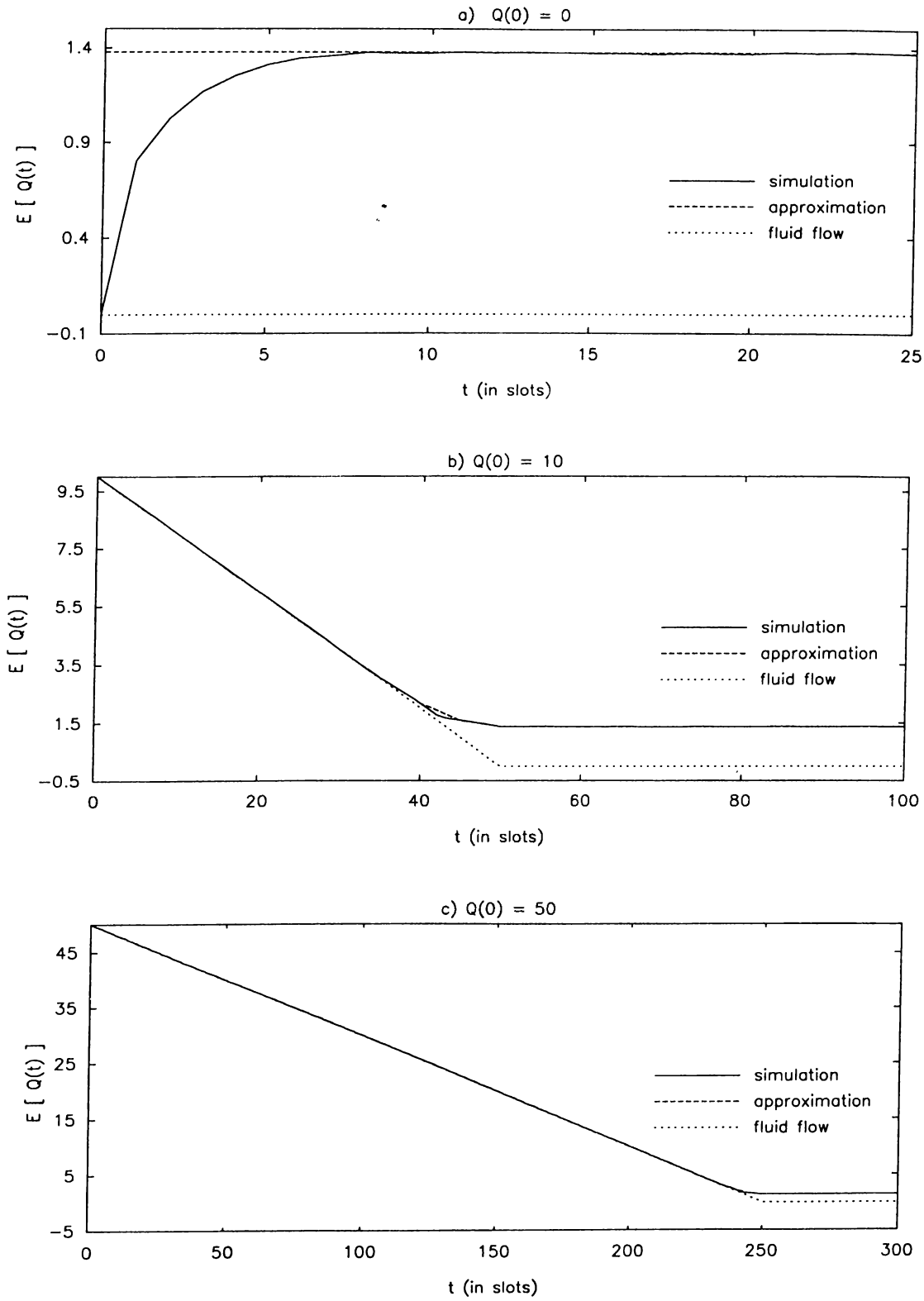
There are two major differences between the expressions (3.4) and (3.5). The first term associated with the short term fluctuations of the queue length is the random variable  $\tilde{Q}_n$  in the periodic model whereas it equals zero in the fluid model. This is in fact why the fluid flow models do not give accurate results in light to moderate traffic when several on/off sources are multiplexed on a common link, as noted by [6],[36]. The second term associated with the dynamical behavior of the queue length in (3.5) is just a linear interpolation of the corresponding term in (3.4).

For the overload states, since the probability that the queue length is zero at some time epoch is negligible, fluid flow approximation gives accurate results in the analysis of the transient response of the queue. Taking (3.4) as our key equality, our approach is mainly based on interpolating the second term as in (3.5) while preserving the first term,  $\tilde{Q}_n$ , which captures the short term fluctuations in the packet layer. In regard of this, we approximate  $Q(t)$  by

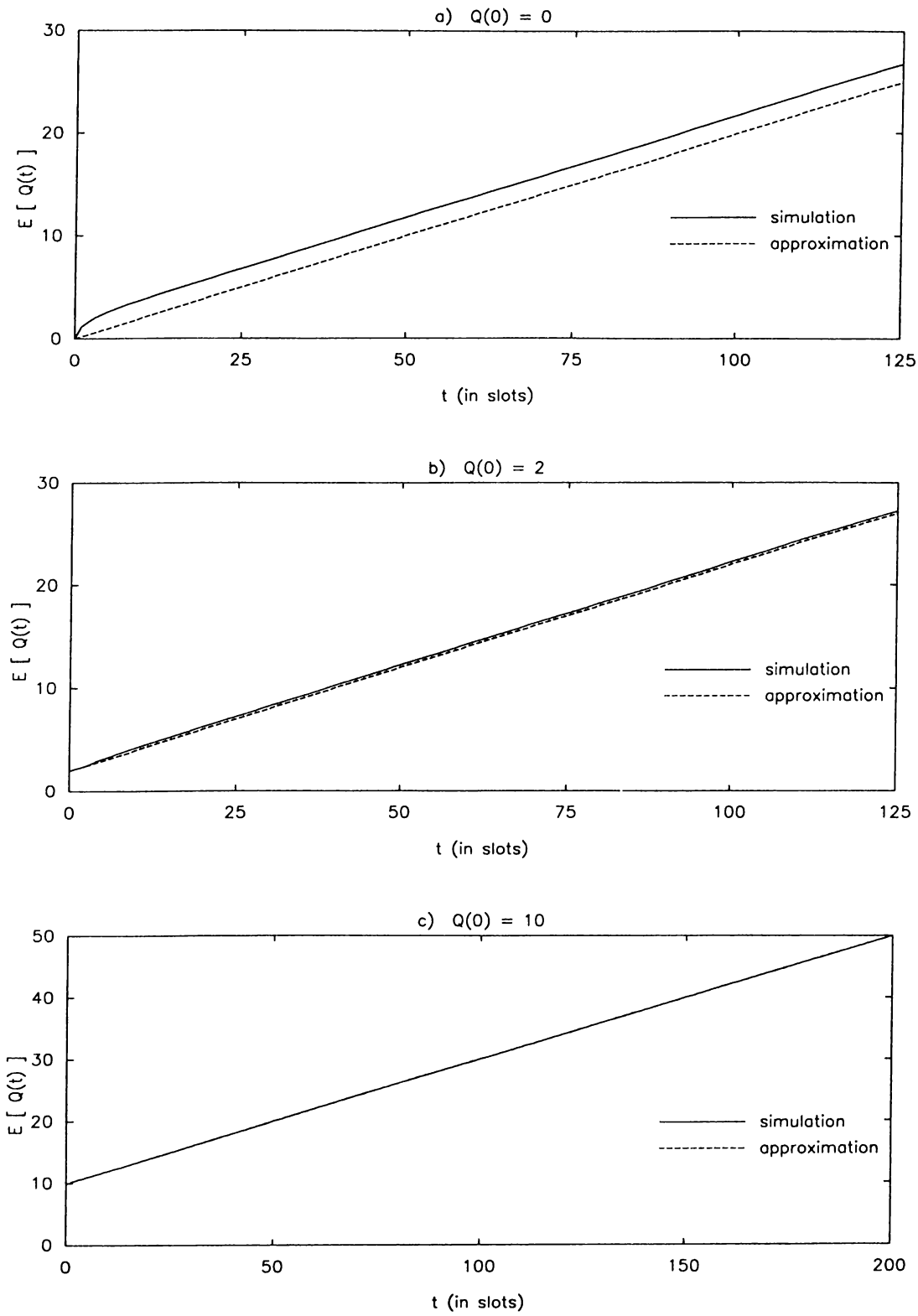
$$Q(t) = \begin{cases} \max(\tilde{Q}_n, Q_0 + (Pn - C)t), & n < R \\ Q_0 + (Pn - C)t, & n \geq R \end{cases} \quad (3.6)$$

The accuracy of this approximation for the  $nD/D/1$  average queue length is examined in Figures 3.1-3.4 and compared with simulation results and fluid flow approximations.

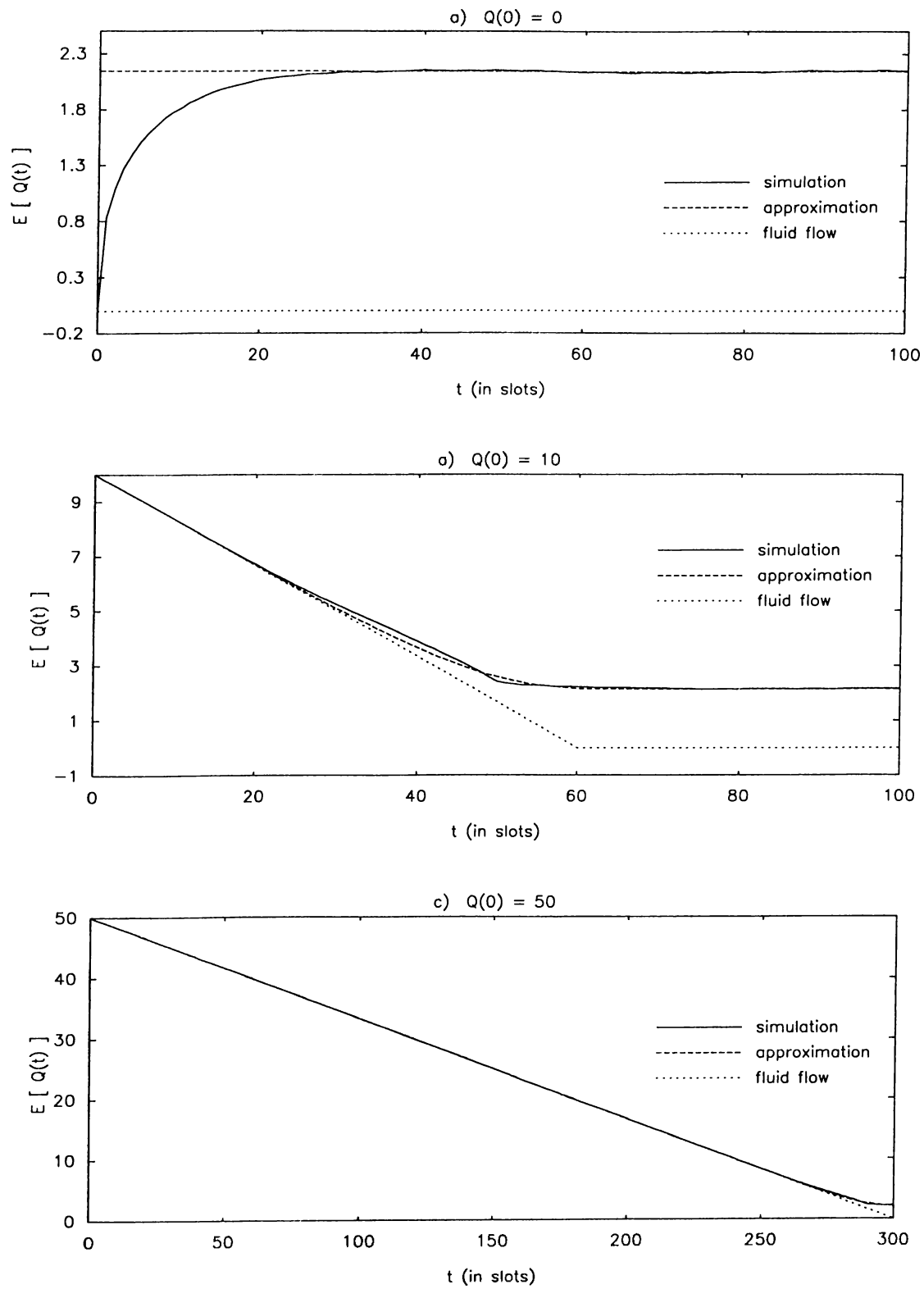
Each of these figures corresponds to a specific  $(R, n)$  pair and the buffer is allowed to



**Figure 3.1:** Comparison of approximations for the expected value of the queue length for the case  $R = 10$  and  $n = 8$  (underload).

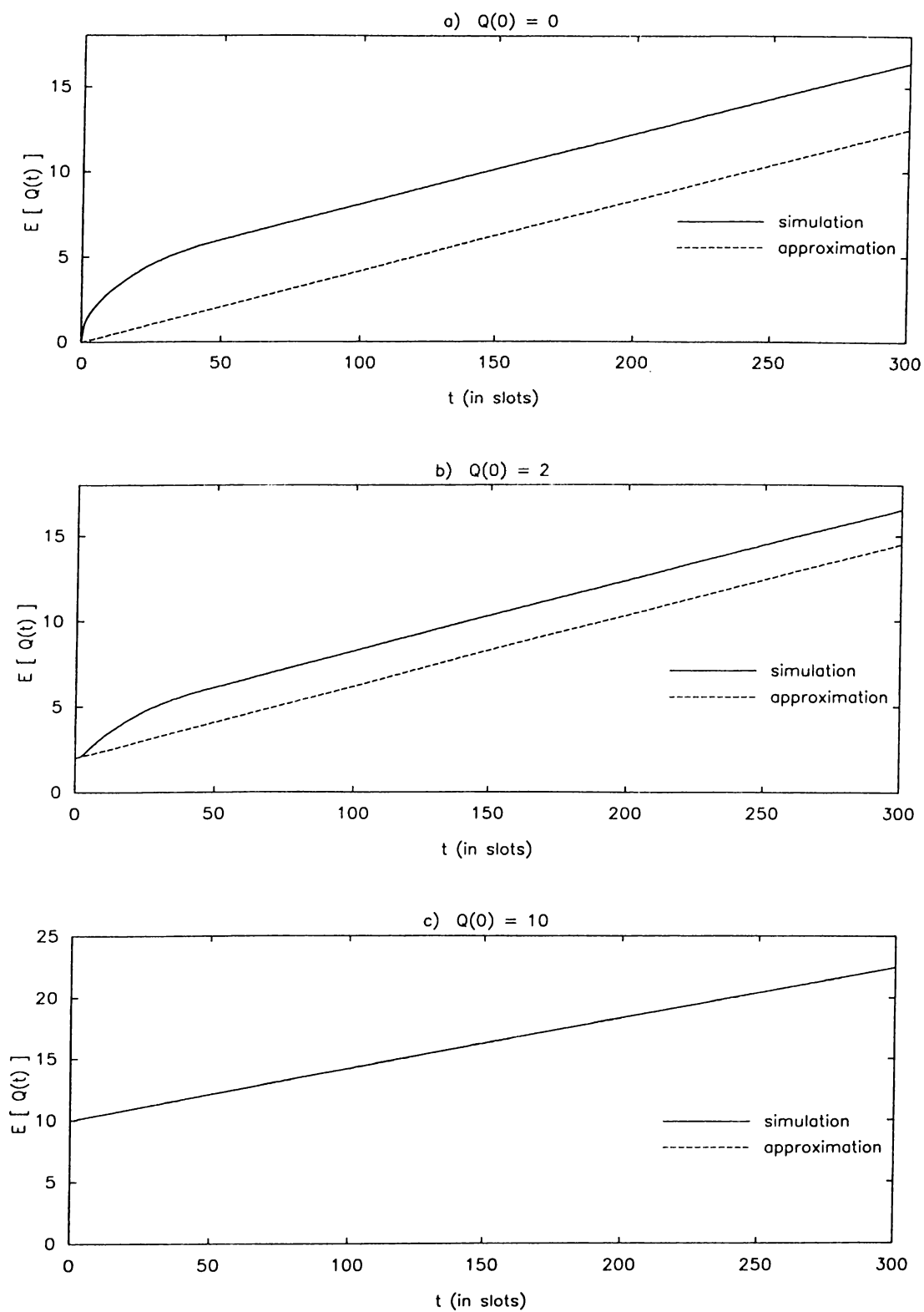


**Figure 3.2:** Performance of the proposed approximation for the expected value of the queue length for the case  $R = 10$  and  $n = 12$  (overload).



**Figure 3.3:** Comparison of approximations for the expected value of the queue length for the case  $R = 48$  and  $n = 40$  (underload).





**Figure 3.4:** Performance of the proposed approximation for the expected value of the queue length for the case  $R = 48$  and  $n = 50$  (overload).

start from different initial conditions. The major observation is that, the approximation (3.6) is very accurate for all the traffic regimes for both underload and overload operation when the initial buffer content is not in the vicinity of zero. Even in this case, the approximation is satisfactory for the underload case and is able to track the simulation curve after the queue length reaches its steady-state distribution. When the queue is in the overload regime, if the queue starts from an almost empty initial length, the proposed approximation underestimates the expected queue length. This is due to the assumption that the server will never be idle, however there are actually occasions which yield empty queues even when the aggregate incoming rate exceeds the multiplex capacity. The probability of the queue being empty at some time epoch decreases rapidly with increasing initial buffer content.

Better results compared with fluid flow approximations are always obtained irrespective of degree of utilization in the system. While choosing the  $(R, n)$  pairs, we let the number of active users  $n$  be close to the period  $R$ , this actually corresponds to the worst-case performance of the approximation (3.6). It is not difficult to visualize that when  $R$  and  $n$  are farther apart, the performance will tend to improve. In the next section, the fundamental approximation (3.6) will be used to derive formulas for the queue length cumulative distribution function when the number of active users ( $n$ ) is modulated by our birth-death model (Figure 1.7).

### 3.3 MMPAP/D/1 Queue

Let us now consider the traffic model in Figure 1.7 and concentrate on a particular state  $(n, 0 \leq n \leq N)$  of the Markov chain. Let  $F_b(n, x)$  and  $F_e(n, x)$  be defined in the same way as in the definitions (2.8) and (2.9), respectively. Recall that, the state holding time at state  $n$  is exponentially distributed with parameter  $\sigma_n$  which equals  $(N - n)\lambda + n\mu$ , which is the total probability flow rate out of state  $n$ . We assume that each time the Markov system changes a state, a complete phase randomization of all the sources is assumed to occur whereas for the original system, an active source's phase is independent of the

other sources' state transitions. With this assumption, the stationary queue length at the moment of state transition to  $n$  and  $\tilde{Q}_n$  become independent.

By exploiting the approximation in (3.6) and with the above assumption one obtains

$$F_e(n, x) = \begin{cases} \tilde{Q}_n(x)F_f(n, x), & n < R \\ F_f(n, x), & n \geq R \end{cases} \quad (3.7)$$

where

$$F_f(n, x) \triangleq \left( \int_0^\infty F_b(n, x + (C - Pn)t) \sigma_n \exp(-\sigma_n t) dt \right) u(x),$$

and the subscript  $f$  denotes the fluid flow term. We can therefore write

$$F_f(n, x) = \begin{cases} F_b(n, x) * \left( \frac{-\sigma_n}{Pn - C} \exp\left(\frac{\sigma_n x}{C - Pn}\right) u(-x) \right), & x \geq 0, \quad n < R \\ F_b(n, x), & n = R \\ F_b(n, x) * \left( \frac{\sigma_n}{Pn - C} \exp\left(\frac{\sigma_n x}{C - Pn}\right) u(x) \right), & x \geq 0, \quad n > R \end{cases} \quad (3.8)$$

Note the analogy between the above expression and the output equation (pertaining to  $F_f(n, x)$ ) of a first order linear system with the input  $F_b(n, x)$ . This analogous linear system is anti-causal in the case  $n < R$  whereas it is causal when  $n > R$ . Writing down the state equations of this system, we now have

$$\frac{d}{dx} F_f(n, x) = \frac{\sigma_n}{C - Pn} F_f(n, x) + \frac{\sigma_n}{Pn - C} F_b(n, x), \quad x \geq 0, \quad n \neq R. \quad (3.9)$$

The balance equations of the continuous-time Markov chain are now employed to relate  $F_b(n, x)$ 's to  $F_e(n, x)$ 's (see the derivation of equation (2.16)):

$$\sigma_n F_b(n, x) = \sum_{m \neq n} p(m, n) F_e(m, x), \quad (3.10)$$

in which  $p(m, n)$  is defined to be the state transition rate from state  $m$  to state  $n$ . Actually, for our birth-death model,

$$\begin{aligned} p(n, n+1) &= (N - n)\lambda, \quad n = 0, 1, \dots, N - 1, \\ p(n, n-1) &= n\mu, \quad n = 1, 2, \dots, N. \end{aligned}$$

Combining (3.7),(3.9) and (3.10), we finally obtain the following differential equations for  $F_f(n, x)$ 's:

$$\begin{aligned} \frac{d}{dx} F_f(n, x) &= \frac{\sigma_n}{C - P_n} F_f(n, x) + \frac{1}{P_n - C} \sum_{m \neq n} p(m, n) \tilde{Q}_m(x) F_f(m, x), \quad n \neq R \\ F_f(R, x) &= \frac{1}{\sigma_R} \sum_{m \neq R} p(m, R) \tilde{Q}_m(x) F_f(m, x). \end{aligned} \quad (3.11)$$

In the above equations,  $\tilde{Q}_m(x) \triangleq 1, \forall x \geq 0, m \geq R$ . If the term  $\tilde{Q}_n(x)$  is further taken as unity  $\forall n, n = 0, 1, \dots, N$ , then the above equations are equivalent to the fluid flow equations [1] up to a similarity transformation. The equation belonging to  $F_f(R, x)$  is algebraic and may be eliminated. This is achieved by substituting the second expression in the first equations of (3.11) so that by defining

$$F_f(x) = \left[ F_f(0, x) \quad F_f(1, x) \quad \cdots \quad F_f(R-1, x) \quad F_f(R+1, x) \quad \cdots \quad F_f(N, x) \right]^T,$$

we finally have

$$\frac{d}{dx} F_f(x) = A(x) F_f(x), \quad x \geq 0. \quad (3.12)$$

Here the  $N \times N$  matrix  $A(x)$  is determined through a suitable arrangement of the differential equations in (3.11). Actually,

$$A(x) = A_i, \quad x \in [i, i+1), \quad i \in \mathcal{Z}_+, \quad 0 \leq i \leq R-2,$$

and

$$A(x) = A, \quad x \in [R-1, \infty)$$

for some appropriate constant matrices  $A_i$ 's and  $A$ , due to the piecewise constant structure of the distributions  $\tilde{Q}_n(\cdot)$ 's. Given the initial condition  $F_f(0)$ , the differential equation (3.12) has a unique continuous solution described by

$$F_f(x) = \exp(A_i(x-i)) F_f(i), \quad x \in [i, i+1], \quad 0 \leq i \leq R-2, \quad (3.13)$$

and

$$F_f(x) = \exp(A(x-(R-1))) F_f(R-1), \quad x \geq R-1. \quad (3.14)$$

In order to find the initial condition, we make use of the following observations:

- 1) For  $n > R$ , the queue is always increasing, so the queue length cannot be zero. Therefore,  $F_f(n, 0) = 0$  for  $n > R$ .
- 2) The matrix  $A$  is, in fact, equivalent to the state matrix in fluid flow models, therefore it is known to have  $R - 1$  positive real eigenvalues,  $N - R$  negative real eigenvalues and an eigenvalue at the origin. In order for the solution not to blow up as  $x \rightarrow \infty$ , the coefficients associated with the positive eigenvalues of  $A$  should be set to zero by the choice of  $F_f(0)$ .
- 3) Defining  $\pi_n$  to be the the equilibrium probability of  $n$  sources being active, we write

$$F_f(n, \infty) = \pi_n, \quad 0 \leq n \leq N.$$

To explain the observation 3),  $F_f(n, \infty) = F_e(n, \infty)$ , which equals  $\pi_n$ , the equilibrium probability of  $n$  sources being active (by definition).

Now, let  $z_i$  be a stable eigenvalue of  $A$  and  $\phi_i$  be its corresponding right eigenvector. Then, by observation 2) and (3.14), the solution to  $F_f(x)$  can be written in the form

$$F_f(x) = F_f(\infty) + \sum_{i=1}^{N-R} \exp(z_i(x - R + 1)) a_i \phi_i, \quad x \geq R - 1$$

which yields

$$F_f(R - 1) = F_f(\infty) + \sum_{i=1}^{N-R} a_i \phi_i, \quad (3.15)$$

where  $a_i$ 's are coefficients to be determined. The relationship between  $F_f(0)$  and  $F_f(R - 1)$  now needs to be established. Using (3.13), one can write

$$F_f(R - 1) = Z F_f(0) \triangleq \left( \prod_{i=0}^{R-2} \exp(A_i) \right) F_f(0). \quad (3.16)$$

Besides, by observation 1),  $F_f(0)$  is in the form

$$F_f(0) = \begin{bmatrix} f \\ 0 \end{bmatrix},$$

where  $f$  is of size  $R \times 1$ . Combining (3.15) and (3.16), one can solve for  $a_i$ 's and  $f$ , and thus the initial condition  $F_f(0)$  through a linear matrix equation of size  $N$ . In fact, if we define

$$\begin{aligned}\Phi &= \begin{bmatrix} \phi_1 & \phi_2 & \cdots & \phi_{N-R} \end{bmatrix} \\ a &= \begin{bmatrix} a_1 & a_2 & \cdots & a_{N-R} \end{bmatrix}^T,\end{aligned}$$

then

$$\Phi a + F_f(\infty) = Z_1 f$$

and

$$\begin{bmatrix} f \\ a \end{bmatrix} = \begin{bmatrix} Z_1 & -\Phi \end{bmatrix}^{-1} F_f(\infty), \quad (3.17)$$

where  $Z_1$  is composed of the first  $R$  columns of  $Z$ . Having found the initial condition, the solutions given in (3.13) and (3.14) complete our description of the stationary queue length distribution through the equation (3.7). The essential difference between the method presented here and computations encountered in solving the fluid flow models is the computation of the linear operator  $Z$  defined in (3.16).

Using equality (3.7), one can evaluate  $F_e(n, x)$ ' from  $F_f(n, x)$ 's so that the overall cdf of queue length is written as the sum of the individual elements  $F_e(n, x)$  (by (2.10)):

$$Pr(\text{queue length} \leq x) = \sum_{n=0}^N F_e(n, x). \quad (3.18)$$

One other goal is actually finding the distribution of the queueing delay rather than the queue length. Queue length can easily be converted to queueing delay by substituting  $Ct$  for  $x$ . However, to form the cdf of the queueing delay, each  $F_e(n, x)$  should be weighted before summation (see for the derivation of equality (2.7)):

$$Pr(\text{delay} \leq t \text{ sec.}) = \frac{1}{\alpha N} \sum_{n=0}^N n F_e(n, Ct) \quad (3.19)$$

where

$$\alpha = \frac{\lambda}{\lambda + \mu}$$

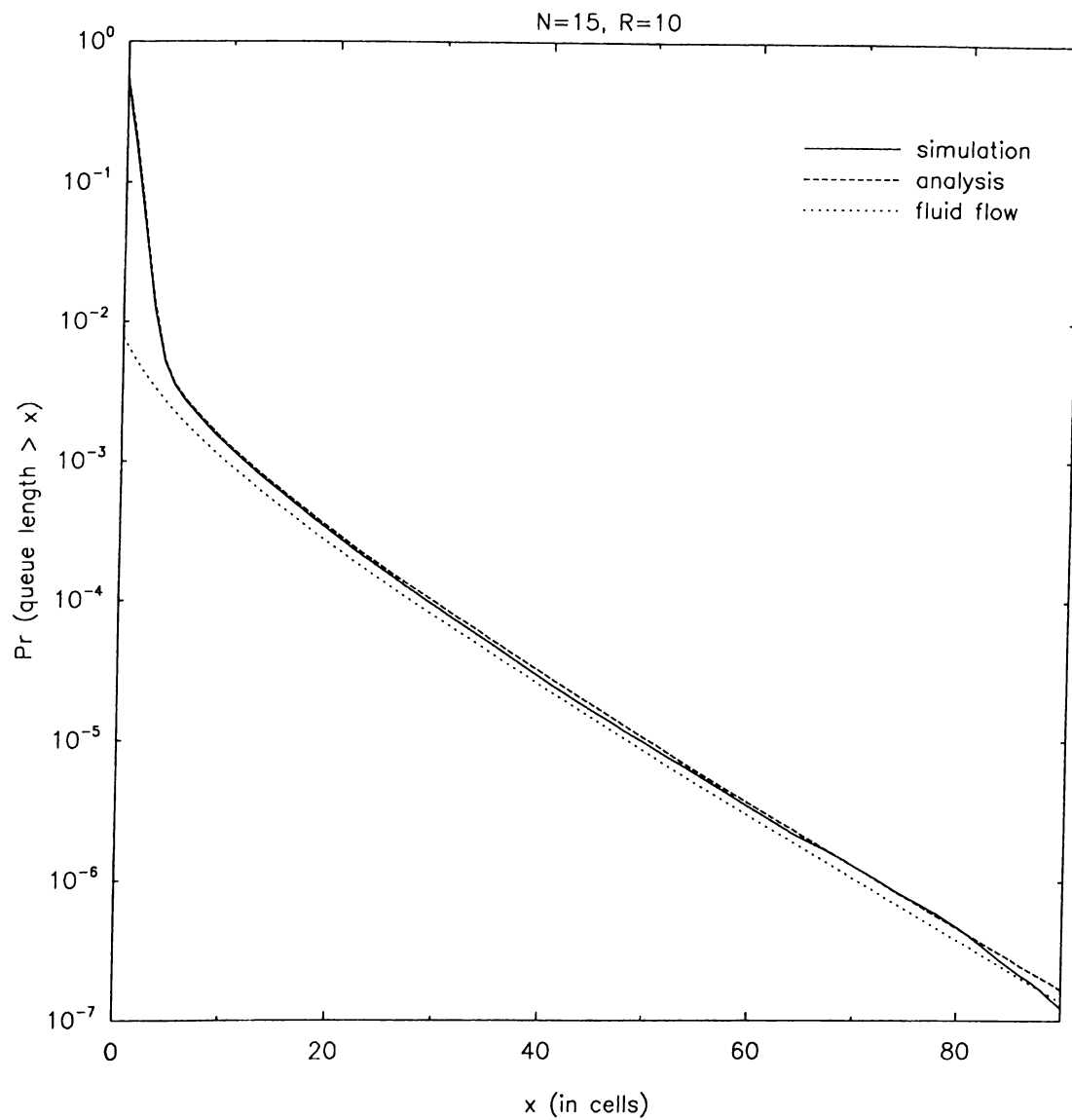
is the average fraction of sources being active.

### 3.3.1 Numerical Examples

We first consider a packetized voice system with line speed 320 kbits/s, voice peak rate 32 kbits/s, mean active period 353 ms and mean silent period 650 ms. The mean number of packets in a talkspurt is approximately 22. The packets are 64 kbytes and the packet transmission time is 1.6 ms. This corresponds to a link rate of 320 kbits/s, which is not the ATM rate. We refer to this example since many authors have concentrated on this packetized voice framework [1],[20] to demonstrate their results for the analysis of statistical multiplexing.

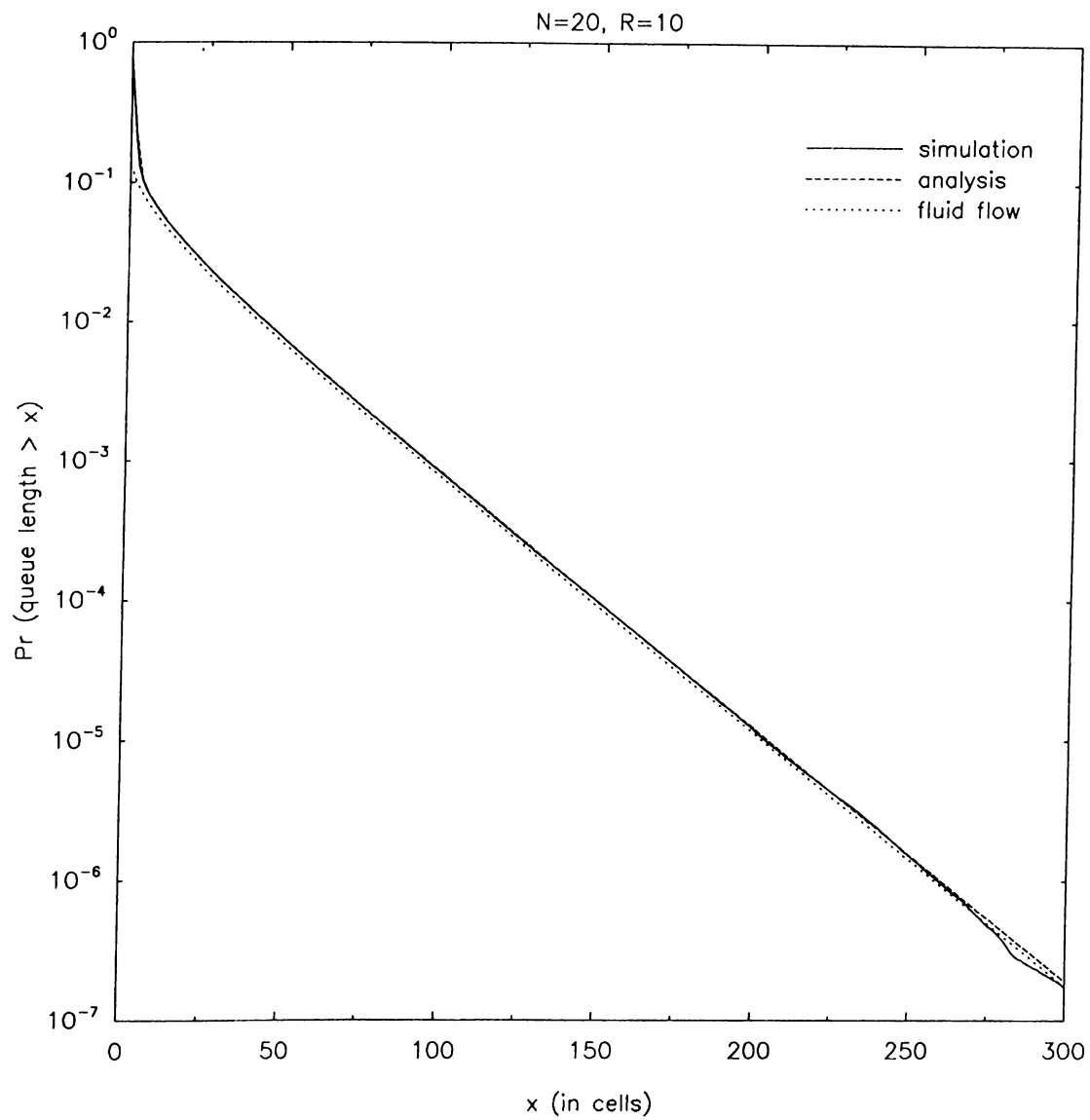
The simulation results are obtained based on the discrete-time queueing scheme as described in Section 3.1. In the simulation, the continuous-time random variables, the active and idle periods, are first chosen from the corresponding exponential random variables and then rounded to the nearest slot-times. Within an active period, the packets from an individual voice source are transmitted in a periodic manner, each source's phase being uniform between 0 and  $R - 1$  ( $R = 10$  for this example). In Table 3.1, the mean waiting time in the queue with respect to the number of voice sources by our analysis method and the fluid flow approximation is given and these values are compared with the simulation results. The analysis method proposed in this paper gives highly accurate results independent of the degree of utilization in the system whereas fluid flow approximation is only satisfactory in the heavy load regime. Figures 3.5 and 3.6 are devoted to the queue length survivor function, which are obtained for the cases  $N = 15$  and 20, respectively. In both cases, the method we propose is able to capture the simulation curve for the buffer survivor function very accurately.

We then extend the example to the case where the link speed = 1.536 Mbits/s and the packet transmission time = 0.333 ms ( $R = 48$  in this case). This is actually the classical packetized voice example found in the literature [36],[54]. To avoid numerical inaccuracies in the calculation of the initial condition  $F_f(0)$ , we approximate  $A_i$  by  $A$  for  $i \geq 9$ . Actually, a closer study of  $\tilde{Q}_n(\cdot)$  shows that this leads to an error no more than 2% for each entry of  $A_i$ . We note that, as can be verified easily, fluid flow approximation



**Figure 3.5:** Comparison of the queue length survivor function for our proposed method with simulation results and the fluid flow approximation ( $N = 15$ ,  $R = 10$ , utilization = 0.52).





**Figure 3.6:** Comparison of the queue length survivor function for our proposed method with simulation results and the fluid flow approximation ( $N = 20$ ,  $R = 10$ , utilization = 0.70).

No. voice sources	simulation results (ms.)	95 % conf. interval	approximations [ms]	
			analysis	fluid flow
4	0.0929	$\pm 0.0021$	0.0948	0.00
6	0.1638	$\pm 0.003$	0.1591	0.00
8	0.2474	$\pm 0.003$	0.2383	0.00
12	0.4716	$\pm 0.0035$	0.4813	0.0023
14	0.6474	$\pm 0.0065$	0.6918	0.0383
16	1.044	$\pm 0.03$	1.136	0.269
18	2.205	$\pm 0.04$	2.311	1.199
20	5.32	$\pm 0.26$	5.46	4.09
22	13.64	$\pm 0.38$	13.61	12.02
24	35.53	$\pm 0.96$	35.16	33.40
25	61.6	$\pm 2.0$	58.8	57.0
26	111.0	$\pm 3.5$	105.6	103.8
27	258.1	$\pm 8.7$	224.6	222.9

**Table 3.1:** Comparison of approximations of the mean waiting time with the simulation results for the case  $R = 10$ .

is equivalent to setting

$$A_i = A, \quad \forall i \geq 0.$$

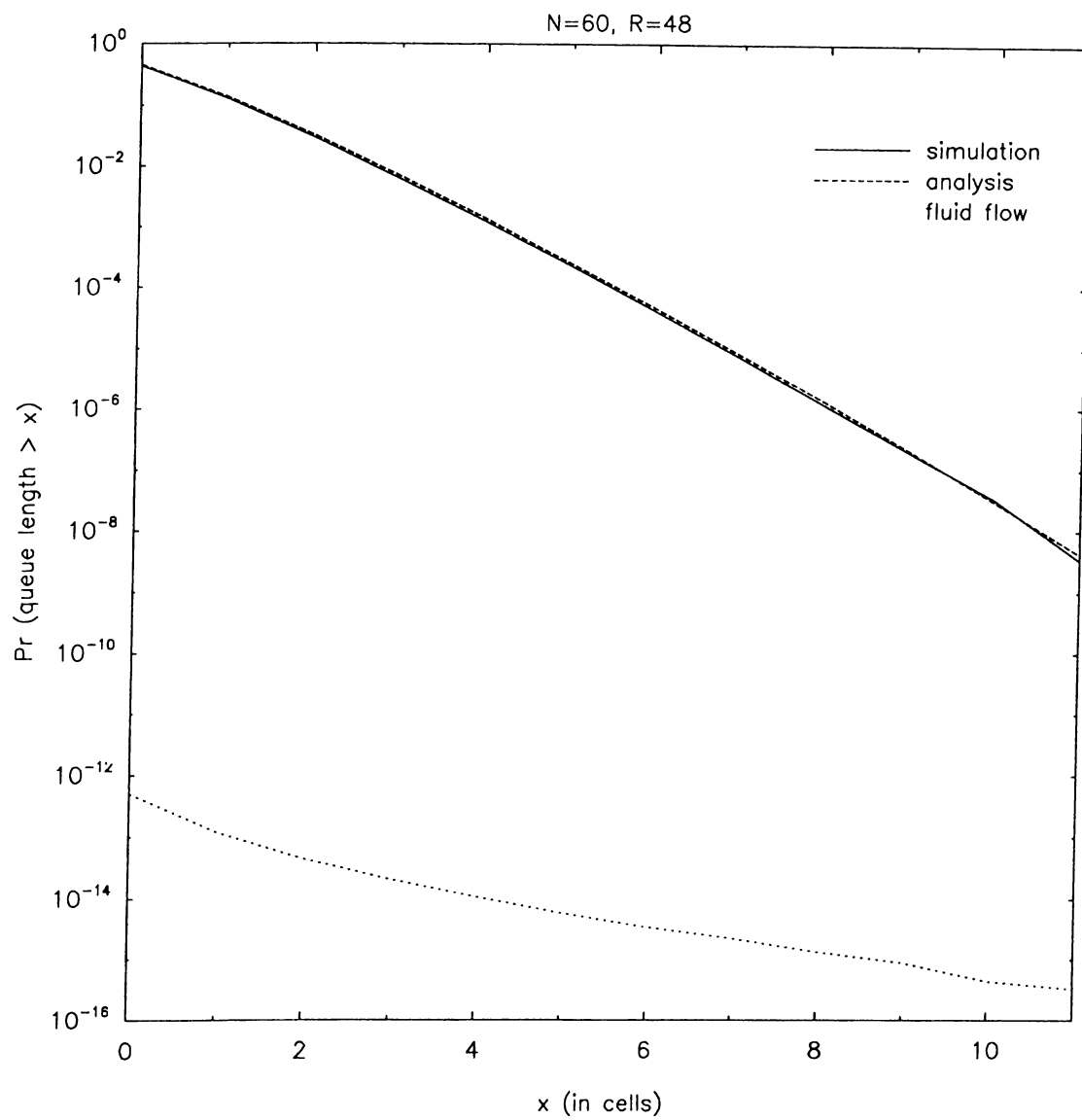
We further use the trapezoidal approximation [55]:

$$\exp(A_i) \approx (I - A_i/2)^{-1}(I + A_i/2) \quad (3.20)$$

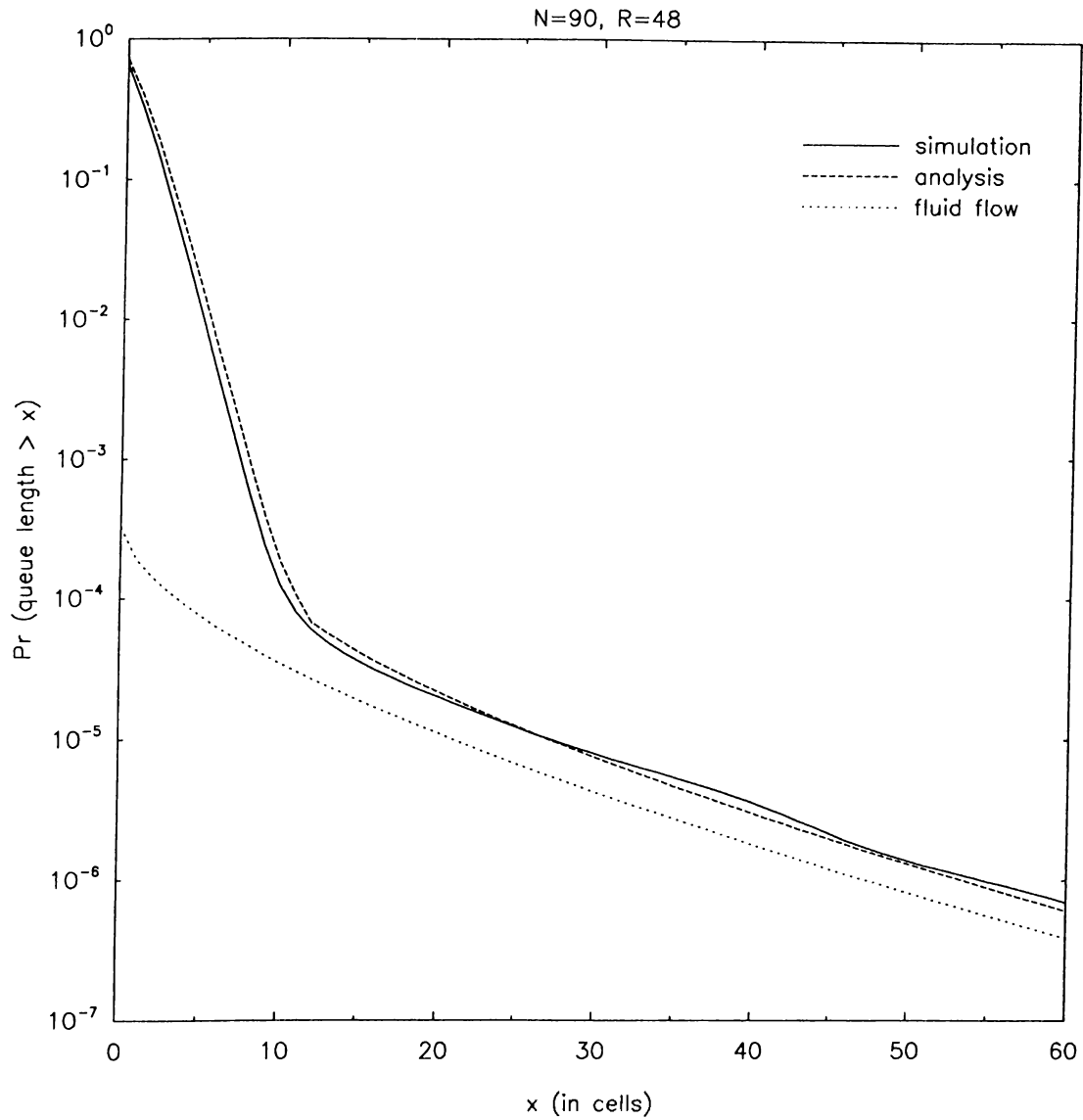
to avoid eigensystem calculations. The results associated with the mean waiting time in the queue is given in Table 3.2.

The queue length survivor functions, for the cases  $N = 60, 90,$  and  $120$  are presented in Figures 3.7-3.9, respectively. Remarkably accurate results are obtained for all the cases compared with the fluid flow approximations in spite of the employment of the above-mentioned approximations. The numerical results provided here demonstrate three significant aspects of our proposed method:

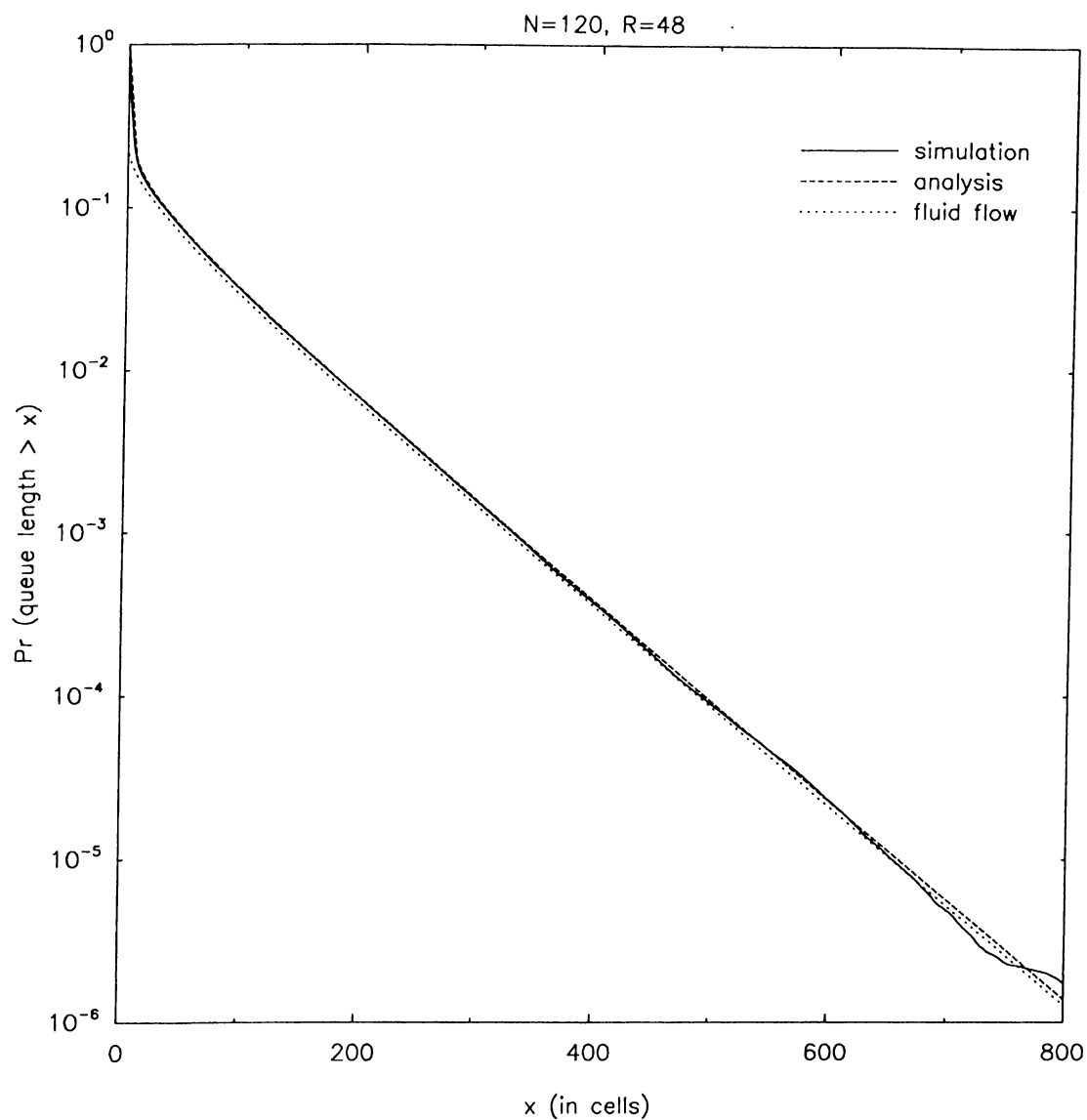
- mathematical formulation is simple and similar to fluid flow models but always



**Figure 3.7:** Comparison of the queue length survivor function for our proposed method with simulation results and the fluid flow approximation ( $N = 60$ ,  $R = 48$ , utilization = 0.44).



**Figure 3.8:** Comparison of the queue length survivor function for our proposed method with simulation results and the fluid flow approximation ( $N = 90$ ,  $R = 48$ , utilization = 0.66).



**Figure 3.9:** Comparison of the queue length survivor function for our proposed method with simulation results and the fluid flow approximation ( $N = 120$ ,  $R = 48$ , utilization = 0.88).

No. voice sources	simulation results (ms.)	95 % conf. interval	approximations [ms]	
			analysis	fluid flow
60	0.1243	$\pm 0.0004$	0.0774	0.00
80	0.222	$\pm 0.001$	0.186	0.00
90	0.298	$\pm 0.003$	0.290	0.00
100	0.431	$\pm 0.016$	0.482	0.02
110	0.976	$\pm 0.077$	1.096	0.36
120	4.52	$\pm 0.2$	4.72	3.49
125	12.51	$\pm 0.7$	12.20	10.14
130	36.91	$\pm 1.15$	36.20	32.60
132	66.0	$\pm 5.2$	64.26	57.90
134	151.5	$\pm 12.7$	142.7	124.6

**Table 3.2:** Comparison of approximations of the mean waiting time with the simulation results for the case  $R = 48$

yields better results,

- the method provides satisfactory results through all traffic regimes,
- besides the averaged performance criteria (e.g., mean queueing delay, mean buffer size), the method is able to capture the whole cdf of the queue length.

Before going through a second example, let us focus our attention on the computational complexity of the proposed algorithm. The main difference between the algorithm presented here and the fluid flow approximations lies under the computation of the operator  $Z$  defined in (3.16). This requires the computation of  $R - 1$  matrix exponentials of size  $N + 1$  and  $R - 2$  matrix multiplications of size  $N + 1$ . It is in fact possible to decrease this number of matrix exponentiations and matrix multiplications without significantly degrading the performance of the algorithm.

Note that, the convergence rate of  $A_i$  to the matrix  $A$  as  $i \rightarrow R - 2$  is fast. This can be concluded from the equation (3.3) where  $Pr(\tilde{Q}_n > q)$  approaches zero quickly as  $q \rightarrow n - 1$ . One can therefore approximate  $A_i$  by  $A$  for  $i > j - 1$ , for some appropriate

$j, j < R - 1$ . The revised version of the algorithm then reduces to computing

$$Z \approx \left( \prod_{i=0}^{j-1} \exp(A_i) \right), \quad (3.21)$$

and solving for the unknowns  $a_i$ 's and  $f$  through the continuity of the solution of the differential equations at  $x = j$  instead of the continuity at  $x = R - 1$  as in the exact version of the algorithm. This revised scheme which we call a  $j^{\text{th}}$ -order approximation, requires the computation of  $j$  matrix exponentials and  $j - 1$  matrix multiplications. Recall that

$$A_i \approx A, \quad \forall i \geq 0,$$

in fluid flow approximations, therefore this technique can be interpreted as a  $0^{\text{th}}$ -order approximation of our proposed analysis scheme in our context.

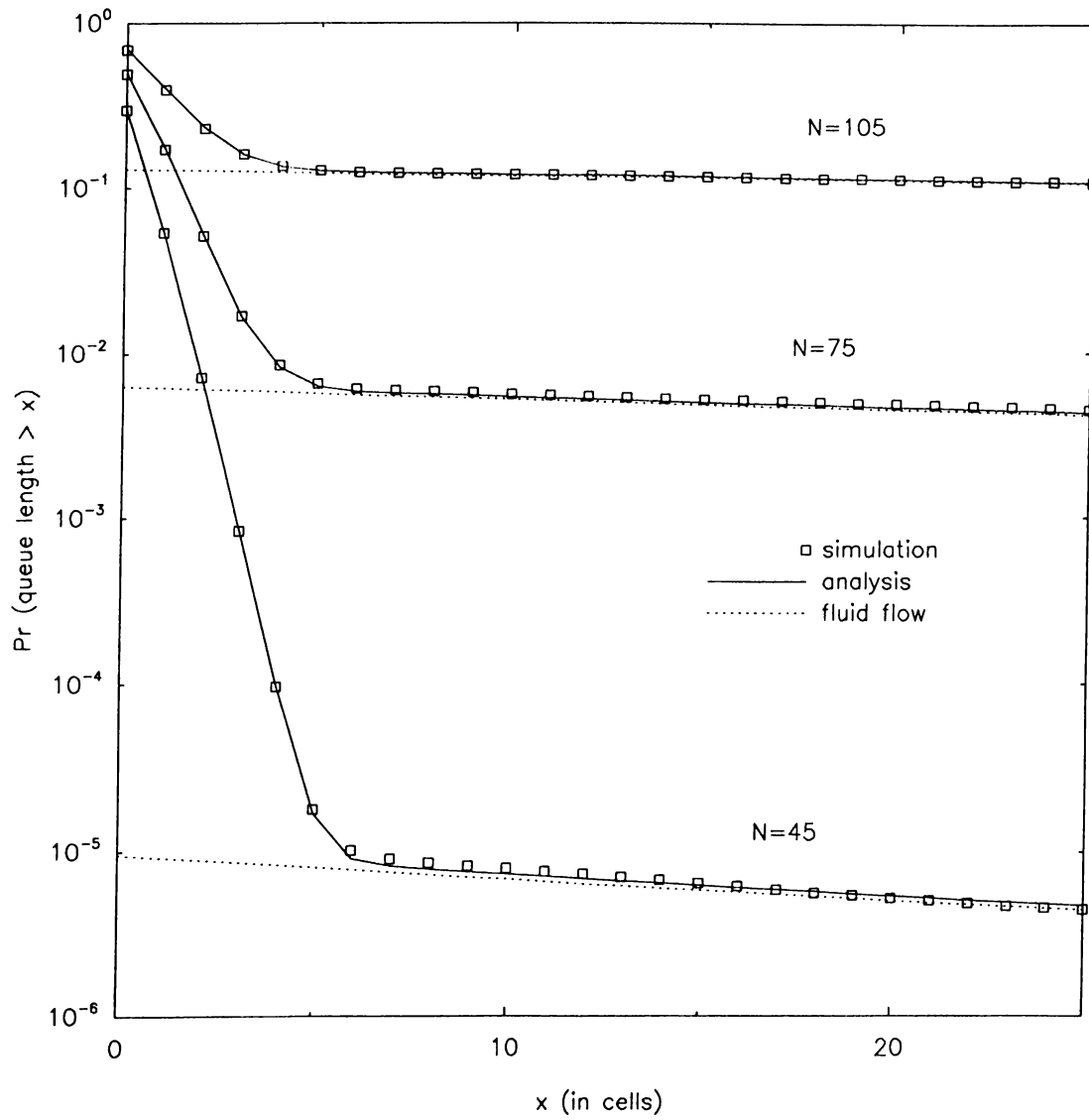
We now investigate the performance of the  $j^{\text{th}}$ -order approximations we have proposed in the following example. We consider an ATM multiplexer that serves LAN (Local Area Network)-generated data streams. The cell emission process for an individual LAN source is widely recognized to be adequately represented by means of an on-off source model [3]. Let us consider a set of  $N$  independent and homogeneous LAN sources characterized by i) the peak rate  $F_p$  ii) the activity factor  $p$ , defined as the ratio between the average bit rate and  $F_p$  iii) mean burst length  $L_b$ . We choose a reference LAN source as in [3] which is characterized by  $F_p = 10$  Mbits/s,  $p = 0.1$ , and  $L_b = 16250$  bytes, where these values are representative of a large class of information flows arising from LAN's accessing to an ATM network. As for the multiplexer, we assume an output capacity equal to 150 Mbits/s (ATM transport rate) and a cell length of 53 bytes. Note that  $R = \text{link rate} / \text{user peak rate} = 15$  for this specific example.

In Figure 3.10, the queue length survivor function computed via our proposed method is plotted along with the simulation results for the cases  $N = 45$  (utilization = 0.3),  $N = 75$  (utilization = 0.5), and  $N = 105$  (utilization = 0.7). For the three particular regimes, the method matches closely with the simulation results whereas the fluid flow approximation only captures the asymptotic behavior. Figure 3.11 addresses the performance issues of  $j^{\text{th}}$ -order approximations when  $N = 75$  and the mean burst length

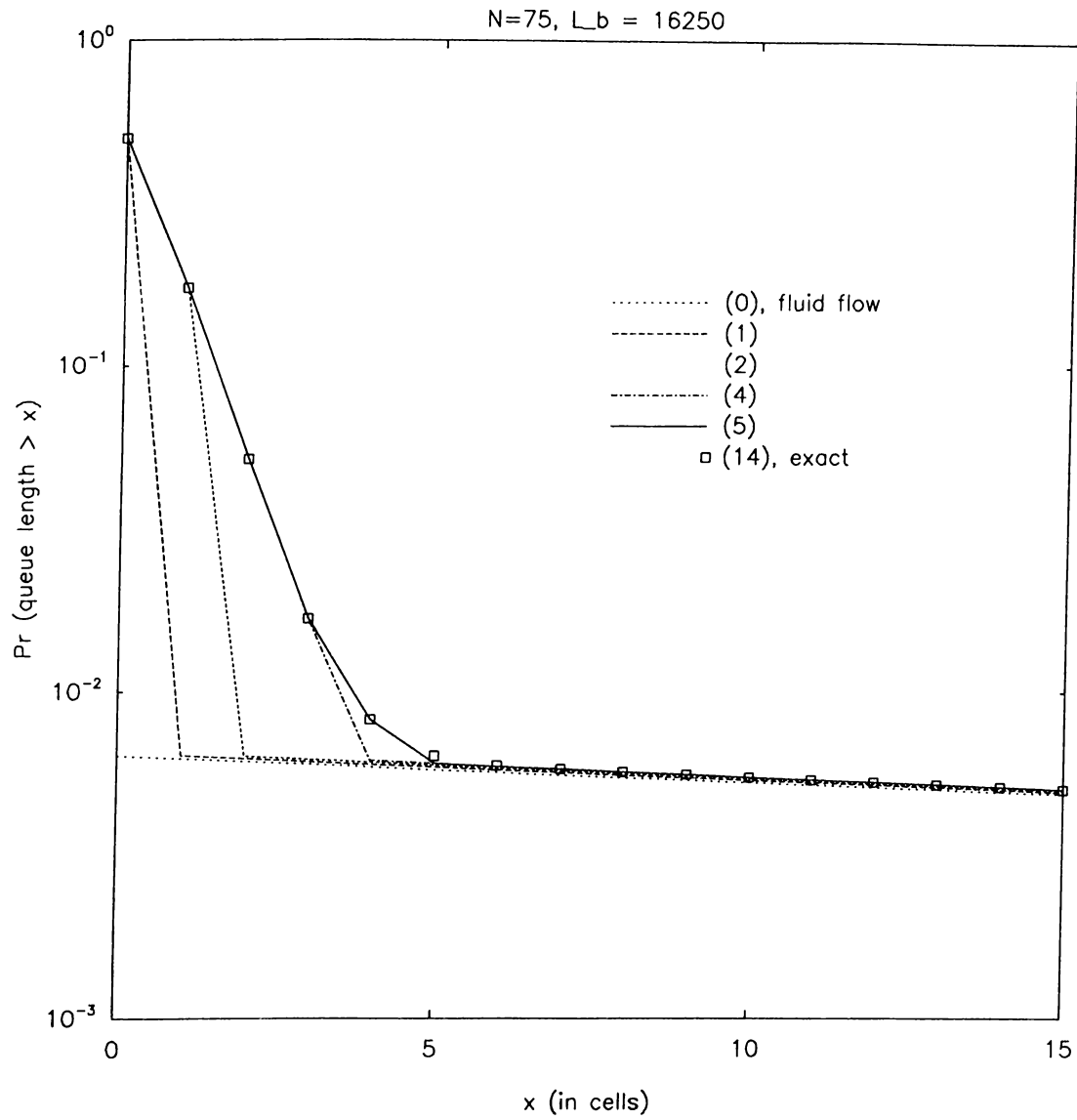
of each user,  $L_b$ , equals 16250 bytes. The notation  $(j)$  in the figure denotes the  $j^{\text{th}}$ -order approximation. We observe that as the degree of approximation increases, the performance of the approximation scheme improves. The performance of the exact version of the algorithm (may also be referred as the  $14^{\text{th}}$ -order approximation, since  $R - 1 = 14$  for this example) and the performance of the  $5^{\text{th}}$ -order approximation are more or less the same. The degree of freedom in choosing the order  $j$  of the approximation scheme gives us a chance to play out with accuracy and computational load. If you pay more (increase  $j$ ), you get more (increase accuracy). In regard of the results presented here, you don't have to pay much (the choice of  $j = 4$  or  $5$  is enough for all practical purposes). For this example, in terms of the asymptotic behavior, all the approximations give almost the same results, but this is not always the case. We demonstrate this fact through the following multiplexer example where we change the mean burst length  $L_b$  to 500 bytes. Figure 3.12 presents the results of the exact version of the algorithm for three different utilizations. What differs from the previous example is that the fluid flow approximation now captures the asymptotic slope but not the asymptotic constant. The asymptotic behavior of the simulation curve (which agrees well with the analysis) is not the same as that of the fluid flow approximations but is only parallel to that. Note that, accuracy of fluid flow approximations deteriorates as the load is decreased (e.g.,  $N=45$ ). Regardless of the utilization of the queueing system, our method gives acceptable results for all  $x$ . When  $N = 75$  and  $L_b = 500$ , we present our results obtained via  $j^{\text{th}}$ -order approximations in Figure 3.13. The results are similar to the previous one (Figure 3.11) except that the increase in accuracy in finding the asymptotic constant by using  $j^{\text{th}}$ -order approximations now becomes significant.

One can further decrease the computational load by using efficient methods to compute the matrix exponential. One of the methods is to use the trapezoidal approximation (3.20), but a further discussion of numerical matrix exponentiation techniques is out of scope of this dissertation. We refer to [38] for a detailed survey of already-existing methods to compute the matrix exponential while keeping in mind that performance investigation of these methods for the analysis of the queueing system of interest might be a future research topic.

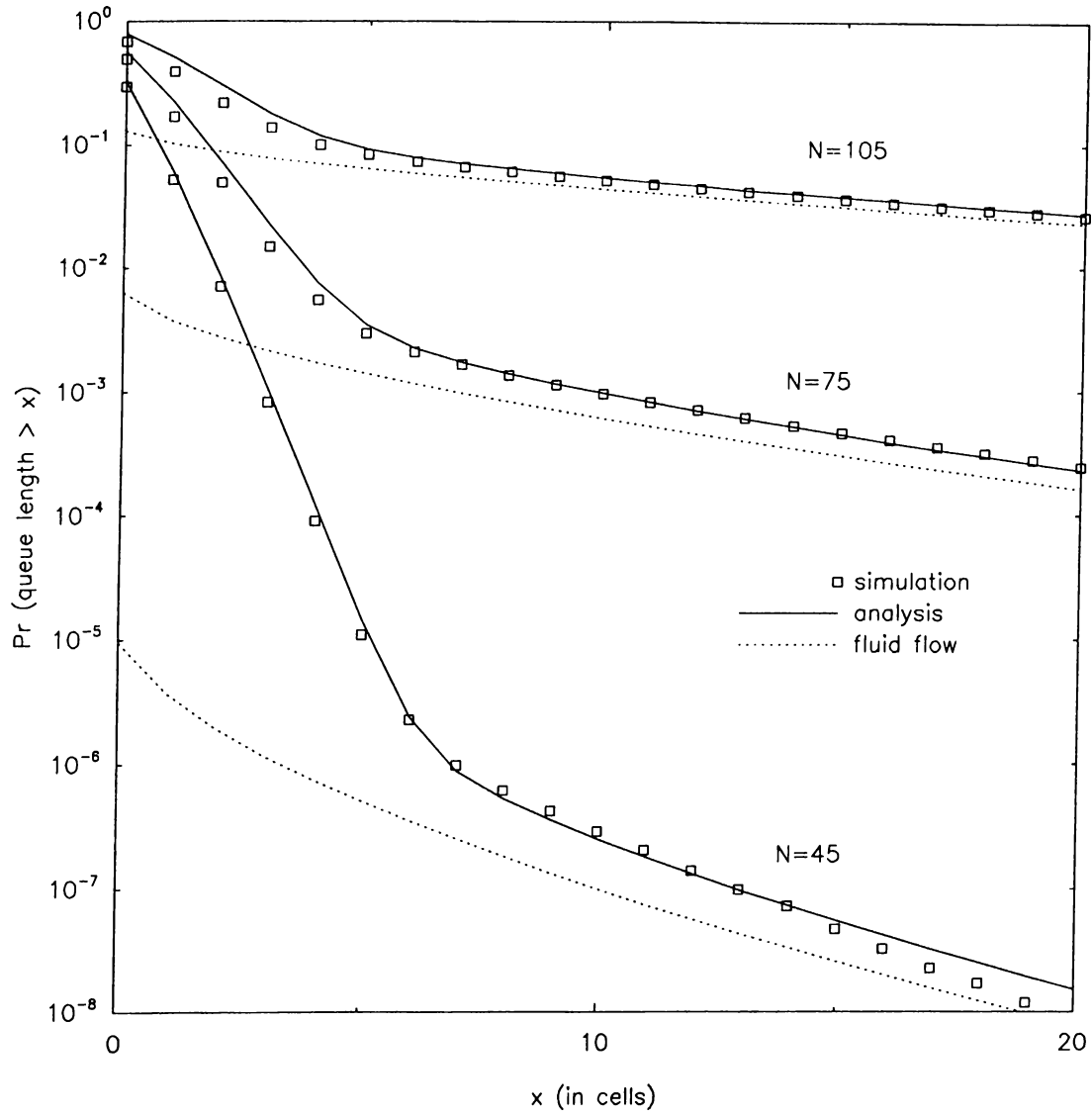




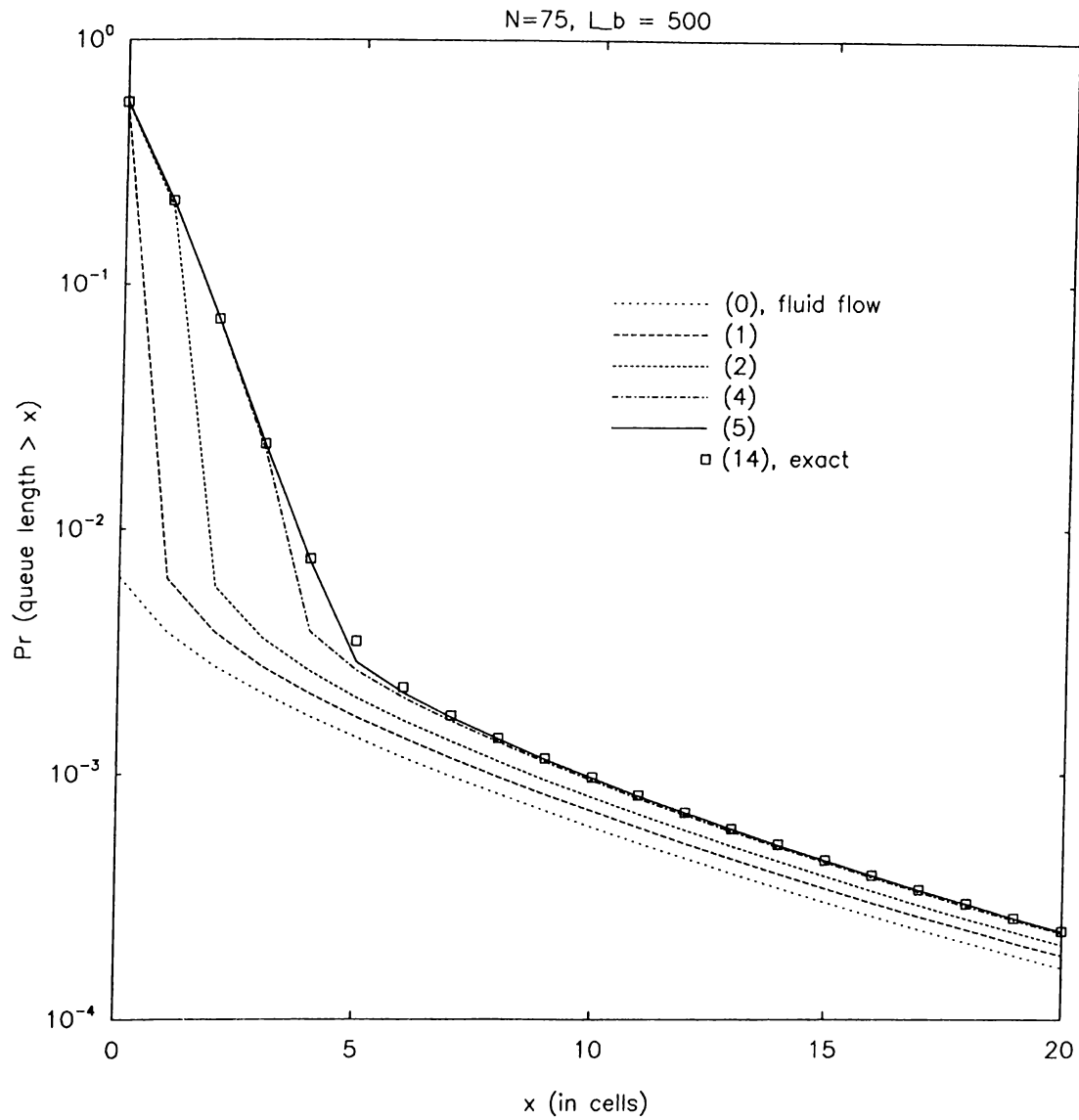
**Figure 3.10:** Queue length survivor function for  $N = 45$ ,  $N = 75$ , and  $N = 105$  when  $L_b = 16250$  bytes.



**Figure 3.11:** Queue length survivor function obtained via  $j^{th}$ -order approximations for  $N = 75$  and  $L_b = 16250$  bytes.



**Figure 3.12:** Queue length survivor function for  $N = 45$ ,  $N = 75$ , and  $N = 105$  when  $L_b = 500$  bytes.



**Figure 3.13:** Queue length survivor function obtained via  $j^{th}$ -order approximations for  $N = 75$  and  $L_b = 500$  bytes.

### 3.4 Finite Buffers

Within this framework, it is also possible to extend the results obtained for the queue length cdf to the case where the buffers are of finite size. This extension is provided based on the formulation of [54] in which fluid flow techniques are successfully applied for a packet-speech multiplexer of finite size.

Let  $K$  be the buffer size in packets and assume that  $K > R$ , which is typically the case. Now, let  $z_i$ 's be the eigenvalues of  $A$  (no stability constraint is now imposed) and  $\phi_i$ 's be the corresponding right eigenvectors. Then,  $F_f(x)$  can be written for large  $x$  as

$$F_f(x) = \sum_{i=1}^N \exp(z_i(x - R + 1)) a_i \phi_i, \quad K > x \geq R - 1, \quad (3.22)$$

where  $a_i$ 's are coefficients to be determined. We also let  $u_n$  to be the probability that the chain is in state  $n$  and the queue is held at its upper limit  $K$ . Defining  $F_f(n, K^-)$  (which equals  $F_e(n, K^-)$ ) as  $\lim_{x \rightarrow K} F_f(n, x)$ ,  $u_n$  is simply the difference between  $F_f(n, K^-)$  and the overall probability of  $n$  active lines

$$u_n = \pi_n - F_f(n, K^-).$$

For the finite buffer case, the observation 1) on page 46 is still valid while for  $n < R$ , the queue is always decreasing, so the queue would never be on its limit. Therefore,

$$u_n = 0, \quad F_f(n, K^-) = \pi_n, \quad n < R. \quad (3.23)$$

Defining

$$\begin{aligned} \Phi &= \begin{bmatrix} \phi_1 & \phi_2 & \cdots & \phi_N \end{bmatrix}, \\ a &= \begin{bmatrix} a_1 & a_2 & \cdots & a_N \end{bmatrix}^T, \\ \bar{\pi} &= \begin{bmatrix} \pi_0 & \pi_1 & \cdots & \pi_{R-1} \end{bmatrix}^T, \\ D &= \text{diag}\{\exp(z_i(K - R + 1))\}, \end{aligned}$$

the expression (3.22) together with the boundary condition (3.23) implies that

$$\Phi D a = \begin{bmatrix} \bar{\pi} \\ * \end{bmatrix}, \quad (3.24)$$

where  $*$  means a don't care vector. The boundary condition 1) and the continuity argument of  $F_f(x)$  at  $x = R - 1$  now suggest that

$$\Phi a = Z_1 f, \quad (3.25)$$

where  $Z_1$  is as defined before. Letting  $\Phi_1$  be composed of the first  $R$  rows of the matrix  $\Phi D$ , we obtain by (3.24) and (3.25) the following expressions for the unknowns  $f$  and  $a$ :

$$\begin{bmatrix} f \\ a \end{bmatrix} = \begin{bmatrix} Z_1 & -\Phi \\ 0 & \Phi_1 \end{bmatrix}^{-1} \begin{bmatrix} 0 \\ \bar{\pi} \end{bmatrix}. \quad (3.26)$$

The initial condition  $f$  together with the coefficients  $a_i$ 's found above completes the analysis for the finite buffer case. The queue length and queueing delay expressions are the same as in (3.18) and (3.19), whereas the rate of the packet loss,  $p_{loss}$ , which occurs when the buffer is full can be determined from  $u_n$ 's as follows [54]:

$$p_{loss} = \frac{1}{\alpha N} \sum_{n=R+1}^N (n - R) u_n.$$

### 3.5 Effective Bandwidth

In the preceding sections, the focus is on the analysis of many on-off sources of a single type onto a communication link. When the uniform arrival and service model (UAS model) is assumed to describe the traffic characteristics of an on-off source, the so-called fluid flow approximation techniques [1] provide efficient algorithms for calculating the probability that the buffer is exceeded in realistically sized systems. Within the same framework, we have examined the case for which the on-off sources generate packets periodically in the active times. The relationship between this model and the UAS model has been investigated in two main perspectives; accuracy in traffic modeling and computational requirements. Two generalizations are natural to consider. Firstly, we would like to be able to consider a single channel offered more than just one type of call which introduces the existence of multiple periods at the input of the link buffer. We note that this requires an approximation to the transient behavior of the  $\sum D_i/D/1$  queue.

Secondly, one would like to be able to consider a network of such channels. Effective bandwidth approximation of a call is an adequate approach for this purpose. To explain, suppose calls of  $m$  different types are offered to our channel and we wish to know how many of each type of call the channel can handle while proving a reasonable grade of service. In other words, we want to find the following region in the  $m$ -dimensional space:

$$\mathcal{A}(B, p) = \{N = (N_1, N_2, \dots, N_m) : P_N\{\text{queue length} \geq B\} \leq p\}, \quad (3.27)$$

where  $B$  is the buffer size and  $p$  is the maximum packet loss probability that each user can tolerate. In connection-oriented networks this problem is called the call admission problem since when a call comes the network checks whether the new user vector  $N$  satisfies (3.27) and decides on accepting or rejecting the new call. It is reasonable to approximate  $\mathcal{A}(B, p)$  by

$$\mathcal{A}(B, p) \approx \{N : \sum_i e_i N_i < C\},$$

where  $e_i$  is the effective bandwidth assigned to type  $i$  and  $C$  is the channel capacity. This approximation enables one to view the channel as a standard circuit-switched link which in turn provides the extension of the model and analysis to a network of channels through approximations such as the Erlang Fixed Point Procedure [29].

When UAS model is used in an ATM network, effective bandwidth approximation is possible in the asymptotic regime  $B \rightarrow \infty$  and  $p \rightarrow 0$  in such a manner that  $\log p/B \rightarrow \xi \in [-\infty, 0]$ . Our aim in this section is to emphasize that the same approximation is also valid for the Markov modulated periodic arrival process in this asymptotic regime. We will not restate the technical details associated with the effective bandwidth approximation for Markov modulated fluid sources in [12] but only the corresponding results in the cited reference together with certain explanations which we find necessary to build up the interconnection between periodic models and UAS models.

We now define

$$F_e(x) = \{F_e(n, x)\}, \quad n = 0, 1, \dots, N$$

$$D = \text{diag}\{Pn - C\}, \quad n = 0, 1, \dots, N$$

$$M = \text{generator matrix for the Markov chain,}$$

so that the governing system of differential equations (3.12) can be rewritten as

$$\frac{d}{dx}DF_e(x) = MF_e(x), \quad x \geq R - 1. \quad (3.28)$$

Since the stationary state distribution has a bounded solution, it has the spectral representation

$$F_e(x) = \sum_{i: \text{Re } z_i < 0} a_i \phi_i \exp(z_i x) + \pi$$

where  $\pi = \{\pi_n\}$ ,  $n = 0, 1, \dots, N$  and the pair  $(z_i, \phi_i)$  is an eigenvalue-eigenvector pair. Such pairs are solutions to the eigenvalue problem

$$zD\phi = M\phi. \quad (3.29)$$

Note that this is the same eigenvalue problem posed for Markov modulated fluid sources. In other words, under the same modulation, fluid sources and periodic sources have the same spectral expansion for large  $x$  except for the coefficients  $a_i$ 's. Indexing now the eigenvalues with negative real parts

$$0 > z_1 \geq \text{Re } z_2 \geq \dots$$

the real eigenvalue  $z_1$  is called the dominant eigenvalue. Let the stationary buffer overflow distribution be given by  $G(x)$ , ( $x \geq R - 1$ ):

$$\begin{aligned} G(x) &= \text{Pr}\{\text{queue length} > x\} \\ &= \sum_i -a_i \left( \sum_j \phi_i(j) \right) \exp(z_i x). \end{aligned}$$

When  $x \rightarrow \infty$ , the eigenvalue  $z_1$  will dominate so that we can write

$$G(x) \sim -a_1 \left( \sum_j \phi_1(j) \right) \exp(z_1 x) \text{ as } x \rightarrow \infty. \quad (3.30)$$

Note that plots of  $\log G(x)$  vs.  $x$  approach linearity as  $x$  increases and the slope approaches  $z_1$  and this slope is the same in both the UAS model and the periodic model.



This is because the models have the same eigenvalue-eigenvector pairs in this regime (see equations (2.4) and (3.29)).

It is now convenient to view  $C$  as a variable parameter and to write down the eigenvalues to be a function of  $C$ ,  $z(C)$ . The inverse problem requires  $C$  to be obtained for given  $z$ . The key fact in this connection is that the inverse problem is also an eigenvalue problem. Writing then  $C = g(z)$ , the equation (3.29) becomes

$$g(z)\phi = A(z)\phi,$$

where

$$A(z) = \Lambda - \frac{1}{z}M, \quad \Lambda \triangleq \text{diag}\{Pn\}.$$

In the above equation,  $g(z)$  is an eigenvalue of the matrix  $A(z)$  in which  $z$  is a parameter.  $g_1(z)$  is called the maximal real eigenvalue; if  $g(z)$  is any other eigenvalue then  $\text{Re } g(z) < g_1(z)$ .

The maximal real eigenvalue has the following properties [12];

- i)  $g_1(0) = \text{mean source rate}$ ,  $g_1(-\infty) = \text{peak source rate}$ ,
- ii)  $g_1'(z) < 0$ , ( $z < 0$ )
- iii) The dominant eigenvalue  $z_1$  is the unique solution in  $(-\infty, 0)$  satisfying  $g_1(z_1) = C$ .

We now consider the admission control problem for an asymptotic regime in which  $B \rightarrow \infty$ ,  $p \rightarrow 0$  in such a manner that  $\log p/B \rightarrow \xi$ . Actually, the graph of  $\log p$  vs.  $B$  has the following description in the prescribed regime:

$$\log p = \xi B + k,$$

for which  $k$  is a finite translation parameter. Then, by (3.30),

$$\frac{G(B)}{p} = k_0 \exp((z_1 - \xi)B), \quad \text{as } B \rightarrow \infty,$$

for some constant  $k_0$ . Notice by property ii) that  $g_1(z)$  decreases as  $z$  increases, hence if  $g_1(\xi) < C$  then  $z_1 < \xi$  and thus  $G(B)/p \rightarrow 0$  as  $B \rightarrow \infty$ . The admission control

criterion is satisfied in this case. Similarly, if  $g_1(\xi) > C$  then  $z_1 > \xi$  and  $G(B)/p \rightarrow \infty$ , so the admission criterion is violated. The above result justifies the use of the term *effective bandwidth* for the quantity  $g_1(\xi)$ . We let  $e = e(M, \Lambda; B, p)$  denote the effective bandwidth of the superposition  $(M, \Lambda)$  in the system for which the admission criterion is  $G(B) \leq p$ . That is,

$$e(M, \Lambda; B, p) = g_1(\xi) \quad (3.31)$$

where  $g_1(\xi)$  is the maximal real eigenvalue of the matrix  $\Lambda - \frac{1}{\xi}M$  and  $\xi = \log p/B$ .

Bad news is that computation of the maximal real eigenvalue of the matrix above turns out to be difficult if the number of states of the underlying Markov chain is large. Good news is that if the incoming traffic is made up of a superposition of independent individual sources then computation of the effective bandwidth of the aggregate source can be made through computation of the effective bandwidths of those simpler individual sources. Actually, the effective bandwidth is an additive quantity in this asymptotic regime. We will give the following powerful result in [12] where all technical results are given for fluid sources. Since the asymptotic slopes of the buffer overflow function of fluid sources and periodic sources are the same, the result is also valid for periodic sources.

**Proposition 3.1.** *Suppose there are  $K$  Markov modulated periodic sources,  $(M^{(k)}, \Lambda^{(k)})$ , ( $1 \leq k \leq K$ ), offered to an ATM multiplexer.  $M^{(k)}$  is the generator matrix of the underlying Markov chain and  $\Lambda^{(k)}$  is the rate matrix (in packets/sec.) associated with source  $k$ . Let the admission control criterion be  $G(B) \leq p$ . Suppose  $B \rightarrow \infty$  and  $p \rightarrow 0$  in such a manner that  $\log p/B \rightarrow \xi \in [-\infty, 0]$ . If*

$$\sum_k g_1^{(k)}(\xi) < C,$$

*then the admission criterion is satisfied. If the inequality sign is reversed then the admission criterion is violated. Here,  $g_1^{(k)}(\xi)$  is the maximal real eigenvalue of  $A^{(k)}(\xi) \triangleq \Lambda^{(k)} - \frac{1}{\xi}M^{(k)}$ .  $\square$*

To give an example, an on-off source with exponentially distributed on and off periods

is obtained by setting

$$M = \begin{bmatrix} -\lambda & \mu \\ \lambda & -\mu \end{bmatrix} \quad \text{and} \quad \Lambda = \begin{bmatrix} 0 & 0 \\ 0 & P \end{bmatrix}.$$

In this case, it can easily be verified that

$$g_1(z) = \frac{Pz + \lambda + \mu - ((Pz + \lambda + \mu)^2 - 4\lambda Pz)^{1/2}}{2z} \quad (3.32)$$

which is a simply computed quantity in terms of  $P$ ,  $\lambda$ , and  $\mu$ . But for sources with more than two states, it is in general difficult to give such explicit expressions for the effective bandwidth.

## Chapter 4

# Padé Approximations in the Analysis of the MMPP/D/1 System

The performance analysis of integrated services networks whose inputs consist of a superposition of different packetized sources has been an intense area of research in the recent years. The Markov Modulated Poisson Process (MMPP) is indeed one of the general input traffic models for ATM networks the performance analysis of which is essential in clarifying the features of statistical multiplexing of bursty sources and developing efficient resource management schemes.

The packet arrival intensity for an MMPP is modulated with respect to the state of a continuous-time  $N$ -state Markov chain. To explain, let the state of the Markov chain be denoted by  $n \in \{1, 2, \dots, N\}$ . The generator matrix of the Markov process is denoted by  $M$ ;  $M(n, m)$  is the transition rate from state  $m$  to state  $n$  and  $M(n, n) = -\sum_{m \neq n} M(m, n)$ . The equilibrium probability of the Markov chain being in state  $n$  is denoted by  $\pi_n$ . Within each state  $n$  of the Markov chain, the packet arrival process is a Poisson process with rate  $\lambda_n$  packets/sec. This arrival process is offered to a server which is specified to be deterministic with rate  $C$  packets/sec. since packet lengths are fixed

in an ATM environment. Such a system is called an MMPP/D/1 queue (see Figure 1.9). We are interested in the queue length distribution function with respect to the utilization in the system given the underlying Markov chain and the associated Poisson rates of packet arrivals.

The deterministic service time is in general hard to tackle. The assumption that packet lengths are exponentially distributed with mean being the deterministic service time of the multiplexer yields an MMPP/M/1 queue which is easier to handle. However, this brings overestimates in the queue length since the behavior of this system over a given period of time is identical to the behavior of an M/M/1 queue and therefore experiences higher frequency fluctuations in system occupancy than does the real queue (M/D/1 queue) over that period of time [6]. Here, we propose a method based on Padé approximations in the transform domain. A Padé approximant is actually the ratio of two polynomials constructed from the coefficients of the Taylor series of the function which needs to be approximated to ensure analytical tractability. The freedom in choosing the degrees of these two polynomials yields different approximations whose accuracy improves as these degrees increase. The number of coefficients in the Taylor series expansion of the function that the Padé approximant can capture is called the order of the approximation. In this respect, fluid flow approximations and MMPP/M/1 queues are shown to be obtained through first order Padé approximations for the service time, the original system being the MMPP/D/1 queue. Motivated by this fact, we use more sophisticated approximations which capture not only the mean but also higher order terms of the service time distribution and whose computational complexities are no worse than the complexity encountered while solving the MMPP/M/1 queue. The authors of [45] use Erlang distributions [30] in time domain to approximate the deterministic service time. Recall that the exponential distribution is the first order Erlang distribution for the deterministic service time. Since any other Erlang distribution will involve an additional complexity in the solution of the MMPP/D/1 queue in our framework, we only concentrate on the particular MMPP/M/1 queueing system to make comparisons.

Another contribution here is a new derivation of the unfinished work distribution

expression in transform domain for the MMPP/G/1 queue. Our proof is simple and based upon the Takács integro-differential equation which describes the transient behavior of the M/G/1 queue. Through the transform domain counterpart of this equation, we obtain transform domain formulas for the unfinished work distribution in an MMPP/G/1 queue. Our formulation is able to include the case when the service time is approximated by Padé approximants. The essential point to note is that a Padé approximant is not necessarily associated with a physical distribution (i.e., may possibly contain higher order derivatives of the Dirac delta function at the origin or may involve negative probabilities).

There is in fact a rich underlying theory developed for Markov modulated fluid sources [49] and the MMPP/M/1 queue [13] which provides efficient algorithms for calculating the stationary state distribution of the corresponding system. This turns out to be especially significant when many identically distributed sources are superposed and fed into a FIFO buffer. In this case, the algebraic theory gives an exact decomposition of the overall system into many smaller subsystems and is therefore able to avoid numerical problems when the number of states in the Markov chain tend to increase. Our proposed method is naturally capable of making use of this theory through the use of Padé approximations. The theory is quite general and ensures an easy treatment of more sophisticated performance evaluation problems in ATM networks.

The organization of this chapter is as follows. We first focus on the transient behavior of the M/G/1 queue from which we obtain transform domain formulas for the steady-state buffer occupancy distribution in an MMPP/G/1 queue. Employing particular Padé approximants of different orders for the deterministic service time distribution, we give a numerical example (for a packetized voice multiplexer) to demonstrate the performance comparison of these approximations in the context of MMPP/D/1 queues. Computational aspects are then discussed and a powerful algorithm is presented which works for the case the input traffic to the buffer is a superposition of many 2-state MMPP's of the same type. We then present a numerically efficient algorithm for the MMPP/D/1/K queue where the buffer size is assumed to be finite. Finally, bandwidth assignment problem to a call of MMPP type is addressed.

## 4.1 Transient Analysis of the M/G/1 Queue

In this section, we take a closer look at the unfinished work  $U(t)$  (time required to empty all the packets present in the system at time  $t$ ) in an M/G/1 queue and derive the forward Kolmogorov equation for its time-dependent behavior. Here we give an account of the solution to  $U(t)$  (see [30] for a detailed analysis) to establish the notation and setting for the statement of our results.

We wish to derive the probability distribution function for  $U(t)$  given its initial value at time zero. Accordingly, we define

$$Q(w, t) \triangleq Pr\{U(t) \leq w | U(0)\}. \quad (4.1)$$

The arrival rate to the queue is Poisson with rate  $\lambda$  packets/sec. Let  $\tilde{x}$  be the random variable associated with the service time. We then define the following functions:

$$\begin{aligned} B(x) &= Pr\{\tilde{x} \leq x\}, \\ b(x) &= \frac{d}{dx}B(x), \\ \hat{B}(s) &= \int_{0^-}^{\infty} \exp(-sx)b(x)dx. \end{aligned}$$

Note that  $b(x)$  and  $B(x)$  are the probability distribution function (pdf) and the cdf of the service time, respectively. We now wish to relate the probability  $Q(w, t + \Delta t)$  to its possible values at time  $t$ . We observe that we can reach this state from  $t$  if, on the one hand, there had been no arrivals during this increment of time (this occurs with probability  $1 - \lambda\Delta t + o(\Delta t)$ ) and the unfinished work was no larger than  $w + \Delta t$  at time  $t$ ; or if, on the other hand, there had been an arrival in this increment of time (with probability  $\lambda\Delta t + o(\Delta t)$ ) such that the unfinished work at time  $t$ , plus the new increment of work brought by this customer do not exceed  $w$ . These observations lead us to the following equation:

$$Q(w, t + \Delta t) = (1 - \lambda\Delta t)Q(w + \Delta t, t) + \lambda\Delta t \int_{x=0}^w B(w - x) \frac{\partial}{\partial x} Q(x, t) dx + o(\Delta t), \quad (4.2)$$

where

$$d_x Q(x, t) \triangleq \frac{\partial}{\partial x} Q(x, t) dx$$

is the pdf for unfinished work at time  $t$ . Expanding our distribution function on its first variable and making some manipulations on (4.2), one finally obtains the Takács integro-differential equation for  $U(t)$ :

$$\frac{\partial Q(w, t)}{\partial t} = \frac{\partial Q(w, t)}{\partial w} - \lambda Q(w, t) + \lambda \int_0^\infty B(w-x) d_x Q(x, t). \quad (4.3)$$

Takacs [50] shows that this equation is good for almost all  $w \geq 0$  and  $t \geq 0$ ; it does not hold at those  $w$  and  $t$  for which  $\partial Q(w, t)/\partial w$  has an accumulation of probability (namely, an impulse). Now defining

$$W^*(s, t) \triangleq \int_{0^-}^\infty \exp(-sw) dQ_w(w, t),$$

we obtain

$$\int_{0^-}^\infty Q(w, t) \exp(-sw) dw = \frac{W^*(s, t)}{s}, \quad (4.4)$$

and, similarly, we have

$$\int_{0^-}^\infty B(w) \exp(-sw) dw = \frac{\hat{B}(s)}{s}.$$

Taking the Laplace transform on the first variable of each side of (4.3), one obtains

$$\frac{1}{s} \frac{\partial W^*(s, t)}{\partial t} = W^*(s, t) - \frac{\lambda W^*(s, t)}{s} + \lambda \frac{W^*(s, t) \hat{B}(s)}{s} - Q_0(t), \quad (4.5)$$

where  $Q_0(t)$  is some suitable function that is to be subtracted to compensate for the impulsive terms since the Takács integro-differential equation does not contain these impulses. The equation (4.5) can then be rewritten as

$$\frac{\partial W^*(s, t)}{\partial t} = (s - \lambda + \lambda \hat{B}(s)) W^*(s, t) - s Q_0(t). \quad (4.6)$$

Takács gives the solution to this equation in [51]. We may now transform on our second variable  $t$  by first defining the double transform

$$Q^{**}(s, r) \triangleq \int_0^\infty \exp(-rt) W^*(s, t) dt.$$

We also need the definition

$$Q_0^*(r) \triangleq \int_0^\infty \exp(-rt) Q_0(t) dt.$$



One may now transform equation (4.6) to obtain

$$Q^{**}(s, r) = \frac{W^*(s, 0) - sQ_0^*(r)}{r - s + \lambda - \lambda\hat{B}(s)} \quad (4.7)$$

where the unknown function  $Q_0^*(r)$  is determined by insisting that the transform  $Q^{**}(s, r)$  be analytic in the region  $Re(s) > 0, Re(r) > 0$ . The expression (4.7) will be the fundamental equation which will lead us to the stationary queue length distribution for the MMPP/G/1 queue in the next section.

## 4.2 MMPP/G/1 Queue

This section is devoted to the analysis of the MMPP/G/1 queue. Consider now the Markov chain that governs the Poisson rate of arrivals with the infinitesimal generator matrix  $M$ . Note that the state holding time at state  $n$  is exponentially distributed with parameter  $\sigma_n = -M(n, n)$ . Let us now focus on a particular phase  $n$  of the MMPP in which the arrival process is Poisson with rate  $\lambda_n$ . We now define  $Q_b(n, w)$  and  $Q_e(n, w)$  to be the equilibrium unfinished work cdf's at the moment of state transition to  $n$  and at the moment of state transition from  $n$ , respectively. Our objective is to find the relation between these two cdf's in transform domain via the use of equation (4.7). For this purpose let  $Q(n, w; t)$  be the unfinished work cdf  $t$  seconds after the chain has made a state transition to  $n$  and that the initial unfinished work has the cdf  $Q_b(n, w)$ . It is now easy to write

$$Q_e(n, w) = \int_0^\infty Q(n, w; t)\sigma_n \exp(-\sigma_n t) dt. \quad (4.8)$$

Also let  $\hat{Q}_b(n, s)$  and  $\hat{Q}_e(n, s)$  be the Laplace transforms of  $Q_b(n, w)$  and  $Q_e(n, w)$ , respectively. At a closer look at the definitions of the preceding section, we then directly make the following substitutions

$$\begin{aligned} Q(n, w; t) &\mapsto Q(w, t) \quad (\text{eq. (4.1)}) \\ \hat{Q}_b(n, s) &\mapsto \frac{W^*(s, 0)}{s} \quad (\text{eq. (4.4)}) \\ \hat{Q}_e(n, s) &\mapsto \frac{\sigma_n Q^{**}(s, \sigma_n)}{s} \quad (\text{eq. (4.7)}) \end{aligned}$$

so that one can write by (4.7)

$$\begin{aligned}\hat{Q}_e(n, s) &= \frac{\sigma_n s \hat{Q}_b(n, s) - s Q_n^*(\sigma_n)}{s \sigma_n - s + \lambda_n - \lambda_n \hat{B}(s)}, \\ &= \frac{\sigma_n (\hat{Q}_b(n, s) - Q_n^*(\sigma_n))}{\sigma_n - s + \lambda_n - \lambda_n \hat{B}(s)},\end{aligned}\quad (4.9)$$

where the unknown constant  $Q_n^*(\sigma_n)$  is chosen by insisting that the transform  $\hat{Q}_e(n, s)$  be analytic in the region  $Re(s) > 0$ .

We define the transform vectors

$$\hat{W}_e(s) = \{\hat{W}_e(n, s) \triangleq \pi_n \hat{Q}_e(n, s)\}, \quad \hat{W}_b(s) = \{\hat{W}_b(n, s) \triangleq \pi_n \hat{Q}_b(n, s)\}. \quad (4.10)$$

Actually,  $\hat{W}_e(s)$  is called the Laplace transform of the virtual delay distribution in an MMPP/G/1 queue [20]. Replacing  $X(t)$  by  $U(t)$  in the derivation of equality (2.16), we immediately obtain

$$\sigma_n \hat{W}_b(n, s) = \sum_{m \neq n} M(n, m) \hat{W}_e(m, s). \quad (4.11)$$

Substituting (4.10) and (4.11) in (4.9), one has  $\forall n = 1, 2, \dots, N$ ,

$$(s - \lambda_n + \lambda_n \hat{B}(s)) \hat{W}_e(n, s) + \sum_{m=1}^N M(n, m) \hat{W}_e(m, s) = \pi_n \sigma_n Q_n^*(\sigma_n). \quad (4.12)$$

One may alternatively rewrite the above equation in matrix form:

$$\hat{W}_e(s) = [sI + M - R + R\hat{B}(s)]^{-1} y_0, \quad (4.13)$$

where

$$R = \text{diag}\{\lambda_1, \lambda_2, \dots, \lambda_N\}$$

and

$$y_0 \triangleq \{\pi_n \sigma_n Q_n^*(\sigma_n), n = 1, 2, \dots, N\}$$

should be determined by insisting that the transform be analytic in the region  $Re(s) > 0$ .

Whenever  $\lim_{s \rightarrow \infty} \frac{\hat{B}(s)}{s} = 0$ , as is typically the case,

$$\begin{aligned}y_0(n) &= \lim_{s \rightarrow \infty} s \hat{W}_e(n, s), \\ &= \pi_n \lim_{s \rightarrow \infty} s \hat{Q}_e(n, s), \\ &= \pi_n Q_e(n, 0),\end{aligned}$$

which then equals to the probability of the system being empty and the chain residing in state  $n$ . In this case, the expression (4.13) is equivalent to the virtual waiting time distribution given in [20] found by matrix analytic methods. Since the Padé approximants of  $\hat{B}(s)$  do not necessarily have this limiting property (will be explained in the following development), the expression (4.13) is more general than the result in [20].

### 4.3 Padé Approximations in the MMPP/D/1 Queue

Let us now focus our attention to the special case of the service time being deterministic, for which

$$b(x) = u_0(x - 1/C),$$

where  $u_0(\cdot)$  is the Dirac delta function and the associated Laplace transform is simply an exponential function of  $s$ :

$$\hat{B}(s) = \exp(-s/C).$$

When the service time has a general distribution, we have used the unfinished work (or virtual waiting time) to describe the queueing behavior in order to conform to the literature. In the deterministic server case, we rather prefer to use the queue length as the entity of interest to facilitate the understanding of the development in this section in view of the other Markov modulated models examined in the preceding chapters.

Let  $F_b(n, x)$  and  $F_e(n, x)$  be defined in the same way as in the equations (2.8) and (2.9), respectively. We similarly define

$$F_e(x) = \{F_e(n, x)\}, \quad \hat{F}_e(s) = L[F_e(x)],$$

where  $L(\cdot)$  denotes the Laplace transform. The queue length distribution is simply related to the virtual waiting time via the following;

$$F_e(x) = W_e(x/C),$$

$$\hat{F}_e(s) = CW_e(sC).$$

The equation (4.13) is therefore reduced to

$$\hat{F}_e(s) = [sCI + M - R + R \exp(-s)]^{-1} f_0, \quad (4.14)$$

where  $f_0 = Cy_0$ . However, the constant vector  $f_0$  is to be determined using the poles of the overall system, which with an irrational term in the denominator is difficult to handle. Therefore, we make use of rational Padé approximants of the irrational transform  $\exp(-s)$  so as to determine a solution for the buffer content distribution. In fact, a Padé approximation with parameters  $n$  and  $m$  is a rational function

$$R_{n,m}(s) = \frac{P_n(s)}{Q_m(s)},$$

where  $P_n(s)$  and  $Q_m(s)$  are polynomials of order  $n$  and  $m$ , respectively, and the first  $(n + m + 1)$  terms of its Taylor series expansion equal to those of the Taylor expansion of  $\exp(-s)$ , or equivalently the first  $(n + m)$  moments of the original service time distribution. A closed form expression for  $R_{n,m}$  exists [55] and is given by

$$R_{n,m}(s) = \frac{\sum_{i=0}^n (m+n-i)! C(n,i) (-1)^i s^i}{\sum_{i=0}^m (m+n-i)! C(m,i) s^i},$$

where

$$C(n,i) = \frac{n!}{i! (n-i)!}.$$

In case  $R_{n,m}(s)$  is used as an approximant for  $\exp(-s)$ , the transform of the queue length cdf,  $\hat{F}_e(s)$  is written as

$$\hat{F}_e(s) = [sCI + M - R + R R_{n,m}(s)]^{-1} f_0. \quad (4.15)$$

An interesting observation is that, if the Padé approximant  $R_{1,0}(s) = 1 - s$  is employed, then  $\hat{F}_e(s)$  becomes

$$\begin{aligned} \hat{F}_e(s) &= [sCI + M - R + R(1 - s)]^{-1} f_0, \\ &= [s(CI - R) + M]^{-1} f_0, \end{aligned}$$

which is in fact equivalent to the expression suggested for the Markov modulated fluid sources [1],[49]. When the approximation  $R_{0,1}(s) = 1/(1+s)$  is imposed, the transform of the virtual waiting time  $W_e(s)$  is rewritten by equation (4.13):

$$\hat{W}_e(s) = [sI + M - R + R(\frac{C}{s+C})]^{-1}y_0,$$

which reduces to the virtual waiting time distribution [20] in an MMPP/M/1 queue, in which the service time has an exponential distribution with mean  $1/C$ . In this respect, fluid flow approximations and MMPP/M/1 queues turn out to be first order Padé approximations for the original MMPP/D/1 queue. These observations are remarkable in their own since it is now possible to unify some of the proposed approximations used for the analysis of statistical multiplexing in ATM networks in this framework.

Our main objective is to use the information on higher order moments of the service time distribution while preserving the degree of the characteristic polynomial

$$\psi(s) = \det((sCI + M - R)Q_m(s) + RP_n(s))$$

same as in the MMPP/M/1 queue (i.e.,  $2N$ ). This determinantal degree denoted by  $d$  is one of the major factors that determine the numerical efficiency of the suggested approximation. Note that

$$d = N \max(m+1, n).$$

Under this degree constraint, there are actually three more rational functions as Padé approximants of  $\exp(-s)$ ;

$$\begin{aligned} R_{1,1}(s) &= \frac{1-s/2}{1+s/2} \\ R_{2,0}(s) &= \frac{1-s+s^2/2}{1} \\ R_{2,1}(s) &= \frac{1-2s/3+s^2/6}{1+s/3} \end{aligned}$$

We now sketch an outline of our solution technique for each approximation  $R_{n,m}(s)$ ;

1) Define  $\pi = [\pi_1 \ \pi_2 \ \cdots \ \pi_N]^T$  and then obtain the spectral expansion of  $F_e(x)$ ;

$$F_e(x) = \pi + \sum_{i: \operatorname{Re} z_i < 0} a_i \exp(z_i x) \phi_i, \quad (4.16)$$

and

$$[(z_i CI + M - R)Q_m(z_i) + RP_n(z_i)]\phi_i = 0.$$

- 2) Use the initial value theorem for Laplace transforms in (4.15) to obtain the linear relationship between  $F_e(0)$  and  $\frac{d}{dx}F_e(x)|_{x=0+}$ .

To explain this stage of the procedure,  $N$  equations are needed to solve for the  $N$  unknown coefficients ( $a_i$ 's). These equations can easily be constructed by the above-mentioned linear relationship which is inherent in the equation (4.15).

- 3) Combine 1) and 2) to find out  $a_i$ 's from a set of  $N$  linear equations.
- 4) The overall cdf of the queue length is the sum of the individual elements of  $F_e(n, x)$ :

$$Pr(\text{queue length} \leq x) = \sum_{n=1}^N F_e(n, x). \quad (4.17)$$

The cdf associated with the queueing delay is a bit different [20]:

$$Pr(\text{delay} \leq t \text{ sec.}) = \frac{1}{\bar{\lambda}} \sum_{n=1}^N \lambda_n F_e(n, Ct) \quad (4.18)$$

where

$$\bar{\lambda} = \sum_{n=1}^N \pi_n \lambda_n$$

is the mean rate of the incoming packet stream.  $\square$

To explain the above procedure neatly, let us assume that  $R_{2,1}(s)$  is imposed as a Padé approximation for the deterministic service time. We then obtain the  $N$  ( $z_i, \phi_i$ ) pairs satisfying

$$[(z_i CI + M - R)(1 + z_i/3) + R(1 - 2z_i/3 + z_i^2/6)]\phi_i = 0,$$

in the first step. For the second step, observe by (4.15) that the following equation holds:

$$(As^2 + Bs + M)\hat{F}_e(s) = (1 + s/3)f_0, \quad (4.19)$$

where

$$\begin{aligned} A &\triangleq \frac{CI}{3} + \frac{R}{6}, \\ B &\triangleq CI - R + \frac{M}{3}. \end{aligned}$$

Now defining

$$F_2(x) = A \frac{d}{dx} F_e(x), \quad \hat{F}_2(s) = L[F_2(x)],$$

we immediately have

$$\begin{bmatrix} sI & -A^{-1} \\ M & sI + BA^{-1} \end{bmatrix} \begin{bmatrix} \hat{F}_e(s) \\ \hat{F}_2(s) \end{bmatrix} = \begin{bmatrix} F_e(0) \\ f_0 - BF_e(0) + s(f_0/3 - AF_e(0)) \end{bmatrix}. \quad (4.20)$$

But by using the initial value theorem on (4.19),

$$\begin{aligned} F_e(0) &= \lim_{s \rightarrow \infty} s \hat{F}_e(s) \\ &= \frac{A^{-1} f_0}{3}, \end{aligned}$$

the term  $s(f_0/3 - AF_e(0))$  vanishes in the second entry of the right-hand side of (4.20).

Taking this fact into account, we obtain by (4.20) that

$$\begin{aligned} \frac{d}{dx} F_e(x)|_{x=0^+} &= A^{-1} F_2(0), \\ &= A^{-1}(f_0 - BF_e(0)), \\ &= (3I - A^{-1}B)F_e(0), \\ &= \Delta F_e(0). \end{aligned} \quad (4.21)$$

This completes the second step of the procedure for  $R_{2,1}(s)$ . Similar procedures apply for the other approximants, but we omit these details in this context. In the third step, we first define

$$\begin{aligned} \Phi &= \begin{bmatrix} \phi_1 & \phi_2 & \cdots & \phi_N \end{bmatrix}, \\ Z &= \text{diag}\{z_1, z_2, \dots, z_N\}, \\ a &= \begin{bmatrix} a_1 & a_2 & \cdots & a_N \end{bmatrix}^T. \end{aligned}$$

Then the initial values for the queue length cdf and the queue length distribution function can be written in state form in terms of the coefficient vector  $a$  via the use of (4.16) as follows:

$$\begin{aligned} F_e(0) &= \pi + \Phi a, \\ \frac{d}{dx} F_e(x)|_{x=0^+} &= \Phi Z a. \end{aligned}$$

This then immediately yields by (4.21) that

$$\Phi Z a = \Delta(\pi + \Phi a),$$

by which

$$a = (\Phi Z - \Delta\Phi)^{-1} \Delta\pi.$$

At this point, we are done with the third step. We complete the procedure by noting that the eigenvalues and the eigenvectors  $z_i$ 's and  $\phi_i$ 's of the system together with the coefficients  $a_i$ 's determine exactly the queue length cdf through the expression (4.16).

### 4.3.1 Numerical Examples

We now demonstrate our results in the following example, in which a superposition of bursty voice sources feeds a deterministic server. The voice sources in their active times are assumed to generate packets with respect to a Poisson process, the rate of which is determined by the peak rate of an individual source (see figure 1.10). This model for packetized voice has indeed been examined in [25]. We consider a packetized voice system with a line speed 10 times the voice peak rate 32 kbits/s, mean active period 353 ms and mean silent period 650 ms. The mean number of packets generated in an active period is 22. The transmission time of a single packet is 1.6 ms.

In Table 4.1, comparison of the average waiting time of the queueing system is given when the above-mentioned Padé approximations are used. The results are obtained by varying the number of voice sources. A parameter pair  $(n, m)$  in this table refers to the approximation associated with  $(\exp(-s) \approx R_{n,m}(s))$ . The simulation results are based



No. sources	simulation results	95 % conf. interval	approximations [ms]				
			(1,0)	(0,1)	(1,1)	(2,0)	(2,1)
4	0.2179	$\pm 0.0007$	0.00	0.4348	0.2181	0.2185	0.2192
6	0.331	$\pm 0.001$	0.00	0.659	0.331	0.331	0.331
8	0.478	$\pm 0.002$	0.00	0.947	0.475	0.477	0.476
10	0.676	$\pm 0.0025$	0.00	1.335	0.674	0.676	0.674
12	0.977	$\pm 0.005$	0.002	1.897	0.970	0.975	0.972
14	1.483	$\pm 0.015$	0.038	2.787	1.475	1.486	1.478
16	2.465	$\pm 0.025$	0.268	4.332	2.449	2.476	2.458
18	4.528	$\pm 0.06$	1.199	7.222	4.497	4.549	4.514
20	9.05	$\pm 0.1$	3.56	12.94	8.96	9.05	8.99
22	19.52	$\pm 0.1$	12.02	24.89	19.00	19.14	19.03
24	43.0	$\pm 0.6$	34.6	52.6	43.3	43.6	43.4
26	120.3	$\pm 1.5$	105.4	136.2	119.4	119.7	119.5
28	1064	$\pm 95$	940	1049	1001	1002	1002

**Table 4.1:** Performance comparison of the Padé approximations in terms of the mean waiting time.

on the discrete-time MMPP/D/1 queue. Note that, while the first two approximations (0,1) and (1,0) match only the first moment (mean), the approximations (1,1), (2,0), and (2,1) match the first two and the first three moments, respectively, of the original service time distribution. Since the second moment of the service time distribution is critical in the mean waiting time expression of the M/G/1 queue in the Pollaczek-Khinchin formula [30], we expect the latter three approximations to give remarkably accurate results. This is in fact the case; all these three approximations find the mean waiting time within 95 % confidence intervals when compared with simulations. However, we did not gain anything in terms of average waiting time estimates by using this third moment information in (2,1) which in a way resembles the Pollaczek-Khinchin formula for the mean waiting time in an M/G/1 queue since mean queueing delay only depends on the mean arrival and service rates and the second moment of the service time. The approximation (0,1) (MMPP/M/1 queueing model) overestimated the packet delay for

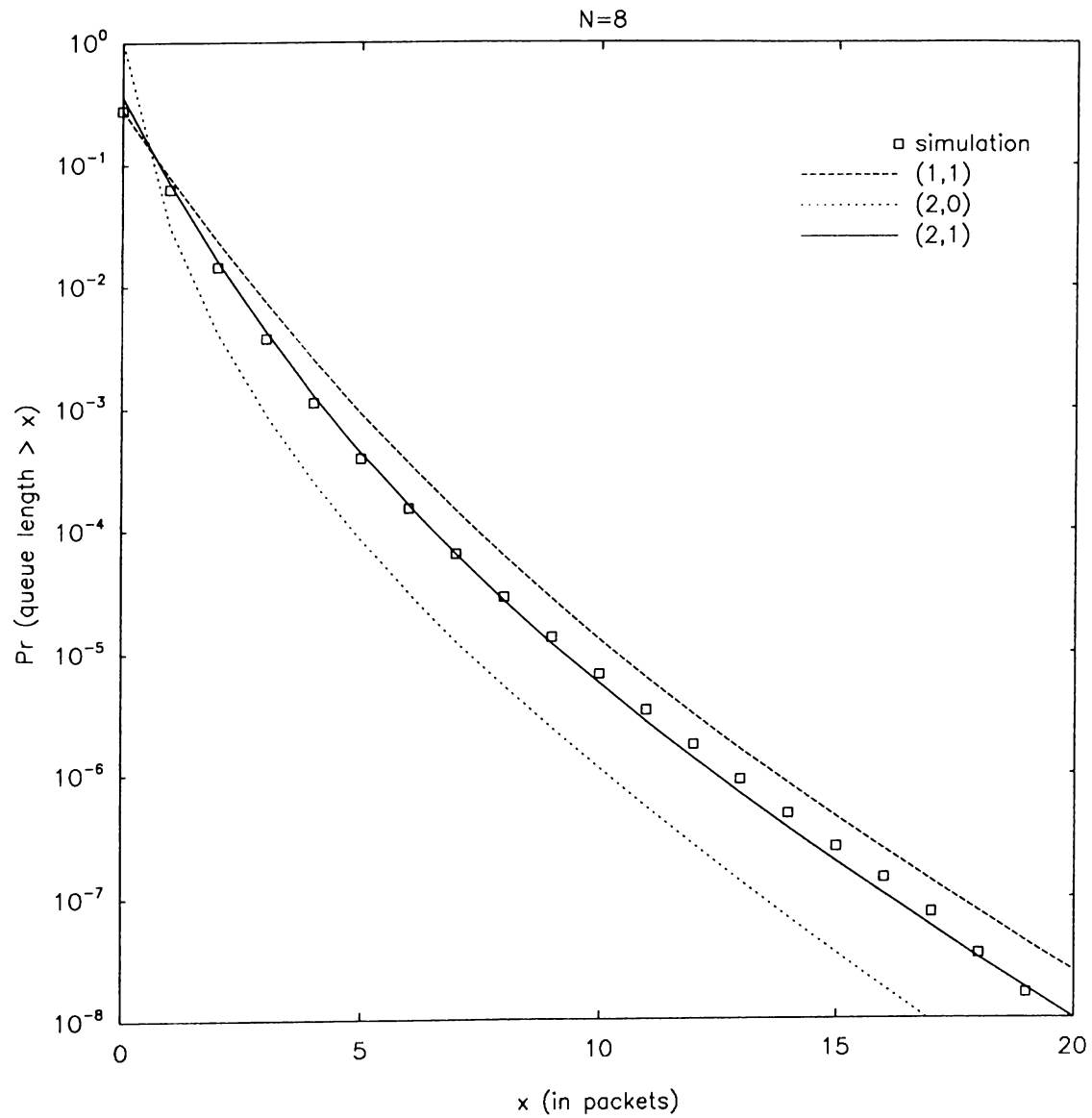
all possible loads, degree of overestimation decreased with increasing load. The fluid model (1,0) does not work for low to moderate loads, but works pretty well in heavy traffic.

In Figures 4.1, 4.2, 4.3, and 4.4, the queue length survivor function  $Pr(\text{Queue length} > x)$  is plotted for the three Padé approximations (1,1), (2,0), and (2,1) for the cases  $N$  (number of voice sources) being 8, 10, 15, and 20, respectively.

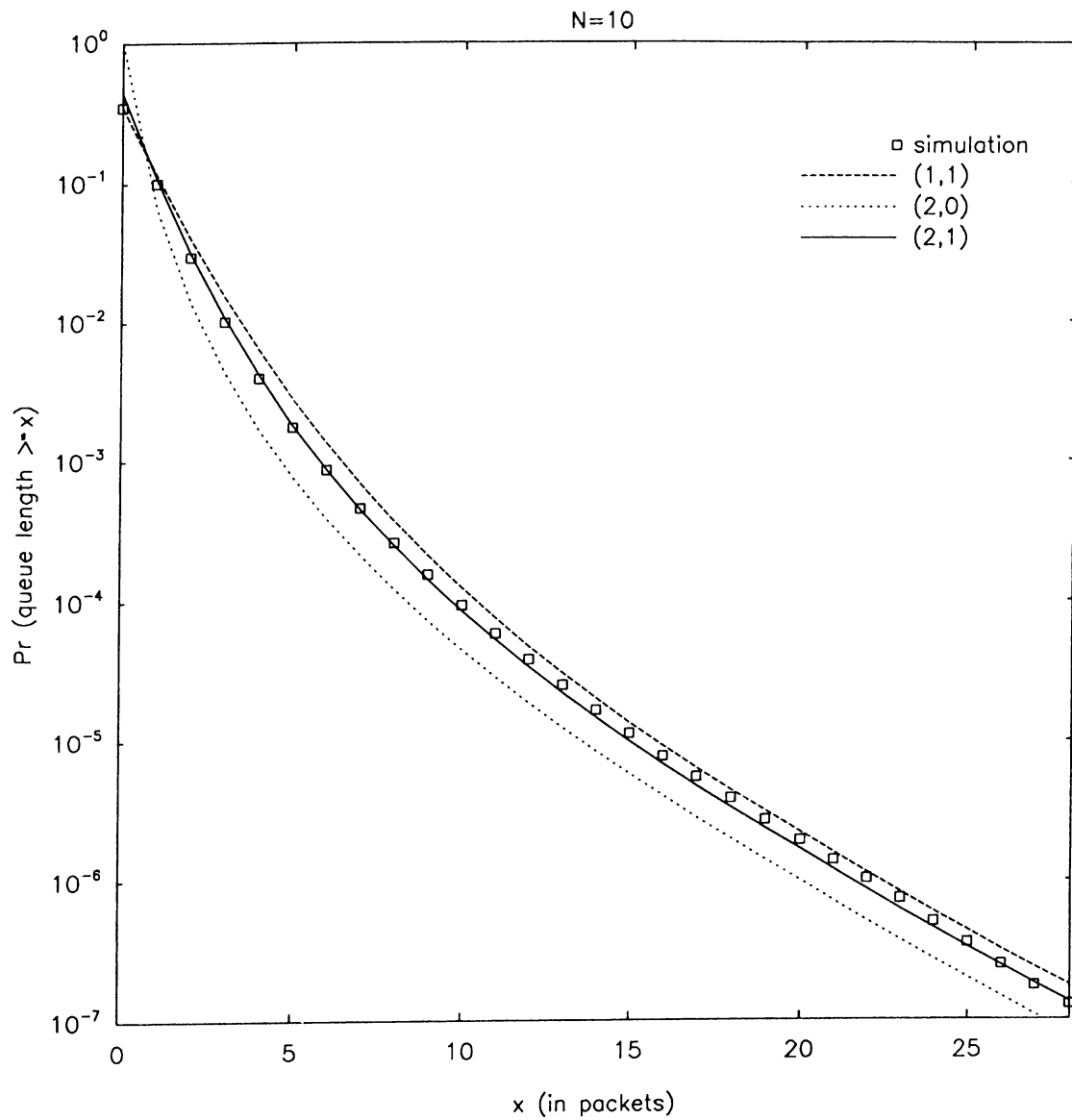
Except for a small deviation from the simulation results for very small  $x$ , all three approximations give satisfactory results, the approximation (2,1) being the most accurate one. When load increases, the performance of the approximations tends to be indistinguishable. In ATM networks, it is believed that the most important performance measures of statistical multiplexers include the tail probabilities  $Pr(\text{queue length} > x)$  for large  $x$  (the asymptotic behavior) as well as the averaged measures (e.g., mean delay, mean delay jitter). The packetized voice example presented here shows us that, these measures can closely be approximated with the use of Padé approximations without a need for cumbersome simulations. This fact has long been known in circuit theory [41]; asymptotic waveform evaluation techniques (or Padé approximations in transform domain terminology) are capable of capturing the asymptotic behavior of the approximated function as well as certain prespecified weighted averages belonging to that particular function. We believe that this particular application of Padé theory to queueing systems of this nature will be significant especially in developing certain congestion control schemes in ATM networks.

## 4.4 Computational Aspects

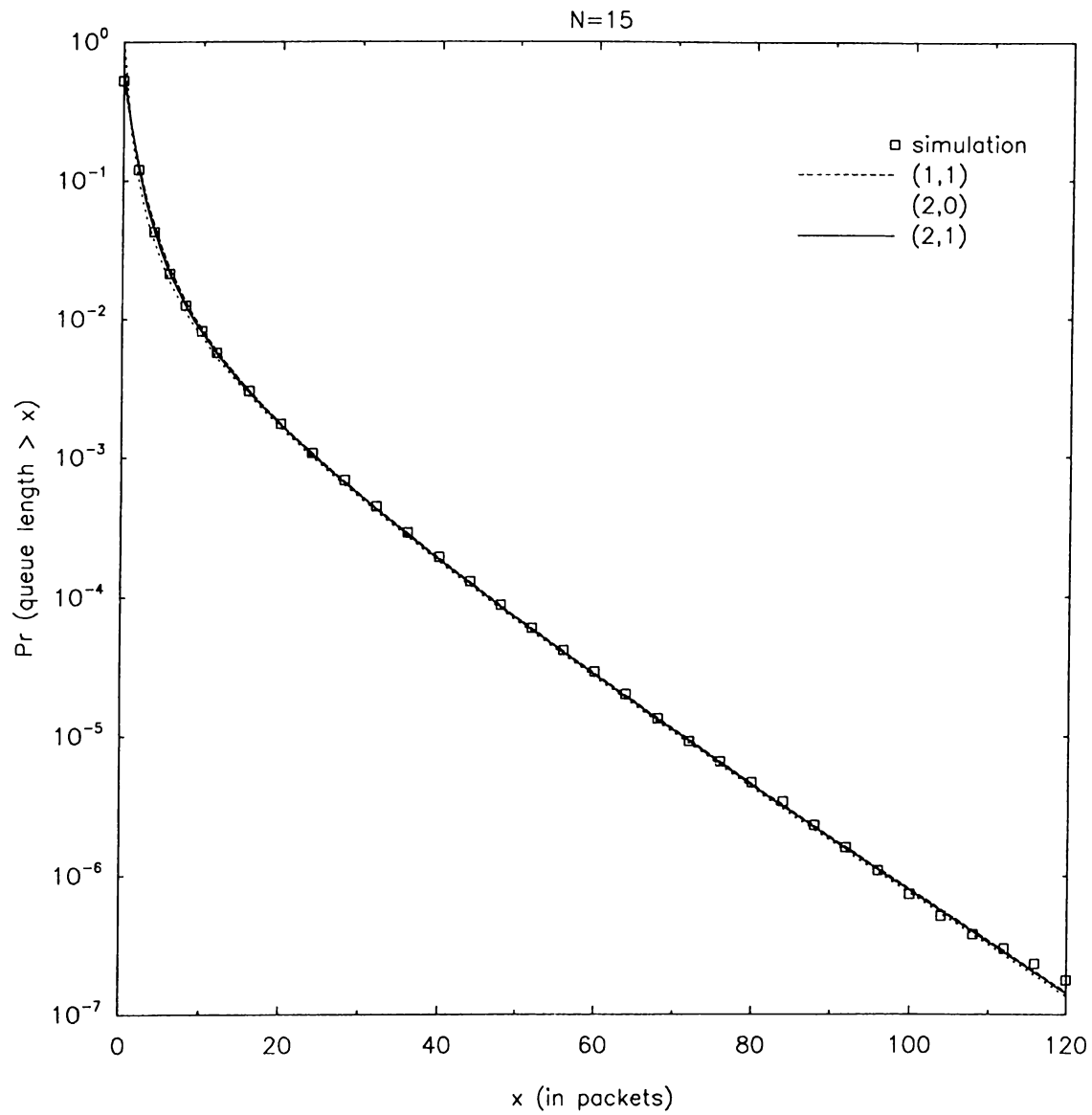
Up to now, we have concentrated on the use of Padé approximations for the deterministic service time distribution and based on this, we presented a novel algorithm to compute the queue length distribution for the MMPP/D/1 queue. We have shown that it is possible to obtain accurate approximations for the queue length distribution via simple



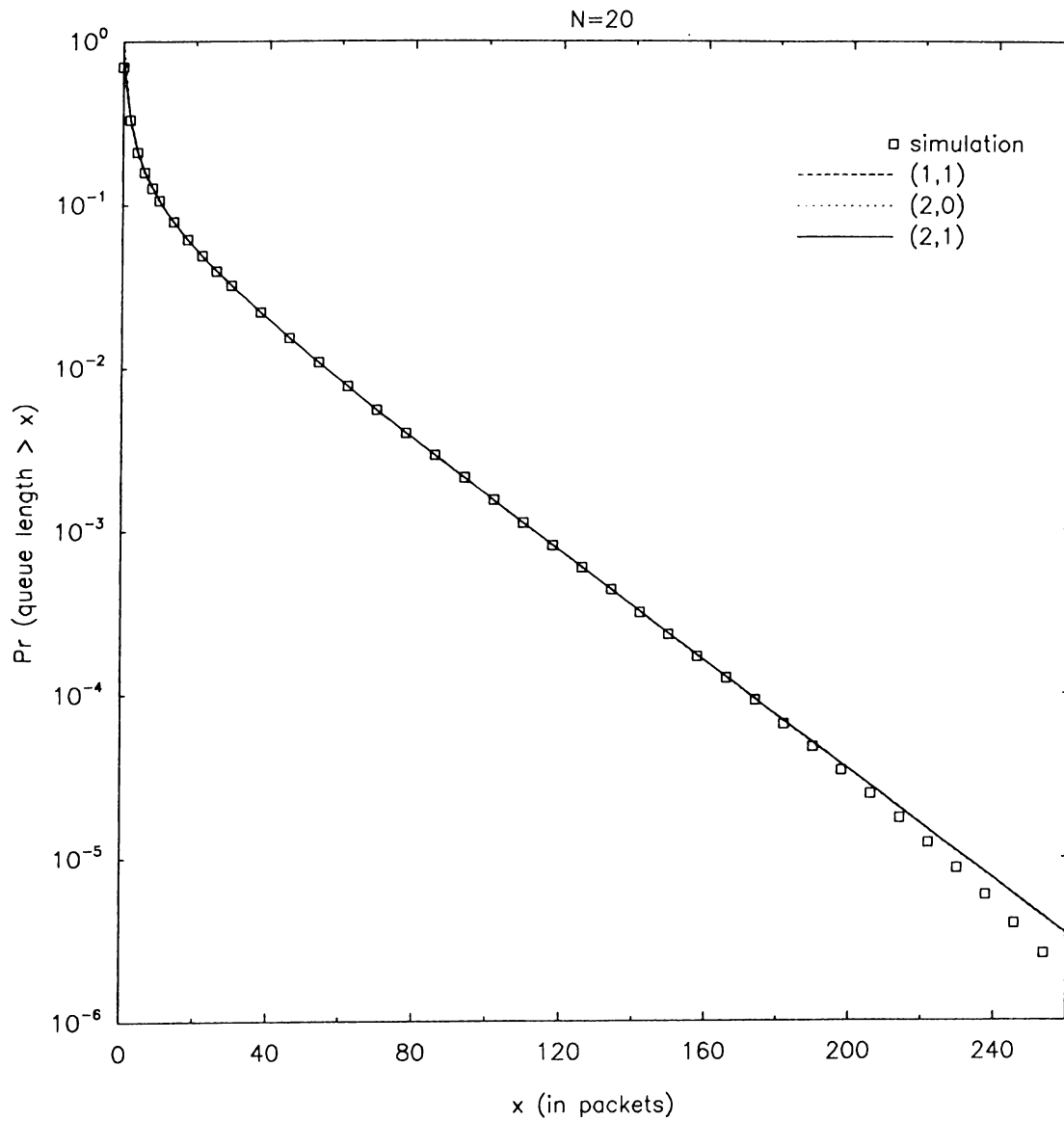
**Figure 4.1:** Performance comparison of the Padé approximations in terms of the queue length survivor function ( $N = 8$ , utilization = 0.28).



**Figure 4.2:** Performance comparison of the Padé approximations in terms of the queue length survivor function ( $N = 10$ , utilization = 0.35).



**Figure 4.3:** Performance comparison of the Padé approximations in terms of the queue length survivor function ( $N = 15$ , utilization = 0.52).



**Figure 4.4:** Performance comparison of the Padé approximations in terms of the queue length survivor function ( $N = 20$ , utilization = 0.70).

Padé approximants (e.g.,  $R_{2,1}(s)$  approximation yielding a determinantal degree  $d = 2N$ ). The computational load of the algorithm consists of an eigenvalue-eigenvector problem of size  $d$  (first step of the algorithm) and a matrix inversion problem of size  $d - N$  (second step of the algorithm). Furthermore, the computational complexity of the first step of the algorithm can be significantly reduced if the incoming MMPP is obtained from a superposition of many MMPP's of smaller dimension. This is in fact the case for ATM networks; individual sources are typically modeled to have an on/off type behavior and can be characterized as in [25] by two-state MMPP's. The aggregate traffic to be statistically multiplexed then turns out to be a superposition of many two-state MMPP's. We examine below in a much general setting how this computational complexity reduction takes place if we consider the case where there are  $K$  independent sources, each represented by an  $N$ -state continuous-time, irreducible generator matrix  $G$ . The rate of the Poisson stream of packets from an individual source in state  $i$  is  $\lambda_i (i = 1, 2, \dots, N)$ . Our objective in this section is to extend for the MMPP/D/1 queue the algebraic theory developed for fluid models [1],[11] and MMPP/M/1 queues [13] that gives the exact decomposition of the eigenvalue problem of the overall system into many small eigenvalue problems.

Although the approach taken here can easily be applied for each Padé approximation  $R_{n,m}(s)$ , we will focus on the approximation  $\exp(-s) \approx R_{2,1}(s)$  since it has been shown to give the best results among the ones that yield the same determinantal degree of  $\psi(s)$ .

Let the state of the source  $i$  be denoted by  $s(i)$  where  $s(i) \in \{1, 2, \dots, N\}$ . The unaggregated state of the sources is given by  $s = \{s(1), s(2), \dots, s(N)\}$ . We let the state space of the unaggregated process by

$$\mathcal{H}_{N,K} \triangleq \{k \mid k \in \mathcal{Z}^K, 1 \leq k(i) \leq N\}.$$

The generator of the unaggregated source process is

$$M = G \oplus G \oplus \dots \oplus G,$$

a  $K$ -fold Kronecker sum on  $G$ , where  $A \otimes B = [a_{ij}B]$ ,  $A \oplus B = A \otimes I + I \otimes B$  [39]. The

associated rate matrix  $R$  is

$$R = \Lambda \oplus \Lambda \oplus \cdots \oplus \Lambda,$$

where  $\Lambda = \text{diag}\{\lambda_1, \lambda_2, \dots, \lambda_N\}$ . The formula for  $\hat{F}_e(s)$  is

$$\hat{F}_e(s) = [(sCI + M - R)(1 + s/3) + R(1 - 2s/3 + s^2/6)]^{-1}(1 + s/3)f_0 \quad (4.22)$$

by which the associated eigenvalue problem is written as

$$[(zCI + M - R)(1 + z/3) + R(1 - 2z/3 + z^2/6)]\phi = 0. \quad (4.23)$$

Due to the Kronecker sum form of  $M$  and  $R$ , it is not difficult to show the equivalence between (4.23) and the existence of a set of  $K$  numbers  $v_1, v_2, \dots, v_K$  satisfying

$$[(zv_i I + G - \Lambda)(1 + z/3) + \Lambda(1 - 2z/3 + z^2/6)]u_i = 0, \quad \sum_{i=1}^K v_i = C. \quad (4.24)$$

If (4.24) holds, then the eigenvector is in the Kronecker product form;

$$\phi = u_1 \otimes u_2 \otimes \cdots \otimes u_K. \quad (4.25)$$

Letting  $g_i(z)$  to be the  $i^{\text{th}}$  eigenvalue ( $i = 1, 2, \dots, N$ ) of the matrix

$$B(z) = \frac{\Lambda(1 - z/6)}{1 + z/3} - \frac{G}{z},$$

by use of (4.24), each  $k \in \mathcal{H}_{N,K}$  gives an equation

$$\sum_{i=1}^K g_{k(i)}(z) = C, \quad (4.26)$$

whose solutions  $z$  are the solutions to the coupled eigenvalue problem in (4.24).

We now consider the aggregated system representation  $\sigma = \{\sigma(1), \sigma(2), \dots, \sigma(N)\}$  where  $\sigma(i)$  is the number of sources in state  $i$ . We also denote the set of all aggregated sources by  $\Sigma$ . Note that for any  $\sigma \in \Sigma$  the aggregated counterpart of (4.26) is

$$\sum_{j=1}^N \sigma(j)g_j(z) = C, \quad (4.27)$$

whose solution  $z$  is also a solution to (4.26) for every  $k \in \mathcal{H}_{N,K}$  that is aggregated to  $\sigma$ . The solutions to (4.27) together with (4.24) and (4.25) gives an exact decomposition



of the eigenvalue problem of the entire aggregated system. At this stage, we avoid presenting technical details of the general algebraic theory corresponding to the eigenvalues and eigenvectors of the aggregated system, the reader may refer to [1],[13], and [11] for a detailed discussion of related issues. We rather prefer to give a procedure without technical proofs for the case  $N = 2$  so as to clarify the concept of extension of the general theory to the case of Padé approximations in the analysis of the MMPP/D/1 queue.

Although the ideas can easily be generalized to multi-state sources as in [13], we rather concentrate on the case  $N = 2$  in which the aggregated state representation is

$$\sigma = (i, j), 0 \leq i, j \leq K, i + j = K.$$

We also let the generator matrix of an individual source to be

$$G = \begin{bmatrix} -\alpha & \beta \\ \alpha & -\beta \end{bmatrix}.$$

In regard of the works done in [13] and [11], an algorithm is developed below to compute the queue length distribution whose Laplace transform is given in (4.22).

**Algorithm.**

1) When  $K$  is odd, for each  $i, 0 \leq i \leq \frac{K-1}{2}$ , define the  $4^{th}$  degree polynomial

$$P(z; i) = (i\tilde{g}_1 + (K - i)\tilde{g}_2 - Cz(1 + z/3))((K - i)\tilde{g}_1 + i\tilde{g}_2 - Cz(1 + z/3)),$$

where  $\tilde{g}_1(z)$  and  $\tilde{g}_2(z)$  are solutions of

$$|tI - \tilde{B}(z)| = 0, \tilde{B}(z) = -\Lambda z^2/6 + z(\Lambda - G/3) - G.$$

Without an explicit knowledge of  $\tilde{g}_i(z)$ , it is also possible to write  $P(z; i)$  as

$$\begin{aligned} P(z; i) &= C^2 z^2 (1 + z/3)^2 - KCz(1 + z/3)tr(\tilde{B}(z)) \\ &+ i(K - i)tr^2(\tilde{B}(z)) + (K^2 - 4i(K - i))|\tilde{B}(z)|. \end{aligned}$$

Here  $tr(\cdot)$  and  $|\cdot|$  denote the trace and the determinant of a matrix, respectively. When the ergodicity condition

$$K(\alpha\lambda_2 + \beta\lambda_1) < C(\alpha + \beta)$$

is satisfied, the polynomial  $P(z; i)$  has two roots in the interval  $(z_2^i < z_1^i)$  in the interval  $(-\infty, 0)$  and two roots in the interval  $[0, \infty)$ . Let

$$z_i \triangleq z_1^i, z_{K-i} \triangleq z_2^i.$$

When  $K$  is even, repeat the procedure described above for  $0 \leq i \leq \frac{K}{2} - 1$ , and then define the  $2^{\text{nd}}$  degree polynomial

$$P(z; K/2) = \frac{K}{2} \text{tr}(\tilde{B}(z)) - zC(1 + z/3),$$

a root  $z_e$  of which lies in  $(-\infty, 0)$ . Now let  $z_{K/2} \triangleq z_e$ . In this stage of the procedure, the eigenvalues of the aggregated system are completely found in terms of the roots of  $4^{\text{th}}$  degree polynomials.

2) For each  $i$ ,  $0 \leq i \leq K$ , solve the eigenvalue problem

$$(-\Lambda z_i^2/6 + z_i(\Lambda - G/3) - G)u = \mu u, \quad (4.28)$$

and define  $u_1$  and  $u_2$  to be the eigenvectors in (4.28) associated with the larger and the smaller eigenvalue, respectively. The entries of the eigenvector  $\phi_i$  associated with the eigenvalue  $z$  is now written as a summation;

$$\phi_i(j) = \sum_{l=\max(0, K-j-i)}^{K-j} \binom{K-i}{l} \binom{i}{K-j-l} u_1(1)^l u_1(2)^{K-i-l} u_2(1)^{K-j-l} u_2(2)^{i-(K-j-l)}.$$

Finally, the queue length cdf  $F_e(j, x)$  is

$$F_e(j, x) = \pi_j + \sum_{i=0}^K a_i \exp(z_i x) \phi_i(j), \quad (4.29)$$

where  $\pi_j$  is the probability of  $j$  users being in the  $2^{\text{nd}}$  state.

3) We use the initial value theorem for the formula (4.22) to show that

$$\frac{d}{dx} F_e(x) \Big|_{x=0^+} = \Delta F_e(0) \quad (4.30)$$

where

$$\Delta = 3I - \left(\frac{CI}{3} + \frac{\bar{R}}{6}\right)^{-1} (CI - \bar{R} + \frac{\bar{M}}{3}).$$

Here,  $\bar{M}$  is the generator matrix for the aggregated process where the state of the aggregated process is the number of sources in the  $2^{nd}$  state and  $\bar{R}$  is the corresponding rate matrix. Defining

$$\begin{aligned}\Phi &= \begin{bmatrix} \phi_0 & \phi_1 & \cdots & \phi_K \end{bmatrix}, \\ Z &= \text{diag}\{z_0, z_1, \dots, z_K\}, \\ a &= \begin{bmatrix} a_0 & a_1 & \cdots & a_K \end{bmatrix}^T, \\ \pi &= \begin{bmatrix} \pi_0 & \pi_1 & \cdots & \pi_K \end{bmatrix}^T,\end{aligned}$$

the equation (4.30) reduces to

$$\Phi Z a = \Delta(\pi + \Phi a),$$

from which one obtains

$$a = (\Phi Z - \Delta \Phi)^{-1} \Delta \pi.$$

The eigenvalues and the eigenvectors  $z_i$ 's and  $\phi_i$ 's of the aggregated system together with the coefficients  $a_i$ 's determine exactly the queue length cdf through the expression (4.29).  $\square$

In regard of the algorithm presented here, let us quantify the reduction obtained in computational complexity. Note that, the eigenvalues of the overall system are obtained through a set of polynomial root finding problems of the form given in the first step of the algorithm. These polynomials are all of  $4^{th}$  degree and there are  $(K + 1)/2$  such polynomials when  $K$  is odd ( $K/2$  polynomials of  $4^{th}$  degree and a polynomial of degree two when  $K$  is even). Since the number of states (number of individual sources in the  $2^{nd}$  state constitute the state of the superposition process) of the superposition process is  $K + 1$ , in the original version of the algorithm, the computation consists of an eigenvalue problem of size  $2(K + 1)$ . In other words, the  $O(K^3)$  computational complexity of the eigenvalue problem in the original version of the algorithm can be made to have a  $O(K)$  complexity in case the MMPP is obtained from a superposition of many identical two-state MMPP's. Although quantification is rather involved, we can still achieve a large computational gain in the general case of a superposition of

a number of heterogeneous multi-state MMPP sources. This can be concluded from equation (4.26) where an eigenvalue of the overall system can be computed from the parameterized eigenvalues of the subsystems through a nonlinear algebraic equation.

In ATM networks, buffering memory in switching nodes is limited. Cells have to be buffered at ATM switches and may have to be dropped (lost) in the case of buffer overflow. The cell loss rate is a key parameter in performance analysis since resource allocation in ATM networks is based on the desired quality of service which is generally expressed in terms of bounds on the loss rate. In the next section, we extend the algorithm proposed for infinite buffers to the case of buffer sizes limited to a finite number  $K$  of cells. This extension yields a buffer size independent computational complexity and performance analysis turns out to be tractable even for large buffer sizes.

## 4.5 MMPP/D/1/K Queue

In this section, we are interested in computing the queue length distribution for the MMPP/D/1/K queue for which the buffer is capable of storing only  $K$  cells. We subdivide the analysis into two parts, first we consider the case of infinite buffers with a state-space realization of the differential equation governing the queue length distribution. At this step, the realization is general so that it covers the use of an arbitrary Padé approximation for the deterministic service time distribution. We then extend the analysis to cover the case of finite queue capacities.

### Infinite Queue Capacity

Let an arbitrary Padé approximation

$$R_{n,m}(s) = \frac{P_n(s)}{Q_m(s)}$$

be imposed as an approximation of  $\exp(-s)$ . The polynomials  $P_n(s)$  and  $Q_m(s)$  are assumed to have degrees  $n$  and  $m$ , respectively. Then the queue length distribution in

(4.15) can be rewritten as

$$\begin{aligned}
\hat{F}_e(s) &= [sCI + M - R + R \frac{P_n(s)}{Q_m(s)}]^{-1} f_0 \\
&= Q_m(s)[(sCI + M - R)Q_m(s) + RP_n(s)]^{-1} f_0 \\
&=: Q_m(s)\hat{H}(s)^{-1} f_0.
\end{aligned} \tag{4.31}$$

The polynomial matrix  $\hat{H}(s) = [(sCI + M - R)Q_m(s) + RP_n(s)]$  has degree

$$k = \max(m + 1, n), \tag{4.32}$$

that is,  $\hat{H}$  can be written as

$$\hat{H}(s) = H_k s^k + H_{k-1} s^{k-1} + \cdots + H_1 s + H_0, \tag{4.33}$$

for some constant matrices  $H_i$ ,  $i = 0, 1, \dots, k$ , with  $H_k$  assumed to be nonsingular. Similarly, the polynomial  $Q_m(s)$  is of the form

$$Q_m(s) = q_{k-1} s^{k-1} + q_{k-2} s^{k-2} + \cdots + q_1 s + q_0, \tag{4.34}$$

since  $\deg(Q_m(s)) = m < k$ . We are now ready to obtain a state-space realization for the transform expression in (4.31). We first define

$$F^i(x) = \frac{d^{i-1}}{dx} F_e(x), \quad i = 1, 2, \dots, k,$$

and

$$F_c(x) = \begin{bmatrix} F^1(x) \\ F^2(x) \\ \vdots \\ F^k(x) \end{bmatrix},$$

where the subscript  $c$  refers to a concatenation of the column vectors  $F^i$ . Using the equation (4.31), one can show that the following ordinary differential equation is valid:

$$\frac{d}{dx} F_c(x) = A_1 F_c(x), \quad x > 0, \tag{4.35}$$

where

$$A_1 = \begin{bmatrix} 0 & I & 0 & 0 \\ 0 & 0 & I & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & I \\ -\bar{H}_0 & -\bar{H}_1 & -\bar{H}_2 & \cdots & -\bar{H}_{k-1} \end{bmatrix},$$

and

$$\bar{H}_i = H_k^{-1} H_i, \quad i = 0, 1, \dots, k-1.$$

On the other hand, the initial condition can be shown to be related to the vector  $f_0$  through the following equalities:

$$\begin{aligned} F^1(0) &= \bar{Q}_{k-1} y =: Z_1 f_0 \\ F^i(0) &= \bar{Q}_{k-i} f_0 - \sum_{j=1}^{i-1} \bar{H}_{k-i+j} F^j(0) =: Z_i f_0, \quad i = 2, 3, \dots, k, \end{aligned}$$

where

$$\bar{Q}_j = q_j H_k^{-1}, \quad j = 0, 1, \dots, k-1.$$

The solution to the linear differential equation (4.35) then takes the form

$$F_c(x) = \exp(A_1 x) \begin{bmatrix} Z_1 \\ Z_2 \\ \vdots \\ Z_k \end{bmatrix} f_0. \quad (4.36)$$

The unknown vector  $f_0$  in the expression (4.36) can be computed using the conventional techniques in [20] based on matrix analytic methods. Another alternative is to use spectral expansion techniques in Section 4.3 and compute  $f_0$  by imposing that no unstable mode of the dynamical system (4.35) is excited. Our concern here is introducing a new mathematical framework for the MMPP/D/1 system that can be extended to the MMPP/D/1/K queues, i.e., to the case of finite buffers, which is considered below.

### Finite Queue Capacity

Let the queue capacity be denoted by  $K$ . When a new arrival finds fewer than  $K$  packets in the queue waiting to be served, it is admitted to the system. Following the scheme in [10],[54], the following differential equation is valid in the interval  $0 < x < K$ :

$$\frac{d}{dx}F_c(x) = A_1F_c(x), \quad 0 < x < K, \quad (4.37)$$

$$F_c(0) = \begin{bmatrix} Z_1 \\ Z_2 \\ \vdots \\ Z_k \end{bmatrix} f_0.$$

Note that the vector  $f_0$  in (4.37) is different from that of equation (4.36).

On the other hand, if an arrival occurs at time  $t$  and the instantaneous queue length at that time is above  $K$ , the packet associated with that arrival is dropped. From the queue length point of view, it is convenient to visualize the incoming MMPP characterized by the matrix pair  $(M, R)$  to change to another MMPP described by the matrix pair  $(M, 0)$  whenever the number of packets in the queue is  $K$ . This is equivalent to assuming that no arrivals will occur and the MMPP will be constituted of only its modulating process. Also note that the queue length cannot exceed  $K + 1$  since there is only one server. Then one can obtain as in (4.35) the following differential equation in the interval  $K \leq x < K + 1$ :

$$\frac{d}{dx}F_c(x) = A_2F_c(x), \quad K \leq x < K + 1. \quad (4.38)$$

In this equation, the matrix  $A_2$  is of the form

$$A_2 = \begin{bmatrix} 0 & I & 0 & \cdots & 0 \\ 0 & 0 & I & \cdots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & 0 & 0 & \cdots & I \\ -\bar{G}_0 & -\bar{G}_1 & -\bar{G}_2 & \cdots & -\bar{G}_{k-1} \end{bmatrix},$$

where

$$\hat{G}(s) = (sCI + M)Q_m(s) = G_k s^k + G_{k-1} s^{k-1} + \cdots + G_1 s + G_0,$$

and  $\tilde{G}_i = G_k^{-1}G_i$ ,  $i = 0, 1, \dots, k - 1$ .

We are now prepared to compute the queue length distribution in a MMPP/D/1/K system except for the boundary conditions. The boundary condition at  $x = K + 1$  is easy to write since i) queue length cannot exceed  $K + 1$ , ii) stationary probability of the queue length being  $K + 1$  is zero, i.e., there may not be a jump in the queue length cdf vector at  $x = K + 1$ . Based on these two observations, one can write

$$F^1(K + 1) = F_e(K + 1) = \pi. \quad (4.39)$$

Making use of the continuity of the solution of the two differential equations (4.37) and (4.38) at  $x = K$ , one can rewrite (4.39) as

$$\begin{bmatrix} I & 0 & \dots & 0 \end{bmatrix} \exp(A_2) \exp(A_1 K) \begin{bmatrix} Z_1 \\ Z_2 \\ \vdots \\ Z_k \end{bmatrix} f_0 = \pi. \quad (4.40)$$

The unknown vector  $f_0$  can be solved through the linear matrix equation (4.40) of size  $Nk$ . At this stage, any algorithm for computing matrix exponentials [38] can be used to compute the left hand side of (4.40). Once the column vector  $f_0$  is computed, the solution to the differential equations for  $F_c$  is easy to write:

$$\begin{aligned} F_c(x) &= \exp(A_1 x) \begin{bmatrix} Z_1 \\ Z_2 \\ \vdots \\ Z_k \end{bmatrix} f_0, \quad 0 \leq x \leq K, \\ F_c(x) &= \exp(A_2(x - K)) F_c(K), \quad K \leq x \leq K + 1. \end{aligned}$$

The stationary queue length cdf  $F_e(x)$  is then expressed as

$$F_e(x) = F^1(x) = \begin{bmatrix} I & 0 & \dots & 0 \end{bmatrix} F_c(x). \quad (4.41)$$

Cell losses occur when arrivals find  $K$  cells waiting in the buffer. The cell loss rate,  $p_{loss}$ , is therefore described by the following expression

$$p_{loss} = \frac{eR(\pi - F_e(K))}{\bar{\lambda}}. \quad (4.42)$$



where  $e$  is a row vector of ones and  $\bar{\lambda}$  is the mean arrival rate.  $\square$

Approximate computation of the queue length distribution in the MMPP/D/1/K queue is shown to be given in terms of the solution of two linear differential equations (4.37) and (4.38). The core of the computation lies in solving the linear equation (4.40) which consists of efficiently computing two matrix exponentials of size  $Nk$  where  $N$  is the number of states of the MMPP and  $k$  is given in (4.32). We now show through numerical examples that a Padé approximation yielding a degree  $k = 3$  will suffice for all practical purposes in order to compute the queue length distribution. Note that the computational complexity of the algorithm is independent of the queue capacity, thus avoiding numerical problems, especially for large buffer sizes.

#### 4.5.1 Numerical Examples

In this section, we investigate the performance of the proposed approximation scheme to compute the cell loss rate in a MMPP/D/1/K queue. We consider an ATM multiplexer that serves LAN (Local Area Network)-generated data streams. The cell emission process for an individual LAN source is widely recognized to be adequately represented by means of an on-off source model [3]. Let us consider a set of  $N$  independent and homogeneous LAN sources characterized by i) the peak rate  $F_p$  ii) the activity factor  $p$ , defined as the ratio between the average bit rate and  $F_p$  iii) mean burst length  $L_b$ . We choose a reference LAN source as in [3] which is characterized by  $F_p = 10$  Mbits/s,  $p = 0.1$ , and  $L_b = 16250$  bytes, where these values are representative of a large class of information flows arising from LAN's accessing to an ATM network. As for the multiplexer, we assume a gross output capacity equal to 150 Mbits/s (ATM transport rate) and a cell length of 53 bytes, 48 of which constitutes the cell payload and the net output capacity is therefore equal to 135.85 Mb/s.

We approximate the superposition of  $N$  such on-off sources by means of a two-state MMPP using the *asymptotic matching* technique proposed in [3]. This technique is shown in [3] to be much more effective than the matching method used in [20] in capturing

the cell loss rate. Our concern is the performance assessment of our proposed algorithm for the MMPP/D/1/K queue rather than examining the performance of the asymptotic matching technique which has already been shown in [3] to be accurate in computation of the loss rate. Therefore, in the performance assessment procedure, our simulations are based on the two-state MMPP model obtained by means of the asymptotic matching technique.

Figure 4.5 is devoted to the computation of the cell loss rate for the case  $N = 40$  via certain Padé approximations. The notation  $(n, m)$  denotes the use of a Padé approximation  $R_{n,m}(s)$  for the deterministic service time. Recall that the key parameter that determines the computational load is the variable  $k$  which equals to  $\max(m + 1, n)$  when a Padé approximation  $R_{n,m}(s)$  is imposed. In Figure 4.5, we show that even with a simple Padé approximant (e.g., (2,2) approximation yielding  $k = 3$ ), high accuracies in cell loss rate computation can be maintained. The Padé approximation (0,1) overestimates the loss rate whereas the approximations (1,1) and (2,1) underestimate this quantity. We do not cover in Figure 4.5 more advanced approximations (e.g., (3,3) approximation) since their performances are quite the same as that of the approximation (2,2) (we observed at most a 2 % difference in the cell loss rate approximations of (2,2) and (3,3) for this example).

The final example is given in Figure 4.6 in which the cell loss rate approximation obtained by using the Padé approximation (2,2) is compared with the simulation results as the number of the LAN sources are varied. Irrespective of the buffer size and the number of users (or equivalently, the load), the Padé approximation (2,2) results in a very accurate approximation of the cell loss rate.

## 4.6 Effective Bandwidth

In this section, we are concerned with the problem of effective bandwidth approximation of a call of MMPP type offered to a deterministic server. The same problem in the

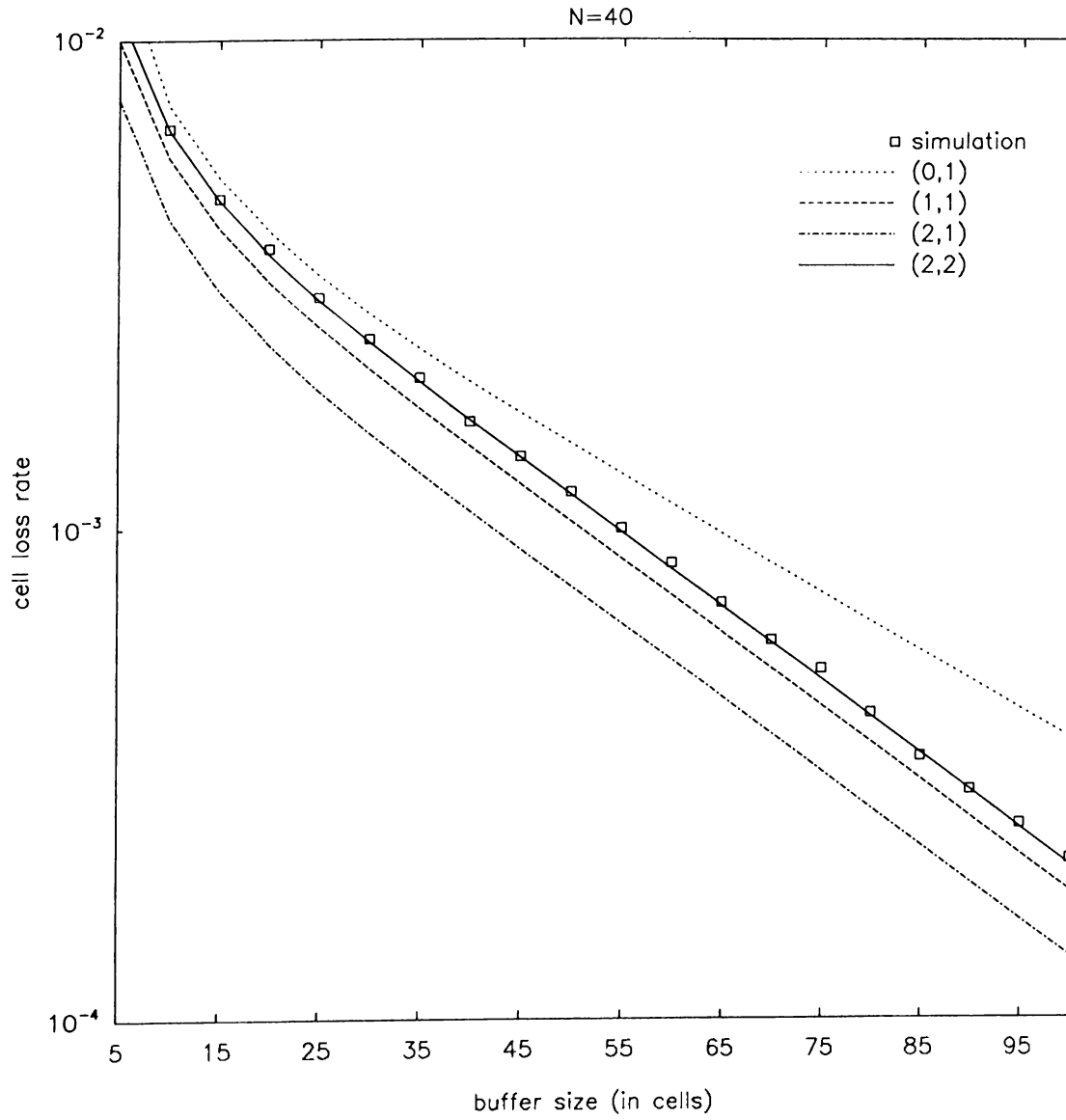
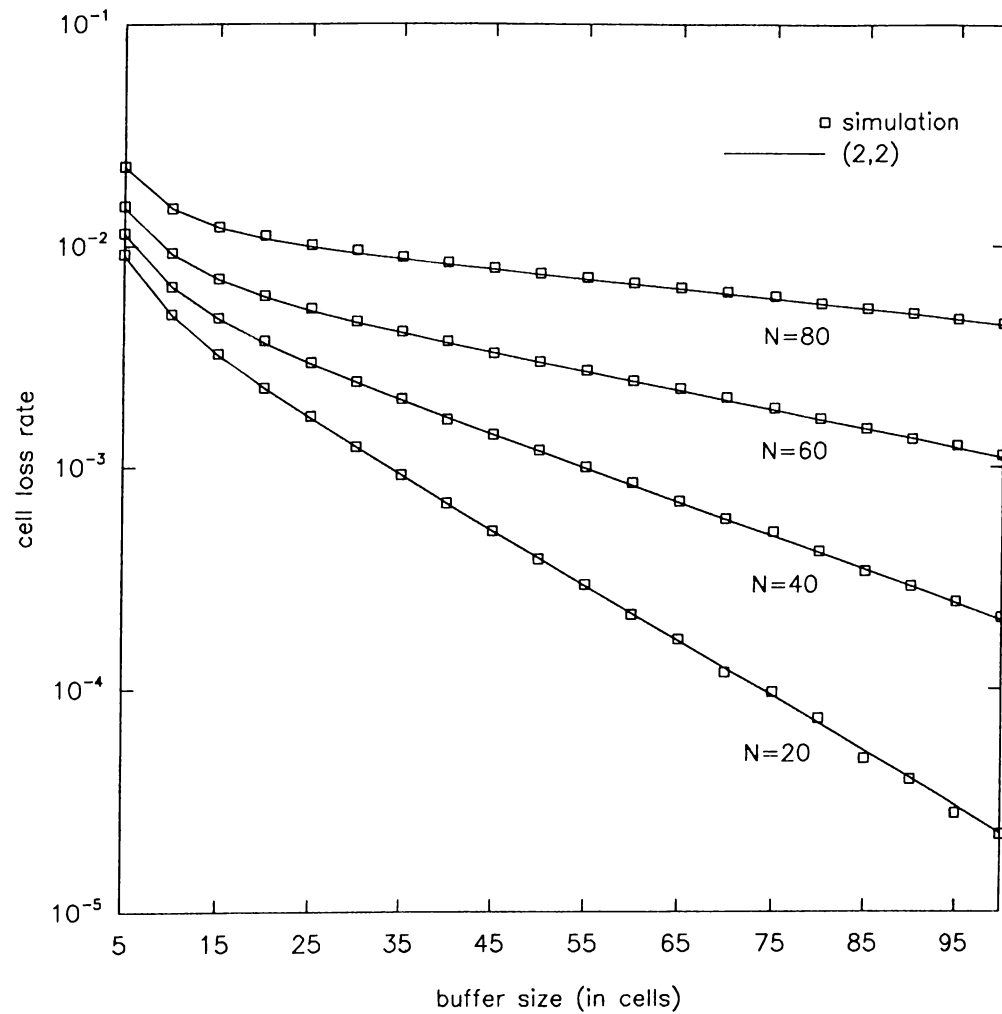


Figure 4.5: Cell loss rate approximations ( $N = 40$ , utilization = 0.29).



**Figure 4.6:** Cell loss rate with respect to the buffer size obtained by Padé approximation (2,2) as  $N$  is varied.

case of Markov modulated periodic arrivals has been considered in Section 3.5. In the exposition that follows, we will extend the results there to the case of Markov modulated Poisson sources and state an analog of Proposition 3.1 as a solution for the call admission problem in the prescribed asymptotic regime of large buffers and small buffer overflow probabilities. In what follows, the Padé approximation  $R_{2,1}(s)$  will be used to approximate the transform of the service time distribution.

We use the same notation in Section 4.2 so that  $M$  is the generator matrix of the underlying Markov chain and  $R$  is the rate matrix associated with that source whose effective bandwidth is our objective. When this process is fed into a statistical multiplexer, the stationary queue length has the following spectral representation

$$F_e(x) = \sum_{i: \text{Re}(z_i) < 0} a_i \phi_i \exp(z_i x) + \pi,$$

where  $\pi = \{\pi_n\}$ ,  $n = 1, 2, \dots, N$  and the pair  $(z_i, \phi_i)$  is an eigenvalue-eigenvector pair. Such pairs are solutions to the generalized eigenvalue problem

$$[(zCI + M - R)(1 + z/3) + R(1 - 2z/3 + z^2/6)]\phi = 0. \quad (4.43)$$

Indexing the eigenvalues with negative real parts

$$0 > z_1 \geq \text{Re}(z_2) \geq \dots$$

we call the real eigenvalue  $z_1$  as the dominant eigenvalue and we also note that

$$\frac{\log G(x)}{x} \rightarrow z_1 \text{ as } x \rightarrow \infty,$$

where  $G(x)$  denotes the buffer survivor function. Writing  $C = g(z)$  in (4.43), we obtain an equivalent expression for (4.43):

$$g(z)\phi = A(z)\phi \triangleq \left[ \frac{R(1 - z/6)}{(1 + z/3)} - \frac{M}{z} \right] \phi. \quad (4.44)$$

Note the additional term  $(1 - z/6)/(1 + z/3)$  contrary to the expression (3.32) suggested for both Markov modulated fluid and periodic sources. Here,  $g(z)$  is an eigenvalue of the matrix  $A(z)$  in which  $z$  is a parameter. Paralleling the development in Section 3.5,

$g_1(z)$  is called the maximal real eigenvalue among the ones that satisfy (4.44) for some  $\phi$ . This particular eigenvalue has the following remarkable properties (proof is omitted due to its length but can be made along the lines followed in [12], principal observation being the essential non-negativity of  $A(z)$  in the interval  $(-3, 0]$ ).

- i)  $g_1(z)$  is a decreasing function of  $z$  in  $(-3, 0]$ . Besides,  $\lim_{z \rightarrow 0^-} g_1(z) = \bar{\lambda}$  (mean source rate) and  $\lim_{z \rightarrow -3^+} g_1(z) = +\infty$ .
- ii) The dominant eigenvalue  $z_1$  is the unique solution in  $(-3, 0]$  satisfying  $g_1(z_1) = C$ .

We next examine the admission criterion  $\{G(B) \leq p\}$  in the asymptotic regime of large buffers  $B$  and small overflow probabilities  $p$ . Our result is

**Proposition 4.1.** *Suppose  $P \rightarrow \infty$  and  $p \rightarrow 0$  in such a manner that  $\log p/B \rightarrow \xi \in (-3, 0]$ . If*

$$g_1(\xi) < C,$$

*then the admission criterion is satisfied. If the inequality sign is reversed then the admission criterion is violated. Here,  $g_1(\xi)$  is the maximal real eigenvalue of*

$$A(\xi) \triangleq \frac{R(1 - \xi/6)}{(1 + \xi/3)} - \frac{M}{\xi}. \quad \square$$

On the basis of the above result, the effective bandwidth of a single source is simply

$$e(M, R; B, p) = g_1(\xi). \quad (4.45)$$

For a two-state MMPP source with the parameter pair  $(M, R)$

$$M = \begin{bmatrix} -\lambda & \mu \\ \lambda & -\mu \end{bmatrix} \quad \text{and} \quad R = \begin{bmatrix} 0 & 0 \\ 0 & P \end{bmatrix},$$

$g_1(z)$  is the following simply computed quantity

$$g_1(z) = \frac{P(z)z + \lambda + \mu - ((P(z)z + \lambda + \mu)^2 - 4\lambda P(z)z)^{1/2}}{2z}$$

in which

$$P(z) \triangleq \frac{P(1 - z/6)}{(1 + z/3)}.$$

We next investigate, as in Section 3.3, the decomposition of the expression in (4.45), when the source  $(M, R)$  is the aggregate of  $K$  sources,  $(M^{(k)}, R^{(k)})$ ,  $(1 \leq k \leq K)$ .

**Proposition 4.2.** *Suppose there are  $K$  Markov modulated Poisson sources,  $(M^{(k)}, R^{(k)})$ ,  $(1 \leq k \leq K)$ , offered to a multiplexing buffer. Let the admission criterion and the asymptotic regime be as in Proposition 4.1. If*

$$\sum_k g_1^{(k)}(\xi) < C,$$

*then the admission criterion is satisfied. If the inequality sign is reversed then the admission criterion is violated. Here,  $g_1^{(k)}(\xi)$  is the maximal real eigenvalue of*

$$A^{(k)}(\xi) \triangleq \frac{R^{(k)}(1 - \xi/6)}{(1 + \xi/3)} - \frac{M}{\xi}. \quad \square$$

# Chapter 5

## Conclusions and Suggestions for Future Work

In this dissertation, we have considered the queueing analysis of an ATM multiplexer fed by sources modeled as Markov modulated processes. We have focused on the MMPAP/D/1 and MMPP/D/1 queueing systems due to their wide-spread use in the performance analysis of ATM networks.

For the MMPAP/D/1 system, an approximation based on the transient behavior of the  $nD/D/1$  queue is proposed which is capable of capturing cell scale fluctuations. The computation encountered is similar to fluid flow approximations except for the determination of the linear operator  $Z$  (see def. (3.16)). Our suggestions for further research are:

- Computation of the operator  $Z$  (see def. (3.16)) is the main part of the overall algorithm. Although certain methods are presented to decrease this computation, there is still work to be done towards obtaining further adequate approximations to ensure tractability particularly when the number of sources to be multiplexed are increased.
- The case of different periods needs to be investigated. In our framework,



this requires the stationary queue length distribution expression in a  $\sum D_i/D/1$  queue for which an exact solution is not available. For this purpose, accurate approximations proposed for the  $\sum D_i/D/1$  system (see [31],[44]) can be made use of.

- Work needs to be done for analysis of queues with priority management. We note that our methodology is suitable to use in partial buffer sharing mechanisms [33] to which fluid flow techniques are proven to apply [11],[63].
- A more general notion of effective bandwidth is required which takes into account of not only the buffer overflow probability in FIFO queues but also in queues with priority management.

For the MMPP/D/1 system, we employ Padé approximations for the deterministic service time distribution in transform domain. We also show that fluid flow approximations and MMPP/M/1 queues are obtained via first order Padé approximations for the MMPP/D/1 system. We obtain particular Padé approximations which involve no more computational complexity than the one encountered in solving the MMPP/M/1 queue. The performance of these approximations is demonstrated in the case of a packetized voice multiplexer. An algorithm is presented for the finite buffered MMPP-driven queue (i.e., MMPP/D/1/K queue) with a computational complexity independent of the buffer size. Our suggestions for further research are:

- Justification of MMPP as a general traffic model in ATM networks is necessary.
- Examining the use of Padé approximations in approximating the performance of general queueing systems rather than the particular MMPP/D/1 system is one direction of future research.
- There is work to be done in order to propose simple schemes to approximate the asymptotic behavior of the cell loss rate with respect to the buffer size which can in turn be employed in call admission control.

We have proposed methods for the teletraffic analysis of an ATM multiplexer offered with a class of Markov modulated processes. We believe that these methods can be made use of to develop traffic control strategies in an ATM network. These strategies possibly include buffer dimensioning, bandwidth allocation and routing. Improved performance analysis schemes will alleviate the dependence on simulations and help propose new congestion control strategies that may turn B-ISDN into a successful reality.

# Bibliography

- [1] D. Anick, D. Mitra, and M. M. Sondhi. Stochastic theory of a data handling system with multiple sources. *Bell Syst. Tech. Jour.*, 61:1871–1894, 1982.
- [2] J. J. Bae and T. Suda. Survey of traffic control schemes and protocols in ATM networks. *IEEE Proceedings*, 79(2):170–189, 1991.
- [3] A. Baiocchi, N. B. Melazzi, M. Listanti, A. Roveri, and R. Winkler. Loss performance analysis of an ATM multiplexer loaded with high-speed ON-OFF sources. *IEEE JSAC*, 9(3):388–393, 1991.
- [4] A. Bhargava, P. Humblet, and M. G. Hluchyj. Queueing analysis of continuous bit-stream transport in packet networks. In *GLOBECOM*, 1989.
- [5] P. T. Brady. A statistical analysis of on-off patterns in 16 conversations. *Bell Syst. Tech. J.*, 47(1):73–91, 1968.
- [6] J. N. Daigle and J. D. Langford. Models for analysis of packet voice communication systems. *IEEE JSAC*, 4(6):1293–1297, 1986.
- [7] M. de Prycker. *Asynchronous Transfer Mode : Solution for Broadband ISDN*. Ellis Horwood Limited, London, 1991.
- [8] A. E. Eckberg. The single server queue with periodic arrival process and deterministic service time. *IEEE Trans. Commun.*, 27:556–562, 1979.

- [9] A. E. Eckberg, D. T. Luan Jr., and D. M. Lucantoni. Meeting the challenge: congestion and flow control strategies for broadband information transport. In *GLOBECOM*, pages 49.3.1–5, 1990.
- [10] A. I. Elwalid and D. Mitra. Statistical multiplexing with loss priorities in rate-based congestion control of high speed networks. to appear in *IEEE Trans. Commun.*
- [11] A. I. Elwalid and D. Mitra. Fluid models for the analysis and design of statistical multiplexing with loss priorities on multiple classes of bursty traffic. In *Proc. IEEE INFOCOM*, 1992.
- [12] A. I. Elwalid and D. Mitra. Effective bandwidth of general Markovian traffic sources and admission control of high speed networks. *IEEE/ACM Trans. on Networking*, 1(3):329–343, 1993.
- [13] A. I. Elwalid, D. Mitra, and T. E. Stern. A theory of statistical multiplexing of Markovian sources: spectral expansions and algorithms. In *Proc. 1st Int. Workshop Numer. Solut. Markov Chains*, 1990.
- [14] B. A. Francis. *A Course in  $H_\infty$  Control Theory*. Springer-Verlag, Berlin, 1987.
- [15] D. P. Gaver and J. P. Lehoczky. Channels that cooperatively service a data stream and voice messages. *IEEE Trans. Commun.*, 30(5):1153–1162, 1982.
- [16] R. J. Gibbens and P. J. Hunt. Effective bandwidth for the multi-type UAS channel. *Queueing Systems*, 9:17–28, 1991.
- [17] H. Gilbert, O. Aboul-Magd, and V. Phung. Developing a cohesive traffic management strategy for ATM networks. *IEEE Commun Mag.*, 29(10):36–45, 1991.
- [18] R. Guerin, H. Ahmadi, and M. Naghshineh. Equivalent capacity and its application to bandwidth allocation in high-speed networks. *IEEE JSAC*, 9:968–981, 1991.
- [19] I. W. Habib and T. N. Saadawi. Multimedia traffic characteristics in broadband networks. *IEEE Commun. Magazine*, 30(7):48–54, 1992.

- [20] H. Heffes and D. M. Lucantoni. A Markov modulated characterization of packetized voice and data traffic and related statistical multiplexer performance. *IEEE JSAC*, 4(6):856–868, 1986.
- [21] M. G. Hluchyj and M. J. Karol. Queueing in high performance switching. *IEEE JSAC*, 6(9):1587–1597, 1988.
- [22] D. Hong and T. Suda. Congestion control and prevention in ATM networks. *IEEE Network Mag.*, 5(4):10–16, 1991.
- [23] J. Y. Hui. Resource allocation for broadband networks. *IEEE JSAC*, 6(9):1598–1608, 1988.
- [24] J. Y. Hui. *Switching and Traffic Theory for Integrated Broadband Networks*. Kluwer, Boston, MA, 1990.
- [25] I. Ide. Superposition of interrupted Poisson processes and its application to packetized voice multiplexers. In *Proc. ITC-12*, 1988.
- [26] ITU-TS. *Proposed recommendation I. 311*, June 1991.
- [27] K. Kawashima and H. Saito. Teletraffic issues in ATM networks. In *Proc. ITC*, pages 17.5.1–8, 1989.
- [28] F. P. Kelly. Effective bandwidths at multi-class queues. *Queueing Systems*, 9:5–16, 1991.
- [29] F. P. Kelly. Loss networks. *The Annals of Applied Probability*, 1(3):319–378, 1991.
- [30] L. Kleinrock. *Queueing Systems Volume 1: Theory*. Wiley-Interscience Publication, 1975.
- [31] K. Nakagawa. Loss and waiting time probability approximation for general queueing. *IEICE Trans. Commun.*, E76-B(11), 1993.

- [32] L. Kosten. Stochastic theory of a data handling system with groups of multiple sources. In *Performance of Computer Communication Systems*, pages 321–331. Elsevier Science Publishers, Amsterdam, 1984.
- [33] H. Kroner, G. Hébuterne, P. Boyer, and A. Gravey. Priority management in ATM switching nodes. *IEEE JSAC*, 9(3):418–427, 1991.
- [34] S. Q. Li. Performance of a non-blocking space division packet switch with correlated input traffic. In *GLOBECOM*, 1990.
- [35] K. Liao and L. G. Mason. A discrete-time single server queue with a two-level modulated input and its applications. In *GLOBECOM*, pages 913–918, 1989.
- [36] K. Liao and L. G. Mason. A heuristic approach for performance analysis of ATM systems. In *GLOBECOM*, pages 1931–1935, 1990.
- [37] B. Maglaris, D. Anastassiou, P. Sen, G. Karlsson, and J. Robbins. Performance models of statistical multiplexing in packet video communications. *IEEE Trans. Commun.*, 36:834–843, 1988.
- [38] C. B. Moler and C. Van Loan. Nineteen dubious ways to compute the matrix exponential. *SIAM Rev.*, 20:801–836, 1978.
- [39] M. F. Neuts. *Matrix Geometric Solutions in Stochastic Models*. Johns Hopkins University Press, Baltimore, MD, 1981.
- [40] I. Norros, J. W. Roberts, A. Simonian, and J. T. Virtamo. The superposition of variable bit rate sources in an ATM multiplexer. *IEEE JSAC*, 9(3):378–387, 1991.
- [41] L. T. Pillage and R. A. Rohrer. Asymptotic waveform evaluation for timing analysis. *IEEE Trans. Computer Aided Design*, 9(4):352–366, 1990.
- [42] J. W. Roberts. Variable-bit-rate traffic control in B-ISDN. *IEEE Commun. Magazine*, 29(9):50–56, 1991.

- [43] J. W. Roberts. Traffic control in the BISDN. *Comp. Networks and ISDN Sys.*, 25:1055–1064, 1993.
- [44] J. W. Roberts and J. T. Virtamo. The superposition of periodic cell arrival processes in an ATM multiplexer. *IEEE Trans. Commun.*, 39(2):298–303, 1991.
- [45] H. Saito, M. Kawarasaki, and H. Yamada. An analysis of statistical multiplexing in an ATM transport network. *IEEE JSAC*, 9(3):359–367, 1991.
- [46] P. Sen, B. Maglaris, N. Rikli, and D. Anastassiou. Models for packet switching of variable-bit-rate sources. *IEEE JSAC*, 7(5):865–869, 1989.
- [47] K. Sriram and W. Whitt. Characterizing superposition arrival processes in packet multiplexers for voice and data. *IEEE JSAC*, 4:833–846, 1986.
- [48] W. Stallings. *ISDN and Broadband ISDN*. Macmillan Publishing Company, New York, 1992.
- [49] T. E. Stern and A. I. Elwalid. Analysis of separable Markov-modulated rate models for information-handling systems. *Adv. Appl. Prob.*, 23:105–139, 1991.
- [50] L. Takács. Investigation of waiting time problems by reduction to Markov processes. *Acta Math Acad. Sci. Hung.*, 6:101–129, 1955.
- [51] L. Takács. A single-server queue with Poisson input. *Operations Research*, 10:388–397, 1962.
- [52] L. Takács. *Combinatorial Methods in the Theory of Stochastic Processes*. Kluwer Academic Publishers, 1967.
- [53] F. A. Tobagi. Fast packet switch architectures for broadband integrated services digital networks. *Proc. IEEE*, 78(1):133–167, 1990.
- [54] R. C. F. Tucker. Accurate method for analysis of a packet-speech multiplexer with limited delay. *IEEE Trans. Commun.*, 36(4):479–483, 1988.

- [55] J. Vlach and K. Singhal. *Computer Methods for Circuit Analysis and Design*. Van Nostrand Reinhold Company, New York, 1983.
- [56] R. W. Wolff. *Stochastic Modeling and the Theory of Queues*. Prentice Hall, Englewood Cliffs, NJ, 1989.
- [57] A. K. Wong. Queueing analysis for ATM switching of continuous-bit-rate traffic - a recursion computation method. In *GLOBECOM*, pages 1438–1444, 1990.
- [58] G. M. Woodruff and R. Kositpaiboon. Multimedia traffic management principles for guaranteed ATM network performance. *IEEE JSAC*, 8(3):437–446, 1990.
- [59] G. M. Woodruff, R. G. Rogers, and P. S. Richards. A congestion control framework for high-speed integrated packetized transport. In *GLOBECOM*, pages 7.1.1–7.1.5, 1988.
- [60] H. Yamada. An analysis of statistical multiplexing characteristics of ATM cells with voice and video inputs. *Trans. IEICE Japan*, J75-B-1(8):509–516, 1992.
- [61] Y. Yatsuzuka. Highly sensitive speech detector and high speed voiceband discriminator in DSI-ADPCM system. *IEEE Trans. Commun.*, 30:739–750, 1982.
- [62] S. Yazid and H. T. Mouftah. Congestion control methods for BISDN. *IEEE Commun. Mag.*, 30(7):42–47, 1992.
- [63] N. Yin, S. Q. Li, and T. E. Stern. Congestion control for packet voice by selective packet discarding. *IEEE Trans. Commun.*, 38(5):674–683, 1990.



# Vita

Nail Akar was born in Çorum, Turkey, in 1967. He received B. Sc. degree from the Middle East Technical University, Ankara, Turkey, and M. Sc. degree from the Bilkent University, Ankara, Turkey, in 1987 and 1989, respectively, both in Electrical and Electronics Engineering. His research interests include performance evaluation of high speed packet-switched networks, traffic control in ATM networks, and applications of linear systems theory and queueing theory to communication networking problems.