

STUDIES
ON
ALGORITHMIC ANALYSIS OF QUEUES
WITH BATCH MARKOVIAN ARRIVAL STREAMS

HIROYUKI MASUYAMA

STUDIES
ON
ALGORITHMIC ANALYSIS OF QUEUES
WITH BATCH MARKOVIAN ARRIVAL STREAMS

HIROYUKI MASUYAMA

STUDIES
ON
ALGORITHMIC ANALYSIS OF QUEUES
WITH BATCH MARKOVIAN ARRIVAL STREAMS

by

HIROYUKI MASUYAMA

Submitted in partial fulfillment of
the requirement for the degree of
DOCTOR OF INFORMATICS
(Applied Mathematics and Physics)

KYOTO UNIVERSITY
KYOTO 606–8501, JAPAN
NOVEMBER 2003

Preface

Algorithmic analysis of queues, pioneered by Neuts [Neut79] in the mid-1970's, focuses on developing numerical algorithms to compute performance measures of interest in queueing models, rather than on obtaining their formal solutions. Its basic idea is to introduce appropriate auxiliary variables and represent queueing models as Markov chains (or Markov processes) whose structures are suitable for algorithmic analysis of their transient and steady-state solutions. Before the era of the high-performance computer, algorithmic analysis did not gain much attention because the resulting numerical algorithms demand considerable computer performance from a practical viewpoint. As computer performance is improving, however, algorithmic analysis becomes important in queueing research.

It is Markovian arrival process (MAP) [Luca90] which must not be forgotten when discussing algorithmic analysis of queues. MAP is a natural extension of Poisson process, in which the arrival rate is governed by a continuous-time Markov chain with finite states. MAP and its extensions (e.g. batch MAP [Luca91] and marked MAP [He96, He01]) have enough flexibilities for practical use and good representations for algorithmic analysis. For the last decade, many researchers have studied queues with MAP and its extensions and made a great contribution to algorithmic analysis in queueing research. In most of those works, however, service times are assumed to be independent and identically distributed (i.i.d.).

This thesis studies queues with batch Markovian arrival streams. The respective arrival streams are batch MAPs (BMAPs), which are extensions of MAP to allow batch arrivals, and they may have different service time distributions. Queues with batch Markovian arrival streams are so flexible that they can represent most of the queues studied in the past as special cases. In chapter 2, we establish a numerical algorithm to compute the stationary joint queue length distribution in a FIFO single-server queue by extending Takine's recent results [Taki01a, Taki01b, Taki01c]. In a very similar way to chapter 2, chapter 3 constructs a numerical algorithm for the stationary joint queue length distribution in a FIFO single-server queue with service interruptions. In chapter 4, we consider an infinite-server queue. Assuming phase-type service time distributions, we obtain an explicit and numerically feasible formula for time-dependent binomial moments of the queue length. Finally, chapter 5 studies a processor-sharing queue, focusing on computing the sojourn time distribution. Although a single arrival stream and exponential services are assumed there, the processor-sharing queue has been recognized to be very hard to obtain results suitable for the computation of the sojourn time distribution.

The main contribution of this thesis is to develop numerical algorithms to compute performance

measures of interest in several queues with stream-dependent service times. Although the results in this thesis is only one small step in such research, the author hopes that they will be helpful to further research in this field.

Hiroyuki Masuyama
November 2003

Acknowledgment

This thesis would not have been completed without the help of many people, whom I would like to acknowledge here.

First of all, I wish to express my gratitude to Associate Professor Tetsuya Takine of Kyoto University for his constant encouragement and instruction since my graduate student days. I had not studied queueing theory until I became a graduate student, but, his seminar on queueing theory and Markov process let me find the pleasure of studying them and decided to pursue research in this field. Also, I deeply appreciate his helpful comments on an earlier draft of this thesis. Without his support, none of this work could have been completed.

I am indebted to Professor Masao Fukushima of Kyoto University for his insightful suggestions and comments. Although I had not many opportunities to receive his direct instruction, his valuable advice greatly improved my research. Further, I wish to thank other members of Fukushima Laboratory, including Assistant Professor Nobuo Yamashita of Kyoto University, for providing a pleasant working environment. I feel happy to have studied at Fukushima Laboratory for five years.

I would also like to express my great appreciation to Asuka Ikeda, who will be my wife. For six years, she has been always encouraging me to overcome difficulties encountered in advancing my research. Without her moral support, I could not complete this thesis.

Finally, my special thanks are due to my parents for their financial assistance and heartfelt encouragement.

Contents

Preface	v
Acknowledgment	vii
List of Figures	xiii
List of Tables	xv
1 Introduction	1
1.1 Markovian Input Process	1
1.1.1 Markovian arrival process	1
1.1.2 Input process with batch Markovian arrival streams	4
1.2 Markov Chains Skip Free to One Direction	5
1.2.1 Markov chains of G/M/1 type	9
1.2.2 Markov chains of M/G/1 type	10
1.3 Markov Processes Skip Free to One Direction	12
1.3.1 Markov processes skip free to the right	12
1.3.2 Markov processes skip free to the left	15
1.4 Queues with Multiple Arrival Streams	19
1.5 Overview of the Thesis	20
2 FIFO Single-Server Queue	23
2.1 Introduction	23
2.2 Input Process	24
2.3 Waiting Time Distribution	25
2.4 Joint Queue Length Distribution	26
2.4.1 Relationship in the joint queue length distributions	26
2.4.2 Joint queue length distribution immediately after departures	27
2.4.3 Recursions for discrete phase-type batch sizes	28
2.5 Implementations of Recursions	34
2.6 Numerical Examples	40
2.6.1 Efficiency of the algorithm	42
2.6.2 Number of customers in Example 1	43

2.6.3	Number of customers in Example 2	45
2.7	Concluding Remarks	46
3	FIFO Single-Server Queue with Service Interruptions	49
3.1	Introduction	49
3.2	Model	50
3.3	Sojourn Time	53
3.4	Joint Queue Length Distribution	55
3.4.1	Joint queue length distribution immediately after departures	56
3.4.2	Recursions for discrete phase-type batch sizes	57
3.5	Numerical Examples	63
3.5.1	Impact of service time dependency	63
3.5.2	Impact of variation of on- and off-periods	64
3.5.3	Impact of correlation in on- and off-periods	65
3.5.4	Impact of correlation between on-off and arrival processes	68
4	Infinite-Server Queue	71
4.1	Introduction	71
4.2	Model	72
4.3	Time-Dependent Joint Distribution of the Number of Customers	72
4.4	Numerically Feasible Formulas for Phase-Type Service Times	75
4.4.1	Time-dependent joint binomial moments	75
4.4.2	Time-dependent formula for phase-type services	78
4.4.3	Limiting formula for phase-type services	82
4.5	Numerical Examples	84
4.5.1	Impact of service time distribution on $\text{Var}[N]$	85
4.5.2	Impact of arrival process	86
4.5.3	Impact of correlation in service time sequence	89
5	Processor-Sharing Queue	93
5.1	Introduction	93
5.2	Model and Known Results	94
5.3	Sojourn Time Distribution	94
5.4	Numerical Examples	98
5.4.1	Numerical procedure and accuracy guarantee	98
5.4.2	Impact of variation in inter-arrival times	99
5.4.3	Impact of correlation in inter-arrival times	99
5.5	Concluding Remarks	101

6 Conclusion	103
6.1 Summary of Results	103
6.2 Future Work	103
A Uniformization	105
B Queue Length Distribution in a BMAP/GI/1 Queue	107
C Total Queue Length Distribution in a FIFO Queue	111
D Proof of Lemma 3.2	115
E Proof of Theorem 3.5	117
F Total Queue Length Distribution in a Queue with Service Interruptions	119
G Proof of Lemma 4.1	121

List of Figures

1.1	The bivariate Markov process with the skip free to the right property.	13
1.2	The bivariate Markov process with the skip free to the left property.	16
2.1	Complementary distribution of total number of customers in Example 1.	44
2.2	Complementary distribution of total number of customers in Example 1.	45
2.3	Complementary distribution of total number of customers in Example 1.	45
3.1	Expected total queue length $E[N]$	64
3.2	99.9 percentile (99.9 PT) of the total queue length.	66
3.3	Expected total queue length $E[N]$	66
3.4	99.9 percentile (99.9 PT) of the total queue length.	67
3.5	Expected total queue length $E[N]$	67
3.6	Conditional expected total queue length.	67
3.7	Conditional expected total queue length.	67
3.8	99.9 percentile (99.9 PT) of the total queue length.	69
3.9	Expected total queue length $E[N]$	69
4.1	Limiting variance of the number of customers.	86
4.2	Limiting variance of the number of customers.	87
4.3	Time-dependent covariance $\text{Cov}[N(t)]$ of the number of customers.	89
4.4	Time-dependent variance $\text{Var}[N(t)]$ of the number of customers	89
4.5	Variance and covariance in Case 1.	91
4.6	Variance and covariance in Case 2.	91
4.7	Variance and covariance in Case 3.	91
5.1	The computational domain of the $\mathbf{h}_{n,k}$	99
5.2	The complementary distribution $\overline{W}(x)$ of the sojourn time.	100
5.3	The complementary distribution $\overline{W}(x)$ of the sojourn time.	100

List of Tables

2.1	Number of computed $\mathbf{F}_m(\mathbf{n})$'s in Example 1.	43
2.2	Number of stored $\check{\mathbf{F}}_m(\mathbf{n})$'s in Example 1.	44
2.3	Joint queue length distribution $\mathbf{p}(n_1, n_2)\mathbf{e}$	46
2.4	Expected total number of customers in Example 1.	46
2.5	Expected total number of customers in Example 2.	47
3.1	Joint queue length distribution $\mathbf{p}(n_1, n_2)\mathbf{e}$	65

Chapter 1

Introduction

This chapter provides materials to be required in the following chapters and the brief survey of previous works on queues fed by multiple arrival streams with different service time distributions. Throughout this thesis, we denote matrices and vectors by bold capital letters and bold small letters, respectively, and the empty sum is defined as zero.

1.1 Markovian Input Process

This section introduces an input process with batch Markovian arrival streams having different service time distributions. In this input process, customer arrivals from each arrival stream follow a BMAP [Luca91], which is an extension of a MAP [Luca90]. Therefore, after discussing MAP and BMAP [Luca91] in subsection 1.1.1, we formally define the input process with batch Markovian arrival streams in subsection 1.1.2.

1.1.1 Markovian arrival process

We begin with the definition of MAP. We consider a time homogeneous, stationary Markov chain with finite states $\mathcal{M} = \{1, \dots, M\}$, which is assumed to be irreducible. The Markov chain is called the underlying Markov chain hereafter.

MAP is defined as follows. The underlying Markov chain stays in state i ($i \in \mathcal{M}$) for an exponential interval of time with mean μ_i^{-1} . When the sojourn time in state i has elapsed, there are two possibilities: (1) With probability $d_{i,j}$, a transition to state j ($j \in \mathcal{M}$) happens with an arrival. (2) With probability $c_{i,j}$ ($j \in \mathcal{M}, j \neq i$), a transition to state j happens without an arrival. Note that

$$\sum_{\substack{j \in \mathcal{M} \\ j \neq i}} c_{i,j} + \sum_{j \in \mathcal{M}} d_{i,j} = 1, \quad \text{for all } i \in \mathcal{M}.$$

For the later use, we introduce two matrices, \mathbf{C} and \mathbf{D} . Let \mathbf{C} denote an $M \times M$ matrix whose (i, j) th ($i, j \in \mathcal{M}$) element $C_{i,j}$ is given by

$$C_{i,j} = \begin{cases} -\mu_i, & \text{if } i = j, \\ \mu_i c_{i,j}, & \text{otherwise.} \end{cases}$$

Let \mathbf{D} denote an $M \times M$ nonnegative matrix whose (i, j) th $(i, j \in \mathcal{M})$ element $D_{i,j}$ is given by

$$D_{i,j} = \mu_i d_{i,j}.$$

Note that the infinitesimal generator of the underlying Markov chain is given by $\mathbf{C} + \mathbf{D}$. We define $\boldsymbol{\pi}$ as the stationary probability vector of the underlying Markov chain. Because of the finite state space \mathcal{M} and the irreducibility of the underlying Markov chain, $\boldsymbol{\pi}$ is uniquely determined. Note that $\boldsymbol{\pi}$ satisfies

$$\boldsymbol{\pi}(\mathbf{C} + \mathbf{D}) = \mathbf{0}, \quad \boldsymbol{\pi}\mathbf{e} = 1, \quad (1.1)$$

where \mathbf{e} denotes a column vector whose elements are all equal to one. The arrival rate λ is then given by

$$\lambda = \boldsymbol{\pi}\mathbf{D}\mathbf{e}.$$

We define $N(t)$ and $S(t)$ as the number of arrivals in time interval $(0, t]$ and the state of the underlying Markov chain at time t , respectively. Further, let $\mathbf{N}(t, n)$ denote an $M \times M$ matrix whose (i, j) th $(i, j \in \mathcal{M})$ element represents

$$\Pr[N(t) = n, S(t) = j \mid N(0) = 0, S(0) = i].$$

$\mathbf{N}(t, n)$'s ($t \geq 0, n = 0, 1, \dots$) satisfy the forward Chapman-Kolmogorov equations:

$$\begin{aligned} \frac{d}{dt}\mathbf{N}(t, n) &= \mathbf{N}(t, n)\mathbf{C} + \mathbf{N}(t, n-1)\mathbf{D}, & t \geq 0, n = 1, 2, \dots, \\ \mathbf{N}(0, 0) &= \mathbf{I}, \end{aligned}$$

where \mathbf{I} denotes an identity matrix with appropriate dimension. Thus, the z -transform $\mathbf{N}^*(t, z)$ of $\mathbf{N}(t, n)$ is given by

$$\mathbf{N}^*(t, z) = \exp[(\mathbf{C} + z\mathbf{D})t].$$

As mentioned above, MAP can be characterized by (\mathbf{C}, \mathbf{D}) .

In what follows, we show some special cases of MAP.

Example 1.1 Poisson process. For $\mathbf{C} = -\lambda$ and $\mathbf{D} = \lambda$, MAP is reduced to a Poisson process with arrival rate λ .

Example 1.2 Phase-type (PH) renewal process. PH renewal process is defined in a very similar way to MAP except that the state of the underlying Markov chain immediately after an arrival is chosen according to a $1 \times M$ probability vector $\boldsymbol{\beta}$, independent of its state immediately before the arrival. Namely, $\mathbf{D} = (-\mathbf{T}\mathbf{e})\boldsymbol{\beta}$, where we use, by convention, \mathbf{T} for \mathbf{C} to distinguish PH renewal process from MAP. Note that the i th $(i \in \mathcal{M})$ element of $M \times 1$ vector $(-\mathbf{T}\mathbf{e})$ represents the arrival (renewal) rate given the underlying Markov chain is in state i . Note also that inter-arrival (renewal) times are i.i.d. according to distribution function $\psi(x)$:

$$\psi(x) = 1 - \boldsymbol{\beta}\exp(\mathbf{T}x)\mathbf{e},$$

which is called the phase-type distribution with representation $(\boldsymbol{\beta}, \mathbf{T})$.

Example 1.3 Markov modulated Poisson process (MMPP). *MMPP is a natural extension of a Poisson process, which has the arrival rate modulated by an M -state irreducible Markov chain. To put it more concretely, the arrival rate takes M values $\lambda_1, \dots, \lambda_M$, and is equal to λ_j whenever the Markov chain is in state j . Let \mathbf{R} denote the infinitesimal generator of the underlying Markov chain. Let $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_M)$. MMPP can then be characterized by \mathbf{R} and $\mathbf{\Lambda}$. Note here that this MMPP is equivalent to a MAP with representation (\mathbf{C}, \mathbf{D}) , where*

$$\mathbf{C} = \mathbf{R} - \mathbf{\Lambda}, \quad \mathbf{D} = \mathbf{\Lambda}.$$

We next discuss BMAP. The definition of BMAP is very similar to that of MAP except that multiple arrivals may happen simultaneously. Given the underlying Markov chain is in state i ($i \in \mathcal{M}$), it changes its state to state j ($j \in \mathcal{M}$) at rate $\mu_i d_{i,j}(n)$ and n arrivals occur at the same time. Besides, transitions with no arrivals are driven by \mathbf{C} , as well as MAP. Note here that $d_{i,j}(n)$'s satisfy

$$\sum_{\substack{j \in \mathcal{M} \\ j \neq i}} c_{i,j} + \sum_{j \in \mathcal{M}} \sum_{n=1}^{\infty} d_{i,j}(n) = 1, \quad \text{for all } i \in \mathcal{M},$$

where $c_{i,j}$ denotes the (i, j) th element of \mathbf{C} . Let $\mathbf{D}(n)$ denote an $M \times M$ matrix whose (i, j) th $(i, j \in \mathcal{M})$ element $D_{i,j}(n)$ is given by

$$D_{i,j}(n) = \mu_i d_{i,j}(n).$$

Note that $\mathbf{C} + \sum_{n=1}^{\infty} n \mathbf{D}(n)$ is the infinitesimal generator of the underlying Markov chain. Note further that the arrival rate λ is given by

$$\lambda = \boldsymbol{\pi} \sum_{n=1}^{\infty} n \mathbf{D}(n) \mathbf{e}, \tag{1.2}$$

where $\boldsymbol{\pi}$ denotes the stationary probability vector of the underlying Markov chain. Note also that $\mathbf{N}(t, n)$'s for BMAP satisfy

$$\begin{aligned} \frac{d}{dt} \mathbf{N}(t, n) &= \mathbf{N}(t, n) \mathbf{C} + \sum_{m=1}^n \mathbf{N}(t, n-m) \mathbf{D}(m), \quad t \geq 0, n = 1, 2, \dots, \\ \mathbf{N}(0, 0) &= \mathbf{I}. \end{aligned}$$

It then follows from the above equations that the z -transform $\mathbf{N}^*(t, z)$ of $\mathbf{N}(t, n)$ is given by

$$\mathbf{N}^*(t, z) = \exp[(\mathbf{C} + \mathbf{D}^*(z)) t],$$

where

$$\mathbf{D}^*(z) = \sum_{n=1}^{\infty} z^n \mathbf{D}(n).$$

Thus, BMAP can be characterized by $(\mathbf{C}, \mathbf{D}(n))$.

1.1.2 Input process with batch Markovian arrival streams

It is assumed that there exist K ($K \geq 1$) arrival streams. Customers arriving from the k th ($k \in \mathcal{K} = \{1, \dots, K\}$) stream is called *class k* customers. Service times of class k customers are assumed to be i.i.d. according to a distribution function $H_k(x)$ ($x \geq 0$) with finite mean h_k .

As well as in subsection 1.1.1, we consider the underlying Markov chain with the finite state space $\mathcal{M} = \{1, \dots, M\}$ and assume its irreducibility. Customer arrivals happen in the following way. The underlying Markov chain stays in state i ($i \in \mathcal{M}$) for an exponential interval of time with mean μ_i^{-1} , and then changes its state to state j ($j \in \mathcal{M}$) with probability $\sigma_{i,j}$, where $\sum_{j \in \mathcal{M}} \sigma_{i,j} = 1$ for all i ($i \in \mathcal{M}$). When a transition from state i to state j happens, n_k ($k \in \mathcal{K}$) class k customers simultaneously arrive at the queue with probability $\sigma_{i,j}(n_1, \dots, n_K) / \sigma_{i,j}$, where $\sigma_{i,i}(0, \dots, 0)$ is assumed to be zero for all i ($i \in \mathcal{M}$), and $\sigma_{i,j}(n_1, \dots, n_K)$ satisfies

$$\sigma_{i,j} = \sum_{n_1=1}^{\infty} \cdots \sum_{n_K=1}^{\infty} \sigma_{i,j}(n_1, \dots, n_K).$$

The assumption $\sigma_{i,i}(0, \dots, 0) = 0$ ($\forall i \in \mathcal{M}$) implies that at least one customer arrives whenever a transition from any state i to itself happens.

We now introduce some notations to describe the above input process. We define \mathbf{C} as an $M \times M$ matrix whose (i, j) th element $C_{i,j}$ is given by

$$C_{i,j} = \begin{cases} -\mu_i, & i = j, \\ \sigma_{i,j}(\mathbf{0})\mu_i, & \text{otherwise,} \end{cases}$$

where $\mathbf{0}$ denotes a vector of zeros with appropriate dimension. We also define \mathcal{Z} and \mathcal{Z}^+ as

$$\begin{aligned} \mathcal{Z} &= \{(n_1, \dots, n_K); n_k = 0, 1, \dots \text{ for all } k \in \mathcal{K}\}, \\ \mathcal{Z}^+ &= \mathcal{Z} - \{\mathbf{0}\}, \end{aligned}$$

respectively. Let \mathbf{n} denote a $1 \times K$ nonnegative vector (n_1, \dots, n_K) . We then define $\mathbf{D}(\mathbf{n})$ ($\mathbf{n} \in \mathcal{Z}^+$) as an $M \times M$ matrix whose (i, j) th element $D_{i,j}(\mathbf{n})$ is given by

$$D_{i,j}(\mathbf{n}) = \sigma_{i,j}(\mathbf{n})\mu_i, \quad \mathbf{n} \in \mathcal{Z}^+.$$

When a state transition driven by $\mathbf{D}(\mathbf{n})$ occurs, n_k customers of each class k simultaneously arrive. On the other hand, when a state transition driven by \mathbf{C} occurs, no customers arrive. Note here that the infinitesimal generator of the underlying Markov chain is given by $\mathbf{C} + \mathbf{D}$, where \mathbf{D} is defined as

$$\mathbf{D} = \sum_{\mathbf{n} \in \mathcal{Z}^+} \mathbf{D}(\mathbf{n}). \tag{1.3}$$

Thus, the counting process of arrivals is characterized by $(\mathbf{C}, \mathbf{D}(\mathbf{n}))$.

Example 1.4 *Suppose that*

$$\mathbf{D}(\mathbf{n}) = \mathbf{O} \quad \text{if } \mathbf{n} \in \mathcal{Z} - \{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_K\},$$

where \mathbf{O} denotes a matrix of zeros and \mathbf{e}_k denotes the k th unit vector:

$$\mathbf{e}_k = \underbrace{(0, \dots, 0, 1, 0, \dots, 0)}_{k\text{th}}. \quad (1.4)$$

This assumption means that at most one arrival from an arrival stream happens at any time. This special case is called marked MAP (MMAP) in [He96].

Finally, we consider the marginal arrival process of the k th ($k \in \mathcal{K}$) stream. We define $\tilde{\mathbf{C}}_k$ and $\tilde{\mathbf{D}}_k(n)$ ($n = 1, 2, \dots$) as

$$\tilde{\mathbf{C}}_k = \mathbf{C} + \sum_{\substack{\mathbf{n} \in \mathcal{Z}^+ \\ n_k = 0}} \mathbf{D}_k(\mathbf{n}), \quad \tilde{\mathbf{D}}_k(n) = \sum_{\substack{\mathbf{n} \in \mathcal{Z}^+ \\ n_k = n}} \mathbf{D}_k(\mathbf{n}),$$

respectively. Note that $\tilde{\mathbf{C}}_k$ and $\tilde{\mathbf{D}}_k(n)$ satisfy

$$\tilde{\mathbf{C}}_k + \sum_{n=1}^{\infty} \tilde{\mathbf{D}}_k(n) = \mathbf{C} + \mathbf{D}.$$

When a transition driven by $\tilde{\mathbf{D}}_k(n)$ happens, n customers of class k arrive at the queue. Besides, a transition driven by $\tilde{\mathbf{C}}_k$ happens without a class k customer. Thus, the marginal arrival process of the k th stream is a BMAP with representation $(\tilde{\mathbf{C}}_k, \tilde{\mathbf{D}}_k(n))$. We define λ_k ($k \in \mathcal{K}$) as the arrival rate of class k customers, i.e.,

$$\lambda_k = \boldsymbol{\pi} \sum_{n=1}^{\infty} n \tilde{\mathbf{D}}_k(n) \mathbf{e}_k,$$

where $\boldsymbol{\pi}$ denotes the invariant probability vector for $\mathbf{C} + \mathbf{D}$. Let ρ_k denote the utilization factor of class k customers, i.e., $\rho_k = \lambda_k h_k$. Note that the overall arrival rate λ and utilization factor ρ are given by

$$\lambda = \sum_{k \in \mathcal{K}} \lambda_k, \quad \rho = \sum_{k \in \mathcal{K}} \rho_k,$$

respectively. Note also that in a work-conserving single-server queue, $\rho < 1$ ensures that the system is stable [Loyn62].

1.2 Markov Chains Skip Free to One Direction

This section summarizes analytical results on time-homogeneous Markov chains with the skip free to one direction and spatial homogeneity properties. Markov chains with these properties are classified into two cases, i.e., the skip free to the left case and to the right case, which are natural extensions of embedded Markov chains in an M/G/1 queue and GI/M/1 queue, respectively. Such Markov chains now play a vital role in algorithmic analysis of the queue length distribution. In what follows, we first introduce Markov chains with the skip free to the right/left properties and briefly discuss the steady-state solutions of these two types of Markov chains in subsections 1.2.1 and 1.2.2.

In a single-server queue, let X_k^- and D_k denote the number of customers immediately before the k th arrival and the number of departures during the k th inter-arrival time, respectively. We then have

$$X_{k+1}^- = \max(X_k^- + 1 - D_{k+1}, 0), \quad k = 0, 1, \dots$$

We now consider a GI/M/1 queue, where inter-arrival times are i.i.d. according to a distribution function $G(x)$ with mean λ^{-1} and service times are exponential with mean μ^{-1} . In such a queue, D_k 's are i.i.d. random variables, and hence the sequence $\{X_k^-; k = 0, 1, \dots\}$ of random variables is a Markov chain with transition probability matrix:

$$\begin{bmatrix} f_0 & e_0 & 0 & 0 & \cdots \\ f_1 & e_1 & e_0 & 0 & \cdots \\ f_2 & e_2 & e_1 & e_0 & \cdots \\ f_3 & e_3 & e_2 & e_1 & \cdots \\ f_4 & e_4 & e_3 & e_2 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix},$$

where

$$e_n = \int_0^\infty e^{-\mu x} \frac{(\mu x)^n}{n!} dG(x), \quad f_n = \sum_{m=n+1}^\infty e_m.$$

If $\rho = \lambda/\mu < 1$, the Markov chain $\{X_k^-; k = 0, 1, \dots\}$ has the stationary distribution. Assuming $\rho < 1$, we define $q_{\text{GI/M/1}}(n)$ ($n = 0, 1, \dots$) as $\lim_{k \rightarrow \infty} \Pr[X_k^- = n]$. Then $q_{\text{GI/M/1}}(n)$ is given by

$$q_{\text{GI/M/1}}(n) = \gamma^n (1 - \gamma), \quad n = 0, 1, \dots, \quad (1.5)$$

where γ is a positive number satisfying $\gamma = \sum_{n=0}^\infty e_n \gamma^n$ and $\gamma < 1$. The assumption $\rho < 1$ ensures that γ is uniquely determined.

On the other hand, in a single-server queue, let X_k^+ and A_k denote the number of customers immediately after the k th departure and the number of arrivals during the k th service time. We then have

$$X_{k+1}^+ = \max(X_k^+ - 1, 0) + A_{k+1}.$$

We consider an M/G/1 queue, where arrivals follow Poisson process with a rate λ and service times are i.i.d. according to a distribution $H(x)$ with mean μ^{-1} . In the M/G/1 queue, A_k 's are i.i.d. random variables, and hence $\{X_k^+; k = 0, 1, \dots\}$ is a Markov chain with transition probability matrix:

$$\begin{bmatrix} a_0 & a_1 & a_2 & a_3 & \cdots \\ a_0 & a_1 & a_2 & a_3 & \cdots \\ 0 & a_0 & a_1 & a_2 & \cdots \\ 0 & 0 & a_0 & a_1 & \cdots \\ 0 & 0 & 0 & a_0 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix},$$

where

$$a_n = \int_0^\infty e^{-\lambda x} \frac{(\lambda x)^n}{n!} dH(x).$$

If $\rho = \lambda/\mu < 1$, the Markov chain $\{X_k^+; k = 0, 1, \dots\}$ has the stationary distribution. Under the assumption $\rho < 1$, we define $q_{M/G/1}(n)$ ($n = 0, 1, \dots$) as $\lim_{k \rightarrow \infty} \Pr[X_k^+ = n]$. Then $q_{M/G/1}(n)$'s are recursively determined in the following way.

$$\begin{aligned} q_{M/G/1}(0) &= 1 - \rho, \\ q_{M/G/1}(n) &= \left[q_{M/G/1}(0) \bar{a}_n + \sum_{m=1}^{n-1} q_{M/G/1}(m) \bar{a}_{n+1-m} \right] (1 - \bar{a}_1)^{-1}, \quad n = 1, 2, \dots, \end{aligned}$$

where

$$\bar{a}_n = \sum_{m=n}^{\infty} a_m.$$

Note that Markov chains $\{X_k^-; k = 0, 1, \dots\}$ and $\{X_k^+; k = 0, 1, \dots\}$ have the following two properties. One is the spatial homogeneity except for the boundary part (e.g., if $X_k^+ \geq 1$, X_k^+ increases by $n - 1$ with probability a_n). The other is the skip free to one direction property. To put it more concretely, the Markov chain $\{X_k^-; k = 0, 1, \dots\}$ (resp. $\{X_k^+; k = 0, 1, \dots\}$) can increase (resp. decrease) at most by one when a transition happens. This is called the *skip free to the right* (resp. left) property.

Next we consider a GI/PH/1 queue and a PH/G/1 queue, where the service time distribution and the inter-arrival time distribution are of phase-type with representation $(\boldsymbol{\beta}_H, \mathbf{T}_H)$ and $(\boldsymbol{\beta}_A, \mathbf{T}_A)$, respectively. Let $G(x)$ and $H(x)$ denote the inter-arrival distribution in the GI/PH/1 queue and the service time distribution in the PH/G/1 queue, respectively. Note that D_k 's in the GI/PH/1 queue and A_k 's in the PH/G/1 queue are non-i.i.d. random variables. However, D_{k+1} (resp. A_{k+1}) is conditionally independent of the past history $\{D_l; l = 1, \dots, k\}$ (resp. $\{A_l; l = 1, \dots, k\}$), given the state J_k^- (resp. J_k^+) of the phase-type distribution $(\boldsymbol{\beta}_H, \mathbf{T}_H)$ (resp. $(\boldsymbol{\beta}_A, \mathbf{T}_A)$) immediately before (resp. after) the k th arrival (resp. departure). Thus, $\{(X_k^-, J_k^-); k = 0, 1, \dots\}$ is a bivariate Markov chain whose transition probability matrix is given by

$$\mathbf{P}_{\text{GI/PH/1}} = \begin{bmatrix} \mathbf{F}_0 & \mathbf{E}_0 & \mathbf{O} & \mathbf{O} & \cdots \\ \mathbf{F}_1 & \mathbf{E}_1 & \mathbf{E}_0 & \mathbf{O} & \cdots \\ \mathbf{F}_2 & \mathbf{E}_2 & \mathbf{E}_1 & \mathbf{E}_0 & \cdots \\ \mathbf{F}_3 & \mathbf{E}_3 & \mathbf{E}_2 & \mathbf{E}_1 & \cdots \\ \mathbf{F}_4 & \mathbf{E}_4 & \mathbf{E}_3 & \mathbf{E}_2 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}, \quad (1.6)$$

where \mathbf{E}_n 's ($n = 0, 1, \dots$) satisfy

$$\sum_{n=0}^{\infty} z^n \mathbf{E}_n = \int_0^\infty \exp[(\mathbf{T}_H + z(-\mathbf{T}_H \mathbf{e}) \boldsymbol{\beta}_H) x] dG(x),$$

and

$$\mathbf{F}_n = \sum_{m=n+1}^{\infty} \mathbf{E}_m e\boldsymbol{\beta}_H, \quad n = 0, 1, \dots$$

Similarly, $\{(X_k^+, J_k^+); k = 0, 1, \dots\}$ is a bivariate Markov chain whose transition probability matrix is given by

$$\mathbf{P}_{\text{PH/G/1}} = \begin{bmatrix} \mathbf{B}_0 & \mathbf{B}_1 & \mathbf{B}_2 & \mathbf{B}_3 & \cdots \\ \mathbf{A}_0 & \mathbf{A}_1 & \mathbf{A}_2 & \mathbf{A}_3 & \cdots \\ \mathbf{O} & \mathbf{A}_0 & \mathbf{A}_1 & \mathbf{A}_2 & \cdots \\ \mathbf{O} & \mathbf{O} & \mathbf{A}_0 & \mathbf{A}_1 & \cdots \\ \mathbf{O} & \mathbf{O} & \mathbf{O} & \mathbf{A}_0 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}, \quad (1.7)$$

where \mathbf{A}_n 's ($n = 0, 1, \dots$) satisfy

$$\sum_{n=0}^{\infty} z^n \mathbf{A}_n = \int_0^{\infty} \exp[(\mathbf{T}_A + z(-\mathbf{T}_A e)\boldsymbol{\beta}_A)x] dH(x),$$

and

$$\mathbf{B}_n = e\boldsymbol{\beta}_A \mathbf{A}_n, \quad n = 0, 1, \dots$$

We can see from (1.6) and (1.7) that these two bivariate Markov chains have the spatial homogeneity (except for the boundary part) and the skip free to one direction property with respect to the first variables, i.e., X_k^- and X_k^+ . As mentioned above, the queue length processes in queues we often encounter are reduced to bivariate Markov chains with these two properties. In general, if a bivariate Markov chain has a transition probability matrix with the same structure as

$$\mathbf{P}_R = \begin{bmatrix} \mathbf{F}_0 & \widehat{\mathbf{E}}_0 & \mathbf{O} & \mathbf{O} & \cdots \\ \mathbf{F}_1 & \mathbf{E}_1 & \mathbf{E}_0 & \mathbf{O} & \cdots \\ \mathbf{F}_2 & \mathbf{E}_2 & \mathbf{E}_1 & \mathbf{E}_0 & \cdots \\ \mathbf{F}_3 & \mathbf{E}_3 & \mathbf{E}_2 & \mathbf{E}_1 & \cdots \\ \mathbf{F}_4 & \mathbf{E}_4 & \mathbf{E}_3 & \mathbf{E}_2 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}, \quad (1.8)$$

it is called a Markov chain of G/M/1 type. Further, a bivariate Markov chain has a transition probability matrix with the same structure as

$$\mathbf{P}_L = \begin{bmatrix} \mathbf{B}_0 & \mathbf{B}_1 & \mathbf{B}_2 & \mathbf{B}_3 & \cdots \\ \widehat{\mathbf{A}}_0 & \mathbf{A}_1 & \mathbf{A}_2 & \mathbf{A}_3 & \cdots \\ \mathbf{O} & \mathbf{A}_0 & \mathbf{A}_1 & \mathbf{A}_2 & \cdots \\ \mathbf{O} & \mathbf{O} & \mathbf{A}_0 & \mathbf{A}_1 & \cdots \\ \mathbf{O} & \mathbf{O} & \mathbf{O} & \mathbf{A}_0 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}, \quad (1.9)$$

it is called a Markov chain of M/G/1 type. These two Markov chains can be considered as extensions of the Markov chains describing the queue length processes in the GI/M/1 queue and the M/G/1 queue, respectively, in the sense that each element of a transition probability matrix is replaced by a block matrix. This is the origin of their names.

Remark 1.1 *Some researchers have studied G/M/1-type and M/G/1-type Markov chains with a tree structure. For G/M/1-type Markov chains with a tree structure, matrix product-form solutions are obtained by Yeung and Sengupta [Yeun94]. As for M/G/1-type Markov chains with a tree structure, He [He00a, He00b, He03a] studied conditions for them to be positive recurrent, null recurrent or transient. A related work is found in [He03b], which studied sufficient conditions for stability or instability of a single-server queue with LCFS preemptive-repeat service discipline. Further, Takine et al. [Taki95] established a paradigm to compute the steady-state solution of the M/G/1-type Markov chain with a tree structure under the assumption that it is irreducible and positive recurrent. That paradigm is an extension of the M/G/1 paradigm [Neut89].*

1.2.1 Markov chains of G/M/1 type

In this subsection, we consider a G/M/1-type Markov chain $\{(X_k^-, J_k^-); k = 0, 1, \dots\}$ with a transition probability matrix \mathbf{P}_R in (1.8), where $\mathbf{E}_n, \mathbf{F}_n$ ($n = 0, 1, \dots$) and $\widehat{\mathbf{E}}_0$ denotes $M_1^- \times M_1^-$, $M_{\min(n,1)}^- \times M_0^-$ and $M_0^- \times M_1^-$ matrices, respectively.

We define \mathbf{E} as

$$\mathbf{E} = \sum_{n=0}^{\infty} \mathbf{E}_n.$$

Note here that \mathbf{E} is a stochastic matrix.

We now assume the following:

Assumption 1.1 *\mathbf{E} is irreducible.*

We then define $\boldsymbol{\nu}_E$ as a $1 \times M_1^-$ vector satisfying

$$\boldsymbol{\nu}_E \mathbf{E} = \boldsymbol{\nu}_E, \quad \boldsymbol{\nu}_E \mathbf{e} = 1.$$

Note that Assumption 1.1 ensures that $\boldsymbol{\nu}_E$ is uniquely determined. We also make the following assumption:

Assumption 1.2

(a) *The G/M/1-type Markov chain $\{(X_k^-, J_k^-); k = 0, 1, \dots\}$ is irreducible, and*

$$(b) \boldsymbol{\nu}_E \sum_{n=1}^{\infty} n \mathbf{E}_n \mathbf{e} > 1.$$

Remark 1.2 *The irreducible G/M/1-type Markov chain $\{(X_k^-, J_k^-); k = 0, 1, \dots\}$ is positive recurrent if and only if Assumption 1.4 (b) holds (See Corollary 3.3 on page 319 in [Asmu03]).*

We define X^- and J^- as generic random variables for X_k^- and J_k^- , respectively, in steady state. Let $\mathbf{x}^- = (\mathbf{x}_0^-, \mathbf{x}_1^-, \dots)$ denote the steady-state solution for the G/M/1-type Markov chain $\{(X_k^-, J_k^-); k = 0, 1, \dots\}$, where \mathbf{x}_n^- denotes a $1 \times M_{\min(n,1)}^-$ vector whose j th element represents $\Pr[X^- = n, J^- = j]$. Let \mathbf{R} denote the minimal nonnegative solution [Neut81] of

$$\mathbf{R} = \sum_{n=0}^{\infty} \mathbf{R}^n \mathbf{E}_n.$$

\mathbf{R} is called a *rate matrix* and has the following probabilistic meanings [Neut81]. We consider a cycle of visits to the state set $\{(n, \nu); \nu = 1, \dots, M_1^-\} (n \geq 1)$. The (i, j) th element of \mathbf{R} then represents the mean number of visits on state $(n+1, j)$ during a cycle that starts with state (n, i) . In terms of rate matrix \mathbf{R} , $\mathbf{x}^- = (\mathbf{x}_0^-, \mathbf{x}_1^-, \dots)$ is given in the following way.

Proposition 1.1 ([Neut81]) \mathbf{x}_0^- and \mathbf{x}_1^- are determined by

$$\begin{aligned} \mathbf{x}_0^- &= \mathbf{x}_0^- \mathbf{F}_0 + \mathbf{x}_1^- \sum_{m=1}^{\infty} \mathbf{R}^{m-1} \mathbf{F}_m, \\ \mathbf{x}_1^- &= \mathbf{x}_0^- \widehat{\mathbf{E}}_0 + \mathbf{x}_1^- \sum_{m=1}^{\infty} \mathbf{R}^{m-1} \mathbf{E}_m, \\ \mathbf{x}_0^- \mathbf{e} + \mathbf{x}_1^- (\mathbf{I} - \mathbf{R})^{-1} \mathbf{e} &= \mathbf{1}, \end{aligned}$$

and for $n = 2, 3, \dots$,

$$\mathbf{x}_n^- = \mathbf{x}_1^- \mathbf{R}^{n-1}. \quad (1.10)$$

The \mathbf{x}_n^- can be considered as a matrix version of the geometric solution (1.5) for the GI/M/1 queue, and it is called the *matrix geometric* solution.

1.2.2 Markov chains of M/G/1 type

In this subsection, we consider an M/G/1-type Markov chain $\{(X_k^+, J_k^+); k = 0, 1, \dots\}$ with a transition probability matrix \mathbf{P}_L in (1.9), where $\mathbf{A}_n, \mathbf{B}_n (n = 0, 1, \dots)$ and $\widehat{\mathbf{A}}_0$ denotes $M_1^+ \times M_1^+$, $M_0^+ \times M_{\min(n,1)}^+$ and $M_1^+ \times M_0^+$ matrices, respectively. We define $\boldsymbol{\mu}_A$ and $\boldsymbol{\mu}_B$ as

$$\boldsymbol{\mu}_A = \sum_{n=1}^{\infty} n \mathbf{A}_n \mathbf{e}, \quad \boldsymbol{\mu}_B = \sum_{n=1}^{\infty} n \mathbf{B}_n \mathbf{e},$$

respectively. We also define \mathbf{A} and \mathbf{B} as

$$\mathbf{A} = \sum_{n=0}^{\infty} \mathbf{A}_n, \quad \mathbf{B} = \sum_{n=1}^{\infty} \mathbf{B}_n,$$

respectively. Note that \mathbf{A} is a stochastic matrix and $\mathbf{B}_0 \mathbf{e} + \mathbf{B} \mathbf{e} = \mathbf{e}$.

We now assume the following:

Assumption 1.3 \mathbf{A} is irreducible.

Let ν_A denote a $1 \times M_1^+$ vector satisfying

$$\nu_A \mathbf{A} = \nu_A, \quad \nu_A \mathbf{e} = 1.$$

Note that Assumption 1.3 ensures that ν_A is uniquely determined. We define ρ_A as

$$\rho_A = \nu_A \boldsymbol{\mu}_A.$$

In what follows, we also make the following assumption:

Assumption 1.4

- (a) The M/G/1-type Markov chain $\{(X_k^+, J_k^+); k = 0, 1, \dots\}$ is irreducible,
- (b) $\boldsymbol{\mu}_B$ is a finite vector, and
- (c) $\rho_A < 1$.

Remark 1.3 The irreducible M/G/1-type Markov chain $\{(X_k^+, J_k^+); k = 0, 1, \dots\}$ is positive recurrent if and only if Assumption 1.4 (b) and (c) hold (See Proposition 3.1 on page 318 in [Asmu03]).

We define X^+ and J^+ as generic random variables for X_k^+ and J_k^+ , respectively, in steady state. Let $\mathbf{x}^+ = (\mathbf{x}_0^+, \mathbf{x}_1^+, \dots)$ denote the steady-state solution for the M/G/1-type Markov chain $\{(X_k^+, J_k^+); k = 0, 1, \dots\}$, where \mathbf{x}_n^+ denotes a $1 \times M_{\min(n,1)}^+$ vector whose j th element represents $\Pr[X^+ = n, J^+ = j]$. The numerical algorithm to compute the steady-state solution $\mathbf{x}^+ = (\mathbf{x}_0^+, \mathbf{x}_1^+, \dots)$ is called the M/G/1 paradigm, where an $M_1^+ \times M_1^+$ matrix \mathbf{G} plays a vital role [Neut89]. The (i, j) th element of \mathbf{G} represents the probability that for $n \geq 2$, the first passage time from state (n, i) to any state in $\{(n-1, 1), \dots, (n-1, M_1^+)\}$ ends with state $(n-1, j)$. The \mathbf{G} is often called the fundamental matrix. It is known that \mathbf{G} is the minimal nonnegative solution of [Neut89]

$$\mathbf{G} = \sum_{n=0}^{\infty} \mathbf{A}_n \mathbf{G}^n.$$

The computation of \mathbf{G} is referred to [Neut89, Luca91]. In terms of \mathbf{G} , we define $M_0^+ \times M_0^+$ matrix \mathbf{K} as

$$\mathbf{K} = \mathbf{B}_0 + \sum_{m=1}^{\infty} \mathbf{B}_m \mathbf{G}^{m-1} \left[\mathbf{I} - \sum_{l=1}^{\infty} \mathbf{A}_l \mathbf{G}^{l-1} \right]^{-1} \hat{\mathbf{A}}_0. \quad (1.11)$$

Further, let $\boldsymbol{\kappa}$ denote a $1 \times M_0^+$ vector satisfying

$$\boldsymbol{\kappa} \mathbf{K} = \boldsymbol{\kappa}, \quad \boldsymbol{\kappa} \mathbf{e} = 1.$$

We now show a recursive formula for the computation of \mathbf{x}_n^+ 's ($n = 0, 1, \dots$).

Proposition 1.2 ([Rama88, Sche90]) \mathbf{x}_0^+ is given by

$$\mathbf{x}_0^+ = \left[1 + \frac{\kappa}{1 - \rho_A} \left\{ \boldsymbol{\mu}_B + \left(\mathbf{B} - \sum_{l=1}^{\infty} \mathbf{B}_l \mathbf{G}^l \right) [\mathbf{I} - \mathbf{A} + \mathbf{e}\nu_A]^{-1} \boldsymbol{\mu}_A \right\} \right]^{-1} \boldsymbol{\kappa},$$

and \mathbf{x}_n^+ 's ($n = 1, 2, \dots$) are recursively determined by

$$\mathbf{x}_n^+ = \left[\mathbf{x}_0^+ \bar{\mathbf{B}}_n + \sum_{l=1}^{n-1} \mathbf{x}_l^+ \bar{\mathbf{A}}_{n-l+1} \right] [\mathbf{I} - \bar{\mathbf{A}}_1]^{-1},$$

where

$$\bar{\mathbf{B}}_\nu = \sum_{m=0}^{\infty} \mathbf{B}_{m+\nu} \mathbf{G}^m, \quad \bar{\mathbf{A}}_\nu = \sum_{m=0}^{\infty} \mathbf{A}_{m+\nu} \mathbf{G}^m.$$

1.3 Markov Processes Skip Free to One Direction

This section considers time-homogeneous Markov processes with the skip free to one direction and spatial homogeneity properties, which are continuous analogs of Markov chains discussed in section 1.2. These Markov processes are also classified into two cases, i.e., the skip free to the left case and to the right case, and they are very important models to analyze the waiting time distribution and the virtual waiting time distribution in single-server queues.

1.3.1 Markov processes skip free to the right

This subsection considers a bivariate Markov process $\{(X_R(t), J_R(t)); t \geq 0\}$ with the skip free to the right property. We formally define it as follows:

- (a) $X_R(t)$ takes values on $[0, \infty)$ and $J_R(t)$ takes integer values on $\mathcal{M} = \{1, \dots, M\}$. Further, $X_R(t)$ is skip free to the right and increases at a linear rate of one if there are no downward jumps.
- (b) Given $(X_R(t), J_R(t)) = (x, i)$ ($x > 0, i \in \mathcal{M}$), the bivariate Markov process can change its state to somewhere between $(x - y, j)$ and $(x - y + dy, j)$ ($0 \leq y < x, j \in \mathcal{M}$) at a rate $dE_{i,j}(y)$. Let $\mathbf{E}(y)$ denote an $M \times M$ matrix whose (i, j) th element is given by $E_{i,j}(y)$. To avoid trivialities, $E_{i,i}(0) = 0$ ($\forall i \in \mathcal{M}$) is assumed. Further, $\mathbf{E}(\infty)$ is assumed to be an irreducible matrix.
- (c) Given $(X_R(t), J_R(t)) = (x, i)$ ($x > 0, i \in \mathcal{M}$), the bivariate Markov process can change its state to state $(0, j)$ ($j \in \mathcal{M}$) at rate $F_{\text{jump},i,j}(x)$. Let $\mathbf{F}_{\text{jump}}(x)$ denote an $M \times M$ matrix whose (i, j) th element is given by $F_{\text{jump},i,j}(x)$. It is assumed that $\mathbf{F}_{\text{jump}}(\infty) = \mathbf{O}$.

We define $\mathbf{E}_{\text{jump}}(x)$ ($x \geq 0$) and $\mathbf{E}_{\text{jump}}(\infty)$ as

$$\begin{aligned} \mathbf{E}_{\text{jump}}(x) &= \mathbf{E}(x) - \mathbf{E}(0), \\ \mathbf{E}_{\text{jump}} &= \mathbf{E}(\infty) - \mathbf{E}(0), \end{aligned}$$

respectively. Note that these two matrices are nonnegative. We also define $\mathbf{E}_{\text{slope}}$ as an $M \times M$ matrix whose (i, j) th element $\mathbf{E}_{\text{slope}, i, j}$ is given by

$$\mathbf{E}_{\text{slope}, i, j} = \begin{cases} E_{i, j}(0), & \text{if } i \neq j, \\ -\sum_{\nu \in \mathcal{M}} E_{i, \nu}(\infty), & \text{otherwise.} \end{cases}$$

Note that $\mathbf{E}_{\text{slope}} + \mathbf{E}_{\text{jump}}$ is the infinitesimal generator of an irreducible Markov chain. Thus, $\{(X_R(t), J_R(t)); t \geq 0\}$ can be characterized by $\mathbf{E}_{\text{slope}}$, $\mathbf{E}_{\text{jump}}(x)$ and $\mathbf{F}_{\text{jump}}(x)$ as shown in Figure 1.1

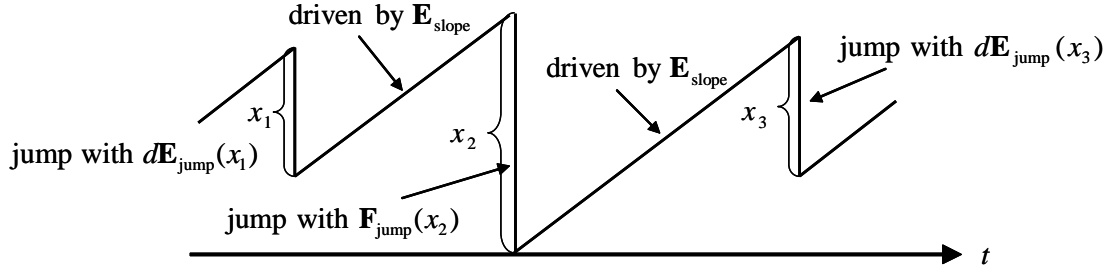


Figure 1.1: The bivariate Markov process with the skip free to the right property.

Proposition 1.3 [Seng89] *A necessary and sufficient condition for the irreducible Markov process $\{(X_R(t), J_R(t)); t \geq 0\}$ to be positive recurrent is*

$$\eta_E \int_0^\infty x d\mathbf{E}_{\text{jump}}(x) \mathbf{e} > 1,$$

where η_E denotes the $1 \times M$ invariant probability vector for $\mathbf{E}_{\text{slope}} + \mathbf{E}_{\text{jump}}$.

In the rest of this subsection, we assume that $\{(X_R(t), J_R(t)); t \geq 0\}$ is positive recurrent and discuss the steady-state density of $\{(X_R(t), J_R(t)); t \geq 0\}$.

Let $\pi_R(x)$ denote a $1 \times M$ vector whose j th element $\pi_{R, j}(x)$ represents

$$\pi_{R, j}(x) dx = \lim_{t \rightarrow \infty} \frac{1}{t} \int_0^t 1\{x \leq X_R(\tau) \leq x + dx, J_R(\tau) = j\} d\tau,$$

where $1\{\chi\}$ denotes an indicator function of event χ . Then, standard flow conservation arguments yield

$$\pi_R(x) = \pi_R(x - \Delta x)(\mathbf{I} + \mathbf{E}_{\text{slope}} \Delta x) + \int_0^\infty \pi_R(x + u) d\mathbf{E}_{\text{jump}}(u) \Delta x + o(\Delta x), \quad x > 0, \quad (1.12)$$

and

$$\pi_R(0) = \int_0^\infty \pi_R(u) d\mathbf{F}_{\text{jump}}(u).$$

Rearranging terms in (1.12), dividing both sides by Δx and taking their limits as $\Delta x \rightarrow 0$, we have

$$\frac{d}{dx} \pi_R(x) = \pi_R(x) \mathbf{E}_{\text{slope}} + \int_0^\infty \pi_R(x + u) d\mathbf{E}_{\text{jump}}(u), \quad x > 0.$$

Further, from the total probability law, we have

$$\int_0^\infty \boldsymbol{\pi}_R(x) \mathbf{e} \, dx = 1.$$

We now introduce an $M \times M$ irreducible matrix $\boldsymbol{\Theta}$, which characterizes $\boldsymbol{\pi}_R(x)$ ($x \geq 0$). We consider the sequence of matrices $\{\boldsymbol{\Theta}_n; n = 0, 1, \dots\}$ with $\boldsymbol{\Theta}_0 = \mathbf{E}_{\text{slope}}$ and

$$\boldsymbol{\Theta}_n = \mathbf{E}_{\text{slope}} + \int_0^\infty \exp(\boldsymbol{\Theta}_{n-1}x) d\mathbf{E}_{\text{jump}}(x), \quad n = 1, 2, \dots$$

It is known [Seng89] that if $\{(X_R(t), J_R(t)); t \geq 0\}$ is positive recurrent, $\boldsymbol{\Theta} = \lim_{n \rightarrow \infty} \boldsymbol{\Theta}_n$ exists and the limit matrix $\boldsymbol{\Theta}$ is the minimal solution of

$$\boldsymbol{\Theta} = \mathbf{E}_{\text{slope}} + \int_0^\infty \exp(\boldsymbol{\Theta}x) d\mathbf{E}_{\text{jump}}(x).$$

Note here that $\boldsymbol{\Theta}$ is non-singular and has strictly negative diagonal elements and nonnegative off-diagonal elements [Seng89].

Proposition 1.4 [Seng89] *The steady-state density $\boldsymbol{\pi}_R(x)$ ($x \geq 0$) is given by*

$$\boldsymbol{\pi}_R(x) = \boldsymbol{\pi}_R(0) \exp(\boldsymbol{\Theta}x), \quad x > 0,$$

where $\boldsymbol{\pi}_R(0)$ satisfies

$$\boldsymbol{\pi}_R(0) = \boldsymbol{\pi}_R(0) \int_0^\infty \exp(\boldsymbol{\Theta}x) \mathbf{F}_{\text{jump}}(x) dx, \quad \boldsymbol{\pi}_R(0)(-\boldsymbol{\Theta})^{-1} \mathbf{e} = 1.$$

We close this subsection by showing an important application of the above result.

We consider the waiting time distribution in a stationary FIFO GI/PH/1 queue. Service times are assumed to be i.i.d. according to a phase-type distribution with representation $(\boldsymbol{\beta}, \mathbf{T})$, where $\boldsymbol{\beta}$ denotes a $1 \times M$ probability vector and \mathbf{T} denotes an $M \times M$ irreducible and defective generator. We define $\boldsymbol{\tau}$ as $\boldsymbol{\tau} = -\mathbf{T}\mathbf{e}$. Let $G(x)$ denote the inter-arrival time distribution. We assume that $G(0+) = 0$. For the stationary GI/PH/1 queue, we define a bivariate process $\{(W(t), S(t)); t \geq 0\}$, where $W(t)$ and $S(t)$ denotes the attained waiting time and the phase of service, respectively, at time t . If no customers are present in the system at time t , $W(t)$ and $S(t)$ are set to zero. Note that the $\{W(t); t \geq 0\}$ increases at a rate of one as long as a customer is in service and jumps downward when a customer in the server departs from the system. The amount of the downward jump is equal to the inter-arrival time between the departing customer and the customer served next to him, unless the departure corresponds to the end of a busy period. Otherwise, $W(t)$ jumps to zero.

We now construct another process by deleting all times when $W(t)$ and $J(t)$ are equal to zero and gluing together all non-zero realizations of $\{(W(t), S(t)); t \geq 0\}$. It is easy to see that the resulting process is equivalent to $\{(X_R(t), J_R(t)); t \geq 0\}$ with

$$\mathbf{E}_{\text{slope}} = \mathbf{T},$$

$$\mathbf{E}_{\text{jump}}(x) = \boldsymbol{\tau} \boldsymbol{\beta} G(x), \tag{1.13}$$

$$\mathbf{F}_{\text{jump}}(x) = \boldsymbol{\tau} \boldsymbol{\beta} (1 - G(x)), \tag{1.14}$$

where $\tilde{\mathbf{T}}$ denote a diagonal matrix whose j th diagonal element is given by the (j, j) th element of \mathbf{T} . Thus, from Proposition 1.4, the steady-state density $\pi_{\mathbf{R}}(x)$ ($x > 0$) is given by

$$\pi_{\mathbf{R}}(x) = -\boldsymbol{\alpha}\boldsymbol{\Theta} \exp(\boldsymbol{\Theta}x), \quad x > 0, \quad (1.15)$$

where $\boldsymbol{\alpha}$ is a invariant probability vector for $(\mathbf{T} + \boldsymbol{\tau}\boldsymbol{\beta})$ and $\boldsymbol{\Theta}$ satisfies

$$\boldsymbol{\Theta} = \mathbf{T} + \int_0^\infty \exp(\boldsymbol{\Theta}x) dG(x) \boldsymbol{\tau}\boldsymbol{\beta}. \quad (1.16)$$

Let t_k denote the time when the k th downward jump happens, i.e., the k th customer departs from the system. Note that the waiting time of $(k + 1)$ th customer is equivalent to $X_{\mathbf{R}}(t_k +)$. It is clear from (1.13) and (1.14) that $X_{\mathbf{R}}(t_k +)$ and $J_{\mathbf{R}}(t_k +)$ are independent. Thus, the waiting time distribution $W_q(x)$ is given by

$$W_q(x) = \int_0^\infty \left(\int_0^{x+u} \frac{\pi_{\mathbf{R}}(y)\boldsymbol{\tau}}{\boldsymbol{\alpha}\boldsymbol{\tau}} dy \right) dG(u), \quad (1.17)$$

Substituting (1.15) and (1.16) into (1.17), we obtain

$$W_q(x) = 1 - \frac{1}{\mu} \exp(\boldsymbol{\Theta}x)(\boldsymbol{\Theta} - \mathbf{T})\mathbf{e},$$

where $\mu = \boldsymbol{\alpha}\boldsymbol{\tau}$ is the reciprocal of the mean service time.

1.3.2 Markov processes skip free to the left

In this subsection, we consider a bivariate Markov process $\{(X_{\mathbf{L}}(t), J_{\mathbf{L}}(t)); t \geq 0\}$ with the skip free to the left property, which is formally defined as follows:

- (a) $X_{\mathbf{L}}(t)$ takes values on $[0, \infty)$ and $J_{\mathbf{L}}(t)$ takes integer values on $\mathcal{M} = \{1, \dots, M\}$. Further, $X_{\mathbf{L}}(t)$ is skip free to the left and decreases at a linear rate of one if there are no upward jumps.
- (b) Given $(X_{\mathbf{L}}(t), J_{\mathbf{L}}(t)) = (x, i)$ ($x > 0, i \in \mathcal{M}$), the bivariate Markov process can change its state to somewhere between $(x + y, j)$ and $(x + y + dy, j)$ ($y \geq 0, j \in \mathcal{M}$) at a rate $dA_{i,j}(y)$. Let $\mathbf{A}(y)$ denote an $M \times M$ matrix whose (i, j) th element is given by $A_{i,j}(y)$. To avoid trivialities, $A_{i,i}(0) = 0$ ($\forall i \in \mathcal{M}$) is assumed. Further, $\mathbf{A}(\infty)$ is assumed to be an irreducible matrix.
- (c) Given $(X_{\mathbf{L}}(t), J_{\mathbf{L}}(t)) = (0, i)$ ($i \in \mathcal{M}$), the bivariate Markov process can change its state to somewhere between (y, j) and $(y + dy, j)$ ($y \geq 0, j \in \mathcal{M}$) with probability $dB_{\text{jump},i,j}(y)$. Let $\mathbf{B}_{\text{jump}}(y)$ denote an $M \times M$ matrix whose (i, j) th element is given by $B_{\text{jump},i,j}(y)$, where $\mathbf{B}_{\text{jump}}(\infty)$ is assumed to be a stochastic matrix. This means that an upward jump happens with probability one when $X_{\mathbf{L}}(t) = 0$.

For later use, we introduce some matrices. We define $\mathbf{A}_{\text{jump}}(x)$ ($x \geq 0$) and \mathbf{A}_{jump} as

$$\begin{aligned} \mathbf{A}_{\text{jump}}(x) &= \mathbf{A}(x) - \mathbf{A}(0), \\ \mathbf{A}_{\text{jump}} &= \mathbf{A}(\infty) - \mathbf{A}(0), \end{aligned}$$

respectively. Note that these two matrices are nonnegative. We also define $\mathbf{A}_{\text{slope}}$ as an $M \times M$ matrix whose (i, j) th element $\mathbf{A}_{\text{slope}, i, j}$ is given by

$$\mathbf{A}_{\text{slope}, i, j} = \begin{cases} A_{i, j}(0), & \text{if } i \neq j, \\ -\sum_{\nu \in \mathcal{M}} A_{i, \nu}(\infty), & \text{otherwise.} \end{cases}$$

Note that $\mathbf{A}_{\text{slope}}$ has strictly negative diagonal elements and nonnegative off-diagonal elements. Note also that $\mathbf{A}_{\text{slope}} + \mathbf{A}_{\text{jump}}$ is an irreducible generator of a Markov chain and hence has the invariant probability vector $\boldsymbol{\eta}_A$. Thus, $\{(X_L(t), J_L(t)); t \geq 0\}$ can be characterized by $\mathbf{A}_{\text{slope}}$, $\mathbf{A}_{\text{jump}}(x)$ and $\mathbf{B}_{\text{jump}}(x)$, as shown in Figure 1.2.

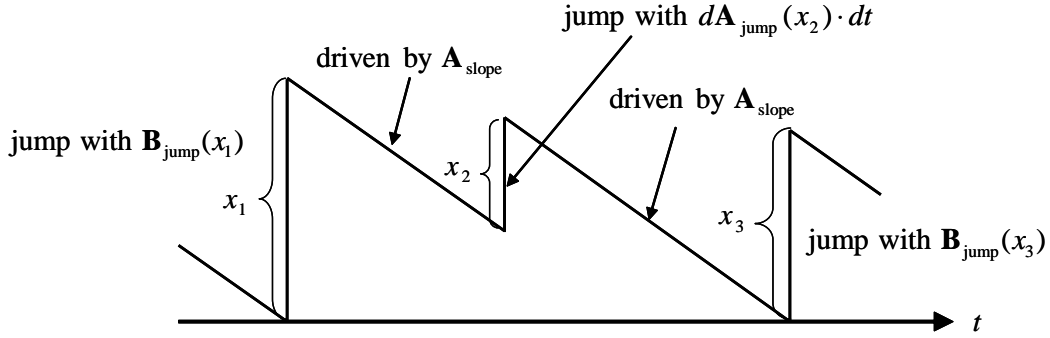


Figure 1.2: The bivariate Markov process with the skip free to the left property.

In what follows, we discuss the steady-state behavior of $\{(X_L(t), J_L(t)); t \geq 0\}$. To do so, we first show a necessary and sufficient condition for its positive recurrence.

Proposition 1.5 [Taki96] *A necessary and sufficient condition for the irreducible Markov process $\{(X_L(t), J_L(t)); t \geq 0\}$ to be positive recurrent is*

$$\boldsymbol{\eta}_A \bar{\mathbf{a}}_{\text{jump}} < \mathbf{1}, \quad \text{and} \quad \left[\int_0^\infty x d\mathbf{B}_{\text{jump}}(x) \right]_i < \infty \quad (i \in \mathcal{M}),$$

where

$$\bar{\mathbf{a}}_{\text{jump}} = \int_0^\infty x d\mathbf{A}_{\text{jump}}(x) \mathbf{e}.$$

Hereafter, we assume that $\{(X_L(t), J_L(t)); t \geq 0\}$ is positive recurrent. Let $\boldsymbol{\pi}_L(x)$ denote a $1 \times M$ vector whose j th element $\pi_{L, j}(x)$ represents

$$\pi_{L, j}(x) dx = \lim_{t \rightarrow \infty} \frac{1}{t} \int_0^t \mathbf{1}\{x \leq X_L(\tau) \leq x + dx, J_L(\tau) = j\} d\tau.$$

We define $\boldsymbol{\Pi}_L(x)$ as

$$\boldsymbol{\Pi}_L(x) = \int_0^x \boldsymbol{\pi}_L(y) dy.$$

Then, standard flow conservation arguments yield

$$\begin{aligned} \mathbf{\Pi}_L(x) &= [\mathbf{\Pi}_L(x + \Delta x) - \mathbf{\Pi}_L(\Delta x)] (\mathbf{I} + \mathbf{A}_{\text{slope}} \Delta x) + \int_0^x \mathbf{\Pi}_L(x + \Delta x - y) d\mathbf{A}_{\text{jump}}(y) \Delta x \\ &\quad + \boldsymbol{\pi}_L(0) \Delta x \mathbf{B}_{\text{jump}}(x) + o(\Delta x). \end{aligned} \quad (1.18)$$

Rearranging terms in (1.18), dividing both sides by Δx and taking their limits as $\Delta x \rightarrow 0$, we have

$$\frac{d}{dx} \mathbf{\Pi}_L(x) + \mathbf{\Pi}_L(x) \mathbf{A}_{\text{slope}} + \int_0^x \mathbf{\Pi}_L(x - y) d\mathbf{A}_{\text{jump}}(y) - \boldsymbol{\pi}_L(0) (\mathbf{I} - \mathbf{B}_{\text{jump}}(x)) = \mathbf{0}.$$

Finally, taking the Laplace-Stieltjes transforms (LSTs) of both sides of the above equation and noting $s\mathbf{I} + \mathbf{A}_{\text{slope}} + \mathbf{A}_{\text{jump}}^*(s)$ is non-singular for $\text{Re}(s) > 0$ [Taki94a], we obtain the following result.

Proposition 1.6 [Taki96] *The LST $\boldsymbol{\pi}_L^*(s)$ of $\boldsymbol{\pi}_L(x)$ is given by*

$$\boldsymbol{\pi}_L^*(s) = \boldsymbol{\pi}_L(0) [\mathbf{I} - \mathbf{B}_{\text{jump}}^*(s)] [s\mathbf{I} + \mathbf{A}_{\text{slope}} + \mathbf{A}_{\text{jump}}^*(s)]^{-1}, \quad \text{Re}(s) > 0, \quad (1.19)$$

where $\mathbf{A}_{\text{jump}}^*(s)$ and $\mathbf{B}_{\text{jump}}^*(s)$ denote the LSTs of $\mathbf{A}_{\text{jump}}(x)$ and $\mathbf{B}_{\text{jump}}(x)$, respectively.

Next, we determine the boundary vector $\boldsymbol{\pi}_L(0)$. To do so, we introduce some notations. Let \mathbf{Q}_L denote an $M \times M$ irreducible matrix representing an infinitesimal generator of a certain Markov chain, which satisfies

$$\mathbf{Q}_L = \mathbf{A}_{\text{slope}} + \int_0^\infty d\mathbf{A}_{\text{jump}}(x) e^{\mathbf{Q}_L x}.$$

Let $\boldsymbol{\kappa}_L$ denote a $1 \times M$ probability vector satisfying $\boldsymbol{\kappa}_L \mathbf{Q}_L = \mathbf{0}$. We define $\mathbf{r}(x)$ as an $M \times 1$ vector whose j th element represents the mean first passage time from state (x, i) to the set of states $\{(0, j); j \in \mathcal{M}\}$. Then $\mathbf{r}(x)$ is given by [Taki94a]

$$\mathbf{r}(x) = (x\mathbf{e}\boldsymbol{\kappa}_L - e^{\mathbf{Q}_L x} + \mathbf{I}) [(e - \bar{\mathbf{a}}_{\text{jump}}) \boldsymbol{\kappa}_L - \mathbf{A}_{\text{slope}} - \mathbf{A}_{\text{jump}}]^{-1} \mathbf{e}.$$

Proposition 1.7 [Taki96] *The boundary vector $\boldsymbol{\pi}_L(0)$ is given by*

$$\boldsymbol{\pi}_L(0) = c \cdot \boldsymbol{\pi}_{L,0},$$

where the $1 \times M$ vector $\boldsymbol{\pi}_{L,0}$ and the constant c are uniquely determined by

$$\boldsymbol{\pi}_{L,0} = \boldsymbol{\pi}_{L,0} \int_0^\infty d\mathbf{B}_{\text{jump}}(x) \exp(\mathbf{Q}_L x), \quad \boldsymbol{\pi}_{L,0} \mathbf{e} = 1,$$

and

$$c = \frac{1}{\boldsymbol{\pi}_{L,0} \int_0^\infty d\mathbf{B}_{\text{jump}}(x) \mathbf{r}(x)},$$

respectively.

We now provide an important application of the above results. We consider the virtual waiting time (i.e., the amount of unfinished work in the system) in a stationary single-server queue with the input process introduced in subsection 1.1.2. Namely, the input process has K batch Markovian arrival streams, which are characterized by \mathbf{C} , $\mathbf{D}(\mathbf{n})$ and K different service time distributions $H_1(x), \dots, H_K(x)$. We define $\overline{\mathbf{D}}(x)$ as

$$\overline{\mathbf{D}}(x) = \sum_{\mathbf{n} \in \mathcal{Z}} \mathbf{D}(\mathbf{n}) H_1^{(n_1)} * \dots * H_K^{(n_K)}(x), \quad (1.20)$$

where $H_k^{(1)}(x) = H_k(x)$ and $H_k^{(n)}(x)$ ($n = 2, 3, \dots$) denotes the n -fold convolution of $H_k(x)$ with itself, and where $H * G(x)$ denotes the convolution of two distribution functions $H(x)$ and $G(x)$ ($x \geq 0$), i.e.,

$$H * G(x) = \int_0^x dH(y)G(x-y).$$

Let $V(t)$ and $S(t)$ denote the virtual waiting time and the state of the underlying Markov chain that governs the input process, respectively, at time t . Let $\mathbf{v}(x)$ denote an $M \times M$ vector whose j th element represents

$$\lim_{t \rightarrow \infty} \frac{1}{t} \int_0^t 1\{V(\tau) \leq x, S(\tau) = j\} d\tau.$$

Also, let $\mathbf{v}^*(s)$ denote the LST of $\mathbf{v}(x)$. We derive $\mathbf{v}^*(s)$ as follows.

We first construct the censored process obtained by observing $\{(V(t), S(t)); t \geq 0\}$ only when $\{V(t) > 0\}$. This censored process is equivalent to $\{(X_L(t), J_L(t)); t \geq 0\}$ with

$$\mathbf{A}_{\text{slope}} = \mathbf{C}, \quad \mathbf{A}_{\text{jump}}(x) = \overline{\mathbf{D}}(x), \quad \mathbf{B}_{\text{jump}}(x) = (-\mathbf{C})^{-1} \overline{\mathbf{D}}(x).$$

Let $\boldsymbol{\kappa}$ denote the invariant probability vector for \mathbf{Q} satisfying

$$\mathbf{Q} = \mathbf{C} + \int_0^\infty d\overline{\mathbf{D}}(x) \exp(\mathbf{Q}x).$$

Note that \mathbf{Q} can be considered as the irreducible generator of the censored Markov chain obtained by observing the underlying Markov chain only when the server is busy [Taki94a]. It then follows from Propositions 1.6 and 1.7 that

$$\boldsymbol{\pi}_L^*(s) = c \boldsymbol{\pi}_{L,0} [\mathbf{I} - (-\mathbf{C})^{-1} \overline{\mathbf{D}}^*(s)] [s\mathbf{I} + \mathbf{C} + \overline{\mathbf{D}}^*(s)]^{-1}, \quad \text{Re}(s) > 0, \quad (1.21)$$

where

$$\boldsymbol{\pi}_{L,0} = \frac{\boldsymbol{\kappa}(-\mathbf{C})}{\boldsymbol{\kappa}(-\mathbf{C})\mathbf{e}}, \quad c = \frac{(1-\rho)\boldsymbol{\kappa}(-\mathbf{C})\mathbf{e}}{\rho}.$$

Note here that $\boldsymbol{\pi}_L^*(s)$ in (1.21) represents the conditional LST of the stationary virtual waiting time distribution, given that the server is busy. Note also that the stationary probability vector of the underlying Markov chain provided that the server is idle is given by

$$\boldsymbol{\kappa} = \frac{\boldsymbol{\pi}_{L,0}(-\mathbf{C})^{-1}}{\boldsymbol{\pi}_{L,0}(-\mathbf{C})^{-1}\mathbf{e}}.$$

Further, since the fraction of time when the server is idle (resp. busy) is given by $1 - \rho$ (resp. ρ), we obtain

$$\mathbf{v}^*(s) = (1 - \rho)\boldsymbol{\kappa} + \rho\boldsymbol{\pi}_L^*(s). \quad (1.22)$$

From (1.21) and (1.22), we obtain the following result.

Proposition 1.8 ([Taki94a]) *$\mathbf{v}^*(s)$ is given by*

$$\mathbf{v}^*(s) \left[s\mathbf{I} + \mathbf{C} + \overline{\mathbf{D}}^*(s) \right] = s(1 - \rho)\boldsymbol{\kappa}, \quad \operatorname{Re}(s) > 0. \quad (1.23)$$

We finally show the recursive formula for derivatives of $\mathbf{v}^*(s)$ evaluated at $s = 0+$. We define $\overline{\mathbf{v}}^{(n)}$ as

$$\overline{\mathbf{v}}^{(n)} = \lim_{s \rightarrow 0+} \frac{(-1)^n}{n!} \frac{d^n}{ds^n} \mathbf{v}^*(s), \quad n = 1, 2, \dots,$$

and $\overline{\mathbf{v}}^{(0)} = \boldsymbol{\pi}$, where $\boldsymbol{\pi}$ is the invariant probability vector for $(\mathbf{C} + \mathbf{D})$.

Proposition 1.9 [Taki94a] *The $\overline{\mathbf{v}}^{(n)}$ is recursively determined by the following recursion:*

$$\begin{aligned} \mathbf{x}^{(1)} &= \left[\boldsymbol{\pi}(\overline{\mathbf{D}}^{(1)} - \mathbf{I}) + (1 - \rho)\boldsymbol{\kappa} \right] (\mathbf{e}\boldsymbol{\pi} - \mathbf{C} - \mathbf{D})^{-1}, \\ \overline{\mathbf{v}}^{(1)}\mathbf{e} &= \frac{1}{1 - \rho} \left[\boldsymbol{\pi}\overline{\mathbf{D}}^{(2)}\mathbf{e} + \mathbf{x}^{(1)}\overline{\mathbf{D}}^{(1)}\mathbf{e} \right], \\ \overline{\mathbf{v}}^{(1)} &= \overline{\mathbf{v}}^{(1)}\mathbf{e}\boldsymbol{\pi} + \mathbf{x}^{(1)}, \end{aligned}$$

and for $n = 2, 3, \dots$,

$$\begin{aligned} \mathbf{x}^{(n)} &= \left[\overline{\mathbf{v}}^{(n-1)}(\overline{\mathbf{D}}^{(1)} - \mathbf{I}) + \sum_{k=2}^n \overline{\mathbf{v}}^{(n-k)}\overline{\mathbf{D}}^{(k)} \right] (\mathbf{e}\boldsymbol{\pi} - \mathbf{C} - \mathbf{D})^{-1}, \\ \overline{\mathbf{v}}^{(n)}\mathbf{e} &= \frac{1}{1 - \rho} \left[\sum_{k=2}^{n+1} \overline{\mathbf{v}}^{(n+1-k)}\overline{\mathbf{D}}^{(k)}\mathbf{e} + \mathbf{x}^{(n)}\overline{\mathbf{D}}^{(1)}\mathbf{e} \right], \\ \overline{\mathbf{v}}^{(n)} &= \overline{\mathbf{v}}^{(n)}\mathbf{e}\boldsymbol{\pi} + \mathbf{x}^{(n)}, \end{aligned}$$

where $\overline{\mathbf{D}}^{(n)}$ is defined as

$$\overline{\mathbf{D}}^{(n)} = \lim_{s \rightarrow 0+} \frac{(-1)^n}{n!} \frac{d^n}{ds^n} \overline{\mathbf{D}}^*(s). \quad n = 1, 2, \dots$$

1.4 Queues with Multiple Arrival Streams

In this section, we survey previous works on queues with multiple non-Poissonian arrival streams having different service time distributions. No papers have analyzed such queues with full generality until Takine and Hasegawa's work [Taki94a], which studied the workload in a MAP/G/1 queue with state-dependent service times. Related studies on the workload can be found in [Asmu91, Zhu91, Mach93, Taki96]. As for the queue length, Regterschot and de Smit [Regt86] derived the probability generating function of the total queue length in the MMPP/G/1-FIFO queue with state-dependent service times. Further, for a MAP/G/1 nonpreemptive priority queue with class-dependent service

times, Takine [Taki99] derived the moments and the mass function of the marginal queue length of each class.

To our best of our knowledge, Takine [Taki01a] performed, for the first time, the full-dress analysis of the joint queue length distribution in queues with multiple non-Poissonian arrival streams having different service time distributions. In [Taki01a], an invariance relationship between the time-average joint queue length distribution and the customer-average joint queue length distribution at departures was shown for queues with the Markovian input process described in subsection 1.1.2. Based on this invariance relationship, a distributional form of Little's law for FIFO queues with simple arrivals (i.e., no batch arrivals happen) was derived. Further, using this law and the result in [Taki94a], recursions to compute the joint queue length distributions in FIFO single-server queues with/without vacations were constructed [Taki01a, Taki01b].

He and Alfa [He98] and Takine [Taki02] studied single-server queues, where customer arrivals follow a MMAP (see Example 1.4) and they are served on a Last-Come-First-Served (LCFS) preemptive basis. In [He98], phase-type service time distributions were assumed. As a result, the queue length process forms a G/M/1-type Markov chain with a tree structure (see Remark 1.1), and hence the stationary joint queue distribution has a matrix product-form solution. On the other hands, in [Taki02], general service time distributions were assumed and the matrix product solution for the stationary joint queue length distribution was also obtained by tracking the dynamics of the string representing classes and remaining service times of respective customers in the system.

As mentioned above, some studies have already been done for single server queues with multiple non-Poissonian arrival streams having different service time distributions. As far as we know, however, no researchers have studied the multi-server case. As for infinite-server queues with non-i.i.d. service times, a few studies have been done in [Smit72, Mach86, Liu91].

1.5 Overview of the Thesis

The main purpose of this thesis is to establish numerical algorithms to compute performance measures (e.g., the queue length distribution and the sojourn time distribution) in queues with batch Markovian arrival streams having different service time distributions.

In chapter 2, we consider a FIFO single-server queue. It had been recognized that the joint queue length distribution for a single-server queue with multiple non-Poissonian arrival streams having different service time distributions is hard to analyze. Recently, for the single arrival case (i.e., no batch arrivals happen), Takine developed recursions to compute the joint queue length distribution by a new approach [Taki01a, Taki01b, Taki01c], which is based on the distributional form of Little's law [Taki01a]. By a very similar approach, we can also establish recursions to compute the joint queue length distribution in the batch arrival case. However, in their implementation, we have to determine several truncations and stopping criteria that are associated with batch arrivals. In addition, the straightforward implementation of those recursions is much more computation-intensive, compared with the single arrival case. The algorithm proposed in chapter 2 solves these problems by efficiently discarding components that make a small contribution to the final results.

Chapter 3 considers a FIFO single-server queue with service interruptions, where the marginal on-off process of the server is a phase-type Markov renewal process, the marginal arrival process is a MMAP (see Example 1.4) with batch arrivals and these two processes are possibly correlated. This model is an extension of one in [Taki97] and it allows batch arrivals and non-i.i.d. service times. To analyze the joint queue length, we first consider a censored process obtained by observing the queue length process only when the server is *on*. This censored process is probabilistically similar to a FIFO single-server queue studied in chapter 2. This fact enables us to establish a numerical algorithm to compute the joint queue length distribution in a very similar way to in chapter 2 [Masu03a].

In chapter 4, we consider an infinite-server queue. Since Takács' work [Taka62], infinite-server queues have been studied by many researchers for forty years. To the best of our knowledge, the infinite-server queue studied by Liu et al. [Liu91] is most general among those studied in the past. In their infinite-server queue, arrival times, batch sizes of arrivals and service time distributions of arriving customers are governed by a Markov renewal process. Although their model is very general, their result is not necessarily suitable for computing performance measures of interest, e.g., the time-dependent mean and variance of the number of customers in the system. Compared with the model studied by Liu et al. [Liu91], the infinite-server queue studied in chapter 4 is more general in the sense that service times in the same batch may be not i.i.d. even though the underlying Markov renewal process is of phase-type. For this infinite-server queue, we derive explicit and numerically feasible formulas to compute the time-dependent and limiting joint binomial moments of the number of customers in the system by considering the time-reversed process of the original queueing process.

Chapter 5 considers processor-sharing queues. In chapter 5, we impose two assumptions on the input process introduced in subsection 1.1.2. First, there exists only one arrival stream. This assumption means that the arrival process is a BMAP (see subsection 1.1.1). Second, service times are i.i.d. according to an exponential distribution. Thus, the processor-sharing queue studied in chapter 5 is denoted by a BMAP/M/1-PS. It is well-known that the stationary queue length distribution in a BMAP/M/1-PS queue are identical to that in a BMAP/M/1-FIFO queue. As for the sojourn time distribution in a BMAP/M/1-PS, however, no studies have been done. In chapter 5, we derive a recursive formula to compute the sojourn time distribution by noting that the sojourn time of a customer is equivalent to the first passage time to an absorbing state in an absorbing Markov chain with a special structure. Further, for a MAP/M/1-PS queue, we develop a numerically feasible algorithm to compute the sojourn time distribution, using the uniformization technique [Tijm94]. The main advantage of the algorithm is that the numerical precision of the result can be given in advance.

Finally, the conclusion and suggestions for future research topics are provided in chapter 6.

The results discussed in chapter 2 is mainly based on [Masu03a], chapter 3 on [Masu03c], chapter 4 on [Masu02] and chapter 5 on [Masu03b].

Chapter 2

FIFO Single-Server Queue

2.1 Introduction

In this chapter, we study the joint queue length distribution in a stationary work-conserving FIFO single-server queue fed by batch Markovian arrival streams. A particular feature of this queue is that service time distributions of customers may be different for different arrival streams.

Most of previous works on FIFO single-server queues with Markovian arrival streams assume that service times of all customers are i.i.d. according to a common distribution function. As a result, the bivariate process of the total number of customers and the state of the Markov chain that governs the arrival process immediately after departures forms a Markov chain of M/G/1 type and the steady-state solution can be computed by well-known M/G/1 paradigm [Neut89].

On the other hand, if service time distributions of customers from respective arrival streams are different from one another, the bivariate process does not have the Markov property [Taki01b], except for queues with a superposition of independent Poisson streams. Thus the queue length analysis of such a queue is not straightforward. Note, however, that the virtual waiting time process in such a queue is characterized by a bivariate Markov process [Asmu91, Taki94a, Taki96, Zhu91], and algorithmic solution methods are known in the literature [Taki94a, Taki96].

Recently, a new approach was developed to characterize the joint queue length distribution in FIFO queues with MMAP (see Example 1.4 in chapter 1) having different service time distributions [Taki01a, Taki01b]. In [Taki01a], the invariant relationship of the joint queue length distributions at a random point in time and at departures was obtained and from this, the distributional form of Little's law was established in [Taki01a]. Further, based on the latter, an algorithmic solution method was developed [Taki01a, Taki01b]. Related works are found in [Mach99, Regt86]. See [Taki01c] for a survey of those developments.

The results in this chapter are considered as an extension of those in [Taki01b], allowing batch arrivals in each arrival stream. Note here that the distributional form of Little's law does not hold for FIFO queues with batch arrivals. Therefore our starting point in analyzing the time-average joint queue length distribution is the invariant relationship of the joint queue length distributions at a random point in time and at departures in [Taki01a]. By doing so, the problem is reduced to find the joint queue length distributions at departures of customers from respective arrival streams.

As you will see, the joint queue length distribution at departures in the FIFO queue is closely related to the virtual waiting time distribution that is readily obtained with the known results. Using these facts, we derive a general formula for the stationary joint queue length distribution at departures in terms of the sojourn time distribution. Further, assuming discrete phase-type batch size distributions, we derive recursions to compute the stationary joint queue length distribution at a random point in time.

The above outline is similar to the single arrival case in [Taki01b]. However, the implementation of some of those recursions is not trivial, because we have to determine several truncation and stopping criteria, which are due to batch arrivals, and their straightforward implementation would require very huge memory space and time-consuming. In this chapter, assuming discrete phase-type batch size distributions, we propose a numerically feasible procedure to compute those recursions, while ensuring the numerical accuracy in the final result. This is the main contribution of this chapter. Note that the procedure is also applicable to FIFO BMAP/G/1 queues with i.i.d. services, when batch sizes follow a discrete phase-type distribution.

The rest of this chapter is divided into six sections. In section 2.2, an input process considered in this chapter is described. In section 2.3, we briefly discuss the virtual and actual waiting time distributions. In section 2.4, we first derive a general formula for the joint queue length distribution, and assuming discrete phase-type batch sizes, we show recursive formulas to compute the joint queue length distribution. In section 2.5, the implementation of the recursions is discussed. In section 2.6, we discuss the efficiency of our algorithm and the qualitative behavior of the queue length through some numerical examples. Finally, concluding remarks are provided in section 2.7.

2.2 Input Process

This chapter imposes the following assumption on the input process introduced in subsection 1.1.2:

Assumption 2.1 *No arrivals from two or more arrival streams happen simultaneously, i.e.,*

$$\begin{aligned} \mathbf{D}(\mathbf{n}) &\geq \mathbf{O}, & \text{if } \mathbf{n} = n_k \mathbf{e}_k, \quad n_k = 1, 2, \dots, \quad k \in \mathcal{K}, \\ \mathbf{D}(\mathbf{n}) &= \mathbf{O}, & \text{otherwise,} \end{aligned}$$

where \mathbf{e}_k is defined in (1.4).

We define $\mathbf{D}_k(n)$ ($n = 1, 2, \dots, k \in \mathcal{K}$) and \mathbf{D}_k ($k \in \mathcal{K}$) as

$$\mathbf{D}_k(n) = \mathbf{D}(n\mathbf{e}_k), \quad \mathbf{D}_k = \sum_{n=1}^{\infty} \mathbf{D}_k(n),$$

respectively. Then, the counting process of arrivals is characterized by \mathbf{C} and $\mathbf{D}_k(n)$ ($k \in \mathcal{K}, n = 1, 2, \dots$). This arrival process is an extension of MMAP (see Example 1.4 in chapter 1) to allow batch arrivals, i.e., a batch MMAP. Roughly speaking, customers arrive in the following way. When a state transition driven by $\mathbf{D}_k(n)$ occurs, n customers of class k arrive simultaneously, and their service times are i.i.d. according to a distribution function $H_k(x)$ with finite mean h_k . On the other

hand, when a state transition driven by \mathbf{C} occurs, no customers arrive. Therefore, the input process is characterized by $(\mathbf{C}, \mathbf{D}_k(n), H_k(x); k \in \mathcal{K})$. Note that the arrival rate λ_k of class k customers is given by

$$\lambda_k = \sum_{n=1}^{\infty} n\pi \mathbf{D}_k(n)\mathbf{e}.$$

In the remainder of this chapter, we assume that utilization factor $\rho = \sum_{k \in \mathcal{K}} \lambda_k h_k < 1$, which ensures that all customers arriving to the system are eventually served [Loyn62].

2.3 Waiting Time Distribution

In this section, we consider the actual waiting time of class k customers in steady state. We define $W_{q,k}(n; m)$ as a generic random variable representing the actual waiting time of a randomly chosen class k customer who is a member of a batch of size n and the m th served customer among members of the same batch. Let $S^{(A_k)}(n)$ denote a generic random variable representing the state of the underlying Markov chain immediately after class k batches of size n arrive. With those, we define $\mathbf{w}_{q,k}(x | n; m)$ as a $1 \times M$ vector whose j th element represents $\Pr[W_{q,k}(n; m) \leq x, S^{(A_k)}(n) = j]$. Note that $\mathbf{w}_{q,k}(x | n; 1)$ ($n \geq 1$) is given by [Taki94a]

$$\mathbf{w}_{q,k}(x | n; 1) = \frac{\mathbf{v}(x)\mathbf{D}_k(n)}{\pi \mathbf{D}_k(n)\mathbf{e}}, \quad (2.1)$$

where $\mathbf{v}(x)$ denotes a $1 \times M$ vector whose j th ($j \in \mathcal{M}$) element represents the stationary joint probability that the total unfinished work is not greater than x and the underlying Markov chain is in state j . Further, for $n \geq 2$,

$$\mathbf{w}_{q,k}(x | n; m) = \int_0^x \mathbf{w}_{q,k}(x - y | n; 1) dH_k^{(m-1)}(y), \quad m = 2, 3, \dots, n. \quad (2.2)$$

Let $W_{q,k}$ and $S^{(A_k)}$ denote generic random variables representing the actual waiting time of class k customers and the state of the underlying Markov chain immediately after arrivals of class k batches, respectively. We then define $\mathbf{w}_{q,k}(x)$ as a $1 \times M$ vector whose j th element represents $\Pr[W_{q,k} \leq x, S^{(A_k)} = j]$. Because a randomly chosen customer of class k is a member of a batch of size n with probability $n\pi \mathbf{D}_k(n)\mathbf{e}/\lambda_k$, we have

$$\mathbf{w}_{q,k}(x) = \sum_{n=1}^{\infty} \frac{n\pi \mathbf{D}_k(n)\mathbf{e}}{\lambda_k} \cdot \frac{1}{n} \sum_{m=1}^n \mathbf{w}_{q,k}(x | n; m). \quad (2.3)$$

Let $\mathbf{w}_{q,k}^*(s)$ denote the LST of $\mathbf{w}_{q,k}(x)$. From (2.1)–(2.3), we have

$$\mathbf{w}_{q,k}^*(s) = \sum_{n=1}^{\infty} \frac{\mathbf{v}^*(s)\mathbf{D}_k(n)}{\lambda_k} \sum_{m=1}^n \{H_k^*(s)\}^{m-1}, \quad \text{Re}(s) > 0, \quad (2.4)$$

where $\mathbf{v}^*(s)$ and $H_k^*(s)$ denote the LSTs of $\mathbf{v}(x)$ and $H_k(x)$, respectively. Note here that $\mathbf{v}^*(s)$ satisfies (1.23), where

$$\overline{\mathbf{D}}^*(s) = \sum_{k \in \mathcal{K}} \sum_{n=1}^{\infty} \mathbf{D}_k(n) \{H_k^*(s)\}^n. \quad (2.5)$$

From (2.4), we can readily obtain the following result.

Theorem 2.1 $\mathbf{w}_{q,k}^*(s)$ ($k \in \mathcal{K}$) is given by

$$\mathbf{w}_{q,k}^*(s) = \frac{\mathbf{v}^*(s)(\mathbf{D}_k - \mathbf{D}_k^*(H_k^*(s)))}{\lambda_k(1 - H_k^*(s))}, \quad \text{Re}(s) > 0,$$

where

$$\mathbf{D}_k^*(z_k) = \sum_{n=1}^{\infty} z_k^n \mathbf{D}_k(n). \quad (2.6)$$

2.4 Joint Queue Length Distribution

This section considers the joint queue length distribution. In subsection 2.4.1, we apply a general relationship between the time-average queue length distribution and the queue length distributions at departures of customers of respective classes [Taki01a] to our specific queue. Then the problem is reduced to characterize the joint queue length distributions at departures of respective classes, which is discussed in subsection 2.4.2. Finally in subsection 2.4.3, assuming discrete phase-type batch size distributions, we derive recursions for some quantities required in computing the joint queue length distribution.

2.4.1 Relationship in the joint queue length distributions

Let N_k ($k \in \mathcal{K}$) and S denote generic random variables representing the number of class k customers and the state of the underlying Markov chain, respectively, in steady state. We define $\mathbf{p}(n_1, \dots, n_K)$ as a $1 \times M$ vector whose j th element represents $\Pr[N_1 = n_1, \dots, N_K = n_K, S = j]$. For simplicity, let \mathbf{n} and \mathbf{z} denote a $1 \times K$ nonnegative integer vector (n_1, \dots, n_K) and a $1 \times K$ complex vector (z_1, \dots, z_K) , respectively. Further we define \mathcal{Z} as

$$\mathcal{Z} = \{(n_1, \dots, n_K); n_k = 0, 1, \dots, \text{ for all } k \in \mathcal{K}\}.$$

We then define $\mathbf{p}^*(\mathbf{z})$ as

$$\mathbf{p}^*(\mathbf{z}) = \sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \mathbf{p}(\mathbf{n}), \quad |z_k| \leq 1 \text{ for all } k \in \mathcal{K}.$$

Note that $\mathbf{p}^*(\mathbf{z})$ denotes the vector generating function of the joint queue length distribution in steady state.

Let $N_\nu^{(D_k)}$ and $S^{(D_k)}$ ($k, \nu \in \mathcal{K}$) denote generic random variables representing the number of class ν customers and the state of the underlying Markov chain, respectively, immediately after departures of class k customers in steady state. We then define $\mathbf{q}_k(\mathbf{n})$ ($k \in \mathcal{K}$, $\mathbf{n} \in \mathcal{Z}$) as a $1 \times M$ vector whose j th element represents $\Pr[N_1^{(D_k)} = n_1, \dots, N_K^{(D_k)} = n_K, S^{(D_k)} = j]$. Further we define $\mathbf{q}_k^*(\mathbf{z})$ ($k \in \mathcal{K}$) as

$$\mathbf{q}_k^*(\mathbf{z}) = \sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \mathbf{q}_k(\mathbf{n}), \quad |z_k| \leq 1 \text{ for all } k \in \mathcal{K}.$$

Note that $\mathbf{q}_k^*(\mathbf{z})$ denotes the vector generating function of the joint queue length distribution immediately after departures of class k customers. Thus, applying Theorem 1 in [Taki01a] to our model, we have the following theorem.

Theorem 2.2 ([Taki01a]) $\mathbf{p}^*(\mathbf{z})$ and $\mathbf{q}_k^*(\mathbf{z})$ are related by

$$\mathbf{p}^*(\mathbf{z}) \left[\mathbf{C} + \sum_{k \in \mathcal{K}} \mathbf{D}_k^*(z_k) \right] = \sum_{k \in \mathcal{K}} \lambda_k (z_k - 1) \mathbf{q}_k^*(\mathbf{z}), \quad k \in \mathcal{K}, \quad (2.7)$$

where $\mathbf{D}_k^*(z_k)$ is given in (2.6).

Further, comparing the coefficient vectors of $z_1^{n_1} \cdots z_K^{n_K}$ on both sides of (2.7), we obtain the following result.

Corollary 2.1 The $\mathbf{p}(\mathbf{n})$ ($\mathbf{n} \in \mathcal{Z}$) is recursively determined by the following way.

$$\mathbf{p}(\mathbf{0}) = \sum_{k \in \mathcal{K}} \lambda_k \mathbf{q}_k(\mathbf{0}) (-\mathbf{C})^{-1}, \quad (2.8)$$

and for $\mathcal{Z}^+ = \mathcal{Z} - \{\mathbf{0}\}$,

$$\mathbf{p}(\mathbf{n}) = \sum_{k \in \mathcal{K}} \left[\lambda_k (\mathbf{q}_k(\mathbf{n}) - \mathbf{q}_k(\mathbf{n} - \mathbf{e}_k)) + \sum_{m_k=1}^{n_k} \mathbf{p}(\mathbf{n} - m_k \mathbf{e}_k) \mathbf{D}_k(m_k) \right] (-\mathbf{C})^{-1}, \quad (2.9)$$

where $\mathbf{q}_k(\mathbf{n}) = \mathbf{0}$ for $\mathbf{n} \notin \mathcal{Z}$ and \mathbf{e}_k ($k \in \mathcal{K}$) is given in (1.4).

This recursion is an extension of the recursion (B.3) and (B.4) for the BMAP/GI/1 queue.

Remark 2.1 If the service discipline is independent of the state of the underlying Markov chain, Theorem 2.2 holds for stationary queues fed by Markovian input process $(\mathbf{C}, \mathbf{D}_k(n), H_k(x); k \in \mathcal{K})$, despite the number of servers, the service discipline and the service time distributions of respective classes. The details are referred to [Taki01a].

2.4.2 Joint queue length distribution immediately after departures

In this subsection, we consider the vector generating function of the joint queue length distribution immediately after departures of each class. We denote, by $C_k(n; m)$ ($k \in \mathcal{K}$, $n = 1, 2, \dots$, $m = 1, 2, \dots, n$), a randomly chosen class k customer who is a member of a batch of size n and the m th served customer among members of the same batch. Let $N_\nu^{(D_k)}(n; m)$ and $S^{(D_k)}(n; m)$ ($k, \nu \in \mathcal{K}$, $n = 1, 2, \dots$, $m = 1, 2, \dots, n$) denote generic random variables representing the number of class ν customers and the state of the underlying Markov chain, respectively, immediately after the departure of customer $C_k(n; m)$ in steady state. We then define $\mathbf{q}_k^*(\mathbf{z} | n; m)$ ($k \in \mathcal{K}$, $n = 1, 2, \dots$, $m = 1, 2, \dots, n$) as a $1 \times M$ vector whose j th element represents

$$\mathbb{E} \left[\prod_{\nu \in \mathcal{K}} z_\nu^{N_\nu^{(D_k)}(n; m)} \mathbf{1}\{S^{(D_k)}(n; m) = j\} \right],$$

where $\mathbf{1}\{\chi\}$ denotes an indicator function of event χ . Because a randomly chosen customer of class k is a member of a batch of size n with probability $n\pi \mathbf{D}_k(n) \mathbf{e} / \lambda_k$, we have

$$\mathbf{q}_k^*(\mathbf{z}) = \sum_{n=1}^{\infty} \frac{n\pi \mathbf{D}_k(n) \mathbf{e}}{\lambda_k} \frac{1}{n} \sum_{m=1}^n \mathbf{q}_k^*(\mathbf{z} | n; m). \quad (2.10)$$

In what follows, we consider $\mathbf{q}_k^*(\mathbf{z} | n; m)$.

We define $W_k(n; m)$ ($k \in \mathcal{K}$, $n = 1, 2, \dots$, $m = 1, 2, \dots, n$) as a generic random variable representing the sojourn time of customer $C_k(n; m)$. Note here that

$$W_k(n; m) = W_{q,k}(n; 1) + H_{k,1} + \dots + H_{k,m},$$

where $W_{q,k}(n; 1)$ denotes the actual waiting time of customer $C_k(n; 1)$, and $H_{k,l}$ ($l = 1, 2, \dots, m$) denotes the service time of customer $C_k(n; l)$. By definition, $W_{q,k}(n; 1)$ depends only on the past history up to the arrival instant of a batch including customer $C_k(n; m)$. On the other hand, the number of customers in the system immediately after the departure of customer $C_k(n; m)$ is equal to the sum of the $n - m$ customers in the same batch and customers who arrived during the sojourn time of customer $C_k(n; m)$. Note here that the latter is conditionally independent of the past history given the length of the sojourn time and the state of the underlying Markov chain immediately after the arrival of the batch. Thus we have

$$\mathbf{q}_k^*(\mathbf{z} | n; m) = z_k^{n-m} \int_0^\infty d\mathbf{w}_{q,k}(x | n; 1) \mathbf{N}^*(x, \mathbf{z}) \left[\int_0^\infty dH_k(y) \mathbf{N}^*(y, \mathbf{z}) \right]^m, \quad (2.11)$$

where

$$\mathbf{N}^*(x, \mathbf{z}) = \exp \left[\left(\mathbf{C} + \sum_{k \in \mathcal{K}} \mathbf{D}_k^*(z_k) \right) x \right]. \quad (2.12)$$

Theorem 2.3 *The vector generating function $\mathbf{q}_k^*(\mathbf{z})$ ($k \in \mathcal{K}$) of the joint queue length distribution immediately after departures of class k customers in the steady state is given by*

$$\mathbf{q}_k^*(\mathbf{z}) = \frac{1}{\lambda_k} \sum_{m=1}^\infty \sum_{l=0}^\infty z_k^l \int_0^\infty d\mathbf{v}(x) \mathbf{D}_k(m+l) \mathbf{N}^*(x, \mathbf{z}) \left[\int_0^\infty dH_k(y) \mathbf{N}^*(y, \mathbf{z}) \right]^m. \quad (2.13)$$

Proof. Using (2.1), (2.10) and (2.11), we have

$$\begin{aligned} \mathbf{q}_k^*(\mathbf{z}) &= \sum_{n=1}^\infty \frac{\pi \mathbf{D}_k(n) \mathbf{e}}{\lambda_k} \sum_{m=1}^n z_k^{n-m} \int_0^\infty \frac{d\mathbf{v}(x) \mathbf{D}_k(n)}{\pi \mathbf{D}_k(n) \mathbf{e}} \mathbf{N}^*(x, \mathbf{z}) \left[\int_0^\infty dH_k(y) \mathbf{N}^*(y, \mathbf{z}) \right]^m \\ &= \frac{1}{\lambda_k} \sum_{n=1}^\infty \sum_{m=1}^n z_k^{n-m} \int_0^\infty d\mathbf{v}(x) \mathbf{D}_k(n) \mathbf{N}^*(x, \mathbf{z}) \left[\int_0^\infty dH_k(y) \mathbf{N}^*(y, \mathbf{z}) \right]^m, \end{aligned}$$

from which (2.13) follows. ■

2.4.3 Recursions for discrete phase-type batch sizes

In this subsection, we develop a recursive formula to compute the vector mass function $\mathbf{q}_k(\mathbf{n})$ of the joint queue length immediately after departures of each class under the following assumption.

Assumption 2.2 *Let α_k denote a $1 \times M_k$ probability vector and \mathbf{P}_k denote an $M_k \times M_k$ substochastic matrix. Then, batch sizes of class k are i.i.d. according to a discrete phase-type distribution with representation (α_k, \mathbf{P}_k) , i.e.,*

$$\mathbf{D}_k(n) = g_k(n) \mathbf{D}_k, \quad (2.14)$$

$$g_k(n) = \alpha_k \mathbf{P}_k^{n-1} (\mathbf{I} - \mathbf{P}_k) \mathbf{e}, \quad n = 1, 2, \dots, \quad (2.15)$$

where $\mathbf{I}(m)$ denotes an $m \times m$ identity matrix (the size m is suppressed when it is clear from the context).

We now define $\mathbf{I}(m)$ as an $m \times m$ identity matrix. When the size of an identity matrix is clear from the context, we suppress (m) .

Lemma 2.1 Under Assumption 2.2, $\mathbf{q}_k^*(\mathbf{z})$ ($k \in \mathcal{K}$) is given by

$$\begin{aligned} \mathbf{q}_k^*(\mathbf{z}) &= \frac{1}{\lambda_k} \int_0^\infty dv(x) \mathbf{D}_k \mathbf{N}^*(x, \mathbf{z}) \left(\boldsymbol{\alpha}_k \otimes \int_0^\infty dH_k(y) \mathbf{N}^*(y, \mathbf{z}) \right) \\ &\quad \cdot \left[\mathbf{I} - \mathbf{P}_k \otimes \int_0^\infty dH_k(y) \mathbf{N}^*(y, \mathbf{z}) \right]^{-1} \left[\{(\mathbf{I} - z_k \mathbf{P}_k)^{-1} (\mathbf{I} - \mathbf{P}_k) \mathbf{e}\} \otimes \mathbf{I}(M) \right]. \end{aligned} \quad (2.16)$$

Proof. Substituting (2.14) and (2.15) into (2.13) and using properties of Kronecker product:

$$a\mathbf{X} = a \otimes \mathbf{X} \text{ for any scalar } a, \text{ and,}$$

$$\begin{aligned} (\mathbf{X}_1 \cdots \mathbf{X}_n) \otimes (\mathbf{Y}_1 \cdots \mathbf{Y}_n) \\ = (\mathbf{X}_1 \otimes \mathbf{Y}_1) \cdots (\mathbf{X}_n \otimes \mathbf{Y}_n) \text{ for any } n = 1, 2, \dots, \end{aligned}$$

we obtain

$$\begin{aligned} \mathbf{q}_k^*(\mathbf{z}) &= \frac{1}{\lambda_k} \int_0^\infty dv(x) \mathbf{D}_k \mathbf{N}^*(x, \mathbf{z}) \cdot \left\{ \boldsymbol{\alpha}_k \sum_{m=1}^\infty \mathbf{P}_k^{m-1} \sum_{l=0}^\infty z_k^l \mathbf{P}_k^l (\mathbf{I} - \mathbf{P}_k) \mathbf{e} \right\} \\ &\quad \otimes \left\{ \left(\int_0^\infty dH_k(y) \mathbf{N}^*(y, \mathbf{z}) \right)^m \right\} \\ &= \frac{1}{\lambda_k} \int_0^\infty dv(x) \mathbf{D}_k \mathbf{N}^*(x, \mathbf{z}) \cdot \left\{ \boldsymbol{\alpha}_k \sum_{m=1}^\infty \mathbf{P}_k^{m-1} (\mathbf{I} - z_k \mathbf{P}_k)^{-1} (\mathbf{I} - \mathbf{P}_k) \mathbf{e} \right\} \\ &\quad \otimes \left\{ \left(\int_0^\infty dH_k(y) \mathbf{N}^*(y, \mathbf{z}) \right)^m \right\} \\ &= \frac{1}{\lambda_k} \int_0^\infty dv(x) \mathbf{D}_k \mathbf{N}^*(x, \mathbf{z}) \cdot \left(\boldsymbol{\alpha}_k \otimes \int_0^\infty dH_k(y) \mathbf{N}^*(y, \mathbf{z}) \right) \\ &\quad \cdot \sum_{m=1}^\infty \left(\mathbf{P}_k \otimes \int_0^\infty dH_k(y) \mathbf{N}^*(y, \mathbf{z}) \right)^{m-1} \left[\{(\mathbf{I} - z_k \mathbf{P}_k)^{-1} (\mathbf{I} - \mathbf{P}_k) \mathbf{e}\} \otimes \mathbf{I}(M) \right] \\ &= \frac{1}{\lambda_k} \int_0^\infty dv(x) \mathbf{D}_k \mathbf{N}^*(x, \mathbf{z}) \cdot \left(\boldsymbol{\alpha}_k \otimes \int_0^\infty dH_k(y) \mathbf{N}^*(y, \mathbf{z}) \right) \\ &\quad \cdot \left[\mathbf{I} - \mathbf{P}_k \otimes \int_0^\infty dH_k(y) \mathbf{N}^*(y, \mathbf{z}) \right]^{-1} \left[\{(\mathbf{I} - z_k \mathbf{P}_k)^{-1} (\mathbf{I} - \mathbf{P}_k) \mathbf{e}\} \otimes \mathbf{I}(M) \right], \end{aligned}$$

which completes the proof. ■

We define $\mathbf{v}_k(\mathbf{n})$ ($k \in \mathcal{K}$, $\mathbf{n} \in \mathcal{Z}$) as a $1 \times M$ vector satisfying

$$\sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \mathbf{v}_k(\mathbf{n}) = \int_0^\infty dv(x) \mathbf{D}_k \mathbf{N}^*(x, \mathbf{z}). \quad (2.17)$$

We also define $\mathbf{A}_k(\mathbf{n})$ and $\mathbf{\Gamma}_k(\mathbf{n})$ ($k \in \mathcal{K}$, $\mathbf{n} \in \mathcal{Z}$) as $M \times M$ and $MM_k \times MM_k$ matrices satisfying

$$\sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \mathbf{A}_k(\mathbf{n}) = \int_0^\infty dH_k(y) \mathbf{N}^*(y, \mathbf{z}), \quad (2.18)$$

$$\sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \mathbf{\Gamma}_k(\mathbf{n}) = \left[\mathbf{I} - \mathbf{P}_k \otimes \int_0^\infty dH_k(y) \mathbf{N}^*(y, \mathbf{z}) \right]^{-1}, \quad (2.19)$$

respectively. Note here that

$$\left\{ (\mathbf{I} - z_k \mathbf{P}_k)^{-1} (\mathbf{I} - \mathbf{P}_k) \mathbf{e} \right\} \otimes \mathbf{I}(M) = \sum_{m=0}^{\infty} z_k^m \{ \mathbf{P}_k^m (\mathbf{I} - \mathbf{P}_k) \mathbf{e} \} \otimes \mathbf{I}(M).$$

Thus (2.16) is rewritten to be

$$\begin{aligned} \mathbf{q}_k^*(\mathbf{z}) &= \frac{1}{\lambda_k} \sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \sum_{m=0}^{n_k} \sum_{\substack{\mathbf{n}_1 + \mathbf{n}_2 + \mathbf{n}_3 \\ = \mathbf{n} - m \mathbf{e}_k}} \mathbf{v}_k(\mathbf{n}_1) \\ &\quad \cdot [\boldsymbol{\alpha}_k \otimes \mathbf{A}_k(\mathbf{n}_2)] \boldsymbol{\Gamma}_k(\mathbf{n}_3) [\{ \mathbf{P}_k^m (\mathbf{I} - \mathbf{P}_k) \mathbf{e} \} \otimes \mathbf{I}(M)], \end{aligned} \quad (2.20)$$

where $\mathbf{n}_j \in \mathcal{Z}$ for $j = 1, 2, 3$. Comparing coefficient vectors of $z_1^{n_1} \cdots z_K^{n_K}$ on both sides of (2.20), we obtain the following result.

Theorem 2.4 *Under Assumption 2.2, $\mathbf{q}_k(\mathbf{n})$ ($k \in \mathcal{K}$, $\mathbf{n} \in \mathcal{Z}$) is given by*

$$\mathbf{q}_k(\mathbf{n}) = \frac{1}{\lambda_k} \sum_{m=0}^{n_k} \sum_{\substack{\mathbf{n}_1 + \mathbf{n}_2 + \mathbf{n}_3 \\ = \mathbf{n} - m \mathbf{e}_k}} \mathbf{v}_k(\mathbf{n}_1) [\boldsymbol{\alpha}_k \otimes \mathbf{A}_k(\mathbf{n}_2)] \boldsymbol{\Gamma}_k(\mathbf{n}_3) [\mathbf{P}_k^m (\mathbf{I} - \mathbf{P}_k) \mathbf{e} \otimes \mathbf{I}(M)],$$

where $\mathbf{n}_j \in \mathcal{Z}$ for $j = 1, 2, 3$.

Theorem 2.4 implies that the computation of $\mathbf{q}_k(\mathbf{n})$ is reduced to those of $\mathbf{v}_k(\mathbf{n})$, $\mathbf{A}_k(\mathbf{n})$ and $\boldsymbol{\Gamma}_k(\mathbf{n})$, which are discussed in the rest of this subsection.

We first consider the $\mathbf{A}_k(\mathbf{n})$. Let θ denote the maximum absolute value of diagonal elements of \mathbf{C} . We define $\mathbf{F}_m(\mathbf{n})$ ($m = 0, 1, \dots$, $\mathbf{n} \in \mathcal{Z}$) as an $M \times M$ matrix that satisfies

$$\sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \mathbf{F}_m(\mathbf{n}) = \left[\mathbf{I} + \theta^{-1} \left(\mathbf{C} + \sum_{k \in \mathcal{K}} \mathbf{D}_k^*(z_k) \right) \right]^m. \quad (2.21)$$

Lemma 2.2 *$\mathbf{A}_k(\mathbf{n})$ is given by*

$$\mathbf{A}_k(\mathbf{n}) = \sum_{m=0}^{\infty} \gamma_k^{(m)}(\theta) \mathbf{F}_m(\mathbf{n}), \quad k \in \mathcal{K}, \mathbf{n} \in \mathcal{Z}, \quad (2.22)$$

where

$$\gamma_k^{(m)}(\theta) = \int_0^{\infty} e^{-\theta y} \frac{(\theta y)^m}{m!} dH_k(y), \quad k \in \mathcal{K}, m = 0, 1, \dots, \quad (2.23)$$

and $\mathbf{F}_m(\mathbf{n})$'s are recursively determined by

$$\mathbf{F}_0(\mathbf{n}) = \begin{cases} \mathbf{I}, & \text{if } \mathbf{n} = \mathbf{0}, \\ \mathbf{O}, & \text{otherwise,} \end{cases} \quad (2.24)$$

and for $m = 0, 1, \dots$,

$$\mathbf{F}_{m+1}(\mathbf{n}) = \mathbf{F}_m(\mathbf{n})(\mathbf{I} + \theta^{-1} \mathbf{C}) + \theta^{-1} \sum_{k \in \mathcal{K}} \sum_{l_k=1}^{n_k} \mathbf{F}_m(\mathbf{n} - l_k \mathbf{e}_k) \mathbf{D}_k(l_k), \quad \mathbf{n} \in \mathcal{Z}. \quad (2.25)$$

Proof. From (2.12), (2.18) and (2.23), we obtain

$$\begin{aligned}
& \sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \mathbf{A}_k(\mathbf{n}) \\
&= \int_0^\infty dH_k(y) \exp \left[\left(\mathbf{C} + \sum_{k \in \mathcal{K}} \mathbf{D}_k^*(z_k) \right) y \right] \\
&= \sum_{m=0}^\infty \int_0^\infty e^{-\theta y} \frac{(\theta y)^m}{m!} dH_k(y) \left[\mathbf{I} + \theta^{-1} \left(\mathbf{C} + \sum_{k \in \mathcal{K}} \mathbf{D}_k^*(z_k) \right) \right]^m \\
&= \sum_{m=0}^\infty \gamma_k^{(m)}(\theta) \left[\mathbf{I} + \theta^{-1} \left(\mathbf{C} + \sum_{k \in \mathcal{K}} \mathbf{D}_k^*(z_k) \right) \right]^m. \tag{2.26}
\end{aligned}$$

Substituting (2.21) into (2.26) and changing the order of summations, we have

$$\sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \mathbf{A}_k(\mathbf{n}) = \sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \sum_{m=0}^\infty \gamma_k^{(m)}(\theta) \mathbf{F}_m(\mathbf{n}). \tag{2.27}$$

Comparing the coefficient matrices of $z_1^{n_1} \cdots z_K^{n_K}$ on both sides of (2.27), we obtain (2.22).

The remaining is to show (2.24) and (2.25). From (2.6) and (2.21), we have for $m = 0, 1, \dots$,

$$\begin{aligned}
& \sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \mathbf{F}_{m+1}(\mathbf{n}) \\
&= \sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \mathbf{F}_m(\mathbf{n}) \left[\mathbf{I} + \theta^{-1} \left(\mathbf{C} + \sum_{k \in \mathcal{K}} \mathbf{D}_k^*(z_k) \right) \right] \\
&= \sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \mathbf{F}_m(\mathbf{n}) \left(\mathbf{I} + \theta^{-1} \mathbf{C} \right) \\
&\quad + \sum_{\mathbf{n} \in \mathcal{Z}} \sum_{k \in \mathcal{K}} \left[\sum_{l_k=1}^\infty z_1^{n_1} \cdots z_{k-1}^{n_{k-1}} z_k^{n_k+l_k} z_{k+1}^{n_{k+1}} \cdots z_K^{n_K} \mathbf{F}_m(\mathbf{n}) \cdot \theta^{-1} \mathbf{D}_k(l_k) \right] \\
&= \sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \mathbf{F}_m(\mathbf{n}) \left(\mathbf{I} + \theta^{-1} \mathbf{C} \right) \\
&\quad + \sum_{\mathbf{n} \in \mathcal{Z}^+} z_1^{n_1} \cdots z_K^{n_K} \theta^{-1} \sum_{k \in \mathcal{K}} \sum_{l_k=1}^{n_k} \mathbf{F}_m(\mathbf{n} - l_k \mathbf{e}_k) \mathbf{D}_k(l_k).
\end{aligned}$$

Comparing the coefficient vectors of $z_1^{n_1} \cdots z_K^{n_K}$ on both sides of the above equation, we obtain (2.25). Finally, (2.24) is clear from the definition. \blacksquare

Next we consider the $\mathbf{\Gamma}_k(\mathbf{n})$ in (2.19).

Lemma 2.3 $\mathbf{\Gamma}_k(\mathbf{n})$ ($k \in \mathcal{K}$, $\mathbf{n} \in \mathcal{Z}$) is determined by the following recursion:

$$\begin{aligned}
\mathbf{\Gamma}_k(\mathbf{0}) &= [\mathbf{I} - \mathbf{P}_k \otimes \mathbf{A}_k(\mathbf{0})]^{-1}, \\
\mathbf{\Gamma}_k(\mathbf{n}) &= \sum_{\substack{\mathbf{0} \leq \mathbf{l} \leq \mathbf{n} \\ \mathbf{l} \neq \mathbf{0}}} \mathbf{\Gamma}_k(\mathbf{n} - \mathbf{l}) [\mathbf{P}_k \otimes \mathbf{A}_k(\mathbf{l})] \mathbf{\Gamma}_k(\mathbf{0}), \quad \mathbf{n} \in \mathcal{Z}^+.
\end{aligned}$$

Proof. Note first that (2.19) is equivalent to

$$\sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \Gamma_k(\mathbf{n}) \left[\mathbf{I} - \mathbf{P}_k \otimes \int_0^\infty dH_k(y) \mathbf{N}^*(y, \mathbf{z}) \right] = \mathbf{I}.$$

Substituting (2.18) into the above equation, we have

$$\sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \Gamma_k(\mathbf{n}) \left[\mathbf{I} - \mathbf{P}_k \otimes \sum_{\mathbf{l} \in \mathcal{Z}} z_1^{l_1} \cdots z_K^{l_K} \mathbf{A}_k(\mathbf{l}) \right] = \mathbf{I},$$

from which it follows that

$$\sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \Gamma_k(\mathbf{n}) - \sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \sum_{\mathbf{0} \leq \mathbf{l} \leq \mathbf{n}} \Gamma_k(\mathbf{n} - \mathbf{l}) [\mathbf{P}_k \otimes \mathbf{A}_k(\mathbf{l})] = \mathbf{I}.$$

Comparing the coefficient matrices of $z_1^{n_1} \cdots z_K^{n_K}$ on both sides of the above equation, we have

$$\begin{aligned} \Gamma_k(\mathbf{0}) - \Gamma_k(\mathbf{0}) [\mathbf{P}_k \otimes \mathbf{A}_k(\mathbf{0})] &= \mathbf{I}, \\ \Gamma_k(\mathbf{n}) - \sum_{\mathbf{0} \leq \mathbf{l} \leq \mathbf{n}} \Gamma_k(\mathbf{n} - \mathbf{l}) [\mathbf{P}_k \otimes \mathbf{A}_k(\mathbf{l})] &= \mathbf{O}, \quad \mathbf{n} \in \mathcal{Z}^+, \end{aligned}$$

or equivalently,

$$\Gamma_k(\mathbf{0}) = [\mathbf{I} - \mathbf{P}_k \otimes \mathbf{A}_k(\mathbf{0})]^{-1},$$

and for $\mathbf{n} \in \mathcal{Z}^+$,

$$\Gamma_k(\mathbf{n}) = \sum_{\substack{\mathbf{0} \leq \mathbf{l} \leq \mathbf{n} \\ \mathbf{l} \neq \mathbf{0}}} \Gamma_k(\mathbf{n} - \mathbf{l}) [\mathbf{P}_k \otimes \mathbf{A}_k(\mathbf{l})] [\mathbf{I} - \mathbf{P}_k \otimes \mathbf{A}_k(\mathbf{0})]^{-1},$$

from which Lemma 2.3 follows. ■

Finally, we consider the $\mathbf{v}_k(\mathbf{n})$ in (2.17). In a very similar way to derive (2.22), we obtain the following lemma.

Lemma 2.4 $\mathbf{v}_k(\mathbf{n})$ ($k \in \mathcal{K}$, $\mathbf{n} \in \mathcal{Z}$) is given by

$$\mathbf{v}_k(\mathbf{n}) = \sum_{m=0}^{\infty} \mathbf{v}^{(m)}(\theta) \mathbf{D}_k \mathbf{F}_m(\mathbf{n}), \quad (2.28)$$

where $\mathbf{F}_m(\mathbf{n})$ is given in (2.24) and (2.25), and

$$\mathbf{v}^{(m)}(\theta) = \int_0^\infty e^{-\theta x} \frac{(\theta x)^m}{m!} d\mathbf{v}(x), \quad m = 0, 1, \dots$$

Thus $\mathbf{v}_k(\mathbf{n})$ is given in terms of the $\mathbf{v}^{(m)}(\theta)$. Because the computation of the $\mathbf{v}^{(m)}(\theta)$ has already been studied in [Taki01b], we summarize the result below. As for the details, readers are referred to Lemma 3 in [Taki01b].

Note first that

$$\sum_{m=0}^{\infty} z^m \mathbf{v}^{(m)}(\theta) = \mathbf{v}^*(\theta - \theta z), \quad (2.29)$$

where $\mathbf{v}^*(s)$ is given in (1.23). Thus, substituting $\theta - \theta z$ for s in (1.23) and using (2.29) yield

$$\sum_{m=0}^{\infty} z^m \mathbf{v}^{(m)}(\theta) \left[(\theta - \theta z) \mathbf{I} + \mathbf{C} + \sum_{m=0}^{\infty} z^m \mathbf{D}^{(m)}(\theta) \right] = (\theta - \theta z)(1 - \rho) \boldsymbol{\kappa}, \quad (2.30)$$

where $\mathbf{D}^{(m)}(\theta)$ denotes

$$\mathbf{D}^{(m)}(\theta) = \int_0^{\infty} e^{-\theta x} \frac{(\theta x)^m}{m!} d\bar{\mathbf{D}}(x).$$

Comparing the coefficient vectors of z^m ($m = 0, 1, \dots$) on both sides of (2.30), we can show that the $\mathbf{v}^{(m)}(\theta)$ is identical to the stationary distribution of a Markov chain of M/G/1 type whose transition probability matrix is given by [Taki01b]

$$\begin{bmatrix} \widetilde{\mathbf{B}}_0 + \widetilde{\mathbf{B}}_1 & \widetilde{\mathbf{B}}_2 & \widetilde{\mathbf{B}}_3 & \widetilde{\mathbf{B}}_4 & \cdots \\ \widetilde{\mathbf{B}}_0 & \widetilde{\mathbf{B}}_1 & \widetilde{\mathbf{B}}_2 & \widetilde{\mathbf{B}}_3 & \cdots \\ \mathbf{O} & \widetilde{\mathbf{B}}_0 & \widetilde{\mathbf{B}}_1 & \widetilde{\mathbf{B}}_2 & \cdots \\ \mathbf{O} & \mathbf{O} & \widetilde{\mathbf{B}}_0 & \widetilde{\mathbf{B}}_1 & \cdots \\ \mathbf{O} & \mathbf{O} & \mathbf{O} & \widetilde{\mathbf{B}}_0 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix},$$

where

$$\widetilde{\mathbf{B}}_0 = \mathbf{I} + \theta^{-1}(\mathbf{C} + \mathbf{D}^{(0)}(\theta)), \quad \widetilde{\mathbf{B}}_m = \theta^{-1} \mathbf{D}^{(m)}(\theta), \quad m \geq 1.$$

Thus applying the general theory of Markov chains of M/G/1 type [Neut89], we can compute the $\mathbf{v}^{(m)}(\theta)$. As for the truncation and stopping criteria in computing the steady-state solution of Markov chains of M/G/1 type, readers are referred to [Neut89, Taki94b].

Let $\mathbf{d}_k^{(m)}(\theta)$ ($k \in \mathcal{K}$, $m = 0, 1, \dots$) denote a $1 \times M_k$ vector which satisfies

$$\sum_{m=0}^{\infty} z^m \mathbf{d}_k^{(m)}(\theta) = H_k^*(\theta - \theta z) \boldsymbol{\alpha}_k (\mathbf{I} - \mathbf{P}_k) [\mathbf{I} - H_k^*(\theta - \theta z) \mathbf{P}_k]^{-1}. \quad (2.31)$$

Lemma 2.5 *Under Assumption 2.2, $\mathbf{D}^{(m)}(\theta)$ is given by*

$$\mathbf{D}^{(m)}(\theta) = \sum_{k \in \mathcal{K}} \mathbf{d}_k^{(m)}(\theta) \mathbf{e} \mathbf{D}_k, \quad m = 0, 1, \dots,$$

where $\mathbf{d}_k^{(m)}(\theta)$'s ($k \in \mathcal{K}$) are recursively determined by

$$\mathbf{d}_k^{(0)}(\theta) = \gamma_k^{(0)}(\theta) \boldsymbol{\alpha}_k (\mathbf{I} - \mathbf{P}_k) [\mathbf{I} - \gamma_k^{(0)}(\theta) \mathbf{P}_k]^{-1},$$

and for $m = 1, 2, \dots$,

$$\mathbf{d}_k^{(m)}(\theta) = \frac{\gamma_k^{(m)}(\theta)}{\gamma_k^{(0)}(\theta)} \mathbf{d}_k^{(0)}(\theta) + \left[\sum_{l=1}^m \gamma_k^{(l)}(\theta) \mathbf{d}_k^{(m-l)}(\theta) \right] \mathbf{P}_k [\mathbf{I} - \gamma_k^{(0)}(\theta) \mathbf{P}_k]^{-1}.$$

Proof. Note first that

$$\sum_{m=0}^{\infty} z^m \mathbf{D}^{(m)}(\theta) = \overline{\mathbf{D}}^*(\theta - \theta z),$$

where $\overline{\mathbf{D}}^*(s)$ is given in (2.5). Thus, substituting $\theta - \theta z$ for s in (2.5) and using (2.14) and (2.15), we have

$$\begin{aligned} \sum_{m=0}^{\infty} z^m \mathbf{D}^{(m)}(\theta) &= \sum_{k \in \mathcal{K}} \sum_{n=1}^{\infty} \alpha_k \mathbf{P}_k^{n-1} (\mathbf{I} - \mathbf{P}_k) e \{H_k^*(\theta - \theta z)\}^n \mathbf{D}_k \\ &= \sum_{k \in \mathcal{K}} H_k^*(\theta - \theta z) \alpha_k (\mathbf{I} - \mathbf{P}_k) [\mathbf{I} - H_k^*(\theta - \theta z) \mathbf{P}_k]^{-1} e \mathbf{D}_k. \end{aligned} \quad (2.32)$$

It then follows from (2.31) and (2.32) that

$$\sum_{m=0}^{\infty} z^m \mathbf{D}^{(m)}(\theta) = \sum_{k \in \mathcal{K}} \sum_{m=0}^{\infty} z^m \mathbf{d}_k^{(m)}(\theta) e \mathbf{D}_k = \sum_{m=0}^{\infty} z^m \sum_{k \in \mathcal{K}} \mathbf{d}_k^{(m)}(\theta) e \mathbf{D}_k. \quad (2.33)$$

Note here that

$$H_k^*(\theta - \theta z) = \sum_{m=0}^{\infty} z^m \gamma_k^{(m)}(\theta), \quad k \in \mathcal{K}. \quad (2.34)$$

Thus from (2.31) and (2.34), we have

$$\sum_{m=0}^{\infty} z^m \mathbf{d}_k^{(m)}(\theta) \left[\mathbf{I} - \sum_{l=0}^{\infty} z^l \gamma_k^{(l)}(\theta) \mathbf{P}_k \right] = \sum_{m=0}^{\infty} z^m \gamma_k^{(m)}(\theta) \alpha_k [\mathbf{I} - \mathbf{P}_k],$$

or equivalently,

$$\sum_{m=0}^{\infty} z^m \mathbf{d}_k^{(m)}(\theta) - \sum_{m=0}^{\infty} z^m \sum_{l=0}^m \mathbf{d}_k^{(m-l)}(\theta) \gamma_k^{(l)}(\theta) \mathbf{P}_k = \sum_{m=0}^{\infty} z^m \gamma_k^{(m)}(\theta) \alpha_k [\mathbf{I} - \mathbf{P}_k].$$

Comparing the coefficient vectors of z^m ($m = 0, 1, \dots$) on both sides of the above equation, we have

$$\mathbf{d}_k^{(0)}(\theta) [\mathbf{I} - \gamma_k^{(0)}(\theta) \mathbf{P}_k] = \gamma_k^{(0)}(\theta) \alpha_k [\mathbf{I} - \mathbf{P}_k], \quad (2.35)$$

and for $m = 1, 2, \dots$,

$$\mathbf{d}_k^{(m)}(\theta) [\mathbf{I} - \gamma_k^{(0)}(\theta) \mathbf{P}_k] - \sum_{l=1}^m \mathbf{d}_k^{(m-l)}(\theta) \gamma_k^{(l)}(\theta) \mathbf{P}_k = \gamma_k^{(m)}(\theta) \alpha_k [\mathbf{I} - \mathbf{P}_k]. \quad (2.36)$$

Lemma 2.5 now follows from (2.33), (2.35) and (2.36). ■

2.5 Implementations of Recursions

In this section, we consider the implementation of recursions for $\mathbf{A}_k(\mathbf{n})$, $\mathbf{v}_k(\mathbf{n})$ and $\mathbf{\Gamma}_k(\mathbf{n})$, derived in the preceding section. At a glance, they would seem to be easy to implement. Contrary to the single arrival case [Taki01a, Taki01b], however, the computation of the $\mathbf{F}_m(\mathbf{n})$ appeared in $\mathbf{A}_k(\mathbf{n})$ and $\mathbf{v}_k(\mathbf{n})$ is not straightforward, because the direct implementation of the recursion requires very huge memory space and time-consuming. In what follows, we construct a numerically feasible procedure

to compute the approximate sequences of $\mathbf{A}_k(\mathbf{n})$ and $\mathbf{v}_k(\mathbf{n})$, avoiding the computation of $\mathbf{F}_m(\mathbf{n})$'s whose contributions to $\mathbf{A}_k(\mathbf{n})$ and $\mathbf{v}_k(\mathbf{n})$ are negligible, and establish the truncation/stopping criteria and error bounds. Further, we propose a computational procedure for the $\mathbf{\Gamma}_k(\mathbf{n})$ and establish the error bound.

We start with $\mathbf{A}_k(\mathbf{n})$ and $\mathbf{v}_k(\mathbf{n})$. Note first that for $k \in \mathcal{K}$,

$$\sum_{\mathbf{n} \in \mathcal{Z}} \mathbf{A}_k(\mathbf{n})\mathbf{e} = \mathbf{e}, \quad \sum_{\mathbf{n} \in \mathcal{Z}} \mathbf{v}_k(\mathbf{n})\mathbf{e} = \lambda_k^{(B)},$$

where

$$\lambda_k^{(B)} = \boldsymbol{\pi} \mathbf{D}_k \mathbf{e}.$$

In numerical computation, we have to stop the computation of those sequences. Thus we develop a numerical procedure to obtain approximations $\check{\mathbf{A}}_k(\mathbf{n})$ and $\check{\mathbf{v}}_k(\mathbf{n})$ to $\mathbf{A}_k(\mathbf{n})$ and $\mathbf{v}_k(\mathbf{n})$, respectively, while ensuring the following error bounds: For a given ε ($0 < \varepsilon < 1$), there exist $n_A(k)$ and $n_v(k)$ such that

$$\sum_{\substack{\mathbf{n} \in \mathcal{Z} \\ |\mathbf{n}| \leq n_A(k)}} \check{\mathbf{A}}_k(\mathbf{n})\mathbf{e} > (1 - \varepsilon)\mathbf{e}, \quad (2.37)$$

$$\sum_{\substack{\mathbf{n} \in \mathcal{Z} \\ |\mathbf{n}| \leq n_v(k)}} \check{\mathbf{v}}_k(\mathbf{n})\mathbf{e} > (1 - \varepsilon)\lambda_k^{(B)}, \quad (2.38)$$

where $|\mathbf{n}| = \sum_{k \in \mathcal{K}} n_k$ for $\mathbf{n} \in \mathcal{Z}$. In what follows, we first show our proposed algorithm and then show that the above error bounds are satisfied.

Numerical algorithm for $\mathbf{A}_k(\mathbf{n})$ and $\mathbf{v}_k(\mathbf{n})$

Input.

Stopping criterion : ε ($0 < \varepsilon < 1$),

Underlying Markov chain : \mathbf{C} , \mathbf{D}_k ($k \in \mathcal{K}$),

Batch size distribution : $\boldsymbol{\alpha}_k$, \mathbf{P}_k ($k \in \mathcal{K}$),

Service time distribution : $H_k(x)$ ($k \in \mathcal{K}$).

Step 1. Choose ε_F ($0 < \varepsilon_F < 1$) such that

$$\frac{\varepsilon_F}{\varepsilon} < \min_{k \in \mathcal{K}} \min \left(\frac{1}{\theta h_k}, \frac{\lambda_k^{(B)}}{\theta \bar{\mathbf{v}}^{(1)} \mathbf{D}_k \mathbf{e}} \right), \quad (2.39)$$

where $\bar{\mathbf{v}}^{(1)} = -\lim_{s \rightarrow 0^+} d\mathbf{v}^*(s)/ds$, whose computational procedure can be found in [Taki94a] (see Proposition 1.9). Then compute the $\gamma_k^{(m)}(\theta)$ and the $\mathbf{v}^{(m)}(\theta)$ until they satisfy

$$\sum_{m=0}^{m_{\gamma}(k)} \gamma_k^{(m)}(\theta)(1 - \varepsilon_F)^m > 1 - \varepsilon, \quad k \in \mathcal{K}, \quad (2.40)$$

$$\sum_{m=0}^{m_v(k)} \mathbf{v}^{(m)}(\theta) \mathbf{D}_k \mathbf{e} (1 - \varepsilon_F)^m > (1 - \varepsilon)\lambda_k^{(B)}, \quad k \in \mathcal{K}, \quad (2.41)$$

for some $m_\gamma(k)$ and $m_v(k)$, respectively. Define m_{\max} as

$$m_{\max} = \max_{k \in \mathcal{K}} \max(m_\gamma(k), m_v(k)).$$

Step 2. Choose ε_g such that $0 < \varepsilon_g < \varepsilon_F$. Then compute $g_k(n)$ ($n = 1, 2, \dots$) by (2.15) until the $g_k(n)$ satisfies

$$\theta^{-1} \sum_{n=1}^{n_g(k)} g_k(n) \mathbf{D}_k \mathbf{e} > \theta^{-1} \mathbf{D}_k \mathbf{e} - \frac{\varepsilon_g}{K} \mathbf{e}, \quad k \in \mathcal{K}, \quad (2.42)$$

for some $n_g(k)$.

Step 3. Compute $\check{\mathbf{A}}_k(\mathbf{n})$ and $\check{\mathbf{v}}_k(\mathbf{n})$ by the following procedure, where the initial values of $\check{\mathbf{A}}_k(\mathbf{n})$ and $\check{\mathbf{v}}_k(\mathbf{n})$ ($\mathbf{n} \in \mathcal{Z}$) are assumed to be \mathbf{O} and $\mathbf{0}$, respectively.

Step (3-a). Set $\check{\mathbf{F}}_0(\mathbf{0}) = \mathbf{I}$ and $n_F^{(0)} = 0$. Also set $\check{\mathbf{A}}_k(\mathbf{0}) = \gamma_k^{(0)}(\theta) \mathbf{I}$ and $\check{\mathbf{v}}_k(\mathbf{0}) = \mathbf{v}^{(0)}(\theta) \mathbf{D}_k$ for all $k \in \mathcal{K}$.

Step (3-b). Set $n_F^{(1)} = \max_{k \in \mathcal{K}} n_g(k)$ and $m = 1$, and compute $\check{\mathbf{F}}_1(\mathbf{n})$'s ($|\mathbf{n}| \leq n_F^{(1)}$) by

$$\check{\mathbf{F}}_1(\mathbf{n}) = \begin{cases} \mathbf{I} + \theta^{-1} \mathbf{C}, & \text{if } \mathbf{n} = \mathbf{0}, \\ \theta^{-1} g_k(n_k) \mathbf{D}_k, & \text{if } \mathbf{n} \in \mathcal{Z}_k(F_1), k \in \mathcal{K}, \\ \mathbf{O}, & \text{otherwise,} \end{cases} \quad (2.43)$$

where

$$\mathcal{Z}_k(F_1) = \{\mathbf{n}; \mathbf{n} = n_k \mathbf{e}_k, n_k = 1, 2, \dots, n_g(k)\}, \quad k \in \mathcal{K}.$$

Step (3-c). For each $k \in \mathcal{K}$, if $m \leq m_\gamma(k)$, add $\gamma_k^{(m)}(\theta) \check{\mathbf{F}}_m(\mathbf{n})$ to $\check{\mathbf{A}}_k(\mathbf{n})$ for all \mathbf{n} ($|\mathbf{n}| \leq n_F^{(m)}$). Also, for each $k \in \mathcal{K}$, if $m \leq m_v(k)$, add $\mathbf{v}^{(m)}(\theta) \mathbf{D}_k \check{\mathbf{F}}_m(\mathbf{n})$ to $\check{\mathbf{v}}_k(\mathbf{n})$ for all \mathbf{n} ($|\mathbf{n}| \leq n_F^{(m)}$).

Step (3-d). If $m \geq m_{\max}$, stop computing, and otherwise, add one to m and go to Step (3-e).

Step (3-e). For each $n = 0, 1, \dots$, compute $\check{\mathbf{F}}_m(\mathbf{n})$'s ($|\mathbf{n}| = n$) by

$$\begin{aligned} \check{\mathbf{F}}_m(\mathbf{n}) &= U\left(n_F^{(m-1)} - |\mathbf{n}|\right) \check{\mathbf{F}}_{m-1}(\mathbf{n})(\mathbf{I} + \theta^{-1} \mathbf{C}) \\ &\quad + \theta^{-1} \sum_{k \in \mathcal{K}} \sum_{l_k=1}^{\min(n_k, n_g(k))} U\left(n_F^{(m-1)} - |\mathbf{n} - l_k \mathbf{e}_k|\right) \\ &\quad \cdot \check{\mathbf{F}}_{m-1}(\mathbf{n} - l_k \mathbf{e}_k) g_k(l_k) \mathbf{D}_k, \end{aligned} \quad (2.44)$$

until $\check{\mathbf{F}}_m(\mathbf{n})$'s satisfy $\sum_{|\mathbf{n}| \leq n^*} \check{\mathbf{F}}_m(\mathbf{n}) \mathbf{e} > (1 - \varepsilon_F)^m \mathbf{e}$ for some n^* , where $U(x)$ denotes a unit step function:

$$U(x) = \begin{cases} 1, & x \geq 0, \\ 0, & x < 0. \end{cases}$$

Let $n_F^{(m)} = n^*$ and go to Step (3-c).

Remark 2.2 Note that $\check{\mathbf{A}}_k(\mathbf{n})$ ($|\mathbf{n}| \leq n_A(k)$) and $\check{\mathbf{v}}_k(\mathbf{n})$ ($|\mathbf{n}| \leq n_v(k)$) obtained by the above algorithm satisfy

$$\begin{aligned}\check{\mathbf{A}}_k(\mathbf{n}) &= \sum_{m=0}^{m_\gamma(k)} U\left(n_F^{(m)} - |\mathbf{n}|\right) \gamma_k^{(m)}(\theta) \check{\mathbf{F}}_m(\mathbf{n}), \\ \check{\mathbf{v}}_k(\mathbf{n}) &= \sum_{m=0}^{m_v(k)} U\left(n_F^{(m)} - |\mathbf{n}|\right) \mathbf{v}^{(m)}(\theta) \mathbf{D}_k \check{\mathbf{F}}_m(\mathbf{n}),\end{aligned}$$

respectively, where

$$\begin{aligned}n_A(k) &= \max\left(n_F^{(m)}; m = 0, 1, \dots, m_\gamma(k)\right), \\ n_v(k) &= \max\left(n_F^{(m)}; m = 0, 1, \dots, m_v(k)\right).\end{aligned}\tag{2.45}$$

Remark 2.3 If we are interested only in the $\mathbf{p}(\mathbf{n})$ ($|\mathbf{n}| \leq N_p$) for some N_p , we do not need to compute $\check{\mathbf{F}}(\mathbf{n})$ for \mathbf{n} such that $|\mathbf{n}| > N_p$. Thus, in this case, $n_g(k)$ is redefined as $\min(n_g(k), N_p)$ and Step (3-e) is replaced by

Step (3-e'). For each n ($n = 0, 1, \dots$), compute $\check{\mathbf{F}}_m(\mathbf{n})$'s ($|\mathbf{n}| = n$, $\mathbf{n} \in \mathcal{Z}$) by (2.44) until $\check{\mathbf{F}}_m(\mathbf{n})$'s satisfy

$$\sum_{|\mathbf{n}| \leq n^*} \check{\mathbf{F}}_m(\mathbf{n}) \mathbf{e} > (1 - \varepsilon_F)^m \mathbf{e},$$

for some n^* , or $n = N_p$, whichever occurs first. Let $n_F^{(m)} = \min(n^*, N_p)$ and go to Step (3-c).

This procedure can save the computational cost, while maintaining the accuracy of the results.

We now provide two lemmas that ensure the above procedure eventually stops.

Lemma 2.6 There exist integers $m_\gamma(k)$ and $m_v(k)$ satisfying (2.40) and (2.41), respectively.

Proof. Substituting $1 - \varepsilon_F$ for z in (2.34), we have

$$\sum_{m=0}^{\infty} \gamma_k^{(m)}(\theta) (1 - \varepsilon_F)^m = H_k^*(\theta \varepsilon_F),\tag{2.46}$$

where $H_k^*(s)$ denotes the LST of $H_k(x)$. Similarly, from (2.29), we have

$$\sum_{m=0}^{\infty} \mathbf{v}^{(m)}(\theta) \mathbf{D}_k \mathbf{e} (1 - \varepsilon_F)^m = \mathbf{v}^*(\theta \varepsilon_F) \mathbf{D}_k \mathbf{e}.\tag{2.47}$$

Note here that

$$H_k^*(\theta \varepsilon_F) > 1 - h_k \cdot (\theta \varepsilon_F), \quad k \in \mathcal{K},\tag{2.48}$$

$$\mathbf{v}^*(\theta \varepsilon_F) \mathbf{D}_k \mathbf{e} > \lambda_k^{(B)} - \bar{\mathbf{v}}^{(1)} \mathbf{D}_k \mathbf{e} \cdot (\theta \varepsilon_F), \quad k \in \mathcal{K},\tag{2.49}$$

because $H_k^*(s)$ and each element of $\mathbf{v}^*(s)$ are convex functions of s . Note also that (2.39) is equivalent to

$$1 - h_k \theta \varepsilon_F > 1 - \varepsilon, \quad k \in \mathcal{K}, \quad (2.50)$$

$$\lambda_k^{(B)} - \bar{\mathbf{v}}^{(1)} \mathbf{D}_k \mathbf{e} \theta \varepsilon_F > (1 - \varepsilon) \lambda_k^{(B)}, \quad k \in \mathcal{K}. \quad (2.51)$$

It then follows from (2.46)–(2.51) that

$$\begin{aligned} \sum_{m=0}^{\infty} \gamma_k^{(m)}(\theta) (1 - \varepsilon_F)^m &> 1 - h_k \theta \varepsilon_F > 1 - \varepsilon, \\ \sum_{m=0}^{\infty} \mathbf{v}^{(m)}(\theta) \mathbf{D}_k \mathbf{e} (1 - \varepsilon_F)^m &> \lambda_k^{(B)} - \bar{\mathbf{v}}^{(1)} \mathbf{D}_k \mathbf{e} \theta \varepsilon_F > (1 - \varepsilon) \lambda_k^{(B)}, \end{aligned}$$

which complete the proof. \blacksquare

Lemma 2.7 *There exists an integer $n_F^{(m)}$ such that*

$$\sum_{\substack{\mathbf{n} \in \mathcal{Z} \\ |\mathbf{n}| \leq n_F^{(m)}}} \check{\mathbf{F}}_m(\mathbf{n}) \mathbf{e} > (1 - \varepsilon_F)^m \mathbf{e}, \quad \forall m = 1, 2, \dots, m_{\max}. \quad (2.52)$$

Proof. We first consider the case $m = 1$. It follows from (2.42) and (2.43) that

$$\begin{aligned} \sum_{\substack{\mathbf{n} \in \mathcal{Z} \\ |\mathbf{n}| \leq n_F^{(1)}}} \check{\mathbf{F}}_1(\mathbf{n}) \mathbf{e} &= \left[\mathbf{I} + \theta^{-1} \mathbf{C} + \theta^{-1} \sum_{k \in \mathcal{K}} \sum_{n_k=1}^{n_g(k)} g_k(n_k) \mathbf{D}_k \right] \mathbf{e} \\ &> \mathbf{e} + \theta^{-1} \mathbf{C} \mathbf{e} + \sum_{k \in \mathcal{K}} \left[\theta^{-1} \mathbf{D}_k \mathbf{e} - \frac{\varepsilon_g}{K} \mathbf{e} \right] \\ &= (1 - \varepsilon_g) \mathbf{e} > (1 - \varepsilon_F) \mathbf{e}, \end{aligned} \quad (2.53)$$

where we use $(\mathbf{C} + \mathbf{D}) \mathbf{e} = \mathbf{0}$ and $\varepsilon_g < \varepsilon_F$.

Suppose that for some m ($1 \leq m \leq m_{\max} - 1$), there exists an integer $n_F^{(m)}$ such that

$$\sum_{\substack{\mathbf{n} \in \mathcal{Z} \\ |\mathbf{n}| \leq n_F^{(m)}}} \check{\mathbf{F}}_m(\mathbf{n}) \mathbf{e} > (1 - \varepsilon_F)^m \mathbf{e}. \quad (2.54)$$

Using (2.43) and (2.44), we have for $\mathbf{n} \in \mathcal{Z}$,

$$\begin{aligned} \check{\mathbf{F}}_{m+1}(\mathbf{n}) &= U \left(n_F^{(m)} - |\mathbf{n}| \right) \check{\mathbf{F}}_m(\mathbf{n}) \check{\mathbf{F}}_1(\mathbf{0}) \\ &\quad + \sum_{k \in \mathcal{K}} \sum_{l_k=1}^{\min(n_k, n_g(k))} U \left(n_F^{(m)} - |\mathbf{n} - l_k \mathbf{e}_k| \right) \check{\mathbf{F}}_m(\mathbf{n} - l_k \mathbf{e}_k) \check{\mathbf{F}}_1(l_k \mathbf{e}_k). \end{aligned} \quad (2.55)$$

It then follows from (2.53), (2.54) and (2.55) that

$$\sum_{\substack{\mathbf{n} \in \mathcal{Z} \\ |\mathbf{n}| \leq n_F^{(m)} + n_F^{(1)}}} \check{\mathbf{F}}_{m+1}(\mathbf{n}) \mathbf{e} = \sum_{\substack{\mathbf{n} \in \mathcal{Z} \\ |\mathbf{n}| \leq n_F^{(m)}}} \check{\mathbf{F}}_m(\mathbf{n}) \sum_{\substack{\mathbf{n} \in \mathcal{Z} \\ |\mathbf{n}| \leq n_F^{(1)}}} \check{\mathbf{F}}_1(\mathbf{n}) \mathbf{e} > (1 - \varepsilon_F)^{m+1} \mathbf{e}. \quad (2.56)$$

Thus we can choose $n_F^{(m+1)}$ in such a way that

$$n_F^{(m+1)} \leq n_F^{(m)} + n_F^{(1)}, \quad m = 1, 2, \dots, m_{\max} - 1,$$

which completes the proof. \blacksquare

Theorem 2.5 For $0 < \varepsilon < 1$, the $\check{\mathbf{A}}_k(\mathbf{n})$ and the $\check{\mathbf{v}}_k(\mathbf{n})$, computed by Step 3, satisfy error bounds (2.37) and (2.38), respectively.

Proof. Using Lemma 2.6, Lemma 2.7 and (2.45), we obtain

$$\begin{aligned} \sum_{\substack{\mathbf{n} \in \mathcal{Z} \\ |\mathbf{n}| \leq n_A(k)}} \check{\mathbf{A}}_k(\mathbf{n})\mathbf{e} &= \sum_{\substack{\mathbf{n} \in \mathcal{Z} \\ |\mathbf{n}| \leq n_A(k)}} \sum_{m=0}^{m_\gamma(k)} U\left(\binom{m}{n_F} - |\mathbf{n}|\right) \gamma_k^{(m)}(\theta) \check{\mathbf{F}}_m(\mathbf{n})\mathbf{e} \\ &= \sum_{m=0}^{m_\gamma(k)} \gamma_k^{(m)}(\theta) \sum_{\substack{\mathbf{n} \in \mathcal{Z} \\ |\mathbf{n}| \leq n_F^{(m)}}} \check{\mathbf{F}}_m(\mathbf{n})\mathbf{e} \\ &> \sum_{m=0}^{m_\gamma(k)} \gamma_k^{(m)}(\theta) (1 - \varepsilon_F)^m \mathbf{e} > (1 - \varepsilon)\mathbf{e}, \quad k \in \mathcal{K}. \end{aligned}$$

In the same way, we can obtain (2.38), so that the proof for the $\check{\mathbf{v}}_k(\mathbf{n})$ is omitted. \blacksquare

Finally, we consider the $\mathbf{\Gamma}_k(\mathbf{n})$. Note here that

$$\sum_{\mathbf{n} \in \mathcal{Z}} \mathbf{\Gamma}_k(\mathbf{n})\mathbf{e} = \left\{ (\mathbf{I} - \mathbf{P}_k)^{-1} \mathbf{e}(M_k) \right\} \otimes \mathbf{e}(M),$$

where $\mathbf{e}(m)$ denotes an $m \times 1$ vector whose elements are all equal to one. Keeping the above equation in mind, we propose to compute an approximation $\check{\mathbf{\Gamma}}_k(\mathbf{n})$ to $\mathbf{\Gamma}_k(\mathbf{n})$ in the following way.

Step 4. For each $k \in \mathcal{K}$, compute $\check{\mathbf{\Gamma}}_k(\mathbf{n})$'s ($|\mathbf{n}| = n$) for $n = 0, 1, \dots$ by

$$\check{\mathbf{\Gamma}}_k(\mathbf{0}) = \left[\mathbf{I} - \mathbf{P}_k \otimes \check{\mathbf{A}}_k(\mathbf{0}) \right]^{-1}, \quad (2.57)$$

$$\check{\mathbf{\Gamma}}_k(\mathbf{n}) = \sum_{\substack{0 \leq l \leq n \\ l \neq 0}} U(n_A(k) - |l|) \check{\mathbf{\Gamma}}_k(\mathbf{n} - l) \left[\mathbf{P}_k \otimes \check{\mathbf{A}}_k(l) \right] \check{\mathbf{\Gamma}}_k(\mathbf{0}), \quad \mathbf{n} \in \mathcal{Z}^+, \quad (2.58)$$

until $\check{\mathbf{\Gamma}}_k(\mathbf{n})$'s satisfy

$$\begin{aligned} \sum_{\substack{\mathbf{n} \in \mathcal{Z} \\ |\mathbf{n}| \leq n_\Gamma(k)}} \check{\mathbf{\Gamma}}_k(\mathbf{n})\mathbf{e} &> \left\{ (\mathbf{I} - \mathbf{P}_k)^{-1} \mathbf{e}(M_k) \right\} \otimes \mathbf{e}(M) \\ &\quad - \varepsilon \left\{ (\mathbf{I} - \mathbf{P}_k)^{-2} \mathbf{P}_k \mathbf{e}(M_k) \right\} \otimes \mathbf{e}(M), \end{aligned} \quad (2.59)$$

for some integer $n_\Gamma(k)$.

Remark 2.4 Let G_k ($k \in \mathcal{K}$) denote a generic random variable representing a batch size of class k . We then have

$$(\boldsymbol{\alpha}_k \otimes \boldsymbol{\pi}) \sum_{\mathbf{n} \in \mathcal{Z}} \mathbf{\Gamma}_k(\mathbf{n})\mathbf{e} = \mathbb{E}[G_k],$$

and if (2.59) satisfies for some $n_\Gamma(k)$,

$$(\boldsymbol{\alpha}_k \otimes \boldsymbol{\pi}) \sum_{\substack{\mathbf{n} \in \mathcal{Z} \\ |\mathbf{n}| \leq n_\Gamma(k)}} \check{\mathbf{\Gamma}}_k(\mathbf{n})\mathbf{e} > \mathbb{E}[G_k] - \frac{1}{2} \mathbb{E}[G_k(G_k - 1)]\varepsilon.$$

Lemma 2.8 Suppose $\check{\mathbf{A}}_k(\mathbf{n})$ satisfies (2.37). Then there exists $n_\Gamma(k)$ satisfying (2.59).

Proof. From (2.57) and (2.58), it can be seen that $\check{\mathbf{A}}_k(\mathbf{n})$ ($|\mathbf{n}| \leq n_A(k)$) and $\check{\mathbf{\Gamma}}_k(\mathbf{n})$ ($\mathbf{n} \in \mathcal{Z}$) are related by

$$\sum_{\mathbf{n} \in \mathcal{Z}} \check{\mathbf{\Gamma}}_k(\mathbf{n}) = \sum_{m=0}^{\infty} \left(\mathbf{P}_k \otimes \sum_{\substack{\mathbf{l} \in \mathcal{Z} \\ |\mathbf{l}| \leq n_A(k)}} \check{\mathbf{A}}_k(\mathbf{l}) \right)^m.$$

Post-multiplying both sides of the above equation by $\mathbf{e} = \mathbf{e}(M_k) \otimes \mathbf{e}(M)$, we have

$$\begin{aligned} \sum_{\mathbf{n} \in \mathcal{Z}} \check{\mathbf{\Gamma}}_k(\mathbf{n}) \mathbf{e} &= \sum_{m=0}^{\infty} \left(\mathbf{P}_k \otimes \sum_{\substack{\mathbf{l} \in \mathcal{Z} \\ |\mathbf{l}| \leq n_A(k)}} \check{\mathbf{A}}_k(\mathbf{l}) \right)^m \cdot [\mathbf{e}(M_k) \otimes \mathbf{e}(M)] \\ &= \sum_{m=0}^{\infty} [\mathbf{P}_k^m \mathbf{e}(M_k)] \otimes \left[\left(\sum_{\substack{\mathbf{l} \in \mathcal{Z} \\ |\mathbf{l}| \leq n_A(k)}} \check{\mathbf{A}}_k(\mathbf{l}) \right)^m \mathbf{e}(M) \right]. \end{aligned}$$

Further, using (2.37), we obtain

$$\begin{aligned} \sum_{\mathbf{n} \in \mathcal{Z}} \check{\mathbf{\Gamma}}_k(\mathbf{n}) \mathbf{e} &> \sum_{m=0}^{\infty} (1 - \varepsilon)^m [\mathbf{P}_k^m \mathbf{e}(M_k)] \otimes \mathbf{e}(M) \\ &> \sum_{m=0}^{\infty} (1 - m\varepsilon) [\mathbf{P}_k^m \mathbf{e}(M_k)] \otimes \mathbf{e}(M) \\ &= \left\{ (\mathbf{I} - \mathbf{P}_k)^{-1} \mathbf{e}(M_k) \right\} \otimes \mathbf{e}(M) \\ &\quad - \varepsilon \left\{ (\mathbf{I} - \mathbf{P}_k)^{-2} \mathbf{P}_k \mathbf{e}(M_k) \right\} \otimes \mathbf{e}(M), \end{aligned}$$

which completes the proof. ■

2.6 Numerical Examples

In this section, we show some numerical examples for queues with two arrival streams. Even though the algorithmic analysis has already been done for the single arrival cases [Taki01a, Taki01b], no numerical examples were shown there. Thus the numerical result provided below is the first report in the literature, as for FIFO queues with multiple Markovian arrival streams having different service time distributions.

In all numerical examples, the counting process of class k ($k = 1, 2$) arrivals follows a batch interrupted Poisson process with geometrically distributed batch sizes with mean g . Namely, the counting process of class k ($k = 1, 2$) is characterized by $(\tilde{\mathbf{C}}_k, \tilde{\mathbf{D}}_k(n))$, where

$$\begin{aligned} \tilde{\mathbf{C}}_k &= \begin{bmatrix} -2\lambda_k g^{-1} - 0.1 & 0.1 \\ 0.1 & -0.1 \end{bmatrix}, \\ \tilde{\mathbf{D}}_k(n) &= (1-p)p^{n-1} \begin{bmatrix} 2\lambda_k g^{-1} & 0 \\ 0 & 0 \end{bmatrix}, \quad n = 1, 2, \dots, \end{aligned}$$

where $p = 1 - 1/g$. Note that the arrival rate of class k is fixed to be λ_k regardless of the mean batch size g .

We now consider three types of the superposition of these two streams.

[Case P]

$$\mathbf{C} = \begin{bmatrix} -2(\lambda_1 + \lambda_2)g^{-1} - 0.1 & 0.1 \\ 0.1 & -0.1 \end{bmatrix},$$

and for $n = 1, 2, \dots$,

$$\begin{aligned} \mathbf{D}_1(n) &= (1-p)p^{n-1} \begin{bmatrix} 2\lambda_1 g^{-1} & 0 \\ 0 & 0 \end{bmatrix}, \\ \mathbf{D}_2(n) &= (1-p)p^{n-1} \begin{bmatrix} 2\lambda_2 g^{-1} & 0 \\ 0 & 0 \end{bmatrix}. \end{aligned}$$

[Case I]

$$\mathbf{C} = \tilde{\mathbf{C}}_1 \oplus \tilde{\mathbf{C}}_2,$$

and for $n = 1, 2, \dots$,

$$\mathbf{D}_1(n) = \tilde{\mathbf{D}}_1(n) \otimes \mathbf{I}(2), \quad \mathbf{D}_2(n) = \mathbf{I}(2) \otimes \tilde{\mathbf{D}}_2(n),$$

where \oplus denotes the Kronecker sum, and

[Case N]

$$\mathbf{C} = \begin{bmatrix} -2\lambda_1 g^{-1} - 0.1 & 0.1 \\ 0.1 & -2\lambda_2 g^{-1} - 0.1 \end{bmatrix},$$

and for $n = 1, 2, \dots$,

$$\begin{aligned} \mathbf{D}_1(n) &= (1-p)p^{n-1} \begin{bmatrix} 2\lambda_1 g^{-1} & 0 \\ 0 & 0 \end{bmatrix}, \\ \mathbf{D}_2(n) &= (1-p)p^{n-1} \begin{bmatrix} 0 & 0 \\ 0 & 2\lambda_2 g^{-1} \end{bmatrix}. \end{aligned}$$

Note that in Case P, two arrival streams are positively correlated, in Case I, they are independent each other and in Case N, they are negatively correlated. As for the service time distributions, we consider two cases, Case GD (class-dependent service times) and Case GI (i.i.d. service times):

[Case GD]

$$H_1 = 1, \text{ with prob. } 1, \quad H_2 = 4, \text{ with prob. } 1,$$

[Case GI]

$$H_k = \begin{cases} 1, & \text{with prob. } \lambda_1/(\lambda_1 + \lambda_2), \\ 4, & \text{with prob. } \lambda_2/(\lambda_1 + \lambda_2), \end{cases} \quad k = 1, 2,$$

where H_k ($k = 1, 2$) denotes a generic random variable for a service time of a class k customer. Note that the overall service time distributions are identical in both cases. We denote the queuing model with Case i ($i = P, I, N$) arrivals and Case j ($j = GD, GI$) services by Case (i, j) .

In what follows, we consider two examples, Examples 1 and 2, within the above settings. In Example 1, we set $\lambda_1 = \lambda_2 = 0.15$, so that $\rho_1 = 0.15$ and $\rho_2 = 0.6$ in Case (i, GD) ($i = P, I, N$), and $\rho_1 = \rho_2 = 0.375$ in Case (i, GI) ($i = P, I, N$). On the other hand, in Example 2, we set $\lambda_1 = 0.4$ and $\lambda_2 = 0.1$, so that $\rho_1 = \rho_2 = 0.4$ in Case (i, GD) ($i = P, I, N$) and that $\rho_1 = 0.64$ and $\rho_2 = 0.16$ in Case (i, GI) ($i = P, I, N$).

2.6.1 Efficiency of the algorithm

Before showing the quantitative behavior of the queue length distribution, we discuss the efficiency of our numerical algorithm for the $\mathbf{F}_m(\mathbf{n})$. It follows from (2.21) that for $m = 0, 1, \dots$,

$$\sum_{\mathbf{n} \in \mathcal{Z}} \mathbf{F}_m(\mathbf{n}) = \left[\mathbf{I} + \theta^{-1} \left(\mathbf{C} + \sum_{k \in \mathcal{K}} \sum_{n_k=1}^{\infty} \mathbf{D}_k(n_k) \right) \right]^m, \quad (2.60)$$

where $\mathbf{I} + \theta^{-1}[\mathbf{C} + \sum_{k \in \mathcal{K}} \sum_{n_k=1}^{\infty} \mathbf{D}_k(n_k)]$ is a stochastic matrix. Thus a straightforward implementation of the recursion for the $\mathbf{F}_m(\mathbf{n})$ in (2.24) and (2.25) would be the following. We first truncate the $\mathbf{D}_k(n_k)$ at $n_k = n'_g(k)$ in such a way that

$$\theta^{-1} \sum_{n_k=1}^{n'_g(k)} \mathbf{D}_k(n_k) \mathbf{e} > \theta^{-1} \mathbf{D}_k \mathbf{e} - \frac{\varepsilon'_g}{K} \mathbf{e},$$

so that

$$\left[\mathbf{I} + \theta^{-1} \left(\mathbf{C} + \sum_{k \in \mathcal{K}} \sum_{n_k=1}^{n'_g(k)} \mathbf{D}_k(n_k) \right) \right] \mathbf{e} > (1 - \varepsilon'_g) \mathbf{e}.$$

We then compute all terms obtained by expanding the right hand side of (2.60) with the truncated $\mathbf{D}_k(n_k)$ ($k \in \mathcal{K}$). Note that if $\varepsilon'_g = \varepsilon_F$, the resulting $\mathbf{F}_m(\mathbf{n})$ satisfies (2.52) in Lemma 2.7, where the summation on the left hand side of (2.52) is taken for all computed $\mathbf{F}_m(\mathbf{n})$'s.

In Table 2.1, we show the numbers of $\mathbf{F}_m(\mathbf{n})$'s computed by our algorithm and the above straightforward implementation, using Example 1, where we set $\varepsilon = 10^{-6}$, $\varepsilon_F = \text{r.h.s. of (2.39)} \times \varepsilon/2$ and $\varepsilon_g = \varepsilon_F/10$. We observe that for unbounded batch size cases (i.e., $g > 1$), the number of the computed $\mathbf{F}_m(\mathbf{n})$'s in our algorithm is less than that in the straightforward algorithm about by three order of magnitude. Thus, compared to the straightforward implementation, our algorithm is very efficient in terms of the computational time when the batch size is unbounded.

Table 2.1: Number of computed $\mathbf{F}_m(\mathbf{n})$'s in Example 1.

		A : Our algorithm		B : Straightforward	
Case		$g = 1$	$g = 2$	$g = 5$	$g = 10$
(P, GD)	A	1.021×10^6	3.918×10^6	2.504×10^7	1.257×10^8
	B	2.453×10^6	2.107×10^9	3.696×10^{10}	4.243×10^{11}
(P, GI)	A	1.021×10^6	2.898×10^6	1.476×10^7	6.732×10^7
	B	2.453×10^6	1.438×10^9	1.812×10^{10}	1.825×10^{11}
(I, GD)	A	6.108×10^5	3.314×10^6	3.012×10^7	1.854×10^8
	B	1.993×10^6	2.919×10^9	9.286×10^{10}	1.649×10^{12}
(I, GI)	A	4.123×10^5	1.859×10^6	1.532×10^7	9.032×10^7
	B	1.253×10^6	1.292×10^9	3.741×10^{10}	6.173×10^{11}
(N, GD)	A	6.657×10^4	8.895×10^5	1.165×10^7	8.095×10^7
	B	1.378×10^5	2.620×10^8	1.066×10^{10}	1.955×10^{11}
(N, GI)	A	1.411×10^4	3.113×10^5	4.813×10^6	3.540×10^7
	B	2.743×10^4	7.158×10^7	3.278×10^9	6.442×10^{10}

We note that a very huge memory space is required to store all $\check{\mathbf{F}}_m(\mathbf{n})$'s in some cases, even using our truncation and stopping criteria. For example, in Case (I, GD) with $g = 10$, the memory space to store all $\check{\mathbf{F}}_m(\mathbf{n})$'s is given by $16 \times 1.854 \times 10^8 \times 8$ bytes ≈ 23.73 Gbytes, because each $\check{\mathbf{F}}_m(\mathbf{n})$ is a 4×4 matrix and one element requires 8 bytes in double precision. Thus in our implementation, every time $\check{\mathbf{F}}_m(\mathbf{n})$'s for each m are obtained, we compute the contributions of $\check{\mathbf{F}}_m(\mathbf{n})$'s to $\check{\mathbf{A}}_k(\mathbf{n})$ and $\check{\mathbf{v}}_k(\mathbf{n})$ in Step (3-c), and discard all $\check{\mathbf{F}}_{m-1}(\mathbf{n})$'s.

Table 2.2 shows the maximum number of $\check{\mathbf{F}}_m(\mathbf{n})$'s stored temporarily in our algorithm, where the ratio of it to the total number of computed $\check{\mathbf{F}}_m(\mathbf{n})$'s is also shown in parenthesis. We observe that in most cases, the number of temporarily stored $\check{\mathbf{F}}_m(\mathbf{n})$'s is a few percent of the total number of computed ones. Thus our implementation is expected to save the required memory space, especially when a large number of $\check{\mathbf{F}}_m(\mathbf{n})$'s should be computed.

2.6.2 Number of customers in Example 1

Figures 2.1–2.3 plot the complementary distributions of the total number N of customers in Case (i , GD) and Case (i , GI) ($i = P, I, N$), where the batch size is fixed to be one, i.e., $g = 1$. Note that the overall input processes in Case (P, GD) and Case (P, GI) are identical, so that the distributions of the total number of customers are also identical, as shown in Figure 2.1. However, as shown in Table 2.3, the joint queue length distributions in these two cases are different. Note also that in Case (P, GI), $\mathbf{p}(n_1, n_2)\mathbf{e} = \mathbf{p}(n_2, n_1)\mathbf{e}$, because the conditional joint distribution $\Pr(N_1 = n_1, N_2 = n_2 \mid N_1 + N_2 = n_1 + n_2)$ follows a binomial distribution with parameter 0.5. We also observe that $\mathbf{p}(n, n)\mathbf{e}$'s in both cases take the same value for each n . Unfortunately, we cannot provide any intuitive explanation of this phenomenon.

Table 2.2: Number of stored $\check{F}_m(\mathbf{n})$'s in Example 1.

Case	$g = 1$	$g = 2$	$g = 5$	$g = 10$
(P, GD)	27225 (2.67%)	90601 (2.31%)	455625 (1.82%)	1651227 (1.31%)
(P, GI)	27225 (2.67%)	75076 (2.59%)	330051 (2.24%)	1125723 (1.67%)
(I, GD)	16641 (2.72%)	68121 (2.06%)	399424 (1.33%)	1548781 (0.84%)
(I, GI)	13110 (3.18%)	47524 (2.56%)	263683 (1.72%)	997003 (1.10%)
(N, GD)	4970 (7.47%)	41209 (4.63%)	324331 (2.78%)	1387686 (1.71%)
(N, GI)	1764 (12.51%)	21171 (6.80%)	187491 (3.90%)	833571 (2.35%)

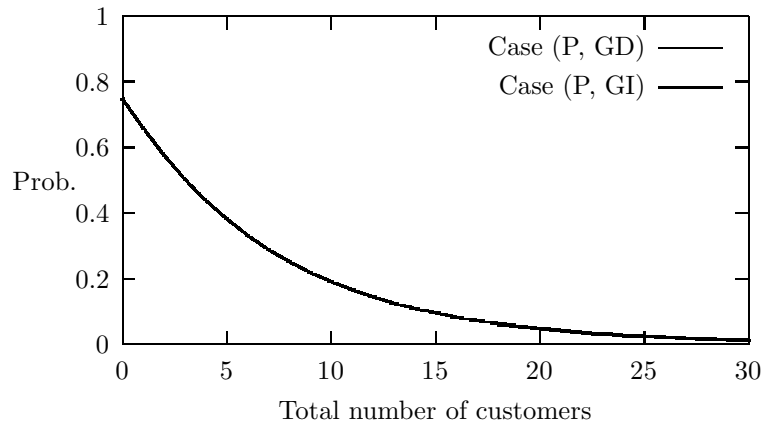


Figure 2.1: Complementary distribution of total number of customers in Example 1.

From Figures 2.2 and 2.3, we observe that class-dependent service times cause longer tails in the total queue length distributions, in these specific examples. We shall explain this phenomenon for Case N. In Case (N, GD), the conditional expected amounts of work brought into the system per unit time given the state of the underlying Markov chain are different, and they are given by 0.3 and 1.2, respectively. Thus in Case (N, GD), the system is overloaded during a half of time. On the other hand, in Case (N, GI), the conditional expected amount of work brought into the system per unit time is fixed to be 0.75, regardless of the state of the underlying Markov chain. Therefore the distribution of the total number of customers in Case (N, GD) has a longer tail than that in Case (N, GI).

Next, we consider the expected total number $E[N]$ of customers as a function of the mean batch

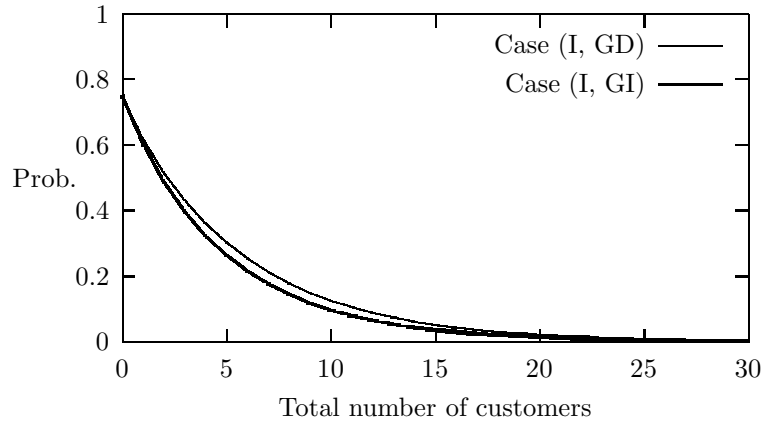


Figure 2.2: Complementary distribution of total number of customers in Example 1.

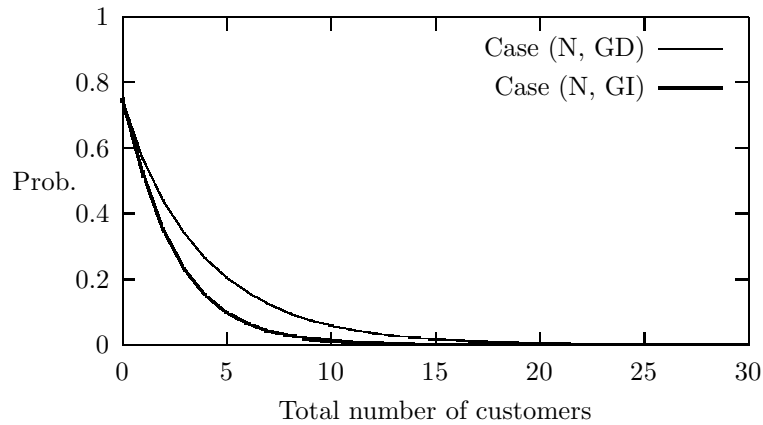


Figure 2.3: Complementary distribution of total number of customers in Example 1.

size g . Table 2.4 shows $E[N]$ for the mean batch size $g = 1, 2, 3, 4, 5$ and 10. We observe that $E[N]$ increases with the mean batch size g in all cases. This phenomenon comes from the fact that the deviation of the amount of work brought into the system per unit time increases with g . We also observe that for a fixed g , the positive correlation in the two streams leads to a larger $E[N]$ in both Cases GD and GI, as expected.

2.6.3 Number of customers in Example 2

Table 2.5 shows the expected total number $E[N]$ of customers for the mean batch size $g = 1, 2, 3, 4, 5$ and 10. We first examine the case of $g = 1$. Contrary to Example 1, we observe that the class-dependent service time (Case GD) decreases the expected total number of customers in Cases I and N. This phenomenon can be explained in a similar way to Example 1. For example, in Case (N, GD), the conditional expected amount of work brought into the system per unit time is fixed to be 0.8, regardless of the state of the underlying Markov chain. On the other hand, in Case (N, GI),

Table 2.3: Joint queue length distribution $\mathbf{p}(n_1, n_2)\mathbf{e}$.
(Upper rows for Case (P, GD) and lower rows for Case (P, GI))

n_1	0	1	2	3
n_2				
0	2.500×10^{-1}	2.472×10^{-2}	8.593×10^{-3}	3.481×10^{-3}
	2.500×10^{-1}	4.501×10^{-2}	2.054×10^{-2}	9.224×10^{-3}
1	6.530×10^{-2}	4.108×10^{-2}	2.193×10^{-2}	1.118×10^{-2}
	4.501×10^{-2}	4.108×10^{-2}	2.767×10^{-2}	1.629×10^{-2}
2	3.249×10^{-2}	3.341×10^{-2}	2.444×10^{-2}	
	2.054×10^{-2}	2.767×10^{-2}	2.444×10^{-2}	
3	1.497×10^{-2}	2.141×10^{-2}		
	9.224×10^{-3}	1.629×10^{-2}		
4	6.630×10^{-3}			
	4.073×10^{-3}			

Table 2.4: Expected total number of customers in Example 1.

Case	$g = 1$	$g = 2$	$g = 3$	$g = 4$	$g = 5$	$g = 10$
(P, GD)	5.8760	9.9815	13.9356	17.8320	21.7001	40.8865
(P, GI)	5.8760	9.1466	12.2898	15.3793	18.4408	33.5873
(I, GD)	4.5417	8.5777	12.4865	16.3524	20.1987	39.3295
(I, GI)	4.0010	7.1857	10.2714	13.3219	16.3555	31.4326
(N, GD)	3.2822	7.2033	11.0527	14.8822	18.7035	37.7739
(N, GI)	2.2800	5.2800	8.2800	11.2800	14.2800	29.2800

the conditional expected amounts of work brought into the system per unit time given the state of the underlying Markov chain are different, and they are given by 1.28 and 0.32, respectively. Thus in Case (N, GI), the system is overloaded during a half of time, so that $E[N]$ in Case (N, GI) is greater than that in Case (N, GD).

We observe that in any case, the expected total number of customer increases with the mean batch size g , as in Example 1, and that $E[N]$ in Case GD eventually becomes greater than $E[N]$ in Case GI. We also observe that for a fixed g , the positive correlation in the two streams leads to a larger $E[N]$ in both Cases GD and GI, as in Example 1.

2.7 Concluding Remarks

We developed a numerically feasible procedure to compute the joint queue length distribution in a FIFO single-server queue with multiple batch Markovian arrival streams having different service time distributions, assuming discrete phase-type batch size distributions. We established several

Table 2.5: Expected total number of customers in Example 2.

Case	$g = 1$	$g = 2$	$g = 3$	$g = 4$	$g = 5$	$g = 10$
(P, GD)	11.5019	17.7712	23.8347	29.8053	35.7261	65.0310
(P, GI)	11.5019	15.9366	20.1933	24.3657	28.4904	48.8117
(I, GD)	7.1517	13.3270	19.3007	25.2050	31.0760	60.2474
(I, GI)	8.7304	13.1052	17.3093	21.4407	25.5333	45.7640
(N, GD)	3.2168	9.0326	14.8425	20.6497	26.4551	55.4705
(N, GI)	6.0892	10.3399	14.4641	18.5407	22.5933	42.7206

truncation and stopping criteria to ensure numerical accuracy in the final result.

Note, however, that the computation of the joint queue length distribution is intensive by nature, especially when the number of classes is large. Even in such a case, the stationary total queue length distribution can be readily computed by modifying our algorithm. For the sake of completeness, we show algorithm steps for the total queue length distribution in Appendix C. Note here that the algorithm to compute $\mathbf{A}_k^{(T)}(n)$ in (C.1) for the number of arrivals in a service time can be also used for the computation of \mathbf{A}_n , which is a matrix for the number of arrivals per service in an ordinary BMAP/GI/1 queue (see Appendix B), because $\mathbf{A}_k^{(T)}(n)$ ($k \in \mathcal{K}$) in the case of $\mathcal{K} = \{1\}$ corresponds to \mathbf{A}_n in the BMAP/GI/1 queue (see (B.2) in Appendix B and (C.2) in Appendix C). To the best of our knowledge, there is no work to consider the truncation and stopping criterion to compute \mathbf{A}_n in the BMAP/GI/1 queue. Thus our development also contributes to the standard algorithm for the BMAP/GI/1 queue.

Chapter 3

FIFO Single-Server Queue with Service Interruptions

3.1 Introduction

This chapter considers a FIFO single-server queue with service interruptions. In such a queue, the state of the server changes *on* and *off* alternately. While in on-state, the server is available for service. On the other hand, while in off-state, the server does not work even if customers are present in the system. Hereafter periods during which the server is in on-state (resp. off-state) are called on-periods (resp. off-periods).

Queues with service interruptions have many applications in the fields of manufacturing, computer and telecommunications systems, and many studies on those queues have been done for a few decades. A detailed survey on queues with service interruptions can be found in the introduction of the paper by Federgruen and Green [Fede86]. They mainly discussed approximation methods for an M/G/1 queue with service interruptions, where on- and off-periods are generally distributed [Fede86]. Further, assuming a phase-type on-period distribution, they established an exact algorithm to compute the steady-state queue length distribution [Fede88].

Recently, more general queues with service interruptions have been studied. Sengupta [Seng90] considered the model where on- and off-period distributions are general, customers arrive according to a Poisson process whose arrival rate depends on the server state, and service times are generally distributed, depending on the server state upon arrival. He showed that the amount of unfinished work in such a queue is closely related to the waiting time in a special GI/G/1 queue. Also, Takine and Sengupta [Taki97] considered a single-server queue with service interruptions, where both the server state and arrival processes are governed by a finite-state Markov chain. Namely, the marginal processes of the server state and customer arrivals form an alternating phase-type Markov renewal process and a MAP (see subsection 1.1.1), respectively, and they may be dependent. For this queue, they obtained the steady-state queue length distribution. The crucial assumption posed in [Taki97] was i.i.d. (independent and identically distributed) service times.

This chapter considers an extension of the results in [Taki97] to allow multiple batch Markovian arrival streams. Thus, the marginal arrival process follows a batch MMAP (see section 2.2).

Further, service times of customers can depend on both their arrival stream and the server state on arrival. As stated in [Taki97], such a queue cannot be analyzed by the conventional M/G/1 paradigm (see subsection 1.2.2). To analyze the extended model, therefore, we use a new approach developed in [Masu03a, Taki01a, Taki01b, Taki01c], which is based on the invariant relationship of the joint queue length distributions at a random point in time and at departures [Taki01a]. We then derive the vector joint generating function of the numbers of customers from respective arrival streams. Further assuming discrete phase-type batch size distributions as in subsection 2.4.3, we provide a computational algorithm for the steady-state joint queue length distribution. We also show some numerical examples and examine the impact of system parameters on the queue length distribution.

The rest of this chapter is divided into four sections. In section 2, the mathematical model is described. Section 3 briefly discusses the sojourn time distribution. In section 4, we first derive a general formula for the joint queue length distribution, and assuming discrete phase-type batch size distributions, we show recursive formulas to compute the joint queue length distribution. Finally, in section 5, we show some numerical examples.

3.2 Model

We consider a FIFO single-server queue with service interruptions. The state of the server changes on and off alternately, and while the server is being on, customers are served successively. On the other hand, services of customers stop temporarily while the server is being off, and interrupted services are restarted in a preemptive-resume manner when the server becomes on again. In what follows, we call the process of the server state the on-off process.

We assume that both the on-off and arrival processes are governed by an underlying finite-state Markov chain that is assumed to be irreducible. Let $\mathcal{M} = \{1, \dots, M\}$ denote the state space of the underlying Markov chain, where $M \geq 2$. It stays in state i ($i \in \mathcal{M}$) for an exponential interval of time with mean μ_i^{-1} , and when the sojourn time in state i has elapsed, the underlying Markov chain changes its state to state j with probability $\sigma_{i,j}$ ($j \in \mathcal{M}$), where

$$\sum_{j \in \mathcal{M}} \sigma_{i,j} = 1.$$

The on-off process of the server is defined in the following way. The state space \mathcal{M} is divided into two disjoint sub-spaces, $\mathcal{M}_{\text{on}} = \{1, \dots, M_{\text{on}}\}$ and $\mathcal{M}_{\text{off}} = \{M_{\text{on}} + 1, \dots, M_{\text{on}} + M_{\text{off}}\}$, where $M_{\text{on}} \geq 1$, $M_{\text{off}} \geq 1$ and $M_{\text{on}} + M_{\text{off}} = M$. The server is assumed to be on (resp. off) while the underlying Markov chain is being in state $i \in \mathcal{M}_{\text{on}}$ (resp. $i \in \mathcal{M}_{\text{off}}$). Thus the on-off process forms a phase-type alternating Markov renewal process.

Next we describe the arrival process of customers. We assume that there are K ($K \geq 1$) arrival streams. Let \mathcal{K} denote a set of class indices, i.e., $\mathcal{K} = \{1, \dots, K\}$. Customers arriving from the k th ($k \in \mathcal{K}$) arrival stream are called class k customers. Given a state transition of the underlying Markov chain from state i to state j ($i, j \in \mathcal{M}$), n ($n = 1, 2, \dots$) customers of class k ($k \in \mathcal{K}$) arrive

in batch with probability $\sigma_{k,i,j}(n)/\sigma_{i,j}$, where

$$\sum_{k \in \mathcal{K}} \sum_{n=1}^{\infty} \sigma_{k,i,j}(n) \leq \sigma_{i,j},$$

for all $i, j \in \mathcal{M}$. Note that customers in the same batch belong to the same class. For later use, we define $\sigma_{i,j}(0)$ ($i, j \in \mathcal{M}$) as

$$\sigma_{i,j}(0) = \sigma_{i,j} - \sum_{k \in \mathcal{K}} \sum_{n=1}^{\infty} \sigma_{k,i,j}(n).$$

Note that $\sigma_{i,j}(0)/\sigma_{i,j}$ represents the conditional probability of no arrivals given that a state transition from state i to state j happens. Without loss of generality, we assume $\sigma_{i,i}(0) = 0$ for all i ($i \in \mathcal{M}$).

In terms of service times, class k customers are further classified into two sub-classes based on the server state on arrival. We call class k customers arriving in on-periods (resp. off-periods) class k -on (resp. k -off) customers. Service times of class k -on (resp. k -off) customers are assumed to be i.i.d. according to a distribution function $H_{k,\text{on}}(x)$ (resp. $H_{k,\text{off}}(x)$) with finite mean $h_{k,\text{on}}$ (resp. $h_{k,\text{off}}$).

In the rest of this chapter, we impose two assumptions on $\sigma_{k,i,j}(n)$. For each k ($k \in \mathcal{K}$), there exists at least one triad (i, j, n) ($i, j \in \mathcal{M}$, $n = 1, 2, \dots$) such that $\sigma_{k,i,j}(n) > 0$. Thus arrivals of class k customers are certain. Further $\sigma_{k,i,j}(n) = 0$ ($k \in \mathcal{K}$) if $i \in \mathcal{M}_{\text{on}}$ and $j \in \mathcal{M}_{\text{off}}$ or if $i \in \mathcal{M}_{\text{off}}$ and $j \in \mathcal{M}_{\text{on}}$. Thus arrivals of customers and changes of the server state never happen simultaneously.

We now introduce some notations. Let \mathbf{C} denote an $M \times M$ matrix whose (i, j) th ($i, j \in \mathcal{M}$) element $C_{i,j}$ is given by

$$C_{i,j} = \begin{cases} -\mu_i, & \text{if } i = j, \\ \sigma_{i,j}(0)\mu_i, & \text{otherwise.} \end{cases}$$

For each $k \in \mathcal{K}$, let $\mathbf{D}_k(n)$ ($n = 1, 2, \dots$) denote an $M \times M$ matrix whose (i, j) th ($i, j \in \mathcal{M}$) element $D_{k,i,j}(n)$ is given by

$$D_{k,i,j}(n) = \begin{cases} \sigma_{k,i,j}(n)\mu_i, & \text{if } i, j \in \mathcal{M}_{\text{on}} \text{ or } i, j \in \mathcal{M}_{\text{off}}, \\ 0, & \text{otherwise.} \end{cases}$$

Then the on-off and arrival processes are characterized by \mathbf{C} and $\mathbf{D}_k(n)$ ($k \in \mathcal{K}, n = 1, 2, \dots$), and the marginal arrival process is a batch MMAP (see section 2.2). Note here that \mathbf{C} and $\mathbf{D}_k(n)$ have the following structure:

$$\mathbf{C} = \begin{bmatrix} \mathbf{C}_{\text{on}} & \mathbf{E}_{\text{on,off}} \\ \mathbf{E}_{\text{off,on}} & \mathbf{C}_{\text{off}} \end{bmatrix}, \quad \mathbf{D}_k(n) = \begin{bmatrix} \mathbf{D}_{k,\text{on}}(n) & \mathbf{O} \\ \mathbf{O} & \mathbf{D}_{k,\text{off}}(n) \end{bmatrix},$$

where \mathbf{C}_{on} and $\mathbf{D}_{k,\text{on}}(n)$ are $M_{\text{on}} \times M_{\text{on}}$ matrices, \mathbf{C}_{off} and $\mathbf{D}_{k,\text{off}}(n)$ are $M_{\text{off}} \times M_{\text{off}}$ matrices, and $\mathbf{E}_{\text{on,off}}$ and $\mathbf{E}_{\text{off,on}}$ are $M_{\text{on}} \times M_{\text{off}}$ and $M_{\text{off}} \times M_{\text{on}}$ matrices, respectively.

We define $\mathbf{D}_{k,\xi}$ ($k \in \mathcal{K}$, $\xi = \text{on, off}$) and $\bar{\mathbf{D}}_\xi$ ($\xi = \text{on, off}$) as

$$\mathbf{D}_{k,\xi} = \sum_{n=1}^{\infty} n \mathbf{D}_{k,\xi}(n), \quad \bar{\mathbf{D}}_\xi = \sum_{k \in \mathcal{K}} \mathbf{D}_{k,\xi},$$

respectively. We also define \mathbf{D} as

$$\mathbf{D} = \begin{bmatrix} \bar{\mathbf{D}}_{\text{on}} & \mathbf{O} \\ \mathbf{O} & \bar{\mathbf{D}}_{\text{off}} \end{bmatrix}.$$

Note that the infinitesimal generator of the underlying Markov chain is given by $\mathbf{C} + \mathbf{D}$. Note also that $(\mathbf{C} + \mathbf{D})\mathbf{e} = \mathbf{0}$, where \mathbf{e} denotes a column vector whose elements are all equal to one. We denote the stationary probability vector of the underlying Markov chain by $\boldsymbol{\pi}$. Because of the finite state space \mathcal{M} and the irreducibility of the underlying Markov chain, $\boldsymbol{\pi}$ is uniquely determined so as to satisfy $\boldsymbol{\pi}(\mathbf{C} + \mathbf{D}) = \mathbf{0}$ and $\boldsymbol{\pi}\mathbf{e} = 1$. Let $\boldsymbol{\pi}_{\text{on}}$ (resp. $\boldsymbol{\pi}_{\text{off}}$) denote a $1 \times M_{\text{on}}$ (resp. $1 \times M_{\text{off}}$) vector representing the conditional stationary probability vector of the underlying Markov chain given that the server is on (resp. off). Note that $\boldsymbol{\pi}_{\text{on}}$ and $\boldsymbol{\pi}_{\text{off}}$ satisfy

$$\begin{aligned} \boldsymbol{\pi}_{\text{on}} \left[\mathbf{C}_{\text{on}} + \bar{\mathbf{D}}_{\text{on}} + \mathbf{E}_{\text{on,off}} \left[-\mathbf{C}_{\text{off}} - \bar{\mathbf{D}}_{\text{off}} \right]^{-1} \mathbf{E}_{\text{off,on}} \right] &= \mathbf{0}, & \boldsymbol{\pi}_{\text{on}}\mathbf{e} &= 1, \\ \boldsymbol{\pi}_{\text{off}} \left[\mathbf{C}_{\text{off}} + \bar{\mathbf{D}}_{\text{off}} + \mathbf{E}_{\text{off,on}} \left[-\mathbf{C}_{\text{on}} - \bar{\mathbf{D}}_{\text{on}} \right]^{-1} \mathbf{E}_{\text{on,off}} \right] &= \mathbf{0}, & \boldsymbol{\pi}_{\text{off}}\mathbf{e} &= 1, \end{aligned}$$

respectively. Let r_{on} and r_{off} denote fractions of time being in on- and off-periods, respectively. We then have

$$r_{\text{on}} = \frac{\bar{I}_{\text{on}}}{\bar{I}_{\text{on}} + \bar{I}_{\text{off}}}, \quad r_{\text{off}} = \frac{\bar{I}_{\text{off}}}{\bar{I}_{\text{on}} + \bar{I}_{\text{off}}},$$

where \bar{I}_{on} and \bar{I}_{off} denote the mean lengths of on- and off-periods, respectively, and they are given by

$$\bar{I}_{\text{on}} = \frac{\boldsymbol{\pi}_{\text{off}} \mathbf{E}_{\text{off,on}}}{\boldsymbol{\pi}_{\text{off}} \mathbf{E}_{\text{off,on}} \mathbf{e}} \left[-\mathbf{C}_{\text{on}} - \bar{\mathbf{D}}_{\text{on}} \right]^{-1} \mathbf{e}, \quad \bar{I}_{\text{off}} = \frac{\boldsymbol{\pi}_{\text{on}} \mathbf{E}_{\text{on,off}}}{\boldsymbol{\pi}_{\text{on}} \mathbf{E}_{\text{on,off}} \mathbf{e}} \left[-\mathbf{C}_{\text{off}} - \bar{\mathbf{D}}_{\text{off}} \right]^{-1} \mathbf{e}.$$

Note here that $\boldsymbol{\pi}$, $\boldsymbol{\pi}_{\text{on}}$ and $\boldsymbol{\pi}_{\text{off}}$ are related by

$$\boldsymbol{\pi} = (r_{\text{on}} \boldsymbol{\pi}_{\text{on}}, r_{\text{off}} \boldsymbol{\pi}_{\text{off}}).$$

We denote the mean arrival rate of class k ($k \in \mathcal{K}$) customers during on- (resp. off-) periods by $\lambda_{k,\text{on}}$ (resp. $\lambda_{k,\text{off}}$):

$$\lambda_{k,\xi} = \boldsymbol{\pi}_\xi \sum_{n=1}^{\infty} n \mathbf{D}_{k,\xi}(n) \mathbf{e}, \quad \xi = \text{on, off}.$$

Let $\lambda_k = r_{\text{on}} \lambda_{k,\text{on}} + r_{\text{off}} \lambda_{k,\text{off}}$ ($k \in \mathcal{K}$) denote the mean arrival rate of class k customers. We define ρ as the offered load, i.e.,

$$\rho = r_{\text{on}} \sum_{k \in \mathcal{K}} \lambda_{k,\text{on}} h_{k,\text{on}} + r_{\text{off}} \sum_{k \in \mathcal{K}} \lambda_{k,\text{off}} h_{k,\text{off}}.$$

Further, we define ρ_{on} as the conditional utilization factor given that the server is on, which is given by

$$\rho_{\text{on}} = r_{\text{on}}^{-1} \rho.$$

In the remainder of this chapter, we assume that $\rho_{\text{on}} < 1$, and the system is in steady state.

3.3 Sojourn Time

This section considers sojourn time. Because sojourn time is closely related to the amount of unfinished work in the system, we first discuss the latter.

Let V and S denote generic random variables representing the amount of unfinished work and the state of the underlying Markov chain, respectively, in steady state. With these, we define $\mathbf{v}(x)$ as a $1 \times M$ vector whose j th element represents $\Pr[V \leq x, S = j]$. Further we define $\mathbf{v}_\xi(x)$ ($\xi = \text{on}, \text{off}$) as a $1 \times M_\xi$ vector whose j th element represents $\Pr[V \leq x, S = j \mid S \in \mathcal{M}_\xi]$. Note here that $\mathbf{v}(x)$ is given in terms of $\mathbf{v}_{\text{on}}(x)$ and $\mathbf{v}_{\text{off}}(x)$:

$$\mathbf{v}(x) = (r_{\text{on}}\mathbf{v}_{\text{on}}(x), r_{\text{off}}\mathbf{v}_{\text{off}}(x)).$$

Thus, we derive the LST $\mathbf{v}_\xi^*(s)$ ($\xi = \text{on}, \text{off}$) of $\mathbf{v}_\xi(x)$ below.

We first introduce some notations. Let $\mathbf{J}_{\text{off,on}}(x)$ denote an $M_{\text{off}} \times M_{\text{on}}$ matrix whose (i, j) th element represents the probability that the amount of work arriving during an off-period is not greater than x and the underlying Markov chain is in state $j \in \mathcal{M}_{\text{on}}$ at the beginning of the next on-period, given that the off-period starts in state $i \in \mathcal{M}_{\text{off}}$. We define $\overline{\mathbf{D}}_\xi(x)$ ($\xi = \text{on}, \text{off}$) as

$$\overline{\mathbf{D}}_\xi(x) = \sum_{k \in \mathcal{K}} \sum_{n=1}^{\infty} \mathbf{D}_{k,\xi}(n) H_{k,\xi}^{(n)}(x),$$

where $H_{k,\xi}^{(1)}(x) = H_{k,\xi}(x)$ and $H_{k,\xi}^{(n)}(x)$ ($n = 2, 3, \dots$) denotes the n -fold convolution of $H_{k,\xi}(x)$ with itself. Note here that the LST $\mathbf{J}_{\text{off,on}}^*(s)$ of $\mathbf{J}_{\text{off,on}}(x)$ is given by

$$\mathbf{J}_{\text{off,on}}^*(s) = \left[-\mathbf{C}_{\text{off}} - \overline{\mathbf{D}}_{\text{off}}^*(s) \right]^{-1} \mathbf{E}_{\text{off,on}},$$

where $\overline{\mathbf{D}}_\xi^*(s)$ ($\xi = \text{on}, \text{off}$) denotes the LST of $\overline{\mathbf{D}}_\xi(x)$. According to [Taki94a], we define an $M_{\text{on}} \times M_{\text{on}}$ matrix \mathbf{Q}_{on} as an infinitesimal generator of the irreducible Markov chain obtained by observing the underlying Markov chain only when the system is idle in on-periods. Note that \mathbf{Q}_{on} satisfies

$$\mathbf{Q}_{\text{on}} = \mathbf{C}_{\text{on}} + \int_0^\infty d\overline{\mathbf{D}}_{\text{on}}(x) \exp(\mathbf{Q}_{\text{on}}x) + \mathbf{E}_{\text{on,off}} \int_0^\infty d\mathbf{J}_{\text{off,on}}(x) \exp(\mathbf{Q}_{\text{on}}x).$$

Because $\rho_{\text{on}} < 1$, \mathbf{Q}_{on} is uniquely determined by the above equation [Taki94a]. Let $\boldsymbol{\kappa}_{\text{on}}$ denote a probability vector satisfying $\boldsymbol{\kappa}_{\text{on}}\mathbf{Q}_{\text{on}} = \mathbf{0}$.

We now derive $\mathbf{v}_{\text{on}}^*(s)$ in such a way as to obtain Proposition 1.8 in subsection 1.3.2. We consider the bivariate process $\{(V(t), S(t)); t \geq 0\}$, where $V(t)$ and $S(t)$ denote the virtual waiting time and the state of the underlying Markov chain, respectively, at time t . We then construct the censored process obtained by observing $\{(V(t), S(t)); t \geq 0\}$ only when $V(t) > 0$ and $S(t) \in \mathcal{M}_{\text{on}}$. It is easy to see that the resulting bivariate process is equivalent to $\{(X_L(t), J_L(t)); t \geq 0\}$ with

$$\begin{aligned} \mathbf{A}_{\text{slope}} &= \mathbf{C}_{\text{on}}, & \mathbf{A}_{\text{jump}}(x) &= \overline{\mathbf{D}}_{\text{on}}(x) + \mathbf{E}_{\text{on,off}}\mathbf{J}_{\text{off,on}}(x), \\ \mathbf{B}_{\text{jump}}(x) &= (-\mathbf{C}_{\text{on}})^{-1} \left[\overline{\mathbf{D}}_{\text{on}}(x) + \mathbf{E}_{\text{on,off}}\mathbf{J}_{\text{off,on}}(x) \right]. \end{aligned}$$

It then follows that the LST $\mathbf{v}_{\text{on}}^*(s)$ of $\mathbf{v}_{\text{on}}(x)$ satisfies [Taki94a]

$$\mathbf{v}_{\text{on}}^*(s) \left[s\mathbf{I} + \mathbf{C}_{\text{on}} + \overline{\mathbf{D}}_{\text{on}}^*(s) + \mathbf{E}_{\text{on,off}} \left[-\mathbf{C}_{\text{off}} - \overline{\mathbf{D}}_{\text{off}}^*(s) \right]^{-1} \mathbf{E}_{\text{off,on}} \right] = s(1 - \rho_{\text{on}})\boldsymbol{\kappa}_{\text{on}}, \quad (3.1)$$

where $\mathbf{I}(m)$ denotes an $m \times m$ identity matrix. We suppress the size m when it is clear from the context. As for the LST $\mathbf{v}_{\text{off}}^*(s)$ of $\mathbf{v}_{\text{off}}(x)$, using the same approach as in [Taki97], we readily obtain

$$\mathbf{v}_{\text{off}}^*(s) = \frac{\mathbf{v}_{\text{on}}^*(s) \mathbf{E}_{\text{on,off}}}{\boldsymbol{\pi}_{\text{on}} \mathbf{E}_{\text{on,off}} \mathbf{e}} \frac{[-\mathbf{C}_{\text{off}} - \overline{\mathbf{D}}_{\text{off}}^*(s)]^{-1}}{\overline{\mathbf{I}}_{\text{off}}}. \quad (3.2)$$

Next we analyze sojourn time. To do so, we first consider *completion time*, which is a time interval from the beginning of a service to its completion, including service interruptions. Let $T_c(u)$ denote a generic random variable representing the completion time of a service of u units. Note that the completion time $T_c(u)$ depends on the state of the underlying Markov chain at the beginning of the service, as well as the amount of the service. We denote the number of class k ($k \in \mathcal{K}$) customers arriving in interval $(0, t]$ by $L_k(t)$. Assuming that the service of u unit commences at time 0, we define $\mathbf{P}^{**}(\mathbf{z}, s | u)$ as an $M_{\text{on}} \times M_{\text{on}}$ matrix whose (i, j) th element represents

$$\mathbb{E} \left[\prod_{k \in \mathcal{K}} z_k^{L_k(T_c(u))} \exp(-s T_c(u)) \mathbf{1}\{S(T_c(u)) = j \in \mathcal{M}_{\text{on}}\} | S(0) = i \in \mathcal{M}_{\text{on}} \right],$$

where \mathbf{z} denotes a $1 \times K$ complex vector (z_1, \dots, z_K) and $\mathbf{1}\{\chi\}$ denotes an indicator function of event χ . Further, we define $\mathbf{D}_{k,\xi}^*(z_k)$ ($k \in \mathcal{K}$, $\xi = \text{on, off}$) as

$$\mathbf{D}_{k,\xi}^*(z_k) = \sum_{n=1}^{\infty} z_k^n \mathbf{D}_{k,\xi}(n). \quad (3.3)$$

Following an approach similar to [Taki97], we obtain $\mathbf{P}^{**}(\mathbf{z}, s | u) = \exp[\mathbf{K}(\mathbf{z}, s)u]$, where

$$\mathbf{K}(\mathbf{z}, s) = \mathbf{C}_{\text{on}} + \sum_{k \in \mathcal{K}} \mathbf{D}_{k,\text{on}}^*(z_k) - s \mathbf{I} + \mathbf{E}_{\text{on,off}} \left[s \mathbf{I} - \mathbf{C}_{\text{off}} - \sum_{k \in \mathcal{K}} \mathbf{D}_{k,\text{off}}^*(z_k) \right]^{-1} \mathbf{E}_{\text{off,on}}.$$

We are now ready to discuss sojourn time.

Let W_k ($k \in \mathcal{K}$) (resp. $W_{k,\xi}$ ($k \in \mathcal{K}$, $\xi = \text{on, off}$)) denote a generic random variable representing the sojourn time of a class k (resp. k - ξ) customer. Also let $W_{k,\xi}(n; m)$ ($k \in \mathcal{K}$, $\xi = \text{on, off}$, $n = 1, 2, \dots$, $m = 1, \dots, n$) denote a generic random variable representing the sojourn time of a randomly chosen class k - ξ customer who is a member of a batch of size n and the m th served customer among members in the same batch. For convenience, we assume that if $\lambda_{k,\xi} = 0$, $W_{k,\xi} = 0$, and if $\mathbf{D}_{k,\xi}(n) = \mathbf{O}$ for some n ($n \geq 1$), $W_{k,\xi}(n; m) = 0$ for all m ($m = 1, \dots, n$). Further, let $w_k^*(s)$, $w_{k,\xi}^*(s)$ and $w_{k,\xi}^*(s | n; m)$ denote the LSTs of the distributions of W_k , $W_{k,\xi}$ and $W_{k,\xi}(n; m)$, respectively. Because a randomly chosen departing customer of class k ($k \in \mathcal{K}$) belongs to class k - ξ ($\xi = \text{on, off}$) with probability $r_\xi \lambda_{k,\xi} / \lambda_k$, we obtain

$$w_k^*(s) = \frac{r_{\text{on}} \lambda_{k,\text{on}}}{\lambda_k} w_{k,\text{on}}^*(s) + \frac{r_{\text{off}} \lambda_{k,\text{off}}}{\lambda_k} w_{k,\text{off}}^*(s).$$

Note here that

$$w_{k,\xi}^*(s) = \sum_{n=1}^{\infty} \frac{n \boldsymbol{\pi}_\xi \mathbf{D}_{k,\xi}(n) \mathbf{e}}{\lambda_{k,\xi}} \cdot \frac{1}{n} \sum_{m=1}^n w_{k,\xi}^*(s | n; m), \quad k \in \mathcal{K}, \xi = \text{on, off}, \quad (3.4)$$

if $\lambda_{k,\xi} > 0$, and otherwise $w_{k,\xi}^*(s) = 1$. Thus, in what follows, we consider $w_{k,\xi}^*(s | n; m)$ ($k \in \mathcal{K}$, $\xi = \text{on, off}$, $n = 1, 2, \dots$, $m = 1, \dots, n$).

Let $H_{k,\xi}(n; m)$ ($k \in \mathcal{K}$, $\xi = \text{on, off}$, $n = 1, 2, \dots$, $m = 1, \dots, n$) denote a generic random variable representing the service time of a randomly chosen class k - ξ customer who is a member of a batch of size n and the m th served customer among members of the same batch. Because $W_{k,\xi}(n; m) = W_{k,\xi}(n; 1) + \sum_{l=2}^m T_c(H_{k,\xi}(n; l))$ for $\xi = \text{on, off}$, $n = 1, 2, \dots$ and $m = 1, \dots, n$, we obtain for $n = 1, 2, \dots$ and $m = 1, \dots, n$,

$$w_{k,\text{on}}^*(s | n; m) = \int_0^\infty \frac{d\mathbf{v}_{\text{on}}(x) \mathbf{D}_{k,\text{on}}(n) \exp[\boldsymbol{\Omega}(s)x]}{\boldsymbol{\pi}_{\text{on}} \mathbf{D}_{k,\text{on}}(n) \mathbf{e}} \left[\int_0^\infty dH_{k,\text{on}}(y) \exp[\boldsymbol{\Omega}(s)y] \right]^m \mathbf{e}, \quad (3.5)$$

$$w_{k,\text{off}}^*(s | n; m) = \int_0^\infty \frac{d\mathbf{v}_{\text{off}}(x) \mathbf{D}_{k,\text{off}}(n) [s\mathbf{I} - \mathbf{C}_{\text{off}} - \overline{\mathbf{D}}_{\text{off}}]^{-1} \mathbf{E}_{\text{off, on}} \exp[\boldsymbol{\Omega}(s)x]}{\boldsymbol{\pi}_{\text{off}} \mathbf{D}_{k,\text{off}}(n) \mathbf{e}} \cdot \left[\int_0^\infty dH_{k,\text{off}}(y) \exp[\boldsymbol{\Omega}(s)y] \right]^m \mathbf{e}, \quad (3.6)$$

respectively, where $\boldsymbol{\Omega}(s) = \mathbf{K}(1, \dots, 1, s)$. Note here that the (i, j) th ($i, j \in \mathcal{M}_{\text{on}}$) element of $\exp[\boldsymbol{\Omega}(s)x]$ represents $\mathbb{E}[\exp(-sT_c(x)) \mathbf{1}\{S(T_c(x)) = j\} \mid \text{a service of } x \text{ units starts at time 0 and } S(0) = i]$. Thus from (3.4) and (3.5), we obtain for $k \in \mathcal{K}$,

$$w_{k,\text{on}}^*(s) = \frac{1}{\lambda_{k,\text{on}}} \sum_{n=1}^{\infty} \int_0^\infty d\mathbf{v}_{\text{on}}(x) \mathbf{D}_{k,\text{on}}(n) \exp[\boldsymbol{\Omega}(s)x] \sum_{m=1}^n \left[\int_0^\infty dH_{k,\text{on}}(y) \exp[\boldsymbol{\Omega}(s)y] \right]^m \mathbf{e},$$

if $\lambda_{k,\text{on}} > 0$, and otherwise $w_{k,\text{on}}^*(s) = 1$. Similarly, from (3.4) and (3.6), we obtain for $k \in \mathcal{K}$,

$$w_{k,\text{off}}^*(s) = \frac{1}{\lambda_{k,\text{off}}} \sum_{n=1}^{\infty} \int_0^\infty d\mathbf{v}_{\text{off}}(x) \mathbf{D}_{k,\text{off}}(n) [s\mathbf{I} - \mathbf{C}_{\text{off}} - \overline{\mathbf{D}}_{\text{off}}]^{-1} \mathbf{E}_{\text{off, on}} \exp[\boldsymbol{\Omega}(s)x] \cdot \sum_{m=1}^n \left[\int_0^\infty dH_{k,\text{off}}(y) \exp[\boldsymbol{\Omega}(s)y] \right]^m \mathbf{e},$$

if $\lambda_{k,\text{off}} > 0$, and otherwise $w_{k,\text{off}}^*(s) = 1$.

3.4 Joint Queue Length Distribution

This section considers the joint queue length distribution. Let N_k ($k \in \mathcal{K}$) denote a generic random variable representing the number of class k customers in the stationary system. We then define $\mathbf{p}(\mathbf{n})$ ($\mathbf{n} \in \mathcal{Z}$) as a $1 \times M$ vector whose j th element represents $\Pr[N_1 = n_1, \dots, N_K = n_K, S = j]$, where \mathbf{n} denotes a $1 \times K$ nonnegative integer vector (n_1, \dots, n_K) and $\mathcal{Z} = \{(n_1, \dots, n_K); n_k = 0, 1, \dots \text{ for all } k \in \mathcal{K}\}$. Further, let $N_\nu^{(\text{D}_k)}$ ($k, \nu \in \mathcal{K}$) and $S^{(\text{D}_k)}$ ($k \in \mathcal{K}$) denote generic random variables that represent the number of class ν customers in the system and the state of the underlying Markov chain, respectively, immediately after departures of class k customers in steady state. We then define $\mathbf{q}_k(\mathbf{n})$ ($k \in \mathcal{K}$, $\mathbf{n} \in \mathcal{Z}$) as a $1 \times M$ vector whose j th element represents $\Pr[N_1^{(\text{D}_k)} = n_1, \dots, N_K^{(\text{D}_k)} = n_K, S^{(\text{D}_k)} = j]$.

Note here that the $\mathbf{p}(\mathbf{n})$ is given in terms of the $\mathbf{q}_k(\mathbf{n})$ (see Corollary 2.1 in chapter 2). Thus, we consider the $\mathbf{q}_k(\mathbf{n})$ in section 4.1. Further, in section 4.2, assuming discrete phase-type batch size distributions, we derive numerically feasible recursions for some quantities required in computing the $\mathbf{q}_k(\mathbf{n})$.

3.4.1 Joint queue length distribution immediately after departures

We define $\mathbf{q}_k^*(\mathbf{z})$ ($k \in \mathcal{K}$) as the vector generating function of the joint queue length distribution immediately after departures of class k customers.

$$\mathbf{q}_k^*(\mathbf{z}) = \sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \mathbf{q}_k(\mathbf{n}), \quad |z_k| \leq 1 \text{ for all } k \in \mathcal{K}.$$

Let $N_\nu^{(D_{k,\xi})}$ and $S^{(D_{k,\xi})}$ ($k, \nu \in \mathcal{K}$, $\xi = \text{on, off}$) denote generic random variables that represent the number of class ν customers in the system and the state of the underlying Markov chain, respectively, immediately after departures of class k - ξ customers in steady state. We define $\mathbf{q}_{k,\xi}(\mathbf{n})$ ($k \in \mathcal{K}$, $\xi = \text{on, off}$) as a $1 \times M_{\text{on}}$ vector whose j th element represents $\Pr[N_1^{(D_{k,\xi})} = n_1, \dots, N_K^{(D_{k,\xi})} = n_K, S^{(D_{k,\xi})} = j]$. We also define $\mathbf{q}_{k,\xi}^*(\mathbf{z})$ ($k \in \mathcal{K}$, $\xi = \text{on, off}$) as

$$\mathbf{q}_{k,\xi}^*(\mathbf{z}) = \sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \mathbf{q}_{k,\xi}(\mathbf{n}).$$

We then have

$$\mathbf{q}_k^*(\mathbf{z}) = \left(\frac{r_{\text{on}} \lambda_{k,\text{on}}}{\lambda_k} \mathbf{q}_{k,\text{on}}^*(\mathbf{z}) + \frac{r_{\text{off}} \lambda_{k,\text{off}}}{\lambda_k} \mathbf{q}_{k,\text{off}}^*(\mathbf{z}), 0, \dots, 0 \right), \quad (3.7)$$

because all departures always occur in on-periods.

In the rest of this subsection, we derive $\mathbf{q}_{k,\text{on}}^*(\mathbf{z})$ and $\mathbf{q}_{k,\text{off}}^*(\mathbf{z})$. We call a randomly chosen class k - ξ ($k \in \mathcal{K}$, $\xi = \text{on, off}$) customer who is a member of a batch of size n and the m th served customer among members in the same batch *the tagged customer*. Further we call the batch to which the tagged customer belongs *the tagged batch*. Let $N_\nu^{(D_{k,\xi})}(n; m)$ and $S^{(D_{k,\xi})}(n; m)$ ($k, \nu \in \mathcal{K}$, $\xi = \text{on, off}$, $n = 1, 2, \dots$, $m = 1, \dots, n$) denote generic random variables that represent the number of class ν customers in the system and the state of the underlying Markov chain, respectively, immediately after the departure of the tagged customer. We then define $\mathbf{q}_{k,\xi}^*(\mathbf{z} | n; m)$ ($k \in \mathcal{K}$, $\xi = \text{on, off}$, $n = 1, 2, \dots$, $m = 1, \dots, n$) as a $1 \times M_{\text{on}}$ vector whose j th element represents $\mathbb{E} \left[\prod_{\nu \in \mathcal{K}} z_\nu^{N_\nu^{(D_{k,\xi})}(n; m)} \mathbf{1}_{\{S^{(D_{k,\xi})}(n; m) = j\}} \right]$. It is easy to see that $\mathbf{q}_{k,\xi}^*(\mathbf{z})$ can be written in terms of $\mathbf{q}_{k,\xi}^*(\mathbf{z} | n; m)$:

$$\mathbf{q}_{k,\xi}^*(\mathbf{z}) = \sum_{n=1}^{\infty} \frac{n \boldsymbol{\pi}_\xi \mathbf{D}_{k,\xi}(n) \mathbf{e}}{\lambda_{k,\xi}} \frac{1}{n} \sum_{m=1}^n \mathbf{q}_{k,\xi}^*(\mathbf{z} | n; m), \quad k \in \mathcal{K}, \xi = \text{on, off}, \quad (3.8)$$

if $\lambda_{k,\xi} > 0$, and otherwise $\mathbf{q}_{k,\xi}^*(\mathbf{z}) = \mathbf{0}$.

Note here that customers who contribute to $\mathbf{q}_{k,\text{on}}^*(\mathbf{z} | n; m)$ can be divided into three types: (i) customers arriving during the completion time of the total unfinished work immediately before the arrival of the tagged batch, (ii) customers arriving during an interval from the beginning of the first service of a member in the tagged batch to the completion of the service of the tagged customer, and (iii) $n - m$ customers who belong to the tagged batch and receive their services after the tagged customer. It then follows that for $k \in \mathcal{K}$, $n = 1, 2, \dots$ and $m = 1, \dots, n$,

$$\mathbf{q}_{k,\text{on}}^*(\mathbf{z} | n; m) = z_k^{n-m} \int_0^\infty \frac{d\mathbf{v}_{\text{on}}(x) \mathbf{D}_{k,\text{on}}(n) \mathbf{N}^*(\mathbf{z} | x)}{\boldsymbol{\pi}_{\text{on}} \mathbf{D}_{k,\text{on}}(n) \mathbf{e}} \left[\int_0^\infty dH_{k,\text{on}}(y) \mathbf{N}^*(\mathbf{z} | y) \right]^m, \quad (3.9)$$

where $\mathbf{N}^*(\mathbf{z} | x) = \mathbf{P}^{**}(\mathbf{z}, 0 | x)$, i.e.,

$$\begin{aligned} \mathbf{N}^*(\mathbf{z} | x) = \exp & \left[\left\{ \mathbf{C}_{\text{on}} + \sum_{k \in \mathcal{K}} \mathbf{D}_{k,\text{on}}^*(z_k) \right. \right. \\ & \left. \left. + \mathbf{E}_{\text{on,off}} \left(-\mathbf{C}_{\text{off}} - \sum_{k \in \mathcal{K}} \mathbf{D}_{k,\text{off}}^*(z_k) \right)^{-1} \mathbf{E}_{\text{off,on}} \right\} x \right]. \end{aligned} \quad (3.10)$$

Note that $\mathbf{N}^*(\mathbf{z} | x)$ denotes the matrix joint generating function for the numbers of arrivals in respective classes during the completion time of a service of x units. Similarly, we have for $k \in \mathcal{K}$, $n = 1, 2, \dots$ and $m = 1, \dots, n$,

$$\begin{aligned} \mathbf{q}_{k,\text{off}}^*(\mathbf{z} | n; m) & \\ = z_k^{n-m} \int_0^\infty \frac{d\mathbf{v}_{\text{off}}(x) \mathbf{D}_{k,\text{off}}(n)}{\boldsymbol{\pi}_{\text{off}} \mathbf{D}_{k,\text{off}}(n) e} & \left[-\mathbf{C}_{\text{off}} - \sum_{k \in \mathcal{K}} \mathbf{D}_{k,\text{off}}^*(z_k) \right]^{-1} \mathbf{E}_{\text{off,on}} \mathbf{N}^*(\mathbf{z} | x) \\ \cdot \left[\int_0^\infty dH_{k,\text{off}}(y) \mathbf{N}^*(\mathbf{z} | y) \right]^m & . \end{aligned} \quad (3.11)$$

Thus, from (3.8), (3.9) and (3.11), we obtain the following theorem.

Theorem 3.1 $\mathbf{q}_{k,\text{on}}^*(\mathbf{z})$ ($k \in \mathcal{K}$) is given by

$$\mathbf{q}_{k,\text{on}}^*(\mathbf{z}) = \frac{1}{\lambda_{k,\text{on}}} \sum_{m=1}^{\infty} \sum_{l=0}^{\infty} z_k^l \int_0^\infty d\mathbf{v}_{\text{on}}(x) \mathbf{D}_{k,\text{on}}(m+l) \mathbf{N}^*(\mathbf{z} | x) \left[\int_0^\infty dH_{k,\text{on}}(y) \mathbf{N}^*(\mathbf{z} | y) \right]^m,$$

if $\lambda_{k,\text{on}} > 0$, and otherwise $\mathbf{q}_{k,\text{on}}^*(\mathbf{z}) = \mathbf{0}$. On the other hand, $\mathbf{q}_{k,\text{off}}^*(\mathbf{z})$ ($k \in \mathcal{K}$) is given by

$$\begin{aligned} \mathbf{q}_{k,\text{off}}^*(\mathbf{z}) = \frac{1}{\lambda_{k,\text{off}}} \sum_{m=1}^{\infty} \sum_{l=0}^{\infty} z_k^l \int_0^\infty d\mathbf{v}_{\text{off}}(x) \mathbf{D}_{k,\text{off}}(m+l) & \left[-\mathbf{C}_{\text{off}} - \sum_{k \in \mathcal{K}} \mathbf{D}_{k,\text{off}}^*(z_k) \right]^{-1} \mathbf{E}_{\text{off,on}} \\ \cdot \mathbf{N}^*(\mathbf{z} | x) \left[\int_0^\infty dH_{k,\text{off}}(y) \mathbf{N}^*(\mathbf{z} | y) \right]^m & , \end{aligned}$$

if $\lambda_{k,\text{off}} > 0$, and otherwise $\mathbf{q}_{k,\text{off}}^*(\mathbf{z}) = \mathbf{0}$.

3.4.2 Recursions for discrete phase-type batch sizes

In this subsection, we develop recursive formulas to compute the joint queue length distribution $\mathbf{q}_k(\mathbf{n})$ immediately after departures of class k customers under the following assumption.

Assumption 3.1 Batch sizes of class k - ξ ($k \in \mathcal{K}$, $\xi = \text{on}, \text{off}$) are i.i.d. according to a discrete phase-type distribution with representation $(\boldsymbol{\alpha}_{k,\xi}, \mathbf{P}_{k,\xi})$, where $\boldsymbol{\alpha}_{k,\xi}$ denotes a $1 \times M_{k,\xi}$ probability vector and $\mathbf{P}_{k,\xi}$ denotes an $M_{k,\xi} \times M_{k,\xi}$ substochastic matrix.

Under Assumption 3.1, $\mathbf{D}_{k,\xi}(n)$ ($k \in \mathcal{K}$, $\xi = \text{on}, \text{off}$) is given by

$$\mathbf{D}_{k,\xi}(n) = g_{k,\xi}(n) \mathbf{D}_{k,\xi}, \quad n = 1, 2, \dots,$$

where $g_{k,\xi}(n)$ denotes the probability mass function of the batch size of class k - ξ :

$$g_{k,\xi}(n) = \boldsymbol{\alpha}_{k,\xi} \mathbf{P}_{k,\xi}^{n-1} (\mathbf{I} - \mathbf{P}_{k,\xi}) \mathbf{e}, \quad n = 1, 2, \dots$$

Thus, Theorem 3.1 is reduced to:

Corollary 3.1 *Under Assumption 3.1, $\mathbf{q}_{k,\text{on}}^*(\mathbf{z})$ ($k \in \mathcal{K}$) is given by*

$$\begin{aligned} \mathbf{q}_{k,\text{on}}^*(\mathbf{z}) &= \frac{1}{\lambda_{k,\text{on}}} \int_0^\infty d\mathbf{v}_{\text{on}}(x) \mathbf{D}_{k,\text{on}} \mathbf{N}^*(\mathbf{z} | x) \\ &\quad \cdot \left(\boldsymbol{\alpha}_{k,\text{on}} \otimes \int_0^\infty dH_{k,\text{on}}(y) \mathbf{N}^*(\mathbf{z} | y) \right) \left[\mathbf{I} - \mathbf{P}_{k,\text{on}} \otimes \int_0^\infty dH_{k,\text{on}}(y) \mathbf{N}^*(\mathbf{z} | y) \right]^{-1} \\ &\quad \cdot \left[\left\{ (\mathbf{I} - z_k \mathbf{P}_{k,\text{on}})^{-1} (\mathbf{I} - \mathbf{P}_{k,\text{on}}) \mathbf{e} \right\} \otimes \mathbf{I}(M_{\text{on}}) \right], \end{aligned} \quad (3.12)$$

if $\lambda_{k,\text{on}} > 0$, and otherwise $\mathbf{q}_{k,\text{on}}^*(\mathbf{z}) = \mathbf{0}$. Similarly, $\mathbf{q}_{k,\text{off}}^*(\mathbf{z})$ ($k \in \mathcal{K}$) is given by

$$\begin{aligned} \mathbf{q}_{k,\text{off}}^*(\mathbf{z}) &= \frac{1}{\lambda_{k,\text{off}}} \int_0^\infty d\mathbf{v}_{\text{off}}(x) \mathbf{D}_{k,\text{off}} \left[-\mathbf{C}_{\text{off}} - \sum_{k \in \mathcal{K}} \mathbf{D}_{k,\text{off}}^*(z_k) \right]^{-1} \mathbf{E}_{\text{off,on}} \mathbf{N}^*(\mathbf{z} | x) \\ &\quad \cdot \left(\boldsymbol{\alpha}_{k,\text{off}} \otimes \int_0^\infty dH_{k,\text{off}}(y) \mathbf{N}^*(\mathbf{z} | y) \right) \left[\mathbf{I} - \mathbf{P}_{k,\text{off}} \otimes \int_0^\infty dH_{k,\text{off}}(y) \mathbf{N}^*(\mathbf{z} | y) \right]^{-1} \\ &\quad \cdot \left[\left\{ (\mathbf{I} - z_k \mathbf{P}_{k,\text{off}})^{-1} (\mathbf{I} - \mathbf{P}_{k,\text{off}}) \mathbf{e} \right\} \otimes \mathbf{I}(M_{\text{on}}) \right], \end{aligned} \quad (3.13)$$

if $\lambda_{k,\text{off}} > 0$, and otherwise $\mathbf{q}_{k,\text{off}}^*(\mathbf{z}) = \mathbf{0}$.

This corollary can be obtained in the same way as Lemma 2.1 in chapter 2, and therefore we omit the proof.

We define $\mathbf{v}_{k,\text{on}}(\mathbf{n})$ and $\mathbf{v}_{k,\text{off}}(\mathbf{n})$ ($k \in \mathcal{K}, \mathbf{n} \in \mathcal{Z}$) as $1 \times M_{\text{on}}$ vectors satisfying

$$\sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \mathbf{v}_{k,\text{on}}(\mathbf{n}) = \int_0^\infty d\mathbf{v}_{\text{on}}(x) \mathbf{D}_{k,\text{on}} \mathbf{N}^*(\mathbf{z} | x), \quad (3.14)$$

$$\sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \mathbf{v}_{k,\text{off}}(\mathbf{n}) = \int_0^\infty d\mathbf{v}_{\text{off}}(x) \mathbf{D}_{k,\text{off}} \left[-\mathbf{C}_{\text{off}} - \sum_{k \in \mathcal{K}} \mathbf{D}_{k,\text{off}}^*(z_k) \right]^{-1} \mathbf{E}_{\text{off,on}} \mathbf{N}^*(\mathbf{z} | x), \quad (3.15)$$

respectively. We also define $\mathbf{A}_{k,\xi}(\mathbf{n})$ and $\boldsymbol{\Gamma}_{k,\xi}(\mathbf{n})$ ($k \in \mathcal{K}, \xi = \text{on}, \text{off}, \mathbf{n} \in \mathcal{Z}$) as $M_{\text{on}} \times M_{\text{on}}$ and $M_{k,\xi} M_{\text{on}} \times M_{k,\xi} M_{\text{on}}$ matrices satisfying

$$\begin{aligned} \sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \mathbf{A}_{k,\xi}(\mathbf{n}) &= \int_0^\infty dH_{k,\xi}(y) \mathbf{N}^*(\mathbf{z} | y), \\ \sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \boldsymbol{\Gamma}_{k,\xi}(\mathbf{n}) &= \left[\mathbf{I} - \mathbf{P}_{k,\xi} \otimes \int_0^\infty dH_{k,\xi}(y) \mathbf{N}^*(\mathbf{z} | y) \right]^{-1}, \end{aligned} \quad (3.16)$$

respectively. Then $\mathbf{q}_{k,\xi}^*(\mathbf{z})$ ($k \in \mathcal{K}, \xi = \text{on}, \text{off}$) in (3.12) and (3.13) are rewritten to be

$$\begin{aligned} \mathbf{q}_{k,\xi}^*(\mathbf{z}) &= \frac{1}{\lambda_{k,\xi}} \sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \sum_{m=0}^{n_k} \sum_{\substack{\mathbf{n}_1, \mathbf{n}_2, \mathbf{n}_3 \in \mathcal{Z} \\ \mathbf{n}_1 + \mathbf{n}_2 + \mathbf{n}_3 \\ = \mathbf{n} - m \mathbf{e}_k}} \mathbf{v}_{k,\xi}(\mathbf{n}_1) [\boldsymbol{\alpha}_{k,\xi} \otimes \mathbf{A}_{k,\xi}(\mathbf{n}_2)] \boldsymbol{\Gamma}_{k,\xi}(\mathbf{n}_3) \\ &\quad \cdot \left[\left\{ \mathbf{P}_{k,\xi}^m (\mathbf{I} - \mathbf{P}_{k,\xi}) \mathbf{e} \right\} \otimes \mathbf{I}(M_{\text{on}}) \right], \end{aligned} \quad (3.17)$$

if $\lambda_{k,\xi} > 0$, and otherwise $\mathbf{q}_{k,\xi}^*(\mathbf{z}) = \mathbf{0}$. Comparing coefficient vectors of $z_1 \cdots z_K$ on both sides of (3.7) and (3.17), respectively, we obtain the following result.

Theorem 3.2 *Under Assumption 3.1, the $\mathbf{q}_k(\mathbf{n})$ is given by*

$$\mathbf{q}_k(\mathbf{n}) = \left(\frac{r_{\text{on}} \lambda_{k,\text{on}}}{\lambda_k} \mathbf{q}_{k,\text{on}}(\mathbf{n}) + \frac{r_{\text{off}} \lambda_{k,\text{off}}}{\lambda_k} \mathbf{q}_{k,\text{off}}(\mathbf{n}), 0, \dots, 0 \right), \quad k \in \mathcal{K}, \mathbf{n} \in \mathcal{Z},$$

where the $\mathbf{q}_{k,\xi}(\mathbf{n})$ ($k \in \mathcal{K}$, $\xi = \text{on}, \text{off}$, $\mathbf{n} \in \mathcal{Z}$) is given by

$$\begin{aligned} \mathbf{q}_{k,\xi}(\mathbf{n}) &= \frac{1}{\lambda_{k,\xi}} \sum_{m=0}^{n_k} \sum_{\substack{\mathbf{n}_1, \mathbf{n}_2, \mathbf{n}_3 \in \mathcal{Z} \\ \mathbf{n}_1 + \mathbf{n}_2 + \mathbf{n}_3 \\ = \mathbf{n} - m \mathbf{e}_k}} \mathbf{v}_{k,\xi}(\mathbf{n}_1) [\boldsymbol{\alpha}_{k,\xi} \otimes \mathbf{A}_{k,\xi}(\mathbf{n}_2)] \boldsymbol{\Gamma}_{k,\xi}(\mathbf{n}_3) \\ &\quad \cdot \left[\left\{ \mathbf{P}_{k,\xi}^m (\mathbf{I} - \mathbf{P}_{k,\xi}) \mathbf{e} \right\} \otimes \mathbf{I}(M_{\text{on}}) \right], \end{aligned}$$

if $\lambda_{k,\xi} > 0$, and otherwise $\mathbf{q}_{k,\xi}(\mathbf{n}) = \mathbf{0}$ for all $\mathbf{n} \in \mathcal{Z}$.

Theorem 3.2 implies that the computation of the $\mathbf{q}_k(\mathbf{n})$ is reduced to those of the $\boldsymbol{\Gamma}_{k,\xi}(\mathbf{n})$, $\mathbf{A}_{k,\xi}(\mathbf{n})$ and $\mathbf{v}_{k,\xi}(\mathbf{n})$ ($\xi = \text{on}, \text{off}$). Note here that the $\boldsymbol{\Gamma}_{k,\xi}(\mathbf{n})$ is given in terms of the $\mathbf{A}_{k,\xi}(\mathbf{n})$.

Lemma 3.1 *The $\boldsymbol{\Gamma}_{k,\xi}(\mathbf{n})$ ($k \in \mathcal{K}$, $\xi = \text{on}, \text{off}$, $\mathbf{n} \in \mathcal{Z}$) is determined by the following recursion:*

$$\begin{aligned} \boldsymbol{\Gamma}_{k,\xi}(\mathbf{0}) &= [\mathbf{I} - \mathbf{P}_{k,\xi} \otimes \mathbf{A}_{k,\xi}(\mathbf{0})]^{-1}, \\ \boldsymbol{\Gamma}_{k,\xi}(\mathbf{n}) &= \sum_{\substack{\mathbf{0} \leq \mathbf{l} \leq \mathbf{n} \\ \mathbf{l} \neq \mathbf{0}}} \boldsymbol{\Gamma}_{k,\xi}(\mathbf{n} - \mathbf{l}) [\mathbf{P}_{k,\xi} \otimes \mathbf{A}_{k,\xi}(\mathbf{l})] \boldsymbol{\Gamma}_{k,\xi}(\mathbf{0}), \quad \mathbf{n} \in \mathcal{Z}^+. \end{aligned}$$

Proof. We can prove this lemma in the same way as Lemma 2.3. ■

The rest of this subsection therefore discusses the computations of the $\mathbf{A}_{k,\xi}(\mathbf{n})$ and the $\mathbf{v}_{k,\xi}(\mathbf{n})$. We first consider the $\mathbf{A}_{k,\xi}(\mathbf{n})$. Let $\bar{\mathbf{F}}_m(\mathbf{n})$ ($m = 0, 1, \dots$, $\mathbf{n} \in \mathcal{Z}$) denote an $M_{\text{on}} \times M_{\text{on}}$ matrix that satisfies

$$\begin{aligned} \sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \bar{\mathbf{F}}_m(\mathbf{n}) &= \left[\mathbf{I} + \theta_{\text{on}}^{-1} \left\{ \mathbf{C}_{\text{on}} + \sum_{k \in \mathcal{K}} \mathbf{D}_{k,\text{on}}^*(z_k) \right. \right. \\ &\quad \left. \left. + \mathbf{E}_{\text{on,off}} \left(-\mathbf{C}_{\text{off}} - \sum_{k \in \mathcal{K}} \mathbf{D}_{k,\text{off}}^*(z_k) \right)^{-1} \mathbf{E}_{\text{off,on}} \right\} \right]^m, \end{aligned} \quad (3.18)$$

where θ_{on} denotes the maximum absolute value of diagonal elements of \mathbf{C}_{on} .

Lemma 3.2 *The $\bar{\mathbf{F}}_m(\mathbf{n})$ is recursively determined by*

$$\bar{\mathbf{F}}_0(\mathbf{n}) = \begin{cases} \mathbf{I}, & \text{if } \mathbf{n} = \mathbf{0}, \\ \mathbf{O}, & \text{otherwise,} \end{cases} \quad (3.19)$$

and for $m = 1, 2, \dots$,

$$\begin{aligned} \bar{\mathbf{F}}_m(\mathbf{n}) &= \bar{\mathbf{F}}_{m-1}(\mathbf{n}) (\mathbf{I} + \theta_{\text{on}}^{-1} \mathbf{C}_{\text{on}}) + \theta_{\text{on}}^{-1} \sum_{k \in \mathcal{K}} \sum_{l_k=1}^{n_k} \bar{\mathbf{F}}_{m-1}(\mathbf{n} - l_k \mathbf{e}_k) \mathbf{D}_{k,\text{on}}(l_k) \\ &\quad + \theta_{\text{on}}^{-1} \left[\sum_{\mathbf{0} \leq \mathbf{l} \leq \mathbf{n}} \bar{\mathbf{F}}_{m-1}(\mathbf{n} - \mathbf{l}) \mathbf{E}_{\text{on,off}} \mathbf{N}_{\text{off}}(\mathbf{l}) \right] \mathbf{E}_{\text{off,on}}, \quad \mathbf{n} \in \mathcal{Z}, \end{aligned} \quad (3.20)$$

where $M_{\text{off}} \times M_{\text{off}}$ matrices $\mathbf{N}_{\text{off}}(\mathbf{n})$'s are determined by the following recursion:

$$\mathbf{N}_{\text{off}}(\mathbf{0}) = (-\mathbf{C}_{\text{off}})^{-1}, \quad (3.21)$$

$$\mathbf{N}_{\text{off}}(\mathbf{n}) = \left[\sum_{k \in \mathcal{K}} \sum_{l_k=1}^{n_k} \mathbf{N}_{\text{off}}(\mathbf{n} - l_k \mathbf{e}_k) \mathbf{D}_{k,\text{off}}(l_k) \right] \mathbf{N}_{\text{off}}(\mathbf{0}), \quad \mathbf{n} \in \mathcal{Z}^+. \quad (3.22)$$

The proof of Lemma 3.2 is given in Appendix D.

The $\mathbf{A}_{k,\xi}(\mathbf{n})$ is given in terms of the $\overline{\mathbf{F}}_m(\mathbf{n})$ in the following way. It follows from (3.10), (3.16) and (3.18) that

$$\begin{aligned} & \sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \mathbf{A}_{k,\xi}(\mathbf{n}) \\ &= \sum_{m=0}^{\infty} \int_0^{\infty} dH_{k,\xi}(y) \frac{(\theta_{\text{on}} y)^m}{m!} e^{-\theta_{\text{on}} y} \left[\mathbf{I} + \theta_{\text{on}}^{-1} \left\{ \mathbf{C}_{\text{on}} + \sum_{k \in \mathcal{K}} \mathbf{D}_{k,\text{on}}^*(z_k) \right. \right. \\ & \quad \left. \left. + \mathbf{E}_{\text{on,off}} \left(-\mathbf{C}_{\text{off}} - \sum_{k \in \mathcal{K}} \mathbf{D}_{k,\text{off}}^*(z_k) \right)^{-1} \mathbf{E}_{\text{off,on}} \right\} \right]^m \\ &= \sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \sum_{m=0}^{\infty} \gamma_{k,\xi}^{(m)}(\theta_{\text{on}}) \overline{\mathbf{F}}_m(\mathbf{n}), \end{aligned} \quad (3.23)$$

where

$$\gamma_{k,\xi}^{(m)}(\theta_{\text{on}}) = \int_0^{\infty} dH_{k,\xi}(y) \frac{(\theta_{\text{on}} y)^m}{m!} e^{-\theta_{\text{on}} y}, \quad k \in \mathcal{K}, \quad \xi = \text{on, off}, \quad m = 0, 1, \dots$$

Comparing coefficient vectors of $z_1^{n_1} \cdots z_K^{n_K}$ on both sides in (3.23), we obtain the following theorem.

Theorem 3.3 *The $\mathbf{A}_{k,\xi}(\mathbf{n})$ is given by*

$$\mathbf{A}_{k,\xi}(\mathbf{n}) = \sum_{m=0}^{\infty} \gamma_{k,\xi}^{(m)}(\theta_{\text{on}}) \overline{\mathbf{F}}_m(\mathbf{n}), \quad k \in \mathcal{K}, \quad \xi = \text{on, off}, \quad \mathbf{n} \in \mathcal{Z},$$

where the $\overline{\mathbf{F}}_m(\mathbf{n})$ is given in Lemma 3.2.

Next we consider the $\mathbf{v}_{k,\text{on}}(\mathbf{n})$ in (3.14) and the $\mathbf{v}_{k,\text{off}}(\mathbf{n})$ in (3.15). Expanding $\mathbf{N}^*(z | x)$ in (3.14) and (3.15), and comparing coefficient vectors of $z_1^{n_1} \cdots z_K^{n_K}$ on both sides of each equation, we obtain the following theorem.

Theorem 3.4 *The $\mathbf{v}_{k,\text{on}}(\mathbf{n})$ ($k \in \mathcal{K}$) and the $\mathbf{v}_{k,\text{off}}(\mathbf{n})$ ($k \in \mathcal{K}$) are given by*

$$\begin{aligned} \mathbf{v}_{k,\text{on}}(\mathbf{n}) &= \sum_{m=0}^{\infty} \mathbf{v}_{\text{on}}^{(m)}(\theta_{\text{on}}) \mathbf{D}_{k,\text{on}} \overline{\mathbf{F}}_m(\mathbf{n}), \quad \mathbf{n} \in \mathcal{Z}, \\ \mathbf{v}_{k,\text{off}}(\mathbf{n}) &= \sum_{m=0}^{\infty} \mathbf{v}_{\text{off}}^{(m)}(\theta_{\text{on}}) \mathbf{D}_{k,\text{off}} \sum_{\mathbf{0} \leq \mathbf{l} \leq \mathbf{n}} \mathbf{N}_{\text{off}}(\mathbf{n} - \mathbf{l}) \mathbf{E}_{\text{off,on}} \overline{\mathbf{F}}_m(\mathbf{l}), \quad \mathbf{n} \in \mathcal{Z}, \end{aligned}$$

respectively, where the $\overline{\mathbf{F}}_m(\mathbf{n})$ and the $\mathbf{N}_{\text{off}}(\mathbf{n})$ are given in Lemma 3.2, and

$$\mathbf{v}_{\xi}^{(m)}(\theta_{\text{on}}) = \int_0^{\infty} d\mathbf{v}_{\xi}(x) \frac{(\theta_{\text{on}} x)^m}{m!} e^{-\theta_{\text{on}} x}, \quad \xi = \text{on, off}, \quad m = 0, 1, \dots$$

Thus the $\mathbf{v}_{k,\xi}(\mathbf{n})$ ($\xi = \text{on}, \text{off}$) is given in terms of the $\mathbf{v}_\xi^{(m)}(\theta_{\text{on}})$ whose computation has already been studied in [Taki01b]. In what follows, we summarize the result. The perceptive reader will have noticed that the computational procedure for $\mathbf{v}_{\text{on}}^{(m)}(\theta_{\text{on}})$ is very similar to that for $\mathbf{v}^{(m)}(\theta)$ in (2.28). Note first that

$$\mathbf{v}_\xi^*(\theta_{\text{on}} - \theta_{\text{on}}z) = \sum_{m=0}^{\infty} z^m \mathbf{v}_\xi^{(m)}(\theta_{\text{on}}), \quad \xi = \text{on}, \text{off}. \quad (3.24)$$

Substituting $\theta_{\text{on}} - \theta_{\text{on}}z$ for s in (3.1) and using (3.24), we have

$$\begin{aligned} \sum_{m=0}^{\infty} z^m \mathbf{v}_{\text{on}}^{(m)}(\theta_{\text{on}}) \left[(\theta_{\text{on}} - \theta_{\text{on}}z) \mathbf{I} + \mathbf{C}_{\text{on}} + \sum_{m=0}^{\infty} z^m \mathbf{D}_{\text{on}}^{(m)}(\theta_{\text{on}}) \right. \\ \left. + \mathbf{E}_{\text{on,off}} \sum_{m=0}^{\infty} z^m \mathbf{J}_{\text{off}}^{(m)}(\theta_{\text{on}}) \mathbf{E}_{\text{off,on}} \right] = (\theta_{\text{on}} - \theta_{\text{on}}z)(1 - \rho_{\text{on}}) \boldsymbol{\kappa}_{\text{on}}, \end{aligned} \quad (3.25)$$

where $\mathbf{D}_\xi^{(m)}(\theta_{\text{on}})$ ($\xi = \text{on}, \text{off}$) and $\mathbf{J}_{\text{off}}^{(m)}(\theta_{\text{on}})$ are matrices satisfying

$$\begin{aligned} \sum_{m=0}^{\infty} z^m \mathbf{D}_\xi^{(m)}(\theta_{\text{on}}) &= \overline{\mathbf{D}}_\xi^*(\theta_{\text{on}} - \theta_{\text{on}}z), \\ \sum_{m=0}^{\infty} z^m \mathbf{J}_{\text{off}}^{(m)}(\theta_{\text{on}}) &= \left[-\mathbf{C}_{\text{off}} - \overline{\mathbf{D}}_{\text{off}}^*(\theta_{\text{on}} - \theta_{\text{on}}z) \right]^{-1}, \end{aligned} \quad (3.26)$$

respectively. The computation of the $\mathbf{D}_\xi^{(m)}(\theta_{\text{on}})$ ($\xi = \text{on}, \text{off}$) has already been studied (see Lemma 2.5 in chapter 2).

Lemma 3.3 *Under Assumption 3.1, the $\mathbf{D}_\xi^{(m)}(\theta_{\text{on}})$ ($\xi = \text{on}, \text{off}$) is given by*

$$\mathbf{D}_\xi^{(m)}(\theta_{\text{on}}) = \sum_{k \in \mathcal{K}} \mathbf{d}_{k,\xi}^{(m)}(\theta_{\text{on}}) \mathbf{e} \mathbf{D}_{k,\xi}, \quad m = 0, 1, \dots,$$

where the $\mathbf{d}_{k,\xi}^{(m)}(\theta_{\text{on}})$ ($k \in \mathcal{K}$, $\xi = \text{on}, \text{off}$) is given by the following recursion:

$$\begin{aligned} \mathbf{d}_{k,\xi}^{(0)}(\theta_{\text{on}}) &= \gamma_{k,\xi}^{(0)}(\theta_{\text{on}}) \boldsymbol{\alpha}_{k,\xi} (\mathbf{I} - \mathbf{P}_{k,\xi}) \left[\mathbf{I} - \gamma_{k,\xi}^{(0)}(\theta_{\text{on}}) \mathbf{P}_{k,\xi} \right]^{-1}, \\ \mathbf{d}_{k,\xi}^{(m)}(\theta_{\text{on}}) &= \frac{\gamma_{k,\xi}^{(m)}(\theta_{\text{on}})}{\gamma_{k,\xi}^{(0)}(\theta_{\text{on}})} \mathbf{d}_{k,\xi}^{(0)}(\theta_{\text{on}}) + \left[\sum_{l=1}^m \gamma_{k,\xi}^{(l)}(\theta_{\text{on}}) \mathbf{d}_{k,\xi}^{(m-l)}(\theta_{\text{on}}) \right] \\ &\quad \cdot \mathbf{P}_{k,\xi} \left[\mathbf{I} - \gamma_{k,\xi}^{(0)}(\theta_{\text{on}}) \mathbf{P}_{k,\xi} \right]^{-1}, \quad m = 1, 2, \dots \end{aligned}$$

Proof. We can prove this lemma in the same way as Lemma 2.5. ■

The recursion for the $\mathbf{J}_{\text{off}}^{(m)}(\theta_{\text{on}})$ can be readily obtained from

$$\sum_{m=0}^{\infty} z^m \mathbf{J}_{\text{off}}^{(m)}(\theta_{\text{on}}) \left[-\mathbf{C}_{\text{off}} - \sum_{m=0}^{\infty} z^m \mathbf{D}_{\text{off}}^{(m)}(\theta_{\text{on}}) \right] = \mathbf{I}.$$

Lemma 3.4 *The $\mathbf{J}_{\text{off}}^{(m)}(\theta_{\text{on}})$ is recursively determined by the following recursion:*

$$\begin{aligned}\mathbf{J}_{\text{off}}^{(0)}(\theta_{\text{on}}) &= \left[-\mathbf{C}_{\text{off}} - \mathbf{D}_{\text{off}}^{(0)}(\theta_{\text{on}})\right]^{-1}, \\ \mathbf{J}_{\text{off}}^{(m)}(\theta_{\text{on}}) &= \left[\sum_{l=0}^{m-1} \mathbf{J}_{\text{off}}^{(l)}(\theta_{\text{on}}) \mathbf{D}_{\text{off}}^{(m-l)}(\theta_{\text{on}})\right] \mathbf{J}_{\text{off}}^{(0)}(\theta_{\text{on}}), \quad m = 1, 2, \dots\end{aligned}$$

The $\mathbf{v}_{\text{on}}^{(m)}(\theta_{\text{on}})$ is computed as follows. Comparing the coefficient vectors of z^m ($m = 0, 1, \dots$) on both sides of (3.25), we can show that the $\mathbf{v}_{\text{on}}^{(m)}(\theta_{\text{on}})$ is identical to the stationary distribution of a Markov chain of M/G/1 type whose transition probability matrix is given by [Taki01b]

$$\begin{bmatrix} \mathbf{B}_0 + \mathbf{B}_1 & \mathbf{B}_2 & \mathbf{B}_3 & \mathbf{B}_4 & \cdots \\ \mathbf{B}_0 & \mathbf{B}_1 & \mathbf{B}_2 & \mathbf{B}_3 & \cdots \\ \mathbf{O} & \mathbf{B}_0 & \mathbf{B}_1 & \mathbf{B}_2 & \cdots \\ \mathbf{O} & \mathbf{O} & \mathbf{B}_0 & \mathbf{B}_1 & \cdots \\ \mathbf{O} & \mathbf{O} & \mathbf{O} & \mathbf{B}_0 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix},$$

where

$$\begin{aligned}\mathbf{B}_0 &= \mathbf{I} + \theta_{\text{on}}^{-1} \left[\mathbf{C}_{\text{on}} + \mathbf{D}_{\text{on}}^{(0)}(\theta_{\text{on}}) + \mathbf{E}_{\text{on,off}} \mathbf{J}_{\text{off}}^{(0)}(\theta_{\text{on}}) \mathbf{E}_{\text{off,on}}\right], \\ \mathbf{B}_m &= \theta_{\text{on}}^{-1} \left[\mathbf{D}_{\text{on}}^{(m)}(\theta_{\text{on}}) + \mathbf{E}_{\text{on,off}} \mathbf{J}_{\text{off}}^{(m)}(\theta_{\text{on}}) \mathbf{E}_{\text{off,on}}\right], \quad m = 1, 2, \dots\end{aligned}$$

Thus applying the general theory of Markov chains of M/G/1 type (see subsection 1.2.2), we can compute the $\mathbf{v}_{\text{on}}^{(m)}(\theta_{\text{on}})$.

On the other hand, the $\mathbf{v}_{\text{off}}^{(m)}(\theta_{\text{on}})$ can be computed by the following theorem whose proof is given in Appendix E.

Theorem 3.5 *The $\mathbf{v}_{\text{off}}^{(m)}(\theta_{\text{on}})$ is determined by the following recursion:*

$$\mathbf{v}_{\text{off}}^{(0)}(\theta_{\text{on}}) = \frac{\mathbf{v}_{\text{on}}^{(0)}(\theta_{\text{on}}) \mathbf{E}_{\text{on,off}} \left[-\mathbf{C}_{\text{off}} - \mathbf{D}_{\text{off}}^{(0)}(\theta_{\text{on}})\right]^{-1}}{\pi_{\text{on}} \mathbf{E}_{\text{on,off}} \mathbf{e} \bar{I}_{\text{off}}}, \quad (3.27)$$

$$\begin{aligned}\mathbf{v}_{\text{off}}^{(m)}(\theta_{\text{on}}) &= \left[\frac{\mathbf{v}_{\text{on}}^{(m)}(\theta_{\text{on}}) \mathbf{E}_{\text{on,off}}}{\pi_{\text{on}} \mathbf{E}_{\text{on,off}} \mathbf{e} \bar{I}_{\text{off}}} + \sum_{l=0}^{m-1} \mathbf{v}_{\text{off}}^{(l)}(\theta_{\text{on}}) \mathbf{D}_{\text{off}}^{(m-l)}(\theta_{\text{on}})\right] \\ &\quad \cdot \left[-\mathbf{C}_{\text{off}} - \mathbf{D}_{\text{off}}^{(0)}(\theta_{\text{on}})\right]^{-1}, \quad m = 1, 2, \dots, \quad (3.28)\end{aligned}$$

where the $\mathbf{D}_{\text{off}}^{(m)}(\theta_{\text{on}})$ is given in Lemma 3.3.

Among the recursions required in computing the joint queue length distribution, Lemma 3.2 for the $\bar{\mathbf{F}}_m(\mathbf{n})$ is the most extensive. In fact, its straightforward implementation will require very huge memory space. Note that an efficient implementation scheme for it is proposed in section 2.5. All other recursions can be readily implemented as they are. See [Taki94b] and chapter 2 in this thesis for details. Note also that from the results in this subsection, we can readily obtain recursions to compute the total queue length distribution, which are much less extensive than those for the joint queue length distribution. The results for the total queue length distribution are summarized in Appendix F.

3.5 Numerical Examples

In this section, we provide some numerical examples using two-class models. Throughout this section, we assume that the marginal arrival process follows a two-state batch MMAP $(\tilde{\mathbf{C}}, \tilde{\mathbf{D}}_1(n), \tilde{\mathbf{D}}_2(n))$:

$$\tilde{\mathbf{C}} = \begin{bmatrix} -\lambda g^{-1} - a & a \\ a & -\lambda g^{-1} - a \end{bmatrix}, \quad (3.29)$$

$$\tilde{\mathbf{D}}_1(n) = g(n) \begin{bmatrix} \lambda g^{-1} & 0 \\ 0 & 0 \end{bmatrix}, \quad \tilde{\mathbf{D}}_2(n) = g(n) \begin{bmatrix} 0 & 0 \\ 0 & \lambda g^{-1} \end{bmatrix}, \quad n = 1, 2, \dots, \quad (3.30)$$

where $\lambda, g, a > 0$ and $g(n) = 1$ if $n = g$, and otherwise $g(n) = 0$.

We also assume that the marginal on-off process follows an alternating Markov renewal process whose infinitesimal generator is given by

$$\begin{bmatrix} \mathbf{S}_{\text{on}} & \mathbf{T}_{\text{on,off}} \\ \mathbf{T}_{\text{off,on}} & \mathbf{S}_{\text{off}} \end{bmatrix},$$

where \mathbf{S}_{on} (resp. \mathbf{S}_{off}) denotes an $m_{\text{on}} \times m_{\text{on}}$ (resp. $m_{\text{off}} \times m_{\text{off}}$) matrix representing an infinitesimal generator that governs transitions in on-periods (resp. off-periods), and $\mathbf{T}_{\text{on,off}}$ (resp. $\mathbf{T}_{\text{off,on}}$) denotes a transition rate matrix from on-states (resp. off-states) to off-states (resp. on-states). When the on-off and arrival processes are independent of each other, the model is characterized as follows:

$$\begin{aligned} \mathbf{C}_{\text{on}} &= \mathbf{S}_{\text{on}} \oplus \tilde{\mathbf{C}}, & \mathbf{C}_{\text{off}} &= \mathbf{S}_{\text{off}} \oplus \tilde{\mathbf{C}}, & \mathbf{E}_{\text{on,off}} &= \mathbf{T}_{\text{on,off}} \otimes \mathbf{I}(2), & \mathbf{E}_{\text{off,on}} &= \mathbf{T}_{\text{off,on}} \otimes \mathbf{I}(2), \\ \mathbf{D}_{1,\text{on}}(n) &= \mathbf{I}(m_{\text{on}}) \otimes \tilde{\mathbf{D}}_1(n), & \mathbf{D}_{1,\text{off}}(n) &= \mathbf{I}(m_{\text{off}}) \otimes \tilde{\mathbf{D}}_1(n), \\ \mathbf{D}_{2,\text{on}}(n) &= \mathbf{I}(m_{\text{on}}) \otimes \tilde{\mathbf{D}}_2(n), & \mathbf{D}_{2,\text{off}}(n) &= \mathbf{I}(m_{\text{off}}) \otimes \tilde{\mathbf{D}}_2(n). \end{aligned}$$

3.5.1 Impact of service time dependency

In this subsection, we discuss the impact of the service time dependency on the queue length. We assume that the on-off and arrival processes are mutually independent. Let $\mathbf{S}_{\text{on}} = \mathbf{S}_{\text{off}} = -\alpha$ and $\mathbf{T}_{\text{on,off}} = \mathbf{T}_{\text{off,on}} = \alpha$, where $\alpha > 0$. Also let $a = 0.1$ and $\lambda = 0.125$ in (3.29) and (3.30). As for the service time, we consider two cases, Case GD (class-dependent service times) and Case GI (i.i.d. service times):

$$\begin{aligned} \text{[Case GD]} & & H_1 &= 1 \text{ with probability } 1, & H_2 &= 5 \text{ with probability } 1, \\ \text{[Case GI]} & & H_k &= \begin{cases} 1, & \text{with probability } 1/2, \\ 5, & \text{with probability } 1/2, \end{cases} & & k = 1, 2, \end{aligned}$$

where H_k ($k = 1, 2$) denotes a generic random variable representing a service time of a class k customer. Note here that the overall service time distributions in both cases are identical and $\rho_{\text{on}} = 6\lambda = 0.75$.

Figure 3.1 plots the expected total queue lengths $E[N]$ in Cases GD and GI as functions of α^{-1} . As α^{-1} goes to 0, the above model gets close to a work-conserving single-server queue (i.e., no service interruptions occur) with the same arrival process and service time distributions, where the

processing speed of the server is reduced by half. We observe that the difference in the expected total queue lengths in the two cases is kept almost constant regardless of the value of α^{-1} and gets large with constant batch size g .

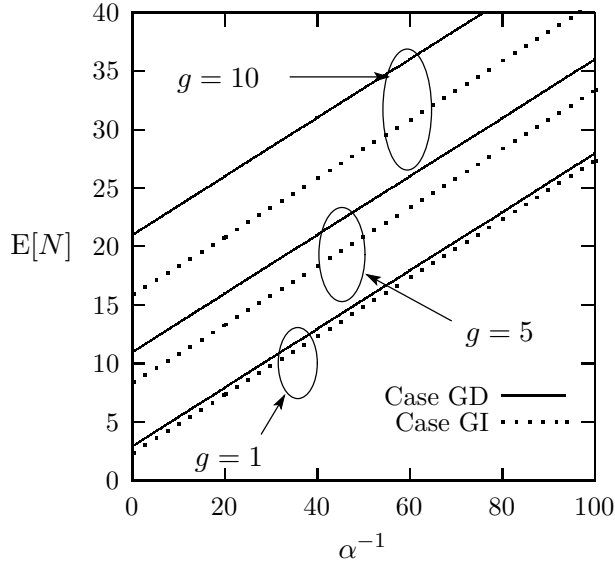


Figure 3.1: Expected total queue length $E[N]$.

Table 3.1 shows the joint queue length distribution for $g = 1$. Let $\mathbf{p}_{\text{GD}}(n_1, n_2)$ (resp. $\mathbf{p}_{\text{GI}}(n_1, n_2)$) denote $\mathbf{p}(n_1, n_2)$ in Case GD (resp. GI). We observe that $\mathbf{p}_{\text{GI}}(n_1, n_2)\mathbf{e} = \mathbf{p}_{\text{GI}}(n_2, n_1)\mathbf{e}$ is due to the symmetry of input parameters of classes 1 and 2 in Case GI. Further, Table 3.1 shows that $\mathbf{p}_{\text{GD}}(n_1, n_2)\mathbf{e} > \mathbf{p}_{\text{GI}}(n_1, n_2)\mathbf{e}$ for $n_1 < n_2$, and vice versa. Thus for m, n such that $m > n$,

$$\mathbf{p}_{\text{GD}}(m, n)\mathbf{e} < \mathbf{p}_{\text{GI}}(m, n)\mathbf{e} = \mathbf{p}_{\text{GI}}(n, m)\mathbf{e} < \mathbf{p}_{\text{GD}}(n, m)\mathbf{e},$$

in this particular example. We conjecture that this phenomenon is caused by the fact that service times of class 2 customers are larger than those of class 1 customer in Case GD, and while a class 2 customer is being served, succeeding class 2 customers are likely to arrive back to back and stay in the system.

3.5.2 Impact of variation of on- and off-periods

Next, we discuss the impact of the variation in on- and off-periods on the total queue length N . We assume that the on-off process follows an alternating renewal process, and the on-off and arrival processes are mutually independent. Let $C_{v,\text{on}}^2$ and $C_{v,\text{off}}^2$ denote the squared coefficients of variation of on-periods and off-periods, respectively. To examine the impact of the variation of on-periods, the off-period distribution is fixed to be exponential with mean 100. For $C_{v,\text{on}}^2 = k^{-1} \leq 1$ ($k = 1, 2, \dots$), on-periods follow a k -stage Erlang distribution with mean 100, and for $C_{v,\text{on}}^2 > 1$, they follow a balanced hyper-exponential distribution $\psi(x)$ with mean 100, where

$$\psi(x) = 1 - p \exp(-0.02px) - (1 - p) \exp[-0.02(1 - p)x], \quad 0 < p < 0.5.$$

Table 3.1: Joint queue length distribution $\mathbf{p}(n_1, n_2)\mathbf{e}$.
(Upper rows for Case GD and lower rows for Case GI)

n_1	0	1	2	5	10	20
n_2						
0	1.34×10^{-1}	1.74×10^{-2}	5.67×10^{-3}	1.38×10^{-3}	1.79×10^{-4}	3.10×10^{-6}
	1.34×10^{-1}	3.41×10^{-2}	1.41×10^{-2}	2.51×10^{-3}	3.44×10^{-4}	6.80×10^{-6}
1	4.67×10^{-2}	9.37×10^{-3}	5.88×10^{-3}	2.12×10^{-3}	4.02×10^{-4}	1.14×10^{-5}
	3.41×10^{-2}	9.25×10^{-3}	6.25×10^{-3}	2.71×10^{-3}	6.02×10^{-4}	2.07×10^{-5}
2	2.18×10^{-2}	6.87×10^{-3}	5.31×10^{-3}	2.60×10^{-3}	6.48×10^{-4}	2.68×10^{-5}
	1.41×10^{-2}	6.25×10^{-3}	5.27×10^{-3}	3.03×10^{-3}	8.78×10^{-4}	4.38×10^{-5}
5	3.68×10^{-3}	3.26×10^{-3}	3.41×10^{-3}	2.87×10^{-3}	1.30×10^{-3}	1.24×10^{-4}
	2.51×10^{-3}	2.71×10^{-3}	3.03×10^{-3}	2.89×10^{-3}	1.50×10^{-3}	1.70×10^{-4}
10	5.02×10^{-4}	7.92×10^{-4}	1.09×10^{-3}	1.68×10^{-3}	1.54×10^{-3}	4.04×10^{-4}
	3.44×10^{-4}	6.02×10^{-4}	8.78×10^{-4}	1.50×10^{-3}	1.55×10^{-3}	4.73×10^{-4}
20	1.04×10^{-5}	2.97×10^{-5}	6.03×10^{-5}	2.15×10^{-4}	5.42×10^{-4}	6.12×10^{-4}
	6.80×10^{-6}	2.07×10^{-5}	4.38×10^{-5}	1.70×10^{-4}	4.73×10^{-4}	6.09×10^{-4}

Note that $C_{v,\text{on}}^2 = 1/\{2p(1-p)\} - 1$ in this case. On the other hand, in examining the impact of the variation of off-periods on the total queue length, the above on- and off-period distributions are exchanged.

As for the arrival process, we set $a = 0.1$, $\lambda = 0.125$ and $g = 1$ in (3.29) and (3.30). Besides, service times of each class are assumed to follow the same service time distribution as in Case GD of the preceding subsection.

Figures 3.2 and 3.3 plot the 99.9 percentile (99.9 PT) and expected value $E[N]$ of the total queue length, respectively, as functions of the squared coefficient of variation $C_{v,\xi}^2$ ($\xi = \text{on, off}$), where the vertical axes are in log-scale. Note that in the case of $C_{v,\xi}^2 = 1$ ($\xi = \text{on, off}$), the two models become identical with exponential on- and off-periods. We observe that both 99.9 PT and $E[N]$ are monotone increasing functions of $C_{v,\xi}^2$ ($\xi = \text{on, off}$) and $C_{v,\text{off}}^2$ has a more impact on the total queue length N than $C_{v,\text{on}}^2$.

3.5.3 Impact of correlation in on- and off-periods

In this subsection, we examine the impact of the correlation in on- and off-periods on the total queue length N . For this purpose, we assume that the on-off and arrival processes are mutually independent and the marginal on-off process is given by

$$\mathbf{S}_{\text{on}} = \mathbf{S}_{\text{off}} = \begin{bmatrix} -1/40 & 0 \\ 0 & -1/160 \end{bmatrix}, \quad \mathbf{T}_{\text{on,off}} = \mathbf{T}_{\text{off,on}} = \begin{bmatrix} p/40 & (1-p)/40 \\ (1-p)/160 & p/160 \end{bmatrix},$$

where $0 < p < 1$. Thus the marginal distributions of on- and off-periods are the same hyper-exponential distribution $\psi(x) = 1 - 0.5 \exp(-x/40) - 0.5 \exp(-x/160)$. Note that parameter p

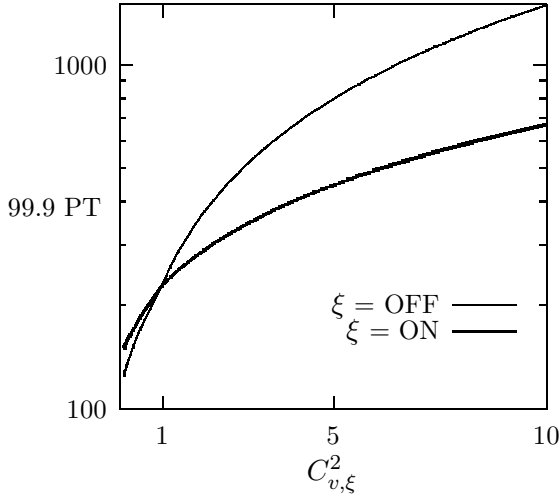


Figure 3.2: 99.9 percentile (99.9 PT) of the total queue length.

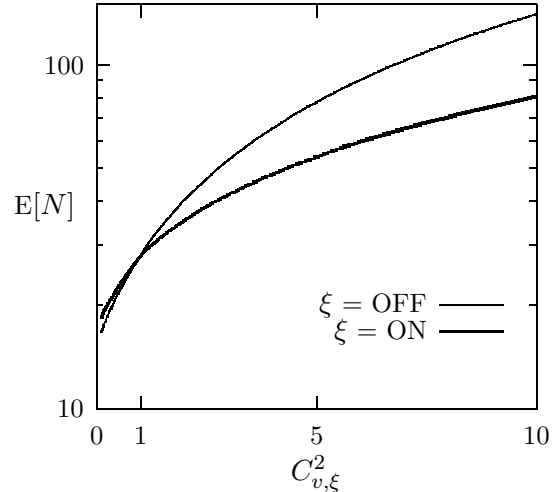


Figure 3.3: Expected total queue length $E[N]$.

controls the correlation in consecutive on- and off-periods. Suppose the on-off process starts with an on-period. Let $I_{\text{on}}(n)$ and $I_{\text{off}}(n)$ ($n = 1, 2, \dots$) denote the lengths of the n th on- and off-periods, respectively. Then, both $\text{Cov}[I_{\text{on}}(n), I_{\text{off}}(n)]$ and $\text{Cov}[I_{\text{off}}(n), I_{\text{on}}(n+1)]$ are negative for $0 < p < 0.5$, equal to zero for $p = 0.5$, and positive for $0.5 < p < 1$. We set $a = 0.1$, $\lambda = 0.125$ and $g = 1$ in (3.29) and (3.30). Service times of each class follow the same distribution as in Case GD of subsection 3.5.1.

Figures 3.4 and 3.5 plot the 99.9 percentile (99.9 PT) and expected value $E[N]$ of the total queue length, respectively, as functions of p . From these figures, we observe the followings. As p goes to zero, both 99.9 PT and $E[N]$ rapidly increase. This phenomenon is due to the fact that once the on-off process is in a long off-period, long off-periods and short on-periods are likely to repeat alternately, and during those intervals, many customers are accumulated in the system. As p becomes large, however, this effect is weakened, and finally, both 99.9 PT and $E[N]$ take their minimums and turn to increase. This implies that there exists some factor to make the queue length increase with p . In what follows, we examine this phenomenon more closely.

Let $\Psi_{\text{short-on}}$, $\Psi_{\text{long-on}}$, $\Psi_{\text{short-off}}$ and $\Psi_{\text{long-off}}$ denote the events that the on-off process is in a short on-period, long on-period, short off-period and long off-period, respectively. Figures 3.6 and 3.7 plot the conditional expected total queue lengths given those events as functions of p . From Figure 3.6, we observe that as expected, $E[N | \Psi_{\text{long-off}}]$ is always larger than $E[N | \Psi_{\text{short-off}}]$, so that the total queue length in an on-period following a long off-period is likely to be larger than that in an on-period following a short off-period, regardless of the value of p . Note here that as p goes to one, the contribution of the total queue length in an on-period following a long off-period to $E[N | \Psi_{\text{long-on}}]$ becomes large, and we conjecture that this factor makes $E[N | \Psi_{\text{long-on}}]$ increase in the region where p is close to one, as shown in Figure 3.7. Moreover, once $E[N | \Psi_{\text{long-on}}]$ turns to increase, this affects the total queue length in the following off-period, and as p goes to one, off-periods following long on-periods are likely to be long off-periods. Thus $E[N | \Psi_{\text{long-off}}]$ turns to increase after $E[N | \Psi_{\text{long-on}}]$ does, as shown in Figures 3.6 and 3.7.

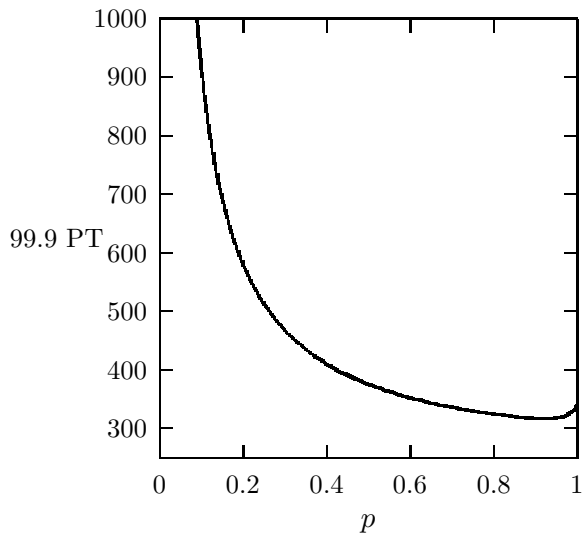


Figure 3.4: 99.9 percentile (99.9 PT) of the total queue length.

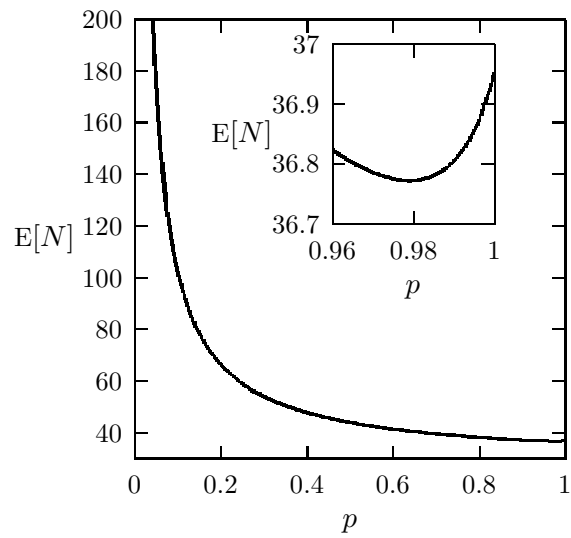


Figure 3.5: Expected total queue length $E[N]$.

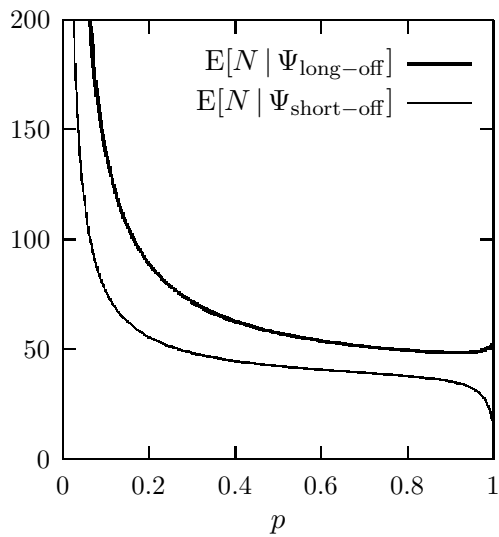


Figure 3.6: Conditional expected total queue length.

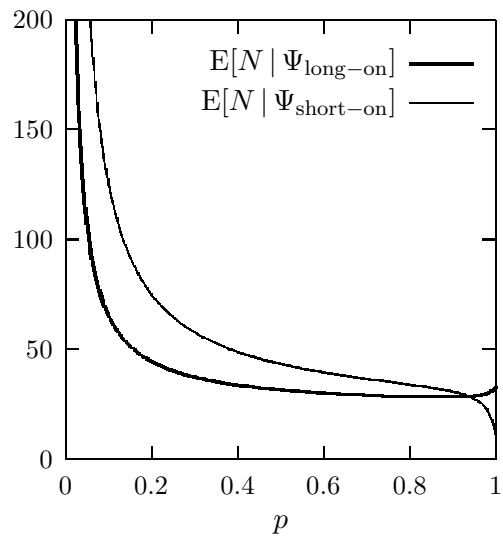


Figure 3.7: Conditional expected total queue length.

Note here that in this particular example,

$$\Pr(\Psi_{\text{long-on}}) = \Pr(\Psi_{\text{long-off}}) = 0.4, \quad \Pr(\Psi_{\text{short-on}}) = \Pr(\Psi_{\text{short-off}}) = 0.1.$$

Therefore the contributions of $E[N \mid \Psi_{\text{long-on}}]$ and $E[N \mid \Psi_{\text{long-off}}]$ to $E[N]$ are four times as large as those of $E[N \mid \Psi_{\text{short-on}}]$ and $E[N \mid \Psi_{\text{short-off}}]$. As a result, $E[N]$ increases for p near one. A similar observation can be applied to the 99.9 percentile, too.

3.5.4 Impact of correlation between on-off and arrival processes

Finally, we examine the impact of the correlation between on-off and arrival processes on the total queue length N . We consider the following three models, where service times of each class follow the same distribution as in Case GD of subsection 3.5.1.

In Model 1, the on-off and arrival processes have a correlation, and they are represented by

$$\mathbf{C} = \left[\begin{array}{c|c} -0.125 - \alpha & \alpha \\ \hline \alpha & -0.125 - \alpha \end{array} \right], \quad \mathbf{D}_1(1) = \left[\begin{array}{c|c} 0.125 & 0 \\ \hline 0 & 0 \end{array} \right], \quad \mathbf{D}_2(1) = \left[\begin{array}{c|c} 0 & 0 \\ \hline 0 & 0.125 \end{array} \right],$$

and $\mathbf{D}_k(n) = \mathbf{O}$ ($k = 1, 2$) for $n = 2, 3, \dots$, where $M_{\text{on}} = M_{\text{off}} = 1$. In Model 2, the on-off and arrival processes are mutually independent. As for the arrival process, we set $a = \alpha$, $\lambda = 0.125$ and $g = 1$ in (3.29) and (3.30). The on-off process is the same as that in subsection 3.5.1. In Model 3, the on-off and arrival processes have a correlation, and they are represented by

$$\mathbf{C} = \left[\begin{array}{c|c} -0.125 - \alpha & \alpha \\ \hline \alpha & -0.125 - \alpha \end{array} \right], \quad \mathbf{D}_1(1) = \left[\begin{array}{c|c} 0 & 0 \\ \hline 0 & 0.125 \end{array} \right], \quad \mathbf{D}_2(1) = \left[\begin{array}{c|c} 0.125 & 0 \\ \hline 0 & 0 \end{array} \right],$$

and $\mathbf{D}_k(n) = \mathbf{O}$ ($k = 1, 2$) for $n = 2, 3, \dots$, where $M_{\text{on}} = M_{\text{off}} = 1$. Note here that the marginal on-off processes in the three models are identical, and so are the marginal arrival processes.

Figures 3.8 and 3.9 show the 99.9 percentile (99.9 PT) and expected value $E[N]$ of the total queue length, respectively, as functions of α^{-1} . Note here that as α^{-1} goes to zero, all the three models get close to a work-conserving single-server queue, where the arrival process follows a Poisson process with rate 0.125, and service times are i.i.d. and take 2 or 10 with equal probability. This is a reason why both 99.9 PTs and $E[N]$'s in all the three models converge the same values, respectively, as $\alpha^{-1} \rightarrow 0$. We also observe that 99.9PT and $E[N]$ in Model 1 (resp. Model 3) are always larger (resp. smaller) than those in Model 2. This is due to the fact the amount of work brought into the system during off-periods in Model 1 (resp. Model 3) is likely to be larger (resp. smaller) than that in Model 2.

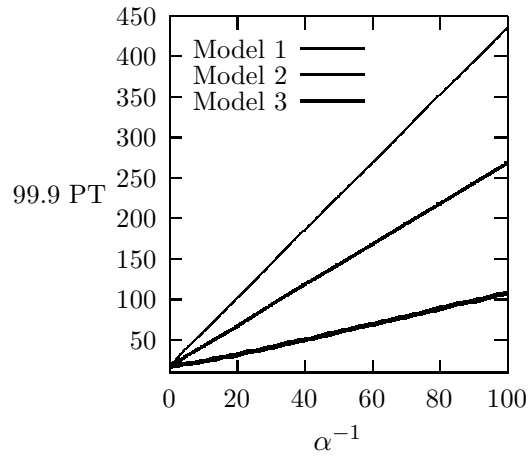


Figure 3.8: 99.9 percentile (99.9 PT) of the total queue length.

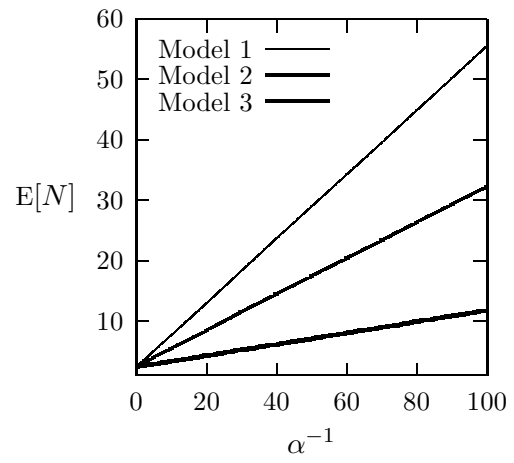


Figure 3.9: Expected total queue length $E[N]$.

Chapter 4

Infinite-Server Queue

4.1 Introduction

For a few decades in the past, several studies have been done on infinite-server queues with batch arrivals, e.g., [Holm82, Holm83, Liu90, Liu91, Shan66]. In particular, Liu and Templeton [Liu91] studied an infinite-server queue, where the arrival process is assumed to be a Markov renewal process, the batch size distribution may depend on the state of the Markov renewal process, and service times of individual customers in the same batch are independent and identically distributed (i.i.d.), given the state of the Markov renewal process on arrival. To the best of our knowledge, this queue is the most general one among those studied in the past. They derived the time-dependent generating function and binomial moments of the number of customers in the system. However, most of those results are given in terms of integrals with respect to a certain Markov renewal function, and therefore numerical computations with those results are not easy to conduct, as stated in [Liu91].

This chapter considers an infinite-server queue fed by the Markovian input process defined in subsection 1.1.2. The input process has multiple arrival streams, and customer arrivals from respective streams occur with a predefined probability when transitions of the underlying Markov chain happen. Further, the service time distribution of customers may be different for different arrival streams. For such a queue, we first derive a system of ordinary differential equations for the time-dependent matrix joint generating function of the number of customers in the system. This is considered as a generalization of the known results for the $N/G/\infty$ queue [Rama78] and for the $PH/G/\infty$ queue [Rama80]. As shown in [Rama80], applying a general-purpose numerical algorithm, we can compute the time-dependent joint distribution and the joint factorial moments of the number of customers in the system. Next we derive the time-dependent joint binomial moments of customers in the system, and assuming phase-type service times, we obtain explicit and numerically feasible expressions of the time-dependent and limiting joint binomial moments. Furthermore, through numerical examples, we reveal the impact of system parameters on the time-dependent and limiting performance.

The remainder of this chapter is organized as follows. In section 2, we mention the input process in this chapter. In section 3, we derive a system of ordinary differential equations for the time-

dependent matrix joint generating function of the number of customers in the system. In section 4, we derive the time-dependent joint binomial moments, and then assuming phase-type service times, we obtain numerically feasible expressions of the time-dependent and limiting joint binomial moments. Finally, in section 5, we show some numerical examples and discuss the impact of system parameters on the performance.

4.2 Model

This chapter considers an infinite-server queue with the Markovian input process introduced in subsection 1.1.2. Recall that the counting process of arrivals is characterized by $(\mathbf{C}, \mathbf{D}(\mathbf{n}))$. Customer arrivals can be viewed as follows. Let $S(t)$ denote the state of the underlying Markov chain (which governs the input process) at time t . Given $S(t) = i$ ($i \in \mathcal{M} = \{1, \dots, M\}$), the conditional joint probability that n_k ($k \in \mathcal{K} = \{1, \dots, K\}$) class k customers simultaneously arrive at the queue during time interval $(t, t + \delta t]$ and $S(t + \delta t) = j$ is given by $D_{i,j}(\mathbf{n})\delta t + o(\delta t)$. Besides, given $S(t) = i$, event $S(t + \delta t) = j$ ($j \neq i$) with no arrivals happens with probability $C_{i,j}\delta t + o(\delta t)$. As for service times, class k ($k \in \mathcal{K}$) customers have a service time distribution $H_k(t)$ ($t \geq 0$) with finite mean h_k .

In what follows, we assume that the arrival rate λ_k ($k \in \mathcal{K}$) of class k is positive and no customers are present in the system at time 0.

4.3 Time-Dependent Joint Distribution of the Number of Customers

In this section, we consider the time-dependent joint generating function of the number of customers of each class in the system. Let T_m denote the m th ($m = 1, 2, \dots$) arrival epoch after time 0, where $0 < T_1 < T_2 < \dots$. Let $X_{m,k}$ ($k \in \mathcal{K}$) denote a random variable representing the number of class k customers arriving at time T_m . We define $N_k(t)$ as a random variable representing the number of class k customers in the system at time t .

When no arrivals happen in time interval $(0, t]$, we have $N_k(t) = 0$ for all k ($k \in \mathcal{K}$). Thus

$$\mathbb{E} \left[z_1^{N_1(t)} \dots z_K^{N_K(t)} \mathbf{1}\{S(t) = j, T_1 > t\} \mid S(0) = i \right] = \left[e^{\mathbf{C}t} \right]_{i,j}, \quad (4.1)$$

where $\mathbf{1}\{\xi\}$ denotes an indicator function of event ξ and $[\mathbf{A}]_{i,j}$ denotes the (i, j) th element of matrix \mathbf{A} .

Next we consider the case $T_\nu \leq t < T_{\nu+1}$ for some $\nu \geq 1$. Let ξ_ν ($\nu = 1, 2, \dots$) denote the event $T_\nu \leq t < T_{\nu+1}$. Given the event ξ_ν happens, we define \mathbf{x}_m ($m = 1, \dots, \nu$) as $(X_{m,1}, \dots, X_{m,K})$. Because customers are served independently, we have for $|z_k| \leq 1$ ($k \in \mathcal{K}$)

$$\begin{aligned} \mathbb{E} \left[z_1^{N_1(t)} \dots z_K^{N_K(t)} \mid \xi_\nu, T_m = t_m \text{ and } \mathbf{x}_m = (n_{m,1}, \dots, n_{m,K}) \text{ for } m = 1, \dots, \nu \right] \\ = \prod_{m=1}^{\nu} \prod_{k_m \in \mathcal{K}} \left[H_{k_m}(t - t_m) + z_{k_m} \overline{H}_{k_m}(t - t_m) \right]^{n_{m,k_m}} \end{aligned}$$

$$= \prod_{m=1}^{\nu} \prod_{k_m \in \mathcal{K}} \left[1 + (z_{k_m} - 1) \overline{H}_{k_m}(t - t_m) \right]^{n_{m,k_m}}, \quad (4.2)$$

where $\overline{H}_k(t) = 1 - H_k(t)$. Because events ξ_ν ($\nu = 1, 2, \dots$) are exclusive, it follows from (4.2) that

$$\begin{aligned} & \mathbb{E} \left[z_1^{N_1(t)} \dots z_K^{N_K(t)} 1\{S(t) = j \text{ and } T_1 < t\} \mid S(0) = i \right] \\ &= \sum_{\nu=1}^{\infty} \sum_{\mathbf{n}_1 \in \mathcal{Z}^+} \sum_{\mathbf{n}_2 \in \mathcal{Z}^+} \dots \sum_{\mathbf{n}_\nu \in \mathcal{Z}^+} \int_0^t dt_1 \int_{t_1}^t dt_2 \dots \int_{t_{\nu-1}}^t dt_\nu \\ & \quad \cdot \prod_{m=1}^{\nu} \prod_{k_m \in \mathcal{K}} \left[1 + (z_{k_m} - 1) \overline{H}_{k_m}(t - t_m) \right]^{n_{m,k_m}} \\ & \quad \cdot \left[e^{\mathbf{C}t_1} \mathbf{D}(\mathbf{n}_1) e^{\mathbf{C}(t_2-t_1)} \mathbf{D}(\mathbf{n}_2) \dots e^{\mathbf{C}(t_\nu-t_{\nu-1})} \mathbf{D}(\mathbf{n}_\nu) e^{\mathbf{C}(t-t_\nu)} \right]_{i,j}, \end{aligned} \quad (4.3)$$

where $\mathcal{Z}^+ = \{(n_1, \dots, n_K); n_k \text{'s } (k \in \mathcal{K}) \text{ are nonnegative integers and at least one of them is strictly positive}\}$ and $\mathbf{n}_m = (n_{m,1}, \dots, n_{m,K})$ for $m = 1, 2, \dots$

We now define $\mathbf{G}^*(t, \mathbf{z})$ as an $M \times M$ matrix whose (i, j) th ($i, j \in \mathcal{M}$) element represents

$$\mathbb{E} \left[z_1^{N_1(t)} \dots z_K^{N_K(t)} 1\{S(t) = j\} \mid S(0) = i \right],$$

where $\mathbf{z} = (z_1, \dots, z_K)$ with $|z_k| \leq 1$ for all k ($k \in \mathcal{K}$). It then follows from (4.1) and (4.3) that

$$\begin{aligned} \mathbf{G}^*(t, \mathbf{z}) &= e^{\mathbf{C}t} + \sum_{\nu=1}^{\infty} \sum_{\mathbf{n}_1 \in \mathcal{Z}^+} \sum_{\mathbf{n}_2 \in \mathcal{Z}^+} \dots \sum_{\mathbf{n}_\nu \in \mathcal{Z}^+} \int_0^t dt_1 \int_{t_1}^t dt_2 \dots \int_{t_{\nu-1}}^t dt_\nu \\ & \quad \cdot \prod_{m=1}^{\nu} \prod_{k_m \in \mathcal{K}} \left[1 + (z_{k_m} - 1) \overline{H}_{k_m}(t - t_m) \right]^{n_{m,k_m}} \\ & \quad \cdot e^{\mathbf{C}t_1} \mathbf{D}(\mathbf{n}_1) e^{\mathbf{C}(t_2-t_1)} \mathbf{D}(\mathbf{n}_2) \dots e^{\mathbf{C}(t_\nu-t_{\nu-1})} \mathbf{D}(\mathbf{n}_\nu) e^{\mathbf{C}(t-t_\nu)}. \end{aligned}$$

Further letting $u_j = t - t_{\nu+1-j}$ ($j = 1, \dots, \nu$) and rearranging terms, we have

$$\begin{aligned} \mathbf{G}^*(t, \mathbf{z}) &= e^{\mathbf{C}t} + \sum_{\nu=1}^{\infty} \sum_{\mathbf{n}_\nu \in \mathcal{Z}^+} \sum_{\mathbf{n}_{\nu-1} \in \mathcal{Z}^+} \dots \sum_{\mathbf{n}_1 \in \mathcal{Z}^+} \int_0^t du_\nu \int_0^{u_\nu} du_{\nu-1} \dots \int_0^{u_2} du_1 \\ & \quad \cdot \prod_{m=1}^{\nu} \prod_{k_m \in \mathcal{K}} \left[1 + (z_{k_m} - 1) \overline{H}_{k_m}(u_m) \right]^{n_{m,k_m}} \\ & \quad \cdot e^{\mathbf{C}(t-u_\nu)} \mathbf{D}(\mathbf{n}_\nu) e^{\mathbf{C}(u_\nu-u_{\nu-1})} \mathbf{D}(\mathbf{n}_{\nu-1}) \dots e^{\mathbf{C}(u_2-u_1)} \mathbf{D}(\mathbf{n}_1) e^{\mathbf{C}u_1} \\ &= e^{\mathbf{C}t} + \sum_{\nu=1}^{\infty} \int_0^t du_\nu \int_0^{u_\nu} du_{\nu-1} \dots \int_0^{u_2} du_1 e^{\mathbf{C}(t-u_\nu)} \mathbf{D}^*(u_\nu, \mathbf{z}) \\ & \quad \cdot e^{\mathbf{C}(u_\nu-u_{\nu-1})} \mathbf{D}^*(u_{\nu-1}, \mathbf{z}) \dots e^{\mathbf{C}(u_2-u_1)} \mathbf{D}^*(u_1, \mathbf{z}) e^{\mathbf{C}u_1}, \end{aligned} \quad (4.4)$$

where

$$\mathbf{D}^*(t, \mathbf{z}) = \sum_{\mathbf{n} \in \mathcal{Z}^+} \prod_{k \in \mathcal{K}} \left[1 + (z_k - 1) \overline{H}_k(t) \right]^{n_k} \mathbf{D}(\mathbf{n}). \quad (4.5)$$

Theorem 4.1 For any closed interval of t where all $\overline{H}_\nu(t)$ are continuous, $\mathbf{G}^*(t, \mathbf{z})$ in (4.4) satisfies the following differential equation.

$$\frac{\partial}{\partial t} \mathbf{G}^*(t, \mathbf{z}) = [\mathbf{C} + \mathbf{D}^*(t, \mathbf{z})] \mathbf{G}^*(t, \mathbf{z}). \quad (4.6)$$

Further (4.6) has a unique continuous solution in $[0, \infty)$ with $\mathbf{G}^*(0, \mathbf{z}) = \mathbf{I}$.

Proof. Pre-multiplying both sides of (4.4) by $\exp(-Ct)$, we have

$$\begin{aligned} e^{-Ct}\mathbf{G}^*(t, \mathbf{z}) &= \mathbf{I} + \sum_{\nu=1}^{\infty} \int_0^t du_{\nu} \int_0^{u_{\nu}} du_{\nu-1} \cdots \int_0^{u_2} du_1 e^{-Cu_{\nu}} \mathbf{D}^*(u_{\nu}, \mathbf{z}) \\ &\quad \cdot e^{C(u_{\nu}-u_{\nu-1})} \mathbf{D}^*(u_{\nu-1}, \mathbf{z}) \cdots \cdots e^{C(u_2-u_1)} \mathbf{D}^*(u_1, \mathbf{z}) e^{Cu_1}. \end{aligned} \quad (4.7)$$

Differentiating both sides of (4.7) with respect to t yields

$$\begin{aligned} e^{-Ct} \frac{\partial}{\partial t} \mathbf{G}^*(t, \mathbf{z}) - e^{-Ct} \mathbf{C} \mathbf{G}^*(t, \mathbf{z}) &= e^{-Ct} \mathbf{D}^*(t, \mathbf{z}) \left[e^{Ct} + \sum_{\nu=2}^{\infty} \int_0^t du_{\nu-1} \int_0^{u_{\nu-1}} du_{\nu-2} \cdots \int_0^{u_2} du_1 e^{C(t-u_{\nu-1})} \mathbf{D}^*(u_{\nu-1}, \mathbf{z}) \right. \\ &\quad \left. \cdot e^{C(u_{\nu-1}-u_{\nu-2})} \mathbf{D}^*(u_{\nu-2}, \mathbf{z}) \cdots \cdots e^{C(u_2-u_1)} \mathbf{D}^*(u_1, \mathbf{z}) e^{Cu_1} \right] \\ &= e^{-Ct} \mathbf{D}^*(t, \mathbf{z}) \mathbf{G}^*(t, \mathbf{z}), \end{aligned} \quad (4.8)$$

where we use (4.4) in the last equality. Rearranging terms in (4.8) and pre-multiplying both sides by $\exp(Ct)$, we obtain (4.6). The uniqueness of the solution of (4.6) with $\mathbf{G}^*(0, \mathbf{z}) = \mathbf{I}$ follows from the well-known result of ordinary differential equations (see [Bell97], p.167, Theorem 1). ■

Remark 4.1 For the $N/G/\infty$ queue and the $PH/G/\infty$ queue, systems of ordinary differential equations for the generating function of the number of customers are derived in [Rama78] and [Rama80], respectively. Theorem 4.1 is considered as a generalization of those results.

We now consider the time-dependent joint distribution of the number of customers of each class in the system. Let $\mathbf{L}(t, \mathbf{n})$ ($\mathbf{n} \in \mathcal{Z}$) denote an $M \times M$ matrix whose (i, j) th element represents

$$\Pr[N_1(t) = n_1, \dots, N_K(t) = n_K, S(t) = j \mid S(0) = i],$$

where $\mathcal{Z} = \{(n_1, \dots, n_K); n_k = 0, 1, \dots \text{ for all } k \in \mathcal{K}\}$. Let $\mathbf{D}(t, \mathbf{n})$ ($\mathbf{n} \in \mathcal{Z}$) denote an $M \times M$ matrix satisfying

$$\mathbf{D}^*(t, \mathbf{z}) = \sum_{\mathbf{n} \in \mathcal{Z}} z^{n_1} \cdots z^{n_K} \mathbf{D}(t, \mathbf{n}),$$

where $\mathbf{D}^*(t, \mathbf{z})$ is given in (4.5). For a given $\mathbf{n} \in \mathcal{Z}$, comparing the coefficient matrices of $z^{n_1} \cdots z^{n_K}$ on both sides of (4.6), we obtain

$$\frac{d}{dt} \mathbf{L}(t, \mathbf{n}) = (\mathbf{C} + \mathbf{D}(t, \mathbf{0})) \mathbf{L}(t, \mathbf{n}) + \sum_{\substack{\mathbf{0} \leq \mathbf{m} \leq \mathbf{n} \\ \mathbf{m} \neq \mathbf{0}}} \mathbf{D}(t, \mathbf{m}) \mathbf{L}(t, \mathbf{n} - \mathbf{m}), \quad \mathbf{n} \in \mathcal{Z}.$$

Therefore, for any given $\mathbf{n} \in \mathcal{Z}$, $\mathbf{L}(t, \mathbf{m})$'s ($\mathbf{m} \in \mathcal{Z}$, $\mathbf{0} \leq \mathbf{m} \leq \mathbf{n}$) can be numerically obtained by solving the above system of ordinary differential equations with a general-purpose numerical algorithm (e.g., see [Rama80]).

Also, from (4.6), we can derive a system of ordinary differential equations for the joint factorial moments for the number of customers of each class in the system at time t . For this purpose, we

define $\Omega(t, \mathbf{n})$ and $\Psi(t, \mathbf{n})$ as

$$\begin{aligned}\Omega(t, \mathbf{n}) &= \lim_{z_1 \rightarrow 1} \cdots \lim_{z_K \rightarrow 1} \frac{\partial^{n_1}}{\partial z_1^{n_1}} \cdots \frac{\partial^{n_K}}{\partial z_K^{n_K}} \mathbf{G}^*(t, \mathbf{z}), \quad \mathbf{n} \in \mathcal{Z}, \\ \Psi(t, \mathbf{n}) &= \lim_{z_1 \rightarrow 1} \cdots \lim_{z_K \rightarrow 1} \frac{\partial^{n_1}}{\partial z_1^{n_1}} \cdots \frac{\partial^{n_K}}{\partial z_K^{n_K}} \mathbf{D}^*(t, \mathbf{z}), \quad \mathbf{n} \in \mathcal{Z},\end{aligned}$$

respectively. By differentiating both sides of (4.6) n_k times with respect to z_k and setting $z_k = 1$ for all k ($k \in \mathcal{K}$), we obtain

$$\frac{\partial}{\partial t} \Omega(t, \mathbf{n}) = (\mathbf{C} + \mathbf{D})\Omega(t, \mathbf{n}) + \sum_{\substack{\mathbf{0} \leq \mathbf{m} \leq \mathbf{n} \\ \mathbf{m} \neq \mathbf{0}}} \prod_{k \in \mathcal{K}} \binom{n_k}{m_k} \Psi(t, \mathbf{m}) \Omega(t, \mathbf{n} - \mathbf{m}), \quad \mathbf{n} \in \mathcal{Z}.$$

Thus, the $\Omega(t, \mathbf{n})$ is also obtained by solving the above system of ordinary differential equations.

4.4 Numerically Feasible Formulas for Phase-Type Service Times

In this section, we consider the joint binomial moments of the number of customers of each class in the system and develop numerically feasible formulas for the joint binomial moments, assuming phase-type service times. We define $\widehat{\mathbf{B}}(t, \mathbf{m})$ ($\mathbf{m} \in \mathcal{Z}^+$) as an $M \times M$ matrix whose (i, j) th element represents

$$[\widehat{\mathbf{B}}(t, \mathbf{m})]_{i,j} = \mathbb{E} \left[\prod_{k \in \mathcal{K}} \binom{N_k(t)}{m_k} \mathbf{1}_{\{S(t) = j\}} \middle| S(0) = i \right].$$

$\widehat{\mathbf{B}}(t, \mathbf{m})$ is called the \mathbf{m} th joint binomial moment matrix of the number of customers of each class in the system at time t .

4.4.1 Time-dependent joint binomial moments

We define $\mathbf{G}_B^*(t, \boldsymbol{\omega})$ as the matrix joint binomial moment generating function of the number of customers of each class in the system at time t , i.e.,

$$\mathbf{G}_B^*(t, \boldsymbol{\omega}) = e^{(\mathbf{C} + \mathbf{D})t} + \sum_{\mathbf{m} \in \mathcal{Z}^+} \omega_1^{m_1} \cdots \omega_K^{m_K} \widehat{\mathbf{B}}(t, \mathbf{m}), \quad (4.9)$$

where $\boldsymbol{\omega} = (\omega_1, \dots, \omega_K)$ and $|\omega_k + 1| \leq 1$ for all k ($k \in \mathcal{K}$). Note here that $\mathbf{G}_B^*(t, \boldsymbol{\omega})$ is given in terms of $\mathbf{G}^*(t, \mathbf{z})$:

$$\mathbf{G}_B^*(t, \boldsymbol{\omega}) = \mathbf{G}^*(t, \omega_1 + 1, \dots, \omega_K + 1). \quad (4.10)$$

Theorem 4.2 *The time-dependent matrix joint binomial moment generating function $\mathbf{G}_B^*(t, \boldsymbol{\omega})$ is given by*

$$\begin{aligned}\mathbf{G}_B^*(t, \boldsymbol{\omega}) &= e^{(\mathbf{C} + \mathbf{D})t} + \sum_{\nu=1}^{\infty} \int_0^t du_\nu \int_0^{u_\nu} du_{\nu-1} \cdots \int_0^{u_2} du_1 e^{(\mathbf{C} + \mathbf{D})(t - u_\nu)} \mathbf{D}_B^*(u_\nu, \boldsymbol{\omega}) \\ &\quad \cdot e^{(\mathbf{C} + \mathbf{D})(u_\nu - u_{\nu-1})} \mathbf{D}_B^*(u_{\nu-1}, \boldsymbol{\omega}) \\ &\quad \cdots \cdots e^{(\mathbf{C} + \mathbf{D})(u_2 - u_1)} \mathbf{D}_B^*(u_1, \boldsymbol{\omega}) e^{(\mathbf{C} + \mathbf{D})u_1},\end{aligned} \quad (4.11)$$

where

$$\mathbf{D}_B^*(t, \boldsymbol{\omega}) = \mathbf{D}^*(t, \omega_1 + 1, \dots, \omega_K + 1) - \mathbf{D}. \quad (4.12)$$

Proof. It is easy to see from (4.6) that $\mathbf{G}_B^*(t, \boldsymbol{\omega})$ in (4.10) satisfies the following differential equation.

$$\frac{\partial}{\partial t} \mathbf{G}_B^*(t, \boldsymbol{\omega}) = [\mathbf{C} + \mathbf{D} + \mathbf{D}_B^*(t, \boldsymbol{\omega})] \mathbf{G}_B^*(t, \boldsymbol{\omega}), \quad (4.13)$$

$$\mathbf{G}_B^*(0, \boldsymbol{\omega}) = \mathbf{I}. \quad (4.14)$$

It is clear that $\mathbf{G}_B^*(t, \boldsymbol{\omega})$ in (4.11) satisfies (4.14). In what follows, we shall show (4.11) is a solution of the differential equation (4.13).

Pre-multiplying both sides of (4.11) by $\exp[-(\mathbf{C} + \mathbf{D})t]$, we have

$$\begin{aligned} e^{-(\mathbf{C}+\mathbf{D})t} \mathbf{G}_B^*(t, \boldsymbol{\omega}) &= \mathbf{I} + \sum_{\nu=1}^{\infty} \int_0^t du_{\nu} \int_0^{u_{\nu}} du_{\nu-1} \cdots \int_0^{u_2} du_1 e^{-(\mathbf{C}+\mathbf{D})u_{\nu}} \mathbf{D}_B^*(u_{\nu}, \boldsymbol{\omega}) \\ &\quad \cdot e^{(\mathbf{C}+\mathbf{D})(u_{\nu}-u_{\nu-1})} \mathbf{D}_B^*(u_{\nu-1}, \boldsymbol{\omega}) \\ &\quad \cdots \cdots e^{(\mathbf{C}+\mathbf{D})(u_2-u_1)} \mathbf{D}_B^*(u_1, \boldsymbol{\omega}) e^{(\mathbf{C}+\mathbf{D})u_1}. \end{aligned} \quad (4.15)$$

Differentiating both sides of (4.15) with respect to t yields

$$\begin{aligned} e^{-(\mathbf{C}+\mathbf{D})t} \frac{\partial}{\partial t} \mathbf{G}_B^*(t, \boldsymbol{\omega}) - e^{-(\mathbf{C}+\mathbf{D})t} (\mathbf{C} + \mathbf{D}) \mathbf{G}_B^*(t, \boldsymbol{\omega}) \\ &= e^{-(\mathbf{C}+\mathbf{D})t} \mathbf{D}_B^*(t, \boldsymbol{\omega}) \left[e^{(\mathbf{C}+\mathbf{D})t} + \sum_{\nu=2}^{\infty} \int_0^t du_{\nu-1} \int_0^{u_{\nu-1}} du_{\nu-2} \cdots \int_0^{u_2} du_1 \right. \\ &\quad \cdot e^{(\mathbf{C}+\mathbf{D})(t-u_{\nu-1})} \mathbf{D}_B^*(u_{\nu-1}, \boldsymbol{\omega}) e^{(\mathbf{C}+\mathbf{D})(u_{\nu-1}-u_{\nu-2})} \mathbf{D}_B^*(u_{\nu-2}, \boldsymbol{\omega}) \\ &\quad \left. \cdots \cdots e^{(\mathbf{C}+\mathbf{D})(u_2-u_1)} \mathbf{D}_B^*(u_1, \boldsymbol{\omega}) e^{(\mathbf{C}+\mathbf{D})u_1} \right] \\ &= e^{-(\mathbf{C}+\mathbf{D})t} \mathbf{D}_B^*(t, \boldsymbol{\omega}) \mathbf{G}_B^*(t, \boldsymbol{\omega}), \end{aligned} \quad (4.16)$$

where we use (4.11) in the last equality. Thus, pre-multiplying both sides of (4.16) by $\exp[(\mathbf{C} + \mathbf{D})t]$, we see that $\mathbf{G}_B^*(t, \boldsymbol{\omega})$ in (4.11) satisfies (4.13). \blacksquare

We now rewrite $\mathbf{G}_B^*(t, \boldsymbol{\omega})$ in (4.11) to be a more appealing form. First, from (4.5) and (4.12), we have

$$\begin{aligned} \mathbf{D}_B^*(t, \boldsymbol{\omega}) &= \sum_{\mathbf{n} \in \mathcal{Z}^+} \prod_{k \in \mathcal{K}} [1 + \omega_k \bar{H}_k(t)]^{n_k} \mathbf{D}(\mathbf{n}) - \mathbf{D} \\ &= \sum_{\mathbf{n} \in \mathcal{Z}^+} \prod_{k \in \mathcal{K}} \left[\sum_{m_k=0}^{n_k} \binom{n_k}{m_k} \omega_k^{m_k} \bar{H}_k^{m_k}(t) \right] \mathbf{D}(\mathbf{n}) - \mathbf{D} \\ &= \sum_{\mathbf{m} \in \mathcal{Z}} \left[\prod_{k \in \mathcal{K}} \omega_k^{m_k} \bar{H}_k^{m_k}(t) \right] \sum_{\substack{\mathbf{n} \geq \mathbf{m} \\ \mathbf{n} \neq \mathbf{0}}} \prod_{k \in \mathcal{K}} \binom{n_k}{m_k} \mathbf{D}(\mathbf{n}) - \mathbf{D} \\ &= \sum_{\mathbf{n} \in \mathcal{Z}^+} \mathbf{D}(\mathbf{n}) + \sum_{\mathbf{m} \in \mathcal{Z}^+} \left[\prod_{k \in \mathcal{K}} \omega_k^{m_k} \bar{H}_k^{m_k}(t) \right] \sum_{\mathbf{n} \geq \mathbf{m}} \prod_{k \in \mathcal{K}} \binom{n_k}{m_k} \mathbf{D}(\mathbf{n}) - \mathbf{D} \end{aligned}$$

$$\begin{aligned}
&= \sum_{\mathbf{m} \in \mathcal{Z}^+} \left[\prod_{k \in \mathcal{K}} \omega_k^{m_k} \overline{H}_k^{m_k}(t) \right] \sum_{\mathbf{n} \geq \mathbf{m}} \prod_{k \in \mathcal{K}} \binom{n_k}{m_k} \mathbf{D}(\mathbf{n}) \\
&= \sum_{\mathbf{m} \in \mathcal{Z}^+} \omega^{(\mathbf{m})} \overline{H}^{(\mathbf{m})}(t) \widehat{\mathbf{D}}(\mathbf{m}), \tag{4.17}
\end{aligned}$$

where for $\mathbf{m} \in \mathcal{Z}^+$,

$$\omega^{(\mathbf{m})} = \prod_{k \in \mathcal{K}} \omega_k^{m_k}, \quad \overline{H}^{(\mathbf{m})}(t) = \prod_{k \in \mathcal{K}} \overline{H}_k^{m_k}(t), \quad \widehat{\mathbf{D}}(\mathbf{m}) = \sum_{\mathbf{n} \geq \mathbf{m}} \prod_{k \in \mathcal{K}} \binom{n_k}{m_k} \mathbf{D}(\mathbf{n}).$$

Thus, with (4.17), $\mathbf{G}_B^*(t, \boldsymbol{\omega})$ in (4.11) is rewritten to be

$$\begin{aligned}
\mathbf{G}_B^*(t, \boldsymbol{\omega}) &= e^{(\mathbf{C}+\mathbf{D})t} + \sum_{\nu=1}^{\infty} \sum_{\mathbf{l}_1 \in \mathcal{Z}^+} \sum_{\mathbf{l}_2 \in \mathcal{Z}^+} \cdots \sum_{\mathbf{l}_\nu \in \mathcal{Z}^+} \omega^{(\mathbf{l}_1 + \cdots + \mathbf{l}_\nu)} \\
&\quad \cdot \int_0^t du_\nu \int_0^{u_\nu} du_{\nu-1} \cdots \int_0^{u_2} du_1 \overline{H}^{(\mathbf{l}_\nu)}(u_\nu) \overline{H}^{(\mathbf{l}_{\nu-1})}(u_{\nu-1}) \cdots \overline{H}^{(\mathbf{l}_1)}(u_1) \\
&\quad \cdot e^{(\mathbf{C}+\mathbf{D})(t-u_\nu)} \widehat{\mathbf{D}}(\mathbf{l}_\nu) e^{(\mathbf{C}+\mathbf{D})(u_\nu-u_{\nu-1})} \widehat{\mathbf{D}}(\mathbf{l}_{\nu-1}) \\
&\quad \quad \quad \cdots \cdots e^{(\mathbf{C}+\mathbf{D})(u_2-u_1)} \widehat{\mathbf{D}}(\mathbf{l}_1) e^{(\mathbf{C}+\mathbf{D})u_1} \\
&= e^{(\mathbf{C}+\mathbf{D})t} + \sum_{\mathbf{m} \in \mathcal{Z}^+} \omega^{(\mathbf{m})} \sum_{\nu=1}^{|\mathbf{m}|} \sum_{\vec{\mathbf{l}}_\nu \in \mathcal{L}_\nu(\mathbf{m})} \\
&\quad \cdot \int_0^t du_\nu \int_0^{u_\nu} du_{\nu-1} \cdots \int_0^{u_2} du_1 \overline{H}^{(\mathbf{l}_\nu)}(u_\nu) \overline{H}^{(\mathbf{l}_{\nu-1})}(u_{\nu-1}) \cdots \overline{H}^{(\mathbf{l}_1)}(u_1) \\
&\quad \cdot e^{(\mathbf{C}+\mathbf{D})(t-u_\nu)} \widehat{\mathbf{D}}(\mathbf{l}_\nu) e^{(\mathbf{C}+\mathbf{D})(u_\nu-u_{\nu-1})} \widehat{\mathbf{D}}(\mathbf{l}_{\nu-1}) \\
&\quad \quad \quad \cdots \cdots e^{(\mathbf{C}+\mathbf{D})(u_2-u_1)} \widehat{\mathbf{D}}(\mathbf{l}_1) e^{(\mathbf{C}+\mathbf{D})u_1}, \tag{4.18}
\end{aligned}$$

where

$$\begin{aligned}
|\mathbf{m}| &= \sum_{k \in \mathcal{K}} m_k, \quad \mathbf{m} \in \mathcal{Z}, \\
\mathbf{l}_j &= (l_{j,1}, \dots, l_{j,K}) \in \mathcal{Z}^+, \\
\vec{\mathbf{l}}_\nu &= \{(\mathbf{l}_1, \dots, \mathbf{l}_\nu) \mid \mathbf{l}_j \in \mathcal{Z}^+, j = 1, \dots, \nu\}, \\
\mathcal{L}_\nu(\mathbf{m}) &= \left\{ \vec{\mathbf{l}}_\nu = (\mathbf{l}_1, \dots, \mathbf{l}_\nu) \mid \mathbf{l}_1 + \cdots + \mathbf{l}_\nu = \mathbf{m}, \widehat{\mathbf{D}}(\mathbf{l}_j) \neq \mathbf{O}, \mathbf{l}_j \in \mathcal{Z}^+, j = 1, \dots, \nu \right\}.
\end{aligned}$$

Therefore comparing (4.9) and (4.18), we obtain the following theorem.

Theorem 4.3 *The \mathbf{m} th ($\mathbf{m} \in \mathcal{Z}^+$) joint binomial moment matrix $\widehat{\mathbf{B}}(t, \mathbf{m})$ is given by*

$$\begin{aligned}
\widehat{\mathbf{B}}(t, \mathbf{m}) &= \sum_{\nu=1}^{|\mathbf{m}|} \sum_{\vec{\mathbf{l}}_\nu \in \mathcal{L}_\nu(\mathbf{m})} \int_0^t du_\nu \int_0^{u_\nu} du_{\nu-1} \cdots \int_0^{u_2} du_1 \\
&\quad \cdot \overline{H}^{(\mathbf{l}_\nu)}(u_\nu) \overline{H}^{(\mathbf{l}_{\nu-1})}(u_{\nu-1}) \cdots \overline{H}^{(\mathbf{l}_1)}(u_1) \\
&\quad \cdot e^{(\mathbf{C}+\mathbf{D})(t-u_\nu)} \widehat{\mathbf{D}}(\mathbf{l}_\nu) e^{(\mathbf{C}+\mathbf{D})(u_\nu-u_{\nu-1})} \widehat{\mathbf{D}}(\mathbf{l}_{\nu-1}) \\
&\quad \quad \quad \cdots \cdots e^{(\mathbf{C}+\mathbf{D})(u_2-u_1)} \widehat{\mathbf{D}}(\mathbf{l}_1) e^{(\mathbf{C}+\mathbf{D})u_1}. \tag{4.19}
\end{aligned}$$

Remark 4.2 Note that $\overline{H}^{(\mathbf{l})}(t)$ ($\mathbf{l} = (l_1, \dots, l_K) \in \mathcal{Z}^+$) has probabilistic meanings. To see this, suppose l_k ($k \in \mathcal{K}$) class k customers simultaneously arrive. Let $H_{k,i}$ ($k \in \mathcal{K}, i = 1, \dots, l_k$) denote a random variable representing the service time of the i th class k customer. $\overline{H}^{(\mathbf{l})}(t)$ then represents the complementary distribution function of $\min_{k \in \mathcal{K}, i=1, \dots, l_k} H_{k,i}$, because

$$\begin{aligned} \overline{H}^{(\mathbf{l})}(t) &= \prod_{k \in \mathcal{K}} \prod_{i=1}^{l_k} \Pr(H_{k,i} > t) \\ &= \Pr\left(\min_{k \in \mathcal{K}, i=1, \dots, l_k} H_{k,i} > t\right). \end{aligned}$$

Remark 4.3 $\widehat{B}(t, \mathbf{m})$ is independent of all $\widehat{D}(\mathbf{l})$ ($\mathbf{l} \geq \mathbf{m}$, $\mathbf{l} \neq \mathbf{m}$). Thus the time-dependent binomial moment of the number of customers in the system depends on the batch size distribution only through its moments of the same and less order.

4.4.2 Time-dependent formula for phase-type services

We now assume that service times of class k customers follow a phase-type distribution with representation $(\boldsymbol{\beta}_k, \mathbf{T}_k)$, i.e.,

$$\overline{H}_k(t) = \boldsymbol{\beta}_k e^{\mathbf{T}_k t} \mathbf{e}, \quad k \in \mathcal{K}.$$

Then using the following properties of Kronecker product \otimes and Kronecker sum \oplus ,

$$\begin{aligned} (\mathbf{U}_1 \mathbf{U}_2 \cdots \mathbf{U}_n) \otimes (\mathbf{V}_1 \mathbf{V}_2 \cdots \mathbf{V}_n) \\ = (\mathbf{U}_1 \otimes \mathbf{V}_1)(\mathbf{U}_2 \otimes \mathbf{V}_2) \cdots (\mathbf{U}_n \otimes \mathbf{V}_n), \quad \forall n \geq 1, \end{aligned} \quad (4.20)$$

$$\exp(\mathbf{U}) \otimes \exp(\mathbf{V}) = \exp(\mathbf{U} \oplus \mathbf{V}), \quad (4.21)$$

we rewrite $\overline{H}^{(\mathbf{l})}(t)$ to be

$$\begin{aligned} \overline{H}^{(\mathbf{l})}(t) &= \prod_{k \in \mathcal{K}} \left(\boldsymbol{\beta}_k e^{\mathbf{T}_k t} \mathbf{e} \right)^{l_k} \\ &= \boldsymbol{\beta}^{<\mathbf{l}>} \exp\left(\mathbf{T}^{[\mathbf{l}]} t\right) \mathbf{e}, \end{aligned}$$

where $\boldsymbol{\beta}^{<\mathbf{l}>}$ and $\mathbf{T}^{[\mathbf{l}]}$ ($\mathbf{l} \in \mathcal{Z}^+$) are given by

$$\begin{aligned} \boldsymbol{\beta}^{<\mathbf{l}>} &= \underbrace{\boldsymbol{\beta}_1 \otimes \cdots \otimes \boldsymbol{\beta}_1}_{l_1} \otimes \underbrace{\boldsymbol{\beta}_2 \otimes \cdots \otimes \boldsymbol{\beta}_2}_{l_2} \otimes \cdots \otimes \underbrace{\boldsymbol{\beta}_K \otimes \cdots \otimes \boldsymbol{\beta}_K}_{l_K}, \\ \mathbf{T}^{[\mathbf{l}]} &= \underbrace{\mathbf{T}_1 \oplus \cdots \oplus \mathbf{T}_1}_{l_1} \oplus \underbrace{\mathbf{T}_2 \oplus \cdots \oplus \mathbf{T}_2}_{l_2} \oplus \cdots \oplus \underbrace{\mathbf{T}_K \oplus \cdots \oplus \mathbf{T}_K}_{l_K}, \end{aligned}$$

respectively. Thus $\widehat{B}(t, \mathbf{m})$ in (4.19) is rewritten to be

$$\widehat{B}(t, \mathbf{m}) = \sum_{\nu=1}^{|\mathbf{m}|} \sum_{\vec{l}_\nu \in \mathcal{L}_\nu(\mathbf{m})} \mathbf{F}_\nu(t, \vec{l}_\nu), \quad (4.22)$$

where

$$\begin{aligned} \mathbf{F}_\nu(t, \vec{l}_\nu) &= \int_0^t du_\nu \int_0^{u_\nu} du_{\nu-1} \cdots \int_0^{u_2} du_1 \prod_{j=1}^\nu \left[\beta^{<l_j>} \exp(\mathbf{T}^{[l_j]} u_j) \mathbf{e} \right] \\ &\quad \cdot e^{(\mathbf{C}+\mathbf{D})(t-u_\nu)} \widehat{\mathbf{D}}(\mathbf{l}_\nu) e^{(\mathbf{C}+\mathbf{D})(u_\nu-u_{\nu-1})} \widehat{\mathbf{D}}(\mathbf{l}_{\nu-1}) \\ &\quad \cdots \cdots e^{(\mathbf{C}+\mathbf{D})(u_2-u_1)} \widehat{\mathbf{D}}(\mathbf{l}_1) e^{(\mathbf{C}+\mathbf{D})u_1}. \end{aligned} \quad (4.23)$$

To obtain a numerically feasible formula for $\mathbf{F}_\nu(t, \vec{l}_\nu)$, we shall rewrite (4.23) by considering the time-reversed arrival process. We define \mathbf{C}^- and $\mathbf{D}^-(\mathbf{n})$ as

$$\mathbf{C}^- = \text{diag}(\boldsymbol{\pi})^{-1} \mathbf{C}^T \text{diag}(\boldsymbol{\pi}), \quad \mathbf{D}^-(\mathbf{n}) = \text{diag}(\boldsymbol{\pi})^{-1} \mathbf{D}(\mathbf{n})^T \text{diag}(\boldsymbol{\pi}),$$

respectively, where $\text{diag}(\boldsymbol{\pi})$ denotes an $M \times M$ diagonal matrix whose i th diagonal element is equal to the i th element of $\boldsymbol{\pi}$. The time-reversed arrival process can then be characterized by \mathbf{C}^- and $\mathbf{D}^-(\mathbf{n})$. Therefore (4.23) is rewritten to be

$$\begin{aligned} \mathbf{F}_\nu(t, \vec{l}_\nu) &= \left[\prod_{\eta=1}^\nu h(\mathbf{l}_\eta) d(\mathbf{l}_\eta) \right] \text{diag}(\boldsymbol{\pi})^{-1} \\ &\quad \cdot \left[\int_0^t du_\nu \int_0^{u_\nu} du_{\nu-1} \cdots \int_0^{u_2} du_1 \prod_{j=1}^\nu \left[\boldsymbol{\kappa}^{(l_j)} \exp(\mathbf{T}^{[l_j]} u_j) (-\mathbf{T}^{[l_j]}) \mathbf{e} \right] \right. \\ &\quad \cdot \exp(\mathbf{Q}^- u_1) \widehat{\mathbf{D}}^-(\mathbf{l}_1) \exp[\mathbf{Q}^-(u_2 - u_1)] \widehat{\mathbf{D}}^-(\mathbf{l}_2) \\ &\quad \left. \cdots \cdots \exp[\mathbf{Q}^-(u_\nu - u_{\nu-1})] \widehat{\mathbf{D}}^-(\mathbf{l}_\nu) \exp[\mathbf{Q}^-(t - u_\nu)] \right]^T \text{diag}(\boldsymbol{\pi}), \end{aligned}$$

where

$$\mathbf{Q}^- = \text{diag}(\boldsymbol{\pi})^{-1} (\mathbf{C} + \mathbf{D})^T \text{diag}(\boldsymbol{\pi}),$$

and for $\mathbf{l} \in \mathcal{Z}^+$,

$$h(\mathbf{l}) = \int_0^\infty dt \beta^{<l>} \exp(\mathbf{T}^{[l]} t) \mathbf{e} = \beta^{<l>} (-\mathbf{T}^{[l]})^{-1} \mathbf{e}, \quad (4.24)$$

$$d(\mathbf{l}) = \max_{i \in \mathcal{M}} \left[\text{diag}(\boldsymbol{\pi})^{-1} \widehat{\mathbf{D}}(\mathbf{l})^T \text{diag}(\boldsymbol{\pi}) \mathbf{e} \right]_i, \quad (4.25)$$

$$\widehat{\mathbf{D}}^-(\mathbf{l}) = d(\mathbf{l})^{-1} \text{diag}(\boldsymbol{\pi})^{-1} \widehat{\mathbf{D}}(\mathbf{l})^T \text{diag}(\boldsymbol{\pi}),$$

$$\boldsymbol{\kappa}^{(l)} = \frac{\beta^{<l>} (-\mathbf{T}^{[l]})^{-1}}{\beta^{<l>} (-\mathbf{T}^{[l]})^{-1} \mathbf{e}}.$$

Remark 4.4 \mathbf{Q}^- represents the infinitesimal generator of the time-reversed process of the underlying Markov chain that governs the arrival process, and satisfies $\boldsymbol{\pi} \mathbf{Q}^- = \mathbf{0}$. Note also that $\widehat{\mathbf{D}}^-(\mathbf{l})$ is a nonnegative matrix whose row sums are all equal to or less than one, i.e., $\widehat{\mathbf{D}}^-(\mathbf{l})$ is a substochastic matrix. Further

$$\boldsymbol{\kappa}^{(l_j)} \exp(\mathbf{T}^{[l_j]} u_j) (-\mathbf{T}^{[l_j]}) \mathbf{e}, \quad j = 1, 2, \dots$$

is considered as the density function of the phase-type distribution with representation $(\boldsymbol{\kappa}^{(l_j)}, \mathbf{T}^{[l_j]})$.

We now define $\mathbf{F}_\nu^-(t, \vec{l}_\nu)$ ($\nu = 1, 2, \dots$) as

$$\begin{aligned} \mathbf{F}_\nu^-(t, \vec{l}_\nu) &= \int_0^t du_\nu \int_0^{u_\nu} du_{\nu-1} \cdots \int_0^{u_2} du_1 \prod_{j=1}^\nu \left[\boldsymbol{\kappa}^{(l_j)} \exp(\mathbf{T}^{[l_j]} u_j) (-\mathbf{T}^{[l_j]}) \mathbf{e} \right] \\ &\quad \cdot \exp(\mathbf{Q}^- u_1) \widehat{\mathbf{D}}^-(l_1) \exp[\mathbf{Q}^-(u_2 - u_1)] \widehat{\mathbf{D}}^-(l_2) \\ &\quad \cdots \cdots \exp[\mathbf{Q}^-(u_\nu - u_{\nu-1})] \widehat{\mathbf{D}}^-(l_\nu) \exp[\mathbf{Q}^-(t - u_\nu)]. \end{aligned}$$

$\mathbf{F}_\nu(t, \vec{l}_\nu)$ can then be rewritten in terms of $\mathbf{F}_\nu^-(t, \vec{l}_\nu)$ as follows:

$$\mathbf{F}_\nu(t, \vec{l}_\nu) = \left[\prod_{\eta=1}^\nu h(l_\eta) d(l_\eta) \right] \text{diag}(\boldsymbol{\pi})^{-1} \left[\mathbf{F}_\nu^-(t, \vec{l}_\nu) \right]^\text{T} \text{diag}(\boldsymbol{\pi}). \quad (4.26)$$

In what follows, we derive a numerically feasible formula for $\mathbf{F}_\nu^-(t, \mathbf{l}_\nu)$.

We first consider $\mathbf{F}_1^-(t, \mathbf{l}_1)$;

$$\begin{aligned} \mathbf{F}_1^-(t, \mathbf{l}_1) &= \int_0^t du_1 \left[\boldsymbol{\kappa}^{(l_1)} \exp(\mathbf{T}^{[l_1]} u_1) (-\mathbf{T}^{[l_1]}) \mathbf{e} \right] \\ &\quad \cdot \exp(\mathbf{Q}^- u_1) \widehat{\mathbf{D}}^-(l_1) \exp[\mathbf{Q}^-(t - u_1)]. \end{aligned} \quad (4.27)$$

Using (4.20) and (4.21), we rewrite (4.27) to be

$$\begin{aligned} \mathbf{F}_1^-(t, \mathbf{l}_1) &= \int_0^t du_1 \left[\boldsymbol{\kappa}^{(l_1)} \cdot \exp(\mathbf{T}^{[l_1]} u_1) \cdot (-\mathbf{T}^{[l_1]}) \mathbf{e} \right] \\ &\quad \cdot \left[\mathbf{I}_Q \cdot \exp(\mathbf{Q}^- u_1) \cdot \widehat{\mathbf{D}}^-(l_1) \right] \exp[\mathbf{Q}^-(t - u_1)] \\ &= \left(\boldsymbol{\kappa}^{(l_1)} \otimes \mathbf{I}_Q \right) \int_0^t du_1 \exp(\mathbf{T}^{[l_1]} \oplus \mathbf{Q}^- u_1) \\ &\quad \cdot \left[(-\mathbf{T}^{[l_1]}) \mathbf{e} \otimes \widehat{\mathbf{D}}^-(l_1) \right] \exp[\mathbf{Q}^-(t - u_1)], \end{aligned}$$

where \mathbf{I}_Q denotes an identity matrix of the same size as \mathbf{Q}^- .

Similarly, $\mathbf{F}_2^-(t, \vec{l}_2)$ is rewritten to be

$$\begin{aligned} \mathbf{F}_2^-(t, \vec{l}_2) &= \int_0^t du_2 \int_0^{u_2} du_1 \\ &\quad \cdot \left[\boldsymbol{\kappa}^{(l_1)} \cdot \exp(\mathbf{T}^{[l_1]} u_1) \cdot (-\mathbf{T}^{[l_1]}) \mathbf{e} \cdot 1 \cdot 1 \cdot 1 \right] \\ &\quad \cdot \left[\boldsymbol{\kappa}^{(l_2)} \cdot \exp(\mathbf{T}^{[l_2]} u_1) \cdot \mathbf{I}_1(l_2) \cdot \exp[\mathbf{T}^{[l_2]}(u_2 - u_1)] \cdot (-\mathbf{T}^{[l_2]}) \mathbf{e} \cdot 1 \right] \\ &\quad \cdot \left[\mathbf{I}_Q \cdot \exp(\mathbf{Q}^- u_1) \cdot \widehat{\mathbf{D}}^-(l_1) \cdot \exp[\mathbf{Q}^-(u_2 - u_1)] \cdot \widehat{\mathbf{D}}^-(l_2) \cdot \exp[\mathbf{Q}^-(t - u_2)] \right] \\ &= \left(\boldsymbol{\kappa}^{(l_1)} \otimes \boldsymbol{\kappa}^{(l_2)} \otimes \mathbf{I}_Q \right) \int_0^t du_2 \int_0^{u_2} du_1 \exp(\mathbf{T}^{[l_1]} \oplus \mathbf{T}^{[l_2]} \oplus \mathbf{Q}^- u_1) \\ &\quad \cdot \left[(-\mathbf{T}^{[l_1]}) \mathbf{e} \otimes \mathbf{I}_1(l_2) \otimes \widehat{\mathbf{D}}^-(l_1) \right] \exp[\mathbf{T}^{[l_2]} \oplus \mathbf{Q}^-(u_2 - u_1)] \\ &\quad \cdot \left[(-\mathbf{T}^{[l_2]}) \mathbf{e} \otimes \widehat{\mathbf{D}}^-(l_2) \right] \exp[\mathbf{Q}^-(t - u_2)], \end{aligned}$$

where $\mathbf{I}_1(l_2)$ denotes an identity matrix of the same size as $\mathbf{T}^{[l_2]}$. Following the same manipulation as in the cases $\nu = 1$ and 2, we obtain for $\nu = 1, 2, \dots$,

$$\mathbf{F}_\nu^-(t, \vec{l}_\nu) = \mathbf{J}_\nu(\vec{l}_\nu) \int_0^t du_\nu \int_0^{u_\nu} du_{\nu-1} \cdots \int_0^{u_2} du_1$$

$$\begin{aligned} & \cdot \exp \left[\mathbf{U}_{\nu,1}(\vec{l}_\nu) u_1 \right] \mathbf{V}_{\nu,1}(\vec{l}_\nu) \exp \left[\mathbf{U}_{\nu,2}(\vec{l}_\nu) (u_2 - u_1) \right] \mathbf{V}_{\nu,2}(\vec{l}_\nu) \\ & \cdots \cdots \exp \left[\mathbf{U}_{\nu,\nu}(\vec{l}_\nu) (u_\nu - u_{\nu-1}) \right] \mathbf{V}_{\nu,\nu}(\vec{l}_\nu) \exp[\mathbf{Q}^-(t - u_\nu)], \end{aligned} \quad (4.28)$$

where

$$\mathbf{J}_\nu(\vec{l}_\nu) = \boldsymbol{\kappa}^{(l_1)} \otimes \cdots \otimes \boldsymbol{\kappa}^{(l_\nu)} \otimes \mathbf{I}_Q, \quad (4.29)$$

$$\mathbf{U}_{\nu,j}(\vec{l}_\nu) = \mathbf{T}^{[l_j]} \oplus \cdots \oplus \mathbf{T}^{[l_\nu]} \oplus \mathbf{Q}^-, \quad j = 1, \dots, \nu, \quad (4.30)$$

$$\mathbf{V}_{\nu,j}(\vec{l}_\nu) = \begin{cases} \left(-\mathbf{T}^{[l_j]} \right) \mathbf{e} \otimes \mathbf{I}_{\nu-j}(l_{j+1}, \dots, l_\nu) \otimes \widehat{\mathbf{D}}^-(l_j), & j = 1, \dots, \nu - 1, \\ \left(-\mathbf{T}^{[l_\nu]} \right) \mathbf{e} \otimes \widehat{\mathbf{D}}^-(l_\nu), & j = \nu, \end{cases} \quad (4.31)$$

and where $\mathbf{I}_{\nu-j}(l_{j+1}, l_{j+2}, \dots, l_\nu)$ ($\nu - j \geq 1$) denotes an identity matrix whose size is equal to that of $\mathbf{T}^{[l_{j+1}]} \oplus \mathbf{T}^{[l_{j+2}]} \oplus \cdots \oplus \mathbf{T}^{[l_\nu]}$.

Lemma 4.1 *Let \mathbf{U}_j ($j = 1, \dots, \nu$) denote an $m_i \times m_i$ matrix and \mathbf{V}_j ($j = 0, \dots, \nu$) denote an $m_j \times m_{j+1}$ matrix. If \mathbf{U}_j 's and \mathbf{V}_j 's are all bounded, the following equation holds for all ν ($\nu = 2, 3, \dots$).*

$$\begin{aligned} & \mathbf{V}_0 \int_0^t du_{\nu-1} \int_0^{u_{\nu-1}} du_{\nu-2} \cdots \int_0^{u_2} du_1 e^{\mathbf{U}_1 u_1} \mathbf{V}_1 e^{\mathbf{U}_2 (u_2 - u_1)} \mathbf{V}_2 \\ & \cdots \cdots e^{\mathbf{U}_{\nu-1} (u_{\nu-1} - u_{\nu-2})} \mathbf{V}_{\nu-1} e^{\mathbf{U}_\nu (t - u_{\nu-1})} \mathbf{V}_\nu \\ & = \left[\mathbf{V}_0 \quad \mathbf{O} \quad \cdots \quad \mathbf{O} \right] \exp \left[\begin{pmatrix} \mathbf{U}_1 & \mathbf{V}_1 & \mathbf{O} & \cdots & \mathbf{O} \\ \mathbf{O} & \mathbf{U}_2 & \mathbf{V}_2 & & \vdots \\ \vdots & & \ddots & \ddots & \mathbf{O} \\ \mathbf{O} & \cdots & \mathbf{O} & \mathbf{U}_{\nu-1} & \mathbf{V}_{\nu-1} \\ \mathbf{O} & \mathbf{O} & \cdots & \mathbf{O} & \mathbf{U}_\nu \end{pmatrix} t \right] \begin{bmatrix} \mathbf{O} \\ \vdots \\ \mathbf{O} \\ \mathbf{V}_\nu \end{bmatrix}. \end{aligned} \quad (4.32)$$

The proof of Lemma 4.1 is given in Appendix G.

We define $\mathbf{T}_\nu(\vec{l}_\nu)$ as

$$\mathbf{T}_\nu(\vec{l}_\nu) = \begin{pmatrix} \mathbf{U}_{\nu,1}(\vec{l}_\nu) & \mathbf{V}_{\nu,1}(\vec{l}_\nu) & \mathbf{O} & \cdots & \mathbf{O} \\ \mathbf{O} & \mathbf{U}_{\nu,2}(\vec{l}_\nu) & \mathbf{V}_{\nu,2}(\vec{l}_\nu) & & \vdots \\ \vdots & & \ddots & \ddots & \mathbf{O} \\ \mathbf{O} & \cdots & \mathbf{O} & \mathbf{U}_{\nu,\nu}(\vec{l}_\nu) & \mathbf{V}_{\nu,\nu}(\vec{l}_\nu) \\ \mathbf{O} & \mathbf{O} & \cdots & \mathbf{O} & \mathbf{Q}^- \end{pmatrix}, \quad (4.33)$$

where $\mathbf{U}_{\nu,j}(\vec{l}_\nu)$ and $\mathbf{V}_{\nu,j}(\vec{l}_\nu)$ ($j = 1, \dots, \nu$) are given in (4.30) and (4.31), respectively. Applying Lemma 4.1 to (4.28) and taking account of (4.22) and (4.26), we have the following theorem.

Theorem 4.4 $\widehat{\mathbf{B}}(t, \mathbf{m})$ is given by

$$\widehat{\mathbf{B}}(t, \mathbf{m}) = \sum_{\nu=1}^{|\mathbf{m}|} \sum_{\vec{l}_\nu \in \mathcal{L}_\nu(\mathbf{m})} \left[\prod_{\eta=1}^{\nu} h(l_\eta) d(l_\eta) \right] \text{diag}(\boldsymbol{\pi})^{-1} \left[\mathbf{F}_\nu^-(t, \vec{l}_\nu) \right]^T \text{diag}(\boldsymbol{\pi}), \quad (4.34)$$

where $h(\mathbf{l}_\eta)$ and $d(\mathbf{l}_\eta)$ ($\mathbf{l}_\eta \in \mathcal{Z}^+$) are given in (4.24) and (4.25), respectively. Further $\mathbf{F}_\nu^-(t, \vec{\mathbf{l}}_\nu)$ ($\nu = 1, 2, \dots$) is given by

$$\mathbf{F}_\nu^-(t, \vec{\mathbf{l}}_\nu) = \begin{bmatrix} \mathbf{J}_\nu(\vec{\mathbf{l}}_\nu) & \mathbf{O} & \cdots & \mathbf{O} \end{bmatrix} \exp[\mathbf{T}_\nu(\vec{\mathbf{l}}_\nu)t] \begin{bmatrix} \mathbf{O} \\ \vdots \\ \mathbf{O} \\ \mathbf{I}_Q \end{bmatrix}, \quad (4.35)$$

where $\mathbf{J}_\nu(\vec{\mathbf{l}}_\nu)$ and $\mathbf{T}_\nu(\vec{\mathbf{l}}_\nu)$ are given in (4.29) and (4.33), respectively.

Remark 4.5 $\mathbf{T}_\nu(\vec{\mathbf{l}}_\nu)$ in (4.33) is considered as the defective infinitesimal generator of an absorbing Markov chain. Namely, $\mathbf{T}_\nu(\vec{\mathbf{l}}_\nu)$ has negative diagonal elements and nonnegative off-diagonal elements, all its row sums are not positive, and at least one row sum is strictly negative. Therefore applying the uniformization technique (see Appendix A), we can readily compute $\exp[\mathbf{T}_\nu(\vec{\mathbf{l}}_\nu)t]$. Further since $\exp[\mathbf{T}_\nu(\vec{\mathbf{l}}_\nu)t]$, $\mathbf{J}_\nu(\vec{\mathbf{l}}_\nu)$, \mathbf{I}_Q , $\text{diag}(\boldsymbol{\pi})$, $\text{diag}(\boldsymbol{\pi})^{-1}$, $h(\mathbf{l}_\eta)$ and $d(\mathbf{l}_\eta)$ are all nonnegative, the computation of $\widehat{\mathbf{B}}(t, \mathbf{m})$ is numerically stable.

4.4.3 Limiting formula for phase-type services

In this subsection, assuming phase-type service times, we derive an explicit and numerically feasible formula for the limit $\widehat{\mathbf{B}}(\mathbf{m})$ of $\widehat{\mathbf{B}}(t, \mathbf{m})$:

$$\widehat{\mathbf{B}}(\mathbf{m}) = \lim_{t \rightarrow \infty} \widehat{\mathbf{B}}(t, \mathbf{m}), \quad \mathbf{m} \in \mathcal{Z}^+.$$

We define $\mathbf{F}_\nu^-(\vec{\mathbf{l}}_\nu)$ as

$$\mathbf{F}_\nu^-(\vec{\mathbf{l}}_\nu) = \lim_{t \rightarrow \infty} \mathbf{F}_\nu^-(t, \vec{\mathbf{l}}_\nu).$$

Theorem 4.5 $\widehat{\mathbf{B}}(\mathbf{m})$ ($\mathbf{m} \in \mathcal{Z}^+$) is given by

$$\widehat{\mathbf{B}}(\mathbf{m}) = \sum_{\nu=1}^{|\mathbf{m}|} \sum_{\vec{\mathbf{l}}_\nu \in \mathcal{L}_\nu(\mathbf{m})} \left[\prod_{\eta=1}^{\nu} h(\mathbf{l}_\eta) d(\mathbf{l}_\eta) \right] \text{diag}(\boldsymbol{\pi})^{-1} \mathbf{F}_\nu^-(\vec{\mathbf{l}}_\nu)^T \text{diag}(\boldsymbol{\pi}), \quad (4.36)$$

where

$$\mathbf{F}_\nu^-(\vec{\mathbf{l}}_\nu) = \mathbf{J}_\nu(\vec{\mathbf{l}}_\nu) \left[-\mathbf{U}_{\nu,1}(\vec{\mathbf{l}}_\nu) \right]^{-1} \mathbf{V}_{\nu,1}(\vec{\mathbf{l}}_\nu) \cdots \cdots \left[-\mathbf{U}_{\nu,\nu}(\vec{\mathbf{l}}_\nu) \right]^{-1} \mathbf{V}_{\nu,\nu}(\vec{\mathbf{l}}_\nu) \mathbf{e}\boldsymbol{\pi}. \quad (4.37)$$

Note that $\mathbf{J}_\nu(\vec{\mathbf{l}}_\nu)$, $\mathbf{U}_{\nu,j}(\vec{\mathbf{l}}_\nu)$ and $\mathbf{V}_{\nu,j}(\vec{\mathbf{l}}_\nu)$ ($j = 1, \dots, \nu$) are given in (4.29), (4.30) and (4.31), respectively.

Proof. Taking the limit of both sides of (4.34) as $t \rightarrow \infty$ yields (4.36). In what follows, we show (4.37). We define a submatrix $\mathbf{T}'_\nu(\vec{\mathbf{l}}_\nu)$ of $\mathbf{T}_\nu(\vec{\mathbf{l}}_\nu)$ as

$$\mathbf{T}'_\nu(\vec{\mathbf{l}}_\nu) = \begin{bmatrix} \mathbf{U}_{\nu,1}(\vec{\mathbf{l}}_\nu) & \mathbf{V}_{\nu,1}(\vec{\mathbf{l}}_\nu) & \mathbf{O} & \cdots & \mathbf{O} \\ \mathbf{O} & \mathbf{U}_{\nu,2}(\vec{\mathbf{l}}_\nu) & \mathbf{V}_{\nu,2}(\vec{\mathbf{l}}_\nu) & & \vdots \\ \vdots & & \ddots & \ddots & \mathbf{O} \\ \mathbf{O} & \cdots & \mathbf{O} & \mathbf{U}_{\nu,\nu-1}(\vec{\mathbf{l}}_\nu) & \mathbf{V}_{\nu,\nu-1}(\vec{\mathbf{l}}_\nu) \\ \mathbf{O} & \mathbf{O} & \cdots & \mathbf{O} & \mathbf{U}_{\nu,\nu}(\vec{\mathbf{l}}_\nu) \end{bmatrix}. \quad (4.38)$$

Then $\mathbf{F}_\nu^-(t, \vec{l}_\nu)$ in (4.35) can be rewritten to be

$$\mathbf{F}_\nu^-(t, \vec{l}_\nu) = \left[\mathbf{J}_\nu(\vec{l}_\nu) \quad \mathbf{O} \quad \cdots \quad \mathbf{O} \mid \mathbf{O} \right] \exp \left(\left(\begin{array}{c|c} \mathbf{T}'_\nu(\vec{l}_\nu) & \begin{bmatrix} \mathbf{O} \\ \vdots \\ \mathbf{O} \\ \mathbf{V}_{\nu,\nu}(\vec{l}_\nu) \end{bmatrix} \\ \hline \mathbf{O} & \mathbf{Q}^- \end{array} \right) t \right) \begin{bmatrix} \mathbf{O} \\ \vdots \\ \mathbf{O} \\ \mathbf{I}_Q \end{bmatrix}.$$

Applying Lemma 4.1 to the above equation yields

$$\mathbf{F}_\nu^-(t, \vec{l}_\nu) = \left[\mathbf{J}_\nu(\vec{l}_\nu) \quad \mathbf{O} \quad \cdots \quad \mathbf{O} \right] \Theta_\nu(t, \vec{l}_\nu), \quad (4.39)$$

where

$$\Theta_\nu(t, \vec{l}_\nu) = \int_0^t du \exp \left[\mathbf{T}'_\nu(\vec{l}_\nu) u \right] \begin{bmatrix} \mathbf{O} \\ \vdots \\ \mathbf{O} \\ \mathbf{V}_{\nu,\nu}(\vec{l}_\nu) \end{bmatrix} \exp[\mathbf{Q}^-(t-u)]. \quad (4.40)$$

Note that $\mathbf{T}'_\nu(\vec{l}_\nu)$ in (4.38) is regarded as the defective infinitesimal generator of an absorbing Markov chain in transient states, and therefore all row sums of $\mathbf{T}'_\nu(\vec{l}_\nu)$ are strictly negative. Thus we have

$$\int_0^\infty du \exp \left[\mathbf{T}'_\nu(\vec{l}_\nu) u \right] = \left[-\mathbf{T}'_\nu(\vec{l}_\nu) \right]^{-1}. \quad (4.41)$$

Further, since \mathbf{Q}^- is the infinitesimal generator of an irreducible Markov chain with the stationary probability vector $\boldsymbol{\pi}$ (see Remark 4.4), we have

$$\lim_{t \rightarrow \infty} \exp(\mathbf{Q}^- t) = \mathbf{e}\boldsymbol{\pi}. \quad (4.42)$$

It then follows from (4.40), (4.41) and (4.42) that

$$\begin{aligned} \Theta_\nu(\vec{l}_\nu) &= \lim_{t \rightarrow \infty} \Theta_\nu(t, \vec{l}_\nu) \\ &= \int_0^\infty du \exp \left[\mathbf{T}'_\nu(\vec{l}_\nu) u \right] \begin{bmatrix} \mathbf{O} \\ \vdots \\ \mathbf{O} \\ \mathbf{V}_{\nu,\nu}(\vec{l}_\nu) \end{bmatrix} \exp(-\mathbf{Q}^- u) \cdot \lim_{t \rightarrow \infty} \exp(\mathbf{Q}^- t) \\ &= \int_0^\infty du \exp \left[\mathbf{T}'_\nu(\vec{l}_\nu) u \right] \cdot \begin{bmatrix} \mathbf{O} \\ \vdots \\ \mathbf{O} \\ \mathbf{V}_{\nu,\nu}(\vec{l}_\nu) \end{bmatrix} \mathbf{e}\boldsymbol{\pi} \\ &= \left[-\mathbf{T}'_\nu(\vec{l}_\nu) \right]^{-1} \begin{bmatrix} \mathbf{O} \\ \vdots \\ \mathbf{O} \\ \mathbf{V}_{\nu,\nu}(\vec{l}_\nu) \end{bmatrix} \mathbf{e}\boldsymbol{\pi}, \end{aligned} \quad (4.43)$$

where we use $\exp(\mathbf{Q}^{-1}t)\mathbf{e} = \mathbf{e}$ ($t \geq 0$). Finally, from (4.39) and (4.43), we obtain

$$\begin{aligned} \mathbf{F}_\nu^-(\vec{l}_\nu) &= \begin{bmatrix} \mathbf{J}_\nu(\vec{l}_\nu) & \mathbf{O} & \cdots & \mathbf{O} \end{bmatrix} \boldsymbol{\Theta}_\nu(\vec{l}_\nu) \\ &= \begin{bmatrix} \mathbf{J}_\nu(\vec{l}_\nu) & \mathbf{O} & \cdots & \mathbf{O} \end{bmatrix} [-\mathbf{T}'_\nu(\vec{l}_\nu)]^{-1} \begin{bmatrix} \mathbf{O} \\ \vdots \\ \mathbf{O} \\ \mathbf{V}_{\nu,\nu}(\vec{l}_\nu) \end{bmatrix} \mathbf{e}\boldsymbol{\pi}, \end{aligned}$$

from which Theorem 4.5 follows. ■

4.5 Numerical Examples

In this section, we show some numerical examples based on Theorems 4.4 and 4.5, and discuss the impact of system parameters on the mean and variance of the number of customers in the system under the assumption that the arrival process is stationary. Note that when the arrival process is stationary, the mean number $\mathbf{E}[N_k(t)]$ ($k \in \mathcal{K}$) of class k customers at time t is given by

$$\mathbf{E}[N_k(t)] = \boldsymbol{\pi} \widehat{\mathbf{B}}(t, \mathbf{e}_k) \mathbf{e} = \lambda_k \int_0^t du \overline{H}_k(u),$$

where λ_k is the arrival rate of class k customers and \mathbf{e}_k is given in (1.4). Further the limit $\mathbf{E}[N_k]$ of $\mathbf{E}[N_k(t)]$ is given by

$$\mathbf{E}[N_k] = \boldsymbol{\pi} \widehat{\mathbf{B}}(\mathbf{e}_k) \mathbf{e} = \lambda_k h_k,$$

where h_k denotes the mean service time of class k customers. Therefore the mean $\mathbf{E}[N(t)]$ of the total number of customers at time t and its limit $\mathbf{E}[N]$ are given by

$$\mathbf{E}[N(t)] = \sum_{k \in \mathcal{K}} \lambda_k \int_0^t du \overline{H}_k(u), \quad \mathbf{E}[N] = \sum_{k \in \mathcal{K}} \lambda_k h_k,$$

respectively.

On the other hand, the variance $\text{Var}[N_k(t)]$ ($k \in \mathcal{K}$) of the number of class k customers at time t and its limit $\text{Var}[N_k]$ are given by

$$\begin{aligned} \text{Var}[N_k(t)] &= 2\boldsymbol{\pi} \widehat{\mathbf{B}}(t, 2\mathbf{e}_k) \mathbf{e} + \mathbf{E}[N_k(t)] - \mathbf{E}[N_k(t)]^2, \\ \text{Var}[N_k] &= 2\boldsymbol{\pi} \widehat{\mathbf{B}}(2\mathbf{e}_k) \mathbf{e} + \mathbf{E}[N_k] - \mathbf{E}[N_k]^2, \end{aligned}$$

respectively. Further the covariance $\text{Cov}[N_k(t), N_{k'}(t)]$ ($k, k' \in \mathcal{K}$) of the numbers of class k and k' customers at time t and its limit $\text{Cov}[N_k, N_{k'}]$ are given by

$$\begin{aligned} \text{Cov}[N_k(t), N_{k'}(t)] &= \boldsymbol{\pi} \widehat{\mathbf{B}}(t, \mathbf{e}_k + \mathbf{e}_{k'}) \mathbf{e} - \mathbf{E}[N_k(t)]\mathbf{E}[N_{k'}(t)], \\ \text{Cov}[N_k, N_{k'}] &= \boldsymbol{\pi} \widehat{\mathbf{B}}(\mathbf{e}_k + \mathbf{e}_{k'}) \mathbf{e} - \mathbf{E}[N_k]\mathbf{E}[N_{k'}], \end{aligned}$$

respectively. The variance $\text{Var}[N(t)]$ of the total number of customers at time t and its limit $\text{Var}[N]$ are given in terms of the marginal variances and covariances:

$$\begin{aligned}\text{Var}[N(t)] &= \sum_{k \in \mathcal{K}} \text{Var}[N_k(t)] + 2 \sum_{k=1}^{K-1} \sum_{k'=k+1}^K \text{Cov}[N_k(t), N_{k'}(t)], \\ \text{Var}[N] &= \sum_{k \in \mathcal{K}} \text{Var}[N_k] + 2 \sum_{k=1}^{K-1} \sum_{k'=k+1}^K \text{Cov}[N_k, N_{k'}],\end{aligned}$$

respectively.

4.5.1 Impact of service time distribution on $\text{Var}[N]$

We first show the impact of the service time distribution on the limiting variance of the number of customers. To do so, we consider the following batch MMPP. There is only one class, i.e., $\mathcal{K} = \{1\}$, and

$$\mathbf{C} = \begin{bmatrix} -9 & 1 \\ 1 & -3 \end{bmatrix}, \quad \mathbf{D} = \begin{bmatrix} 8 & 0 \\ 0 & 2 \end{bmatrix},$$

$$\mathbf{D}(n) = g(n)\mathbf{D}, \quad g(n) = \begin{cases} 1/2, & \text{if } n = 1, \\ 1/2, & \text{if } n = 3, \\ 0, & \text{otherwise.} \end{cases}$$

Note that the mean arrival rate λ_1 is equal to 10. As for the service time distribution, we fix the mean service time to be one (i.e., $\text{E}[N] = 10$), and consider the following three distributions: [k-stage Erlang distribution]

$$\bar{H}_1(t) = \left[1 \underbrace{0 \dots 0}_{k-1} \right] \exp \left[\begin{pmatrix} -k & k & 0 & \dots & 0 \\ 0 & -k & k & \dots & 0 \\ \vdots & & \ddots & \ddots & \vdots \\ 0 & 0 & \dots & -k & k \\ 0 & 0 & \dots & 0 & -k \end{pmatrix} t \right] \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \\ 1 \end{bmatrix} \quad k = 2, 3, \dots, \quad (4.44)$$

[Exponential distribution]

$$\bar{H}_1(t) = e^{-t}, \quad (4.45)$$

[Two-state balanced hyper-exponential distribution]

$$\bar{H}_1(t) = \begin{bmatrix} p & 1-p \end{bmatrix} \exp \left[\begin{pmatrix} -2p & 0 \\ 0 & -2(1-p) \end{pmatrix} t \right] \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad 0 < p < 0.5. \quad (4.46)$$

Let C_v denote the coefficient of variation of the service time distribution. Note that $C_v = 1/\sqrt{k}$ ($k = 2, 3, \dots$) for k -stage Erlang distribution, $C_v = 1$ for exponential distribution and $C_v = \sqrt{1/\{2p(1-p)\} - 1}$ ($0 < p < 0.5$) for two-state balanced hyper-exponential distribution. Figure 4.1 plots the limiting variance $\text{Var}[N]$ of the number of customers in the system as a function of C_v . We observe that $\text{Var}[N]$ decreases as C_v increases.

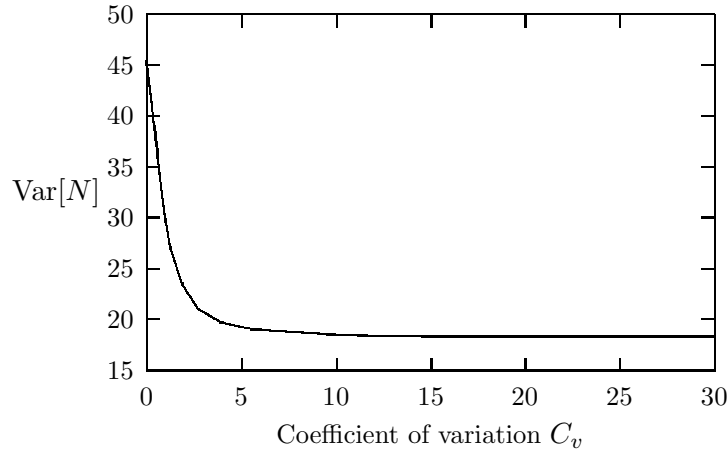


Figure 4.1: Limiting variance of the number of customers.

Remark 4.6 For the $M(t)/G/\infty$ queue with a sinusoidal arrival rate, Eick et al. [Eick93] show that as C_v increases, the deviation of the arrival process has a less impact on the mean number of customers in the system. Thus our observation coincides with that in [Eick93].

4.5.2 Impact of arrival process

We consider the impact of the correlation in the arrival process on the variance of the number of customers. For this purpose, we assume that batches arrive according to the following two-state Markov modulated Poisson process:

$$\mathbf{C} = \begin{bmatrix} -\lambda_1 - c & c \\ c & -\lambda_2 - c \end{bmatrix}, \quad \mathbf{D} = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}, \quad c > 0,$$

$$\mathbf{D}(n) = g(n)\mathbf{D}, \quad n = 1, 2, \dots,$$

where $\{g(n); n = 1, 2, \dots\}$ denotes the batch size distribution, i.e., $\sum_{n=1}^{\infty} g(n) = 1$. Let G denote a random variable representing the batch size. We define $A(t, \tau)$ ($0 \leq t < \tau$) as the number of customers arriving in time interval $(t, \tau]$. We then have

$$\text{Cov}[A(t, t+s), A(t+s+\tau, t+2s+\tau)] = \frac{\mathbb{E}[G]^2 \left(\frac{\lambda_1 - \lambda_2}{2}\right)^2 (1 - e^{-2cs})^2}{(2c)^2} e^{-2c\tau}.$$

Note that the time-correlation of the arrival process increases with the mean sojourn time $1/c$ in each state of the underlying Markov chain. As for the service time distribution, we consider three cases: the two-stage Erlang distribution in (4.44) with $k = 2$, the exponential distribution in (4.45), and the two-state hyper-exponential distribution in (4.46) with $p = 0.25$.

Figure 4.2 plots the limiting variance $\text{Var}[N]$ of the number of customers in the system as a function of $1/2c$, where $\lambda_1 = 8$, $\lambda_2 = 2$ and

$$g(n) = \begin{cases} 1/2, & \text{if } n = 1, \\ 1/2, & \text{if } n = 3, \\ 0, & \text{otherwise.} \end{cases}$$

We observe that $\text{Var}[N]$ increases with the correlation in arrivals.

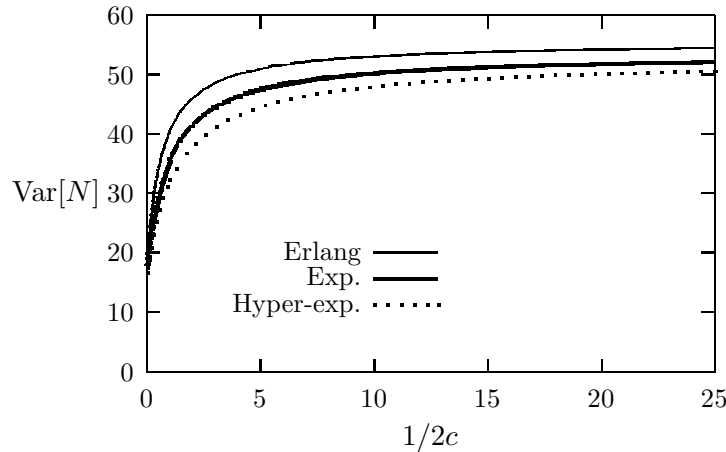


Figure 4.2: Limiting variance of the number of customers.

Next we consider the time-dependent variance $\text{Var}[N(t)]$ with two stationary arrival streams. We fix the marginal characteristics of each arrival stream, which is represented by

$$\mathbf{C} = \begin{bmatrix} -6 & 1 \\ 1 & -1 \end{bmatrix}, \quad \mathbf{D} = \begin{bmatrix} 5 & 0 \\ 0 & 0 \end{bmatrix},$$

$$g(n) = \begin{cases} 1/2, & \text{if } n = 1, \\ 1/2, & \text{if } n = 3, \\ 0, & \text{otherwise,} \end{cases} \quad \mathbf{D}(n) = g(n)\mathbf{D},$$

$$\overline{\mathbf{H}}(t) = \begin{bmatrix} 0.5 & 0.5 \end{bmatrix} \exp \left[\begin{pmatrix} -0.75 & 0 \\ 0 & -1.5 \end{pmatrix} t \right] \begin{bmatrix} 1 \\ 1 \end{bmatrix}.$$

where $\overline{\mathbf{H}}(t)$ denotes the complementary distribution of service times. Under this restriction, we consider the following three cases.

Case N: Negatively correlated arrival streams.

Two arrival streams are negatively correlated. Namely, the two arrival streams are governed by a single two-state underlying Markov chain and when it is in state k ($k = 1, 2$), only class k customers can arrive. More precisely, the overall arrival process is characterized by

$$\mathbf{C} = \begin{bmatrix} -6 & 1 \\ 1 & -6 \end{bmatrix}, \quad \mathbf{D}_1 = \begin{bmatrix} 5 & 0 \\ 0 & 0 \end{bmatrix}, \quad \mathbf{D}_2 = \begin{bmatrix} 0 & 0 \\ 0 & 5 \end{bmatrix},$$

$$\begin{aligned}
g(n) &= \begin{cases} 1/2, & \text{if } n = 1, \\ 1/2, & \text{if } n = 3, \\ 0, & \text{otherwise.} \end{cases} \\
\mathbf{D}(n_1, n_2) &= \begin{cases} g(n_1)\mathbf{D}_1, & \text{if } n_1 \geq 1 \text{ and } n_2 = 0, \\ g(n_2)\mathbf{D}_2, & \text{if } n_1 = 0 \text{ and } n_2 \geq 1, \\ \mathbf{O} & \text{otherwise,} \end{cases} \\
\bar{H}_k(t) &= \begin{bmatrix} 0.5 & 0.5 \end{bmatrix} \exp \left[\begin{pmatrix} -0.75 & 0 \\ 0 & -1.5 \end{pmatrix} t \right] \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad k = 1, 2.
\end{aligned}$$

Case I: Superposition of two independent arrival streams.

Two arrival streams are independent of each other. Thus the overall arrival process is characterized by

$$\begin{aligned}
\mathbf{C} &= \begin{bmatrix} -6 & 1 \\ 1 & -1 \end{bmatrix} \oplus \begin{bmatrix} -6 & 1 \\ 1 & -1 \end{bmatrix} = \begin{bmatrix} -12 & 1 & 1 & 0 \\ 1 & -7 & 0 & 1 \\ 1 & 0 & -7 & 1 \\ 0 & 1 & 1 & -2 \end{bmatrix}, \\
\mathbf{D}_1 &= \begin{bmatrix} 5 & 0 \\ 0 & 0 \end{bmatrix} \otimes \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 5 & 0 & 0 & 0 \\ 0 & 5 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}, \\
\mathbf{D}_2 &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \otimes \begin{bmatrix} 5 & 0 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 5 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 5 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}, \\
g(n) &= \begin{cases} 1/2, & \text{if } n = 1, \\ 1/2, & \text{if } n = 3, \\ 0, & \text{otherwise,} \end{cases} \quad \mathbf{D}(n_1, n_2) = \begin{cases} g(n_1)\mathbf{D}_1, & \text{if } n_1 \geq 1 \text{ and } n_2 = 0, \\ g(n_2)\mathbf{D}_2, & \text{if } n_1 = 0 \text{ and } n_2 \geq 1, \\ \mathbf{O}, & \text{otherwise,} \end{cases} \\
\bar{H}_k(t) &= \begin{bmatrix} 0.5 & 0.5 \end{bmatrix} \exp \left[\begin{pmatrix} -0.75 & 0 \\ 0 & -1.5 \end{pmatrix} t \right] \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad k = 1, 2.
\end{aligned}$$

Case P: Positively correlated arrival streams.

Two arrival streams are positively correlated. Namely, arrivals from both streams occur simultaneously with probability one. More precisely, the overall arrival process is characterized by

$$\begin{aligned}
\mathbf{C} &= \begin{bmatrix} -6 & 1 \\ 1 & -1 \end{bmatrix}, \quad \mathbf{D} = \begin{bmatrix} 5 & 0 \\ 0 & 0 \end{bmatrix}, \\
g(n) &= \begin{cases} 1/2, & \text{if } n = 1, \\ 1/2, & \text{if } n = 3, \\ 0, & \text{otherwise,} \end{cases} \quad \mathbf{D}(n_1, n_2) = g(n_1)g(n_2)\mathbf{D},
\end{aligned}$$

$$\bar{H}_k(t) = \begin{bmatrix} 0.5 & 0.5 \end{bmatrix} \exp \left[\begin{pmatrix} -0.75 & 0 \\ 0 & -1.5 \end{pmatrix} t \right] \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad k = 1, 2.$$

Because the marginal characteristics of each stream in these three cases are the same, the time-dependent and limiting marginal distributions of the numbers of customers in respective classes are identical among these three cases. However, the covariances of the numbers of customers in respective classes are different. Figure 4.3 shows the time-dependent covariance $\text{Cov}[N_1(t), N_2(t)]$ of the numbers of customers in respective classes as a function of t .

We observe that positive (resp. negative) correlation between two arrival streams leads to positive (resp. negative) covariance at any time t . Thus the positive (resp. negative) correlation in arrivals leads to large (resp. small) variance, compared with the superposition of independent streams, as shown in Figure 4.4.

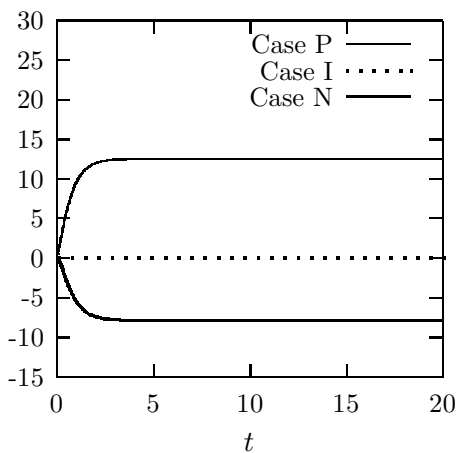


Figure 4.3: Time-dependent covariance $\text{Cov}[N(t)]$ of the number of customers.

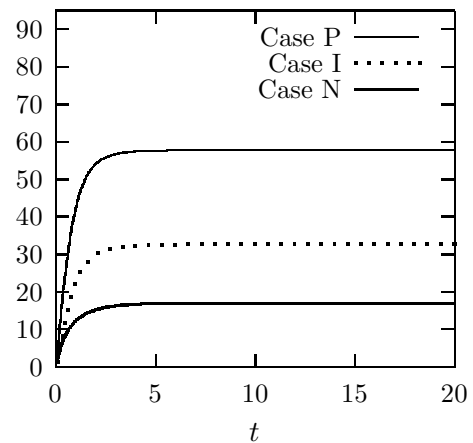


Figure 4.4: Time-dependent variance $\text{Var}[N(t)]$ of the number of customers

4.5.3 Impact of correlation in service time sequence

Finally we consider the $M^X/SM/\infty$ queue to investigate the correlation in the service time sequence. We assume that batches arrive according to a Poisson process with rate 5, and batch sizes are i.i.d. according to the following two-points distribution: The size of the batch is equal to one with probability 1/2 and to 3 with probability 1/2. Let $H(t, n)$ denote the service time distribution of customers in the n th batch. We assume that the sequence $\{H(t, n); n = 1, 2, \dots\}$ of service time distributions constitutes a semi-Markov process whose kernel $\mathbf{H}(t)$ is given by

$$\mathbf{H}(t) = \begin{pmatrix} pH_1(t) & (1-p)H_2(t) \\ (1-p)H_1(t) & pH_2(t) \end{pmatrix}, \quad 0 \leq p < 1,$$

where $H_1(t)$ and $H_2(t)$ ($t \geq 0$) are distribution functions. Note that the sequence has negative correlation for $0 \leq p < 1/2$, it becomes i.i.d. for $p = 1/2$, and it has positive correlation for

$1/2 < p < 1$. We regard customers whose service times follow a distribution function $H_k(t)$ ($k = 1, 2$) as class k customers. Formally this $M^X/SM/\infty$ queue is characterized by

$$\begin{aligned} \mathbf{C} &= \begin{bmatrix} -5 & 0 \\ 0 & -5 \end{bmatrix}, & \mathbf{D}_1 &= \begin{bmatrix} 5p & 0 \\ 5(1-p) & 0 \end{bmatrix}, & \mathbf{D}_2 &= \begin{bmatrix} 0 & 5(1-p) \\ 0 & 5p \end{bmatrix}, \\ g(n) &= \begin{cases} 1/2, & \text{if } n = 1, \\ 1/2, & \text{if } n = 3, \\ 0, & \text{otherwise,} \end{cases} \\ \mathbf{D}(n_1, n_2) &= \begin{cases} g(n_1)\mathbf{D}_1, & \text{if } n_1 \geq 1 \text{ and } n_2 = 0, \\ g(n_2)\mathbf{D}_2, & \text{if } n_1 = 0 \text{ and } n_2 \geq 1, \\ \mathbf{O}, & \text{otherwise,} \end{cases} \end{aligned}$$

where $0 \leq p < 1$. Note that $\lambda_1 = \lambda_2 = 5$ in this formulation. As for service time distributions $H_k(t)$ ($k = 1, 2$), we fix $\lambda_1 h_1 + \lambda_2 h_2 = 10$ (i.e., $E[N] = 10$), and consider the following three cases:

Case 1:

$$H_1(t) = 1 - 0.25e^{-0.5t} - 0.75e^{-1.5t}, \quad H_2(t) = 1 - e^{-t}.$$

Case 2:

$$H_1(t) = 1 - e^{-2t}, \quad H_2(t) = 1 - e^{-\frac{2}{3}t}.$$

Case 3:

$$H_1(t) = 1 - e^{-5t}, \quad H_2(t) = 1 - e^{-\frac{5}{9}t}.$$

In Case 1, the mean service time of each class is equal to one. Thus the autocovariance of the service time sequence is equal to zero for all p . On the other hand, in Case 2 and Case 3, the mean service times are different, so that the autocovariance of the service time sequence is positive for $p > 1/2$ and negative for $p < 1/2$. Note also that the difference between mean service times in Case 2 is equal to one, while that in Case 3 is equal to $8/5$. Thus the absolute value of the autocovariance in Case 3 is larger than in Case 2 for $p \neq 0$.

Figures 4.5, 4.6 and 4.7 plot the marginal variances $\text{Var}[N_1]$ and $\text{Var}[N_2]$, the covariance $\text{Cov}[N_1, N_2]$ and the variance $\text{Var}[N]$ of the total number of customers as functions of parameter p in Cases 1, 2 and 3, respectively. We observe that when the autocovariance of the service time sequence is equal to zero (Case 1), $\text{Var}[N]$ is almost insensitive to parameter p . On the other hand, as shown in Figures 4.6 and 4.7, $\text{Var}[N]$ becomes larger with the correlation in the service time sequence when the autocovariance is not equal to zero. Further, comparing Figure 4.7 with Figure 4.6, we observe that a stronger correlation leads to a larger value of $\text{Var}[N]$.

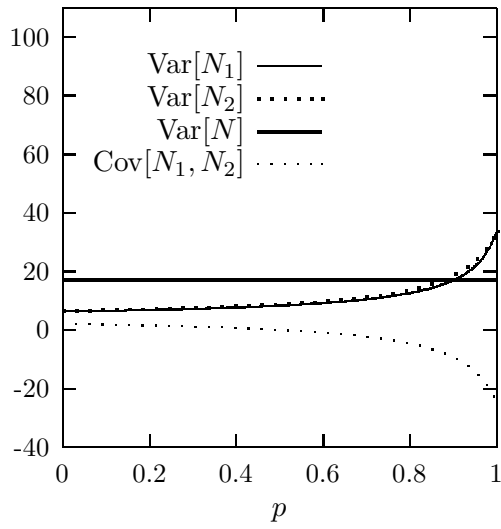


Figure 4.5: Variance and covariance in Case 1.

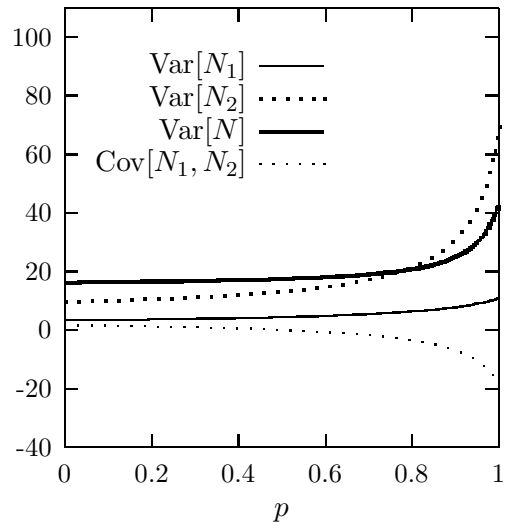


Figure 4.6: Variance and covariance in Case 2.

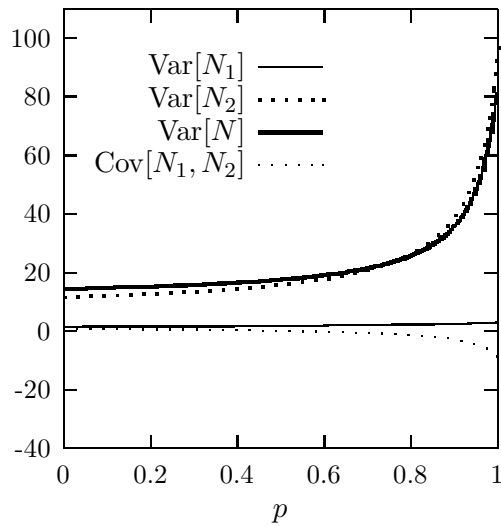


Figure 4.7: Variance and covariance in Case 3.

Chapter 5

Processor-Sharing Queue

5.1 Introduction

Processor-sharing queues are important models in computer engineering. Under the processor-sharing service discipline, the server is shared equally by all customers in the system. To put it more concretely, when n customers are present in the system, each customer receives service at rate $1/n$.

Many researchers have studied processor-sharing queues for a few decades. For the stationary M/M/1 processor-sharing queue (M/M/1-PS), Coffman, Muntz and Trotter [Coff70] derived the LST of the sojourn time distribution. By inverting this LST, Morrison [Morr85] obtained an integral representation for the complementary distribution of the sojourn time, which was also derived by Guillemin and Boyer with the spectral theory [Guil01]. Further Sengupta and Jagerman [Seng85] studied moments of the sojourn time conditioned on the number of customers found by arriving customers.

For the stationary GI/M/1-PS queue, the first two moments and the LST of the sojourn time were derived in Ramaswami [Rama84] and Jagerman and Sengupta [Jage91], respectively. Also, for the stationary M/G/1-PS queue, Ott [Ott84] and Yashkov [Yash83] derived the LST and first two moments of the sojourn time conditioned on the service time of an arriving customer. Besides, Grishechkin [Gris94] and Sengupta [Seng92] studied GI/G/1-PS queues.

To the best of our knowledge, all of previous studies on processor-sharing queues assume that inter-arrival times are independent and identically distributed (i.i.d.). Further there are no works focusing on the computation of the sojourn time distribution. Motivated by these facts, we consider the sojourn time distribution in a stationary processor-sharing queue with BMAP arrivals (see subsection 1.1.1) and exponential service times, i.e., a BMAP/M/1-PS queue.

The main contribution of this chapter is to provide a recursive formula for the stationary sojourn time distribution in a BMAP/M/1-PS queue. The derivation is very simple and the result is easy to be interpreted probabilistically. Moreover, in the single arrival case (i.e., MAP/M/1-PS queue), we discuss numerical procedure and accuracy guarantee.

The rest of this chapter is divided into four sections. In section 2, the mathematical model and some known results are described. In section 3, we derive a recursive formula for the complementary

distribution of the sojourn time in a BMAP/M/1-PS queue. In section 4, for a MAP/M/1-PS queue, we discuss the numerical procedure to control the absolute error in numerical results and show some numerical examples. Finally, in section 5, we provide concluding remarks on how to utilize the results in this chapter to obtain the waiting time distribution in a random order service MAP/M/1 queue.

5.2 Model and Known Results

This chapter considers a stationary BMAP/M/1 processor-sharing (BMAP/M/1-PS) queue. Service times are i.i.d. according to an exponential distribution with finite mean μ_i^{-1} . Customer arrivals follow a BMAP $(\mathbf{C}, \mathbf{D}(n))$ (see subsection 1.1.1). Broadly speaking, customer arrivals occur as follows. There is an irreducible Markov chain with finite state space $\mathcal{M} = \{1, \dots, M\}$, which governs the BMAP. When a state transition driven by $M \times M$ matrix $\mathbf{D}(n)$ happens, n customers arrive simultaneously. On the other hand, when a state transition driven by $M \times M$ matrix \mathbf{C} happens, no customers arrive.

Note that the infinitesimal generator of the underlying Markov chain is given by $\mathbf{C} + \mathbf{D}$, where $\mathbf{D} = \sum_{n=1}^{\infty} \mathbf{D}(n)$. Note also that the mean arrival rate λ is given by

$$\lambda = \boldsymbol{\pi} \sum_{n=1}^{\infty} n \mathbf{D}(n) \mathbf{e},$$

where $\boldsymbol{\pi}$ denotes the invariant probability vector for $(\mathbf{C} + \mathbf{D})$. To avoid trivialities, we assume $\mathbf{D} \neq \mathbf{O}$. Thus $\lambda > 0$. Let ρ denote the utilization factor, i.e., $\rho = \lambda \mu^{-1}$. In the remainder of this chapter, we assume that $\rho < 1$ and the system is in steady state.

We make a comment on the stationary queue length distribution in the BMAP/M/1 queue. Let N and S denote the number of customers in the system and the state of the underlying Markov chain, respectively, in steady state. We then define $\boldsymbol{\pi}_n$ ($n = 0, 1, \dots$) as a $1 \times M$ vector whose j th ($j \in \mathcal{M}$) element represents $\Pr[N = n, S = j]$. Owing to i.i.d. exponential services, the stationary queue length distribution in the BMAP/M/1-PS queue is identical to that in the BMAP/M/1-FCFS. Thus, $\boldsymbol{\pi}_n$'s ($n = 0, 1, \dots$) can be computed by the M/G/1 paradigm (see Appendix B.1).

5.3 Sojourn Time Distribution

In this section, we consider the sojourn time distribution in steady state. We denote, by customer C_n ($n \geq 0$), a randomly chosen customer who is a member of a batch that increases the queue length to $n + 1$ on arrival. Let W_{C_n} denote a random variable representing the sojourn time of customer C_n . Let S_{C_n} denote a random variable representing the state of the underlying Markov chain immediately after the arrival of customer C_n . We then define $\bar{\mathbf{w}}_n(x)$ ($x \geq 0$, $n = 0, 1, \dots$) as an $M \times 1$ vector whose j th ($j \in \mathcal{M}$) element represents $\Pr[W_{C_n} > x \mid S_{C_n} = j]$.

We now consider the sojourn time of a customer C_n arriving at $t = 0$. Conditioning the event at $t = \Delta x$, we readily obtain

$$\bar{\mathbf{w}}_n(x + \Delta x) = \frac{n\mu}{n+1} \Delta x \bar{\mathbf{w}}_{n-1}(x) + \{\mathbf{I} + (\mathbf{C} - \mu\mathbf{I})\Delta x\} \bar{\mathbf{w}}_n(x)$$

$$+ \sum_{m=1}^{\infty} \mathbf{D}(m) \Delta x \bar{\mathbf{w}}_{n+m}(x) + o(\Delta x), \quad n = 0, 1, \dots, \quad (5.1)$$

where $\bar{\mathbf{w}}_{-1}(x) = \mathbf{0}$ for all $x \geq 0$. It then follows from (5.1) that

$$\frac{d}{dx} \bar{\mathbf{w}}_n(x) = \frac{n\mu}{n+1} \bar{\mathbf{w}}_{n-1}(x) + (\mathbf{C} - \mu \mathbf{I}) \bar{\mathbf{w}}_n(x) + \sum_{m=1}^{\infty} \mathbf{D}(m) \bar{\mathbf{w}}_{n+m}(x), \quad n = 0, 1, \dots, \quad (5.2)$$

where $\bar{\mathbf{w}}_{-1}(x) = \mathbf{0}$ for all $x \geq 0$. Note that the differential-difference equation (5.2) is an extension of equation (9.10) on page 111 in [Asmu03].

We now define $\bar{\mathbf{w}}(x)$ and \mathbf{T} as

$$\bar{\mathbf{w}}(x) = \begin{bmatrix} \bar{\mathbf{w}}_0(x) \\ \bar{\mathbf{w}}_1(x) \\ \bar{\mathbf{w}}_2(x) \\ \vdots \end{bmatrix}, \quad (5.3)$$

$$\mathbf{T} = \begin{bmatrix} \mathbf{C} - \mu \mathbf{I} & \mathbf{D}(1) & \mathbf{D}(2) & \mathbf{D}(3) & \mathbf{D}(4) & \cdots \\ \frac{\mu \mathbf{I}}{2} & \mathbf{C} - \mu \mathbf{I} & \mathbf{D}(1) & \mathbf{D}(2) & \mathbf{D}(3) & \cdots \\ \mathbf{O} & \frac{2\mu \mathbf{I}}{3} & \mathbf{C} - \mu \mathbf{I} & \mathbf{D}(1) & \mathbf{D}(2) & \cdots \\ \mathbf{O} & \mathbf{O} & \frac{3\mu \mathbf{I}}{4} & \mathbf{C} - \mu \mathbf{I} & \mathbf{D}(1) & \cdots \\ \mathbf{O} & \mathbf{O} & \mathbf{O} & \frac{4\mu \mathbf{I}}{5} & \mathbf{C} - \mu \mathbf{I} & \cdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix},$$

respectively. (5.2) is then rewritten to be

$$\frac{d}{dx} \bar{\mathbf{w}}(x) = \mathbf{T} \bar{\mathbf{w}}(x). \quad (5.4)$$

Note here that \mathbf{T} can be considered as the defective infinitesimal generator of a continuous-time Markov chain with infinite states (n, j) ($n = 0, 1, \dots, j \in \mathcal{M}$). Therefore the conditional sojourn time of customer C_n given $S_{C_n} = j$ is equivalent to the first passage time from state (n, j) to an implicit absorbing state.

Note that $\bar{\mathbf{w}}_n(0) = \mathbf{e}$ for all $n = 0, 1, \dots$. Then, the formal solution of (5.4) is given by

$$\bar{\mathbf{w}}(x) = \exp(\mathbf{T}x) \mathbf{e}.$$

Thus, applying the uniformization technique (see Appendix A), we obtain

$$\bar{\mathbf{w}}(x) = \sum_{k=0}^{\infty} e^{-(\theta+\mu)x} \frac{(\theta+\mu)^k x^k}{k!} \left[\mathbf{I} + \frac{1}{\theta+\mu} \mathbf{T} \right]^k \mathbf{e}, \quad (5.5)$$

where θ denotes the maximum absolute value of diagonal elements of \mathbf{C} . Let $\mathbf{h}_{n,k}$ ($n = 0, 1, \dots, k = 0, 1, \dots$) denote an $M \times 1$ vector that satisfies

$$\begin{bmatrix} \mathbf{h}_{0,k} \\ \mathbf{h}_{1,k} \\ \vdots \end{bmatrix} = \left[\mathbf{I} + \frac{1}{\theta+\mu} \mathbf{T} \right]^k \mathbf{e}. \quad (5.6)$$

Note here that the j th ($j \in \mathcal{M}$) element of $\mathbf{h}_{n,k}$ represents the probability that the number of transitions in transient states of the uniformized Markov chain with transition matrix $\mathbf{I} + (\theta + \mu)^{-1}\mathbf{T}$ is not less than k during the first passage time from state (n, j) to the implicit absorbing state. Therefore $\mathbf{h}_{n,k}$ ($n = 0, 1, \dots, k = 0, 1, \dots$) satisfies

$$\mathbf{0} < \mathbf{h}_{n,k+1} < \mathbf{h}_{n,k} \leq \mathbf{e}, \quad (5.7)$$

where $\mathbf{h}_{n,0} = \mathbf{e}$ for $n = 0, 1, \dots$

We define $\bar{W}(x) = \Pr[W > x]$ for $x \geq 0$, where W denotes the sojourn time of a randomly chosen customer.

Theorem 5.1 $\bar{w}_n(x)$ ($x \geq 0, n = 0, 1, \dots$) and $\bar{W}(x)$ ($x \geq 0$) are given by

$$\bar{w}_n(x) = \sum_{k=0}^{\infty} e^{-(\theta+\mu)x} \frac{(\theta+\mu)^k x^k}{k!} \mathbf{h}_{n,k}, \quad (5.8)$$

$$\bar{W}(x) = \sum_{n=0}^{\infty} \sum_{m=0}^n \frac{(m+1)\boldsymbol{\pi}_{n-m}\mathbf{D}(m+1)}{\lambda} \sum_{k=0}^{\infty} e^{-(\theta+\mu)x} \frac{(\theta+\mu)^k x^k}{k!} \mathbf{h}_{n,k}, \quad (5.9)$$

respectively, where the $\mathbf{h}_{n,k}$ is determined by the following recursion.

$$\mathbf{h}_{n,0} = \mathbf{e}, \quad n = 0, 1, \dots,$$

and for $k = 0, 1, \dots$,

$$\mathbf{h}_{n,k+1} = \frac{1}{\theta + \mu} \left[\frac{n\mu}{n+1} \mathbf{h}_{n-1,k} + (\theta\mathbf{I} + \mathbf{C})\mathbf{h}_{n,k} + \sum_{m=1}^{\infty} \mathbf{D}(m)\mathbf{h}_{n+m,k} \right], \quad n = 0, 1, \dots, \quad (5.10)$$

where $\mathbf{h}_{-1,k} = \mathbf{0}$.

Proof. It is clear from (5.3), (5.5) and (5.6) that $\bar{w}_n(x)$ satisfies (5.8). Further it follows from (5.6) that

$$\begin{bmatrix} \mathbf{h}_{0,k+1} \\ \mathbf{h}_{1,k+1} \\ \vdots \end{bmatrix} = \left[\mathbf{I} + \frac{1}{\theta + \mu} \mathbf{T} \right] \begin{bmatrix} \mathbf{h}_{0,k} \\ \mathbf{h}_{1,k} \\ \vdots \end{bmatrix}, \quad k = 0, 1, \dots,$$

which shows that the $\mathbf{h}_{n,k}$ satisfies (5.10).

We can show (5.9) as follows. Note that $\lambda^{-1} \sum_{m=0}^n (m+1)\boldsymbol{\pi}_{n-m}\mathbf{D}(m+1)$ ($n = 0, 1, \dots$) is a $1 \times M$ vector whose the j th ($j \in \mathcal{M}$) element represents the joint probability that a customer C_n arrives at the stationary system and the underlying Markov chain is in state j immediately after his arrival [Neut89]. Thus, applying the law of total probabilities to (5.8), we obtain (5.9). ■

In what follows, we consider two special cases, MAP/M/1-PS and M/M/1-PS queues. In the former case, $\mathbf{D}(1) = \mathbf{D}$ and $\mathbf{D}(n) = \mathbf{O}$ for all $n = 2, 3, \dots$. Further, $\boldsymbol{\pi}_n$'s are given by (see Corollary B.1)

$$\begin{aligned} \boldsymbol{\pi}_0 &= \boldsymbol{\pi}(\mathbf{I} - \mathbf{R}), \\ \boldsymbol{\pi}_n &= \boldsymbol{\pi}_0 \mathbf{R}^n, \quad n = 1, 2, \dots, \end{aligned}$$

where \mathbf{R} is the minimal nonnegative solution of

$$\mathbf{D} + \mathbf{R}(\mathbf{C} - \mu\mathbf{I}) + \mu\mathbf{R}^2 = \mathbf{O}.$$

Thus, the following corollary is readily obtained from Theorem 5.1.

Corollary 5.1 *For the MAP/M/1-PS queue, $\bar{w}_n(x)$ ($x \geq 0$, $n = 0, 1, \dots$) is given by (5.8) and $\bar{W}(x)$ ($x \geq 0$) is given by*

$$\bar{W}(x) = \frac{1}{\lambda} \sum_{n=0}^{\infty} \pi_0 \mathbf{R}^n \mathbf{D} \sum_{k=0}^{\infty} e^{-(\theta+\mu)x} \frac{(\theta+\mu)^k x^k}{k!} \mathbf{h}_{n,k}, \quad (5.11)$$

where the $\mathbf{h}_{n,k}$ is determined by the following recursion.

$$\mathbf{h}_{n,0} = \mathbf{e}, \quad n = 0, 1, \dots,$$

and for $k = 0, 1, \dots$,

$$\mathbf{h}_{n,k+1} = \frac{1}{\theta + \mu} \left[\frac{n\mu}{n+1} \mathbf{h}_{n-1,k} + (\theta\mathbf{I} + \mathbf{C})\mathbf{h}_{n,k} + \mathbf{D}\mathbf{h}_{n+1,k} \right], \quad n = 0, 1, \dots, \quad (5.12)$$

where $\mathbf{h}_{-1,k} = \mathbf{0}$.

Next, we discuss an M/M/1-PS queue with arrival rate λ and service rate μ . Note that, in this queue, π_n and $\mathbf{h}_{n,k}$ defined in (5.6) are scalars. To make things clear, we here denote them by π_n and $h_{n,k}$, respectively. Further, we define $\bar{w}_n(x)$ ($x \geq 0$, $n = 0, 1, \dots$) as $\Pr[W_{C_n} > x]$. The following result is the direct conclusion of Corollary 5.1.

Corollary 5.2 *For the M/M/1-PS queue, $\bar{w}_n(x)$ ($x \geq 0$, $n = 0, 1, \dots$) and $\bar{W}(x)$ ($x \geq 0$) are given by*

$$\begin{aligned} \bar{w}_n(x) &= \sum_{k=0}^{\infty} e^{-(\lambda+\mu)x} \frac{(\lambda+\mu)^k x^k}{k!} h_{n,k}, \\ \bar{W}(x) &= \sum_{n=0}^{\infty} (1-\rho)\rho^n \sum_{k=0}^{\infty} e^{-(\lambda+\mu)x} \frac{(\lambda+\mu)^k x^k}{k!} h_{n,k}, \end{aligned}$$

respectively, where the $h_{n,k}$ is determined by the following recursion.

$$h_{n,0} = 1, \quad n = 0, 1, \dots,$$

and for $k = 0, 1, \dots$,

$$h_{n,k+1} = \frac{n}{n+1} \frac{\mu}{\lambda+\mu} h_{n-1,k} + \frac{\lambda}{\lambda+\mu} h_{n+1,k}, \quad n = 0, 1, \dots,$$

where $h_{-1,k} = 0$.

Remark 5.1 *For the M/M/1-PS queue, Asmussen obtained a recursive formula similar to Corollary 5.2 (see Theorem 9.5 on page 112 in [Asmu03]). However, it is numerically unstable because subtraction of positive numbers of similar magnitude is involved. On the other hand, our recursion (5.12) is constructed only by sums and products of nonnegative vectors and matrices, and therefore it is numerically stable.*

5.4 Numerical Examples

This section considers MAP/M/1-PS queues. We first discuss the numerical procedure for $\bar{W}(x)$ and the absolute error in numerical results. Next, we discuss the impact of the arrival process on the sojourn time distribution through some numerical examples.

5.4.1 Numerical procedure and accuracy guarantee

For a fixed x ($x > 0$), we compute an approximation to $\bar{W}(x)$ as follows. We first set some ε ($0 < \varepsilon \ll 1$) and obtain a minimum nonnegative integer $N(\varepsilon)$ satisfying

$$\frac{1}{\lambda} \sum_{n=0}^{N(\varepsilon)} \pi_0 \mathbf{R}^n \mathbf{D} \mathbf{e} > 1 - \varepsilon. \quad (5.13)$$

Note that $N(\varepsilon)$ is uniquely determined because

$$\sum_{n=0}^{\infty} \pi_0 \mathbf{R}^n \mathbf{D} \mathbf{e} = \lambda. \quad (5.14)$$

Next, we set some ε' ($0 < \varepsilon' \ll 1$) and obtain two nonnegative integers $L(\varepsilon', x)$ and $R(\varepsilon', x)$ satisfying $L(\varepsilon', x) \leq R(\varepsilon', x)$ and

$$\sum_{k=L(\varepsilon', x)}^{R(\varepsilon', x)} \frac{(\theta + \mu)^k x^k}{k!} e^{-(\theta + \mu)x} \geq 1 - \varepsilon'. \quad (5.15)$$

One of algorithms to find $L(\varepsilon', x)$ and $R(\varepsilon', x)$ is proposed by Fox and Glynn [Fox88]. We use their algorithm for numerical examples in the following subsections.

Finally, we compute an approximation $\bar{W}_{\text{comp}}(x)$ to $\bar{W}(x)$:

$$\bar{W}_{\text{comp}}(x) = \frac{1}{\lambda} \sum_{n=0}^{N(\varepsilon)} \pi_0 \mathbf{R}^n \mathbf{D} \sum_{k=L(\varepsilon', x)}^{R(\varepsilon', x)} \frac{(\theta + \mu)^k x^k}{k!} e^{-(\theta + \mu)x} \mathbf{h}_{n,k}. \quad (5.16)$$

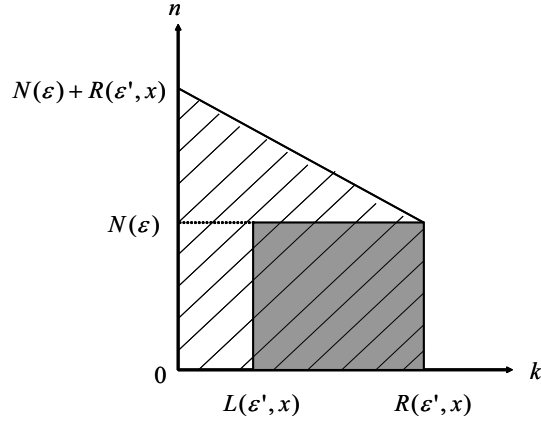
Note here that we have to compute $\mathbf{h}_{n,k}$'s in the trapezoidal domain in Figure 5.1, although $\bar{W}_{\text{comp}}(x)$ requires $\mathbf{h}_{n,k}$'s only in the shaded rectangular domain in Figure 5.1 (see (5.12)).

The bound for the difference between $\bar{W}(x)$ and $\bar{W}_{\text{comp}}(x)$ can be obtained as follows. For simplicity in description, $L(\varepsilon', x)$ and $R(\varepsilon', x)$ are denoted by L and R , respectively. From (5.11) and (5.16), we have

$$\begin{aligned} \bar{W}(x) - \bar{W}_{\text{comp}}(x) &= \frac{1}{\lambda} \sum_{n=N(\varepsilon)+1}^{\infty} \pi_0 \mathbf{R}^n \mathbf{D} \sum_{k=0}^{\infty} \frac{(\theta + \mu)^k x^k}{k!} e^{-(\theta + \mu)x} \mathbf{h}_{n,k} \\ &\quad + \frac{1}{\lambda} \sum_{n=0}^{N(\varepsilon)} \pi_0 \mathbf{R}^n \mathbf{D} \sum_{k < L, k > R} \frac{(\theta + \mu)^k x^k}{k!} e^{-(\theta + \mu)x} \mathbf{h}_{n,k}. \end{aligned}$$

Applying (5.7), (5.13), (5.14) and (5.15) to the above equation, we obtain

$$\bar{W}(x) - \bar{W}_{\text{comp}}(x) < \frac{1}{\lambda} \sum_{n=N(\varepsilon)+1}^{\infty} \pi_0 \mathbf{R}^n \mathbf{D} \mathbf{e} \sum_{k=0}^{\infty} \frac{(\theta + \mu)^k x^k}{k!} e^{-(\theta + \mu)x}$$

Figure 5.1: The computational domain of the $\mathbf{h}_{n,k}$.

$$\begin{aligned}
 & + \frac{1}{\lambda} \sum_{n=0}^{N(\varepsilon)} \pi_0 \mathbf{R}^n \mathbf{D} e \sum_{k < L, k > R} \frac{(\theta + \mu)^k x^k}{k!} e^{-(\theta + \mu)x} \\
 & \leq \varepsilon + \varepsilon'.
 \end{aligned}$$

For numerical examples in sections 5.4.2 and 5.4.3, we set $\varepsilon = 10^{-11}$ and $\varepsilon' = 10^{-10}$.

5.4.2 Impact of variation in inter-arrival times

In this subsection, we discuss the impact of variation of inter-arrival times on the sojourn time. For this purpose, we assume that inter-arrival times are i.i.d. and consider the following three types of inter-arrival time distributions with the same mean λ^{-1} : four-stage Erlang distribution, exponential distribution and two-stage balanced hyper-exponential distribution $\psi(x) = 1 - pe^{-2p\lambda x} - (1 - p)e^{-2(1-p)\lambda x}$ with $p = 0.5 + \sqrt{15}/10$. Let C_v denote the coefficient of variation of the inter-arrival time distribution. Note that $C_v = 0.5$ for the four-stage Erlang distribution, $C_v = 1$ for the exponential distribution and $C_v = 2$ for the above two-stage balanced hyper-exponential distribution.

We set $\lambda = 0.8$ and $\mu = 1$, so that $\rho = 0.8$ for all three cases. Figure 5.2 plots the complementary distributions for these cases. We observe that the sojourn time distribution becomes larger stochastically with the increase of C_v .

5.4.3 Impact of correlation in inter-arrival times

Next we discuss the impact of correlation in inter-arrival times on the sojourn time. For this purpose, we set $\mu = 1$ and assume that the arrival process follows a stationary two-state MAP:

$$\mathbf{C} = \begin{bmatrix} -2 & 0 \\ 0 & -0.5 \end{bmatrix}, \quad \mathbf{D} = \begin{bmatrix} 2p & 2(1-p) \\ 0.5(1-p) & 0.5p \end{bmatrix},$$

where $0 \leq p < 1$. Note here that $\lambda = 0.8$ and hence $\rho = 0.8$.

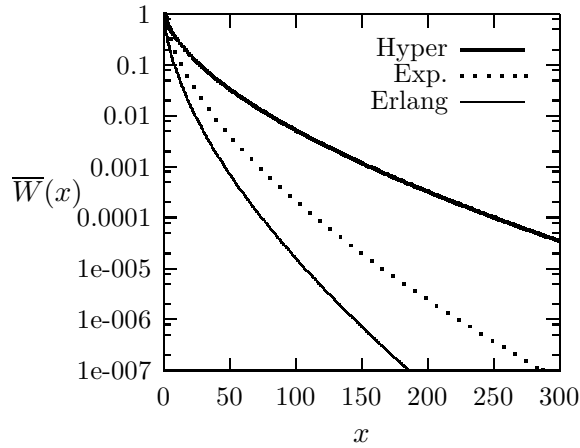


Figure 5.2: The complementary distribution $\overline{W}(x)$ of the sojourn time.

Let Y_n ($n = 1, 2, \dots$) denote a random variable representing the n th inter-arrival time. Note that the marginal distribution of Y_n is independent of p and follows a hyper-exponential distribution $\psi(x) = 1 - 0.5e^{-2x} - 0.5e^{-0.5x}$. Note also that the covariance of Y_n and Y_{n+1} is given by $9(2p-1)/16$. Therefore inter-arrival times are positively (resp. negatively) correlated for $0.5 < p < 1$ (resp. $0 \leq p < 0.5$), and are i.i.d. for $p = 0.5$.

Figure 5.3 plots the complementary distributions for the three cases, $p = 0.1, 0.5$ and 0.9 . We observe that the positive (resp. negative) correlation of inter-arrival times makes the sojourn time stochastically larger (resp. smaller) than that for i.i.d. inter-arrival times.

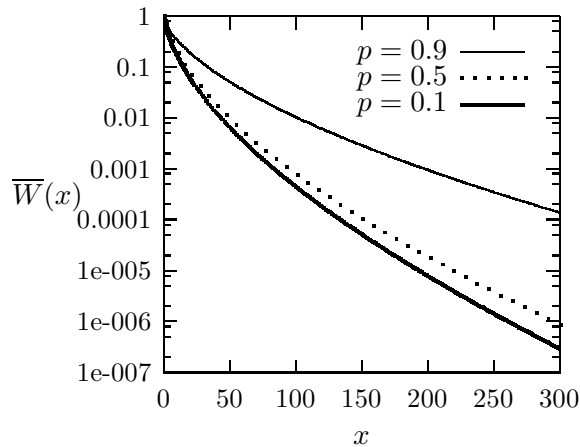


Figure 5.3: The complementary distribution $\overline{W}(x)$ of the sojourn time.

5.5 Concluding Remarks

In this chapter, we derived a numerically feasible formula to compute the sojourn time distribution in a MAP/M/1-PS queue. It is known that in the GI/M/1 queue, the sojourn time distribution for PS is identical to the conditional waiting time distribution for random order of service (ROS) given the waiting time is positive [Bors03, Cohe84]. Note that in the stationary MAP/M/1 queue, the sojourn time distribution for PS is closely related to the conditional waiting time distribution for ROS. To see this, we define $W_{C_n}^{\text{ROS}}$ as a generic random variable representing the waiting time of customers finding n customers in the ROS system on arrival. Let $\bar{\mathbf{w}}_n^{\text{ROS}}(x)$ denote an $M \times 1$ vector whose j th ($j \in \mathcal{M}$) element represents $\Pr[W_{C_n}^{\text{ROS}} > x \mid S_{C_n} = j]$, where S_{C_n} denotes the state of the underlying Markov chain immediately after arrivals finding n customers in the system. It then follows from a very similar reasoning in section 2 of [Bors03] that for $n = 0, 1, \dots$,

$$\bar{\mathbf{w}}_{n+1}^{\text{ROS}}(x) = \bar{\mathbf{w}}_n(x), \quad x \geq 0,$$

where $\bar{\mathbf{w}}_n(x)$ is given in Corollary 5.1. Thus the complementary distribution $\bar{W}^{\text{ROS}}(x)$ of waiting time in a MAP/M/1-ROS queue is given by

$$\bar{W}^{\text{ROS}}(x) = \frac{1}{\lambda} \sum_{n=0}^{\infty} \boldsymbol{\pi}_0 \mathbf{R}^{n+1} \mathbf{D} \bar{\mathbf{w}}_n(x), \quad x \geq 0,$$

which means that our results can be also used to obtain the waiting time distribution in a MAP/M/1-ROS queue.

Chapter 6

Conclusion

6.1 Summary of Results

This thesis studied queues with batch Markovian arrival streams, focusing on the computation of performance measures of interest. The summary of the results is described below.

- (a) We analyzed FIFO single-server queues with/without service interruptions by the approach based on [Taki01a, Taki01b] and established a numerical algorithm to compute the joint queue length distribution. Further, we examined the impact of system parameters on the queue length distribution through some numerical examples.
- (b) Extremely general infinite-server queues were analyzed. By simple probabilistic consideration, we obtained the system of linear differential equations for the time-dependent matrix joint generating function of the queue length. Further, assuming phase-type service times, we obtained explicit and numerically feasible expressions of the time-dependent and limiting joint binomial moments.
- (c) The sojourn time distribution in a BMAP/M/1-PS queue was studied. Based on the uniformization technique (see Appendix A), the recursive formula for the sojourn time distribution was obtained. In addition, for a MAP/M/1-PS queue, we established a numerically feasible algorithm to compute the sojourn time distribution. The advantage of the algorithm is that the accuracy of computational results can be estimated in advance.

6.2 Future Work

This thesis studies some types of queues with batch Markovian arrival streams. However, the study of those queues has just started. The author describes some future works below.

- (a) As for the single-server queue considered chapter 2, we can obtain the mean sojourn time of customers in each class by using Theorem 2.1 and the results on moments of the virtual waiting time distribution [Taki94a]. Further with Little's law, we can readily derive the mean

queue length of each class. We have, however, no way to compute higher-order moments such as a variance. This is a problem to be solved.

- (b) The computational procedures proposed in chapters 2 and 3 are two-step algorithms to obtain the stationary joint queue length distribution. The first step is to obtain the joint queue length distribution immediately after departures, and the second step is to derive the joint queue length distribution at a random point in time, using the relationship of these two joint queue length distributions (see (2.8) and (2.9) in Corollary 2.1). The first step is numerically feasible, but the second step is not necessarily so because we have to subtract positive numbers of similar magnitude (see (2.9)). Therefore, to establish a numerically feasible algorithm directly to compute the stationary joint queue length distribution is a challenging problem.
- (c) A multi-server queue with batch Markovian arrival streams is an interesting model. Such a queue is so flexible that it can cover most of the queues we can now think of. Thus, if we obtained analytic results suitable for computing performance measures of interest, we could establish unified numerical algorithms for almost all of the queues we often encounter.

Appendix A

Uniformization

We consider a continuous-time Markov chain $\{X(t); t \geq 0\}$ with state space \mathcal{L} and infinitesimal generator $\tilde{\mathbf{P}}$. Assuming that diagonal elements of $\tilde{\mathbf{P}}$ are bounded from below, we choose a positive number θ such that

$$\sup_{j \in \mathcal{L}} |[\tilde{\mathbf{P}}]_{j,j}| \leq \theta < \infty,$$

where $[\tilde{\mathbf{P}}]_{i,j}$ denotes the (i, j) th element of $\tilde{\mathbf{P}}$. Let $\tilde{\mathbf{\Pi}}(t)$ denote a matrix whose (i, j) th element represents $\Pr[X(t) = j \mid X(0) = i]$. Then, $\tilde{\mathbf{\Pi}}(t)$ is given by

$$\tilde{\mathbf{\Pi}}(t) = \exp(\tilde{\mathbf{P}}t) = \sum_{k=0}^{\infty} e^{-\theta t} \frac{(\theta t)^k}{k!} [\mathbf{I} + \theta^{-1} \tilde{\mathbf{P}}]^k = \sum_{k=0}^{\infty} e^{-\theta t} \frac{(\theta t)^k}{k!} \mathbf{P}^k, \quad (\text{A.1})$$

where

$$\mathbf{P} = \mathbf{I} + \theta^{-1} \tilde{\mathbf{P}}.$$

Note that \mathbf{P} is a nonnegative matrix because non-diagonal elements of $\tilde{\mathbf{P}}$ are all nonnegative. Note also that $\mathbf{P}\mathbf{e} = \mathbf{e}$ owing to $\tilde{\mathbf{P}}\mathbf{e} = \mathbf{0}$. These two facts mean that \mathbf{P} is a stochastic matrix. Thus, if the state space \mathcal{L} is finite, we can compute the transient distribution of the continuous-time Markov chain $\{X(t); t \geq 0\}$, using (A.1).

We now introduce a discrete-time Markov chain $\{X_k; k = 0, 1, \dots\}$ with the same state space \mathcal{L} as $\{X(t); t \geq 0\}$, whose transition probability matrix is given by \mathbf{P} . Noting that

$$\exp(\tilde{\mathbf{P}}t) = \sum_{k=0}^{\infty} e^{-\theta t} \frac{(\theta t)^k}{k!} \mathbf{P}^k,$$

we see that if a continuous-time Markov chain $\{X(t); t \geq 0\}$ is irreducible, a discrete-time Markov chain $\{X_k; k = 0, 1, \dots\}$ is also irreducible, and vice versa. Assuming that $\{X(t); t \geq 0\}$ is positive recurrent, let $\tilde{\boldsymbol{\pi}}$ denote the stationary distribution of $\{X(t); t \geq 0\}$, i.e., $\tilde{\boldsymbol{\pi}}\tilde{\mathbf{P}} = \mathbf{0}$. We then have

$$\tilde{\boldsymbol{\pi}}\mathbf{P} = \tilde{\boldsymbol{\pi}}(\mathbf{I} + \theta^{-1}\tilde{\mathbf{P}}) = \tilde{\boldsymbol{\pi}}.$$

Thus, $\{X_k; k = 0, 1, \dots\}$ is also positive recurrent and has the same stationary distribution $\tilde{\boldsymbol{\pi}}$.

We call, *uniformization*, the above technique that converts a continuous-time Markov chain to a discrete-time Markov chain with the same stationary distribution [Tijm94]. In general, discrete-time Markov chains are more manageable than continuous-time Markov chains, because transition probability matrices, which characterize discrete-time Markov chains, are nonnegative. Therefore, the uniformization technique is very often used in algorithmic analysis of queues.

Appendix B

Queue Length Distribution in a BMAP/GI/1 Queue

We discuss a stationary FIFO single-server queue with a single batch Markovian arrival stream, i.e., a BMAP/GI/1 queue. Customer arrivals follow a BMAP characterized by \mathbf{C} and $\mathbf{D}(n)$ (see subsection 1.1.1), and service times are i.i.d. according to a distribution function $H(x)$. Let $\mathbf{q}_{\text{BMAP/GI/1}}(n)$ ($n = 0, 1, \dots$) denote a vector whose j th element represents the joint probability that the queue length is equal to n and the underlying Markov chain is in state j immediately after departures. Then the $\mathbf{q}_{\text{BMAP/GI/1}}(n)$ is identical to the steady-state solution for the M/G/1-type Markov chain whose transition probability matrix is given by

$$\begin{bmatrix} \mathbf{B}_0 & \mathbf{B}_1 & \mathbf{B}_2 & \mathbf{B}_3 & \cdots \\ \mathbf{A}_0 & \mathbf{A}_1 & \mathbf{A}_2 & \mathbf{A}_3 & \cdots \\ \mathbf{O} & \mathbf{A}_0 & \mathbf{A}_1 & \mathbf{A}_2 & \cdots \\ \mathbf{O} & \mathbf{O} & \mathbf{A}_0 & \mathbf{A}_1 & \cdots \\ \mathbf{O} & \mathbf{O} & \mathbf{O} & \mathbf{A}_0 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}, \quad (\text{B.1})$$

where \mathbf{A}_n and \mathbf{B}_n satisfy

$$\begin{aligned} \sum_{n=0}^{\infty} z^n \mathbf{A}_n &= \int_0^{\infty} \exp \left[\left(\mathbf{C} + \sum_{m=0}^{\infty} z^m \mathbf{D}(m) \right) x \right] dH(x), \\ \mathbf{B}_n &= (-\mathbf{C})^{-1} \mathbf{A}_n. \end{aligned} \quad (\text{B.2})$$

Thus, $\mathbf{q}_{\text{BMAP/GI/1}}(n)$'s can be computed by the M/G/1 paradigm (see subsection 1.2.2).

Remark B.1 *The recursion of \mathbf{A}_n 's in the BMAP/GI/1 queue is given in Lemma 2.2 with $K = 1$, and a remark on its implementation is found in section 2.7.*

We now define $\mathbf{p}_{\text{BMAP/GI/1}}(n)$ ($n = 0, 1, \dots$) as a vector whose j th element represents the stationary joint probability that the queue length is equal to n and the underlying Markov chain is in state j at a random point in time. It is known that $\mathbf{p}_{\text{BMAP/GI/1}}(n)$ and $\mathbf{q}_{\text{BMAP/GI/1}}(n)$ have

the following relationship:

$$\mathbf{p}_{\text{BMAP/GI/1}}(0) = \lambda \mathbf{q}_{\text{BMAP/GI/1}}(0) (-\mathbf{C})^{-1}, \quad (\text{B.3})$$

$$\begin{aligned} \mathbf{p}_{\text{BMAP/GI/1}}(n) = & \left[\lambda \left\{ \mathbf{q}_{\text{BMAP/GI/1}}(n) - \mathbf{q}_{\text{BMAP/GI/1}}(n-1) \right\} \right. \\ & \left. + \sum_{m=1}^n \mathbf{p}_{\text{BMAP/GI/1}}(n-m) \mathbf{D}(m) \right] (-\mathbf{C})^{-1}, \end{aligned} \quad (\text{B.4})$$

where λ is given by (1.2). Thus, $\mathbf{p}_{\text{BMAP/GI/1}}(n)$'s can be obtained in terms of $\mathbf{q}_{\text{BMAP/GI/1}}(n)$'s. Note here that the above recursion is not always stable because we have to execute the subtraction $\mathbf{q}_{\text{BMAP/GI/1}}(n) - \mathbf{q}_{\text{BMAP/GI/1}}(n-1)$.

Recently, Takine established a stable and direct recursion for the $\mathbf{p}_{\text{BMAP/GI/1}}(n)$ [Taki00]. Takine's recursion is based on the fact that the $\mathbf{p}_{\text{BMAP/GI/1}}(n)$ is identical to the steady-state solution for a certain embedded Markov chain in the corresponding BMAP/GI/1 queue with multiple vacations and exhaustive services. In conclusion, it is shown that the $\mathbf{p}_{\text{BMAP/GI/1}}(n)$ is identical to the steady-state solution for the M/G/1-type Markov chain with the following transition probability matrix [Taki00]:

$$\begin{bmatrix} \mathbf{A}_0 & \mathbf{A}_1 & \mathbf{A}_2 & \mathbf{A}_3 & \cdots \\ \mathbf{A}_0 & \mathbf{A}_1 & \mathbf{A}_2 & \mathbf{A}_3 & \cdots \\ \mathbf{O} & \mathbf{A}_0 & \mathbf{A}_1 & \mathbf{A}_2 & \cdots \\ \mathbf{O} & \mathbf{O} & \mathbf{A}_0 & \mathbf{A}_1 & \cdots \\ \mathbf{O} & \mathbf{O} & \mathbf{O} & \mathbf{A}_0 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}, \quad (\text{B.5})$$

where \mathbf{A}_n 's satisfy (B.2). A related work is found in [Lee00], which considered an N/G/1 queue. Note that \mathbf{K} in (1.11) is identical to fundamental matrix \mathbf{G} for this special Markov chain of M/G/1 type (see subsection 1.2.2). Let \mathbf{g} denote a probability vector satisfying $\mathbf{g}\mathbf{G} = \mathbf{g}$. Then, Proposition 1.2 is reduced to the following result.

Proposition B.1 ([Taki00]) *$\mathbf{p}_{\text{BMAP/GI/1}}(n)$'s ($n = 0, 1, \dots$) are recursively determined in the following way.*

$$\mathbf{p}_{\text{BMAP/GI/1}}(0) = (1 - \rho)\mathbf{g},$$

where ρ denotes the utilization factor. For $n = 1, 2, \dots$,

$$\mathbf{p}_{\text{BMAP/GI/1}}(n) = \left[\mathbf{p}_{\text{BMAP/GI/1}}(0) \bar{\mathbf{A}}(n) + \sum_{l=1}^{n-1} \mathbf{p}_{\text{BMAP/GI/1}}(l) \bar{\mathbf{A}}(n-l+1) \right] [\mathbf{I} - \bar{\mathbf{A}}(1)]^{-1}.$$

We next consider a special case of a BMAP/GI/1 queue, where the arrival process is a MAP characterized by (\mathbf{C}, \mathbf{D}) (see subsection 1.1.1) and service times follow a phase-type distribution $(\boldsymbol{\beta}, \mathbf{T})$. This queue is denoted by MAP/PH/1. Let $\mathbf{p}_{\text{MAP/PH/1}}(0)$ denote a vector whose j th element represents the stationary probability that the system is idle and the underlying Markov chain is in state j . Also let $p_{\text{MAP/PH/1},j,\nu}(n)$ ($n = 1, 2, \dots$) denote the stationary joint probability that

the queue length is equal to n , the underlying Markov chain is in state j and the phase of service is ν . We construct a vector $\mathbf{p}_{\text{MAP/PH/1}}(n)$ by arranging $p_{\text{MAP/PH/1},j,\nu}(n)$'s in the lexicographical order, where the first index is the state j of the underlying Markov chain and the second index is the phase ν of service. Then the $\mathbf{p}_{\text{MAP/PH/1}}(n)$ is identical to the steady-state solution for the continuous-time M/G/1-type Markov chain whose infinitesimal generator is given by

$$\begin{bmatrix} \mathbf{C} & \mathbf{D} \otimes \boldsymbol{\beta} & \mathbf{O} & \mathbf{O} & \cdots \\ \mathbf{I}_{\text{MAP}} \otimes (-\mathbf{T})\mathbf{e} & \mathbf{C} \oplus \mathbf{T} & \mathbf{D} \otimes \mathbf{I}_{\text{PH}} & \mathbf{O} & \cdots \\ \mathbf{O} & \mathbf{I}_{\text{MAP}} \otimes (-\mathbf{T})\mathbf{e}\boldsymbol{\beta} & \mathbf{C} \oplus \mathbf{T} & \mathbf{D} \otimes \mathbf{I}_{\text{PH}} & \cdots \\ \mathbf{O} & \mathbf{O} & \mathbf{I}_{\text{MAP}} \otimes (-\mathbf{T})\mathbf{e}\boldsymbol{\beta} & \mathbf{C} \oplus \mathbf{T} & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}, \quad (\text{B.6})$$

where \mathbf{I}_{MAP} and \mathbf{I}_{PH} denote identity matrices with the same dimensions as \mathbf{C} and \mathbf{T} , respectively. Note that with the uniformization technique (see Appendix A), the above continuous-time Markov chain of G/M/1-type can be converted into a discrete-time Markov chain of G/M/1 type, whose transition probability matrix has the same structure as \mathbf{P}_{R} in (1.8), setting $\mathbf{E}_n = \mathbf{O}$ for $n \geq 3$. Thus, from Proposition 1.1, we can see that the $\mathbf{p}_{\text{MAP/PH/1}}(n)$ is a matrix geometric solution. This result has an advantage over Takine's recursion (see Proposition B.1) if the number of phases of service is considerably small. In the case of exponential service, we recommend computing the stationary queue length distribution with the matrix geometric solution.

Finally, we show the matrix geometric solution of a MAP/M/1 queue. Let $\mathbf{p}_{\text{MAP/M/1}}(n)$ ($n = 0, 1, \dots$) denote a vector whose j th element represents the stationary probability that the queue length is equal to n and the underlying Markov chain is in state j .

Corollary B.1 $\mathbf{p}_{\text{MAP/M/1}}(n)$'s are given by

$$\begin{aligned} \mathbf{p}_{\text{MAP/M/1}}(0) &= \boldsymbol{\pi}(\mathbf{I} - \mathbf{R}), \\ \mathbf{p}_{\text{MAP/M/1}}(n) &= \mathbf{p}_{\text{MAP/M/1}}(0)\mathbf{R}^n, \quad n = 1, 2, \dots, \end{aligned}$$

where $\boldsymbol{\pi}$ is a probability vector satisfying $\boldsymbol{\pi}(\mathbf{C} + \mathbf{D}) = \mathbf{0}$ and \mathbf{R} is the minimal nonnegative solution of

$$\mathbf{D} + \mathbf{R}(\mathbf{C} - \mu\mathbf{I}) + \mu\mathbf{R}^2 = \mathbf{O}.$$

Appendix C

Total Queue Length Distribution in a FIFO Queue

We show a numerical algorithm to compute the stationary total queue length distribution, by modifying our algorithm for the joint queue length distribution. We define $\mathbf{p}^{(\text{T})}(n)$ ($n = 0, 1, \dots$) and $\mathbf{q}_k^{(\text{T})}(n)$ ($k \in \mathcal{K}$, $n = 0, 1, \dots$) as

$$\mathbf{p}^{(\text{T})}(n) = \sum_{\substack{\mathbf{n} \in \mathcal{Z} \\ |\mathbf{n}|=n}} \mathbf{p}(\mathbf{n}), \quad \mathbf{q}_k^{(\text{T})}(n) = \sum_{\substack{\mathbf{n} \in \mathcal{Z} \\ |\mathbf{n}|=n}} \mathbf{q}_k(\mathbf{n}),$$

respectively. Corollary 2.1 is then reduced to

Corollary C.1 *The $\mathbf{p}^{(\text{T})}(n)$ is recursively determined in the following way:*

$$\mathbf{p}^{(\text{T})}(0) = \sum_{k \in \mathcal{K}} \lambda_k \mathbf{q}_k^{(\text{T})}(0) (-\mathbf{C})^{-1},$$

and for $n = 1, 2, \dots$,

$$\mathbf{p}^{(\text{T})}(n) = \sum_{k \in \mathcal{K}} \left[\lambda_k \left(\mathbf{q}_k^{(\text{T})}(n) - \mathbf{q}_k^{(\text{T})}(n-1) \right) + \sum_{m=1}^n \mathbf{p}^{(\text{T})}(n-m) \mathbf{D}_k(m) \right] (-\mathbf{C})^{-1}.$$

Further, under Assumption 2.2, Theorem 2.4 is reduced to

$$\mathbf{q}_k^{(\text{T})}(n) = \frac{1}{\lambda_k} \sum_{\substack{m_1+m_2+m_3 \\ +m_4=n}} \mathbf{v}_k^{(\text{T})}(m_1) [\boldsymbol{\alpha}_k \otimes \mathbf{A}_k^{(\text{T})}(m_2)] \boldsymbol{\Gamma}_k^{(\text{T})}(m_3) [\{\mathbf{P}_k^{m_4} (\mathbf{I} - \mathbf{P}_k) \mathbf{e}\} \otimes \mathbf{I}(M)],$$

where

$$\begin{aligned} \mathbf{A}_k^{(\text{T})}(n) &= \sum_{\substack{\mathbf{n} \in \mathcal{Z} \\ |\mathbf{n}|=n}} \mathbf{A}_k(\mathbf{n}), \\ \mathbf{v}_k^{(\text{T})}(n) &= \sum_{\substack{\mathbf{n} \in \mathcal{Z} \\ |\mathbf{n}|=n}} \mathbf{v}_k(\mathbf{n}), \quad \boldsymbol{\Gamma}_k^{(\text{T})}(n) = \sum_{\substack{\mathbf{n} \in \mathcal{Z} \\ |\mathbf{n}|=n}} \boldsymbol{\Gamma}_k(\mathbf{n}). \end{aligned} \tag{C.1}$$

Thus the $\mathbf{p}^{(\text{T})}(n)$ is obtained if we compute the $\mathbf{A}_k^{(\text{T})}(n)$, the $\mathbf{v}_k^{(\text{T})}(n)$ and the $\boldsymbol{\Gamma}_k^{(\text{T})}(n)$.

Note here that $\mathbf{A}_k^{(\text{T})}(n)$, $\mathbf{v}_k^{(\text{T})}(n)$ and $\mathbf{\Gamma}_k^{(\text{T})}(n)$ satisfy

$$\begin{aligned} \sum_{n=0}^{\infty} z^n \mathbf{A}_k^{(\text{T})}(n) &= \int_0^{\infty} dH_k(x) \exp \left[\left(\mathbf{C} + \sum_{k \in \mathcal{K}} \mathbf{D}_k^*(z) \right) x \right], \\ \sum_{n=0}^{\infty} z^n \mathbf{v}_k^{(\text{T})}(n) &= \int_0^{\infty} d\mathbf{v}(x) \mathbf{D}_k \exp \left[\left(\mathbf{C} + \sum_{k \in \mathcal{K}} \mathbf{D}_k^*(z) \right) x \right], \\ \sum_{n=0}^{\infty} z^n \mathbf{\Gamma}_k^{(\text{T})}(n) &= \left[\mathbf{I} - \mathbf{P}_k \otimes \int_0^{\infty} dH_k(x) \exp \left[\left(\mathbf{C} + \sum_{k \in \mathcal{K}} \mathbf{D}_k^*(z) \right) x \right] \right]^{-1}, \end{aligned} \quad (\text{C.2})$$

respectively. Thus $\mathbf{A}_k^{(\text{T})}(n)$ and $\mathbf{v}_k^{(\text{T})}(n)$ can be written to be

$$\begin{aligned} \mathbf{A}_k^{(\text{T})}(n) &= \sum_{m=0}^{\infty} \gamma_k^{(m)}(\theta) \mathbf{F}_m^{(\text{T})}(n), \\ \mathbf{v}_k^{(\text{T})}(n) &= \sum_{m=0}^{\infty} \mathbf{v}^{(m)}(\theta) \mathbf{D}_k \mathbf{F}_m^{(\text{T})}(n), \end{aligned}$$

respectively, where $\mathbf{F}_m^{(\text{T})}(n)$ denotes an $M \times M$ matrix which satisfies

$$\sum_{n=0}^{\infty} z^n \mathbf{F}_m^{(\text{T})}(n) = \left[\mathbf{I} + \theta^{-1} \left(\mathbf{C} + \sum_{k \in \mathcal{K}} \mathbf{D}_k^*(z) \right) \right]^m.$$

Further the $\mathbf{\Gamma}_k^{(\text{T})}(n)$ ($k \in \mathcal{K}$, $n \geq 0$) is determined by the following recursion:

$$\mathbf{\Gamma}_k^{(\text{T})}(0) = \left[\mathbf{I} - \mathbf{P}_k \otimes \mathbf{A}_k^{(\text{T})}(0) \right]^{-1},$$

and for $n = 1, 2, \dots$,

$$\mathbf{\Gamma}_k^{(\text{T})}(n) = \sum_{l=1}^n \mathbf{\Gamma}_k^{(\text{T})}(n-l) \left[\mathbf{P}_k \otimes \mathbf{A}_k^{(\text{T})}(l) \right] \mathbf{\Gamma}_k^{(\text{T})}(0).$$

Thus we can compute $\mathbf{A}_k^{(\text{T})}(n)$, $\mathbf{v}_k^{(\text{T})}(n)$ and $\mathbf{\Gamma}_k^{(\text{T})}(n)$ by replacing Steps 3 and 4 with the followings.

Step 3. Compute $\check{\mathbf{A}}_k^{(\text{T})}(n)$ and $\check{\mathbf{v}}_k^{(\text{T})}(n)$ by the following procedure, where the initial values of $\check{\mathbf{A}}_k^{(\text{T})}(n)$ and $\check{\mathbf{v}}_k^{(\text{T})}(n)$ ($n \geq 0$) are assumed to be \mathbf{O} and $\mathbf{0}$, respectively.

Step (3-a). Set $\check{\mathbf{F}}_0^{(\text{T})}(0) = \mathbf{I}$ and $n_F^{(0)} = 0$. Also set

$$\check{\mathbf{A}}_k^{(\text{T})}(0) = \gamma_k^{(0)}(\theta) \mathbf{I}, \quad \check{\mathbf{v}}_k^{(\text{T})}(0) = \mathbf{v}^{(0)}(\theta) \mathbf{D}_k, \quad \forall k \in \mathcal{K}.$$

Step (3-b). Set $n_F^{(1)} = \max_{k \in \mathcal{K}} n_g(k)$ and $m = 1$, and compute $\check{\mathbf{F}}_1^{(\text{T})}(n)$ by the following recursion:

$$\check{\mathbf{F}}_1^{(\text{T})}(0) = \mathbf{I} + \theta^{-1} \mathbf{C},$$

and for $n = 1, 2, \dots, n_F^{(1)}$,

$$\check{\mathbf{F}}_1^{(\text{T})}(n) = \theta^{-1} \sum_{k \in \mathcal{K}} U(n_g(k) - n) g_k(n) \mathbf{D}_k.$$

Step (3-c). For each $k \in \mathcal{K}$, if $m \leq m_\gamma(k)$, add $\gamma_k^{(m)}(\theta)\check{\mathbf{F}}_m^{(T)}(n)$ to $\check{\mathbf{A}}_k^{(T)}(n)$ for all $n \leq n_F^{(m)}$. Also, for each $k \in \mathcal{K}$, if $m \leq m_v(k)$, add $\mathbf{v}^{(m)}(\theta)\mathbf{D}_k\check{\mathbf{F}}_m^{(T)}(n)$ to $\check{\mathbf{v}}_k^{(T)}(n)$ for all $n \leq n_F^{(m)}$.

Step (3-d). If $m \geq m_{\max}$, stop computing, and otherwise, add one to m and go to Step (3-e).

Step (3-e). For each $n = 0, 1, \dots$, compute $\check{\mathbf{F}}_m^{(T)}(n)$ by

$$\begin{aligned} \check{\mathbf{F}}_m^{(T)}(n) &= U\left(n_F^{(m-1)} - n\right)\check{\mathbf{F}}_{m-1}^{(T)}(n)(\mathbf{I} + \theta^{-1}\mathbf{C}) \\ &\quad + \sum_{l=1}^{\min(n, n_F^{(1)})} U\left(n_F^{(m-1)} - n + l\right)\check{\mathbf{F}}_{m-1}^{(T)}(n-l)\check{\mathbf{F}}_1^{(T)}(l), \end{aligned}$$

until $\check{\mathbf{F}}_m^{(T)}(n)$'s satisfy $\sum_{n \leq n^*} \check{\mathbf{F}}_m^{(T)}(n)\mathbf{e} > (1 - \varepsilon_F)^m \mathbf{e}$ for some n^* . Let $n_F^{(m)} = n^*$ and go to Step (3-c).

Step 4. Set

$$n_A(k) = \max\left(n_F^{(m)}; m = 0, 1, \dots, m_\gamma(k)\right),$$

and for each $k \in \mathcal{K}$, compute $\check{\mathbf{\Gamma}}_k^{(T)}(n)$ by the following recursion:

$$\check{\mathbf{\Gamma}}_k^{(T)}(0) = \left[\mathbf{I} - \mathbf{P}_k \otimes \check{\mathbf{A}}_k^{(T)}(0)\right]^{-1},$$

and for $n = 1, 2, \dots$,

$$\check{\mathbf{\Gamma}}_k^{(T)}(n) = \sum_{l=1}^n U\left(n_A(k) - l\right)\check{\mathbf{\Gamma}}_k^{(T)}(n-l) \left[\mathbf{P}_k \otimes \check{\mathbf{A}}_k^{(T)}(l)\right] \check{\mathbf{\Gamma}}_k^{(T)}(0),$$

until $\check{\mathbf{\Gamma}}_k^{(T)}(n)$'s satisfy

$$\begin{aligned} \sum_{n=0}^{n_\Gamma(k)} \check{\mathbf{\Gamma}}_k^{(T)}(n)\mathbf{e} &> \left\{(\mathbf{I} - \mathbf{P}_k)^{-1} \mathbf{e}(M_k)\right\} \otimes \mathbf{e}(M) \\ &\quad - \varepsilon \left\{(\mathbf{I} - \mathbf{P}_k)^{-2} \mathbf{P}_k \mathbf{e}(M_k)\right\} \otimes \mathbf{e}(M), \end{aligned}$$

for some integer $n_\Gamma(k)$.

Remark C.1 *The above algorithm ensures that*

$$\sum_{n=0}^{n_A(k)} \check{\mathbf{A}}_k^{(T)}(n)\mathbf{e} > (1 - \varepsilon)\mathbf{e}, \quad \sum_{n=0}^{n_v(k)} \check{\mathbf{v}}_k^{(T)}(n)\mathbf{e} > (1 - \varepsilon)\lambda_k^{(B)},$$

respectively, where $n_v(k)$ is given by

$$n_v(k) = \max\left(n_F^{(m)}; m = 0, 1, \dots, m_v(k)\right).$$

Further $\check{\mathbf{\Gamma}}_k^{(T)}(n)$ satisfies

$$(\boldsymbol{\alpha}_k \otimes \boldsymbol{\pi}) \sum_{n=0}^{n_\Gamma(k)} \check{\mathbf{\Gamma}}_k^{(T)}(n)\mathbf{e} > \mathbb{E}[G_k] - \frac{1}{2}\mathbb{E}[G_k(G_k - 1)]\varepsilon.$$

Appendix D

Proof of Lemma 3.2

From the definition (3.18) of $\bar{\mathbf{F}}_m(\mathbf{n})$, (3.19) is clearly satisfied. (3.20) is proved in the following way. Let $\mathbf{N}_{\text{off}}(\mathbf{n})$ denote an $M_{\text{off}} \times M_{\text{off}}$ matrix satisfying

$$\sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \mathbf{N}_{\text{off}}(\mathbf{n}) = \left(-\mathbf{C}_{\text{off}} - \sum_{k \in \mathcal{K}} \mathbf{D}_{k,\text{off}}^*(z_k) \right)^{-1}. \quad (\text{D.1})$$

From (3.3), (3.18) and (D.1), we have for $m = 1, 2, \dots$,

$$\begin{aligned} & \sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \bar{\mathbf{F}}_m(\mathbf{n}) \\ &= \sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \bar{\mathbf{F}}_{m-1}(\mathbf{n}) \left[\mathbf{I} + \theta_{\text{on}}^{-1} \left\{ \mathbf{C}_{\text{on}} + \sum_{k \in \mathcal{K}} \sum_{n_k=1}^{\infty} z_k^{n_k} \mathbf{D}_{k,\text{on}}(n_k) \right. \right. \\ & \quad \left. \left. + \mathbf{E}_{\text{on,off}} \sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \mathbf{N}_{\text{off}}(\mathbf{n}) \mathbf{E}_{\text{off,on}} \right\} \right] \\ &= \sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \bar{\mathbf{F}}_{m-1}(\mathbf{n}) (\mathbf{I} + \theta_{\text{on}}^{-1} \mathbf{C}_{\text{on}}) \\ & \quad + \sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \theta_{\text{on}}^{-1} \sum_{k \in \mathcal{K}} \sum_{l_k=1}^{n_k} \bar{\mathbf{F}}_{m-1}(\mathbf{n} - l_k \mathbf{e}_k) \mathbf{D}_{k,\text{on}}(l_k) \\ & \quad + \sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \theta_{\text{on}}^{-1} \sum_{\mathbf{0} \leq \mathbf{l} \leq \mathbf{n}} \bar{\mathbf{F}}_{m-1}(\mathbf{n} - \mathbf{l}) \mathbf{E}_{\text{on,off}} \mathbf{N}_{\text{off}}(\mathbf{l}) \mathbf{E}_{\text{off,on}}. \end{aligned}$$

Comparing coefficient matrices of $z_1^{n_1} \cdots z_K^{n_K}$ on both sides of the above equation, we obtain (3.20).

Next, we show (3.21) and (3.22). From (3.3) and (D.1), we have

$$\sum_{\mathbf{n} \in \mathcal{Z}} z_1^{n_1} \cdots z_K^{n_K} \mathbf{N}_{\text{off}}(\mathbf{n}) \left(-\mathbf{C}_{\text{off}} - \sum_{k \in \mathcal{K}} \sum_{n_k=1}^{\infty} z_k^{n_k} \mathbf{D}_{k,\text{off}}(n_k) \right) = \mathbf{I}.$$

Comparing coefficient matrices of $z_1^{n_1} \cdots z_K^{n_K}$ on both sides of the above equation yields

$$\begin{aligned} \mathbf{N}_{\text{off}}(\mathbf{0}) (-\mathbf{C}_{\text{off}}) &= \mathbf{I}, \\ \mathbf{N}_{\text{off}}(\mathbf{n}) (-\mathbf{C}_{\text{off}}) - \sum_{k \in \mathcal{K}} \sum_{l_k=1}^{n_k} \mathbf{N}_{\text{off}}(\mathbf{n} - l_k \mathbf{e}_k) \mathbf{D}_{k,\text{off}}(l_k) &= \mathbf{O}, \end{aligned}$$

from which (3.21) and (3.22) follow. ■

Appendix E

Proof of Theorem 3.5

Post-multiplying both sides of (3.2) by $-\mathbf{C}_{\text{off}} - \overline{\mathbf{D}}_{\text{off}}^*(s)$ and substituting $\theta_{\text{on}} - \theta_{\text{on}}z$ for s , we obtain

$$\sum_{m=0}^{\infty} z^m \mathbf{v}_{\text{off}}^{(m)}(\theta_{\text{on}}) \left[-\mathbf{C}_{\text{off}} - \sum_{m=0}^{\infty} z^m \mathbf{D}_{\text{off}}^{(m)}(\theta_{\text{on}}) \right] = \frac{\sum_{m=0}^{\infty} z^m \mathbf{v}_{\text{on}}^{(m)}(\theta_{\text{on}}) \mathbf{E}_{\text{on,off}}}{\boldsymbol{\pi}_{\text{on}} \mathbf{E}_{\text{on,off}} \mathbf{e}} \frac{1}{\overline{I}_{\text{off}}}, \quad (\text{E.1})$$

where we use (3.24) and (3.26). Comparing coefficient vectors of z^m ($m = 0, 1, \dots$) on both sides of (E.1) yields

$$\mathbf{v}_{\text{off}}^{(0)}(\theta_{\text{on}}) \left[-\mathbf{C}_{\text{off}} - \mathbf{D}_{\text{off}}^{(0)}(\theta_{\text{on}}) \right] = \frac{\mathbf{v}_{\text{on}}^{(0)}(\theta_{\text{on}}) \mathbf{E}_{\text{on,off}}}{\boldsymbol{\pi}_{\text{on}} \mathbf{E}_{\text{on,off}} \mathbf{e}} \frac{1}{\overline{I}_{\text{off}}},$$

and for $m = 1, 2, \dots$,

$$\mathbf{v}_{\text{off}}^{(m)}(\theta_{\text{on}}) \left[-\mathbf{C}_{\text{off}} - \mathbf{D}_{\text{off}}^{(0)}(\theta_{\text{on}}) \right] - \sum_{l=0}^{m-1} \mathbf{v}_{\text{off}}^{(l)}(\theta_{\text{on}}) \mathbf{D}_{\text{off}}^{(m-l)}(\theta_{\text{on}}) = \frac{\mathbf{v}_{\text{on}}^{(m)}(\theta_{\text{on}}) \mathbf{E}_{\text{on,off}}}{\boldsymbol{\pi}_{\text{on}} \mathbf{E}_{\text{on,off}} \mathbf{e}} \frac{1}{\overline{I}_{\text{off}}}.$$

(3.27) and (3.28) follow the above two equations. ■

Appendix F

Total Queue Length Distribution in a Queue with Service Interruptions

This appendix summarizes the recursions for the total queue length distribution. Because they are readily derived from the results in Section 3.4, we omit the proofs.

We define $\mathbf{p}^{(\text{T})}(n)$ ($n = 0, 1, \dots$) and $\mathbf{q}_k^{(\text{T})}(n)$ ($k \in \mathcal{K}, n = 0, 1, \dots$) as the stationary total queue length distributions at a random point in time and at immediately after departures of class k , respectively.

$$\mathbf{p}^{(\text{T})}(n) = \sum_{\substack{\mathbf{n} \in \mathcal{Z} \\ |\mathbf{n}|=n}} \mathbf{p}(\mathbf{n}), \quad \mathbf{q}_k^{(\text{T})}(n) = \sum_{\substack{\mathbf{n} \in \mathcal{Z} \\ |\mathbf{n}|=n}} \mathbf{q}_k(\mathbf{n}),$$

where $|\mathbf{n}| = |n_1| + \dots + |n_K|$. Since $\mathbf{p}^{(\text{T})}(n)$'s can be obtained in terms of $\mathbf{q}_k^{(\text{T})}(n)$'s (see Corollary C.1), we hereafter consider the total queue length distribution $\mathbf{q}_k^{(\text{T})}(n)$ immediately after departures.

We first introduce the following notations: For $k \in \mathcal{K}$, $\xi = \text{on, off}$, $n = 0, 1, \dots$ and $m = 0, 1, \dots$,

$$\begin{aligned} \mathbf{q}_{k,\xi}^{(\text{T})}(n) &= \sum_{\substack{\mathbf{n} \in \mathcal{Z} \\ |\mathbf{n}|=n}} \mathbf{q}_{k,\xi}(\mathbf{n}), & \mathbf{\Gamma}_{k,\xi}^{(\text{T})}(n) &= \sum_{\substack{\mathbf{n} \in \mathcal{Z} \\ |\mathbf{n}|=n}} \mathbf{\Gamma}_{k,\xi}(\mathbf{n}), & \mathbf{A}_{k,\xi}^{(\text{T})}(n) &= \sum_{\substack{\mathbf{n} \in \mathcal{Z} \\ |\mathbf{n}|=n}} \mathbf{A}_{k,\xi}(\mathbf{n}), \\ \mathbf{v}_{k,\xi}^{(\text{T})}(n) &= \sum_{\substack{\mathbf{n} \in \mathcal{Z} \\ |\mathbf{n}|=n}} \mathbf{v}_{k,\xi}(\mathbf{n}), & \mathbf{\bar{F}}_m^{(\text{T})}(n) &= \sum_{\substack{\mathbf{n} \in \mathcal{Z} \\ |\mathbf{n}|=n}} \mathbf{\bar{F}}_m(\mathbf{n}). \end{aligned}$$

From Theorem 3.2, we have the following corollary.

Corollary F.1 *Under Assumption 3.1, the $\mathbf{q}_k^{(\text{T})}(n)$ ($k \in \mathcal{K}, n = 0, 1, \dots$) is given by*

$$\mathbf{q}_k^{(\text{T})}(n) = \left(\frac{r_{\text{on}} \lambda_{k,\text{on}}}{\lambda_k} \mathbf{q}_{k,\text{on}}^{(\text{T})}(n) + \frac{r_{\text{off}} \lambda_{k,\text{off}}}{\lambda_k} \mathbf{q}_{k,\text{off}}^{(\text{T})}(n), 0, \dots, 0 \right),$$

where the $\mathbf{q}_{k,\xi}^{(\text{T})}(n)$ ($k \in \mathcal{K}$, $\xi = \text{on, off}$, $n = 0, 1, \dots$) is given by

$$\begin{aligned} \mathbf{q}_{k,\xi}^{(\text{T})}(n) &= \frac{1}{\lambda_{k,\xi}} \sum_{\substack{m_1+m_2+m_3 \\ +m_4=n}} \mathbf{v}_{k,\xi}^{(\text{T})}(m_1) [\boldsymbol{\alpha}_{k,\xi} \otimes \mathbf{A}_{k,\xi}^{(\text{T})}(m_2)] \mathbf{\Gamma}_{k,\xi}^{(\text{T})}(m_3) \\ &\quad \cdot \left[\left\{ \mathbf{P}_{k,\xi}^{m_4} (\mathbf{I} - \mathbf{P}_{k,\xi}) \mathbf{e} \right\} \otimes \mathbf{I}(M_{\text{on}}) \right], \end{aligned}$$

if $\lambda_{k,\xi} > 0$, and otherwise $\mathbf{q}_{k,\xi}^{(\text{T})}(n) = \mathbf{0}$.

The following corollary for the $\mathbf{\Gamma}_{k,\xi}^{(\text{T})}(n)$ can be readily obtained from Lemma 3.1.

Corollary F.2 *The $\mathbf{\Gamma}_{k,\xi}^{(\text{T})}(n)$ ($k \in \mathcal{K}$, $\xi = \text{on}, \text{off}$, $n = 0, 1, \dots$) is determined by the following recursion:*

$$\begin{aligned}\mathbf{\Gamma}_{k,\xi}^{(\text{T})}(0) &= \left[\mathbf{I} - \mathbf{P}_{k,\xi} \otimes \mathbf{A}_{k,\xi}^{(\text{T})}(0) \right]^{-1}, \\ \mathbf{\Gamma}_{k,\xi}^{(\text{T})}(n) &= \sum_{l=1}^n \mathbf{\Gamma}_{k,\xi}^{(\text{T})}(n-l) \left[\mathbf{P}_{k,\xi} \otimes \mathbf{A}_{k,\xi}^{(\text{T})}(l) \right] \mathbf{\Gamma}_{k,\xi}^{(\text{T})}(0), \quad n = 1, 2, \dots\end{aligned}$$

As for the $\mathbf{A}_{k,\xi}^{(\text{T})}(n)$, the following recursions can be obtained from Lemma 3.2 and Theorem 3.3.

Corollary F.3 *The $\mathbf{A}_{k,\xi}^{(\text{T})}(n)$ is given by*

$$\mathbf{A}_{k,\xi}^{(\text{T})}(n) = \sum_{m=0}^{\infty} \gamma_{k,\xi}^{(m)}(\theta_{\text{on}}) \overline{\mathbf{F}}_m^{(\text{T})}(n), \quad k \in \mathcal{K}, \xi = \text{on}, \text{off}, n = 0, 1, \dots,$$

where the $\overline{\mathbf{F}}_m^{(\text{T})}(n)$ is recursively determined by

$$\overline{\mathbf{F}}_0^{(\text{T})}(n) = \begin{cases} \mathbf{I}, & \text{if } n = 0, \\ \mathbf{O}, & \text{otherwise,} \end{cases}$$

and for $m = 1, 2, \dots$,

$$\begin{aligned}\overline{\mathbf{F}}_m^{(\text{T})}(n) &= \overline{\mathbf{F}}_{m-1}^{(\text{T})}(n) (\mathbf{I} + \theta_{\text{on}}^{-1} \mathbf{C}_{\text{on}}) + \theta_{\text{on}}^{-1} \sum_{l=1}^n \overline{\mathbf{F}}_{m-1}^{(\text{T})}(n-l) \sum_{k \in \mathcal{K}} \mathbf{D}_{k,\text{on}}(l) \\ &\quad + \theta_{\text{on}}^{-1} \left[\sum_{l=0}^n \overline{\mathbf{F}}_{m-1}^{(\text{T})}(n-l) \mathbf{E}_{\text{on},\text{off}} \mathbf{N}_{\text{off}}^{(\text{T})}(l) \right] \mathbf{E}_{\text{off},\text{on}}, \quad n = 0, 1, \dots,\end{aligned}$$

with the $\mathbf{N}_{\text{off}}^{(\text{T})}(n)$, which is given by the following recursion:

$$\begin{aligned}\mathbf{N}_{\text{off}}^{(\text{T})}(0) &= (-\mathbf{C}_{\text{off}})^{-1}, \\ \mathbf{N}_{\text{off}}^{(\text{T})}(n) &= \left[\sum_{l=1}^n \mathbf{N}_{\text{off}}^{(\text{T})}(n-l) \sum_{k \in \mathcal{K}} \mathbf{D}_{k,\text{off}}(l) \right] \mathbf{N}_{\text{off}}^{(\text{T})}(0), \quad n = 1, 2, \dots\end{aligned}$$

Finally, from Theorem 3.4, we obtain the following result.

Corollary F.4 *The $\mathbf{v}_{k,\text{on}}^{(\text{T})}(n)$ ($k \in \mathcal{K}$) and the $\mathbf{v}_{k,\text{off}}^{(\text{T})}(n)$ ($k \in \mathcal{K}$) are determined by*

$$\begin{aligned}\mathbf{v}_{k,\text{on}}^{(\text{T})}(n) &= \sum_{m=0}^{\infty} \mathbf{v}_{\text{on}}^{(m)}(\theta_{\text{on}}) \mathbf{D}_{k,\text{on}} \overline{\mathbf{F}}_m^{(\text{T})}(n), \quad n = 0, 1, \dots, \\ \mathbf{v}_{k,\text{off}}^{(\text{T})}(n) &= \sum_{m=0}^{\infty} \mathbf{v}_{\text{off}}^{(m)}(\theta_{\text{on}}) \mathbf{D}_{k,\text{off}} \sum_{l=0}^n \mathbf{N}_{\text{off}}^{(\text{T})}(n-l) \mathbf{E}_{\text{off},\text{on}} \overline{\mathbf{F}}_m^{(\text{T})}(l), \quad n = 0, 1, \dots,\end{aligned}$$

respectively.

Appendix G

Proof of Lemma 4.1

By mathematical induction, we prove Lemma 4.1. We first consider the case of $\nu = 2$. We define $\mathbf{G}_1(t, \mathbf{U}_1, \mathbf{U}_2, \mathbf{V}_1)$ as

$$\mathbf{G}_1(t, \mathbf{U}_1, \mathbf{U}_2, \mathbf{V}_1) = \int_0^t du e^{\mathbf{U}_1 u} \mathbf{V}_1 e^{\mathbf{U}_2(t-u)}. \quad (\text{G.1})$$

By differentiating (G.1) with respect to t , we have

$$\frac{d}{dt} \mathbf{G}_1(t, \mathbf{U}_1, \mathbf{U}_2, \mathbf{V}_1) - \mathbf{G}_1(t, \mathbf{U}_1, \mathbf{U}_2, \mathbf{V}_1) \mathbf{U}_2 = e^{\mathbf{U}_1 t} \mathbf{V}_1. \quad (\text{G.2})$$

On the other hand, by straightforward calculation based on the definition of matrix exponential, we obtain

$$\exp \left[\begin{pmatrix} \mathbf{U}_1 & \mathbf{V}_1 \\ \mathbf{O} & \mathbf{U}_2 \end{pmatrix} t \right] = \begin{bmatrix} e^{\mathbf{U}_1 t} & \tilde{\mathbf{G}}_1(t, \mathbf{U}_1, \mathbf{U}_2, \mathbf{V}_1) \\ \mathbf{O} & e^{\mathbf{U}_2 t} \end{bmatrix}, \quad (\text{G.3})$$

where

$$\tilde{\mathbf{G}}_1(t, \mathbf{U}_1, \mathbf{U}_2, \mathbf{V}_1) = \sum_{n=1}^{\infty} \frac{t^n}{n!} \sum_{j=0}^{n-1} \mathbf{U}_1^j \mathbf{V}_1 \mathbf{U}_2^{n-1-j}.$$

It is easy to verify that $\tilde{\mathbf{G}}_1(t, \mathbf{U}_1, \mathbf{U}_2, \mathbf{V}_1)$ satisfies

$$\frac{d}{dt} \tilde{\mathbf{G}}_1(t, \mathbf{U}_1, \mathbf{U}_2, \mathbf{V}_1) - \tilde{\mathbf{G}}_1(t, \mathbf{U}_1, \mathbf{U}_2, \mathbf{V}_1) \mathbf{U}_2 = e^{\mathbf{U}_1 t} \mathbf{V}_1. \quad (\text{G.4})$$

Taking the difference between (G.2) and (G.4), we have

$$\begin{aligned} \frac{d}{dt} \left[\mathbf{G}_1(t, \mathbf{U}_1, \mathbf{U}_2, \mathbf{V}_1) - \tilde{\mathbf{G}}_1(t, \mathbf{U}_1, \mathbf{U}_2, \mathbf{V}_1) \right] \\ = \left[\mathbf{G}_1(t, \mathbf{U}_1, \mathbf{U}_2, \mathbf{V}_1) - \tilde{\mathbf{G}}_1(t, \mathbf{U}_1, \mathbf{U}_2, \mathbf{V}_1) \right] \mathbf{U}_2, \end{aligned}$$

from which it follows that

$$\mathbf{G}_1(t, \mathbf{U}_1, \mathbf{U}_2, \mathbf{V}_1) - \tilde{\mathbf{G}}_1(t, \mathbf{U}_1, \mathbf{U}_2, \mathbf{V}_1) = \mathbf{K} e^{\mathbf{U}_2 t},$$

for some constant matrix \mathbf{K} . Note here that

$$\tilde{\mathbf{G}}_1(0, \mathbf{U}_1, \mathbf{U}_2, \mathbf{V}_1) = \mathbf{G}_1(0, \mathbf{U}_1, \mathbf{U}_2, \mathbf{V}_1) = \mathbf{O},$$

so that $\mathbf{K} = \mathbf{O}$. Thus we have

$$\tilde{\mathbf{G}}_1(t, \mathbf{U}_1, \mathbf{U}_2, \mathbf{V}_1) = \mathbf{G}_1(t, \mathbf{U}_1, \mathbf{U}_2, \mathbf{V}_1), \quad t \geq 0,$$

and therefore (G.3) is rewritten to be

$$\exp \left[\begin{pmatrix} \mathbf{U}_1 & \mathbf{V}_1 \\ \mathbf{O} & \mathbf{U}_2 \end{pmatrix} t \right] = \begin{bmatrix} e^{\mathbf{U}_1 t} & \int_0^t du e^{\mathbf{U}_1 u} \mathbf{V}_1 e^{\mathbf{U}_2(t-u)} \\ \mathbf{O} & e^{\mathbf{U}_2 t} \end{bmatrix}. \quad (\text{G.5})$$

As a result, we have

$$\begin{aligned} \mathbf{V}_0 \int_0^t du e^{\mathbf{U}_1 u} \mathbf{V}_1 e^{\mathbf{U}_2(t-u)} \mathbf{V}_2 &= \mathbf{V}_0 \begin{bmatrix} \mathbf{I} & \mathbf{O} \end{bmatrix} \exp \left[\begin{pmatrix} \mathbf{U}_1 & \mathbf{V}_1 \\ \mathbf{O} & \mathbf{U}_2 \end{pmatrix} t \right] \begin{bmatrix} \mathbf{O} \\ \mathbf{I} \end{bmatrix} \mathbf{V}_2 \\ &= \begin{bmatrix} \mathbf{V}_0 & \mathbf{O} \end{bmatrix} \exp \left[\begin{pmatrix} \mathbf{U}_1 & \mathbf{V}_1 \\ \mathbf{O} & \mathbf{U}_2 \end{pmatrix} t \right] \begin{bmatrix} \mathbf{O} \\ \mathbf{V}_2 \end{bmatrix}, \end{aligned} \quad (\text{G.6})$$

which shows (4.32) holds for $\nu = 2$.

Suppose that (4.32) holds for some $\nu = m$ ($m \geq 2$), i.e.,

$$\begin{aligned} &\mathbf{V}_0 \int_0^t du_{m-1} \int_0^{u_{m-1}} du_{m-2} \cdots \int_0^{u_2} du_1 e^{\mathbf{U}_1 u_1} \mathbf{V}_1 e^{\mathbf{U}_2(u_2-u_1)} \mathbf{V}_2 \\ &\quad \cdots \cdots e^{\mathbf{U}_{m-1}(u_{m-1}-u_{m-2})} \mathbf{V}_{m-1} e^{\mathbf{U}_m(t-u_{m-1})} \mathbf{V}_m \\ &= \begin{bmatrix} \mathbf{V}_0 & \mathbf{O} & \cdots & \mathbf{O} \end{bmatrix} \exp \left[\begin{pmatrix} \mathbf{U}_1 & \mathbf{V}_1 & \mathbf{O} & \cdots & \mathbf{O} \\ \mathbf{O} & \mathbf{U}_2 & \mathbf{V}_2 & & \vdots \\ \vdots & & \ddots & \ddots & \mathbf{O} \\ \mathbf{O} & \cdots & \mathbf{O} & \mathbf{U}_{m-1} & \mathbf{V}_{m-1} \\ \mathbf{O} & \mathbf{O} & \cdots & \mathbf{O} & \mathbf{U}_m \end{pmatrix} t \right] \begin{bmatrix} \mathbf{O} \\ \vdots \\ \mathbf{O} \\ \mathbf{V}_m \end{bmatrix}. \end{aligned}$$

We then have

$$\begin{aligned} &\int_0^t du_m \left[\mathbf{V}_0 \int_0^{u_m} du_{m-1} \int_0^{u_{m-1}} du_{m-2} \cdots \int_0^{u_2} du_1 e^{\mathbf{U}_1 u_1} \mathbf{V}_1 e^{\mathbf{U}_2(u_2-u_1)} \mathbf{V}_2 \right. \\ &\quad \left. \cdots \cdots e^{\mathbf{U}_m(u_m-u_{m-1})} \mathbf{V}_m \right] e^{\mathbf{U}_{m+1}(t-u_m)} \mathbf{V}_{m+1} \\ &= \begin{bmatrix} \mathbf{V}_0 & \mathbf{O} & \cdots & \mathbf{O} \end{bmatrix} \\ &\quad \cdot \int_0^t du_m \exp \left[\begin{pmatrix} \mathbf{U}_1 & \mathbf{V}_1 & \mathbf{O} & \cdots & \mathbf{O} \\ \mathbf{O} & \mathbf{U}_2 & \mathbf{V}_2 & & \vdots \\ \vdots & & \ddots & \ddots & \mathbf{O} \\ \mathbf{O} & \cdots & \mathbf{O} & \mathbf{U}_{m-1} & \mathbf{V}_{m-1} \\ \mathbf{O} & \mathbf{O} & \cdots & \mathbf{O} & \mathbf{U}_m \end{pmatrix} u_m \right] \begin{bmatrix} \mathbf{O} \\ \vdots \\ \mathbf{O} \\ \mathbf{V}_m \end{bmatrix} e^{\mathbf{U}_{m+1}(t-u_m)} \mathbf{V}_{m+1}. \end{aligned} \quad (\text{G.7})$$

Thus, applying (G.6) to (G.7), we obtain

$$\int_0^t du_m \left[\mathbf{V}_0 \int_0^{u_m} du_{m-1} \int_0^{u_{m-1}} du_{m-2} \cdots \int_0^{u_2} du_1 e^{\mathbf{U}_1 u_1} \mathbf{V}_1 e^{\mathbf{U}_2(u_2-u_1)} \mathbf{V}_2 \right.$$

$$\begin{aligned}
& \dots e^{U_m(u_m - u_{m-1})} \mathbf{V}_m \Big] e^{U_{m+1}(t - u_m)} \mathbf{V}_{m+1} \\
= & \left[\mathbf{V}_0 \quad \mathbf{O} \quad \dots \quad \mathbf{O} \quad \Big| \quad \mathbf{O} \right] \\
& \cdot \exp \left(\left(\left[\begin{array}{ccccc} U_1 & V_1 & \mathbf{O} & \dots & \mathbf{O} \\ \mathbf{O} & U_2 & V_2 & & \vdots \\ \vdots & & \ddots & \ddots & \mathbf{O} \\ \mathbf{O} & \dots & \mathbf{O} & U_{m-1} & V_{m-1} \\ \mathbf{O} & \mathbf{O} & \dots & \mathbf{O} & U_m \end{array} \right] \Big| \left[\begin{array}{c} \mathbf{O} \\ \vdots \\ \mathbf{O} \\ V_m \end{array} \right] \right) t \Bigg] \left[\begin{array}{c} \mathbf{O} \\ \vdots \\ \mathbf{O} \\ \hline V_{m+1} \end{array} \right],
\end{aligned}$$

which shows that (4.32) holds for $\nu = m + 1$, too. ■

Bibliography

- [Asmu91] S. Asmussen, Ladder Heights and the Markov-Modulated M/G/1 Queue, *Stochastic Processes and their Applications*, 37 (1991) 313–326.
- [Asmu93] S. Asmussen and G. Koole, Marked Point Processes as Limits of Markovian Arrival Streams, *Journal of Applied Probability*, 30 (1993) 365–372.
- [Asmu03] S. Asmussen, *Applied Probability and Queues*, 2nd ed. (Springer Verlag, New York, 2003).
- [Bell97] R. Bellman, *Introduction to Matrix Analysis (Classics in Applied Mathematics)*, 2nd ed. (SIAM, Philadelphia, 1997).
- [Bors03] S. C. Borst, O. J. Boxma, J. A. Morrison and R. Núñez Queija, The Equivalence between Processor Sharing and Service in Random Order, *Operations Research Letters*, 31 (2003) 254–262.
- [Coff70] E. G. Coffman, Jr., R. R. Muntz and H. Trotter, Waiting Time Distributions for Processor-Sharing Systems, *Journal of the Association for Computing Machinery*, 17 (1970) 123–130.
- [Cohe84] J. W. Cohen, On Processor Sharing and Random Service, *Journal of Applied Probability*, 21 (1984) 937.
- [Eick93] S. G. Eick, W. A. Massey and W. Whitt, $M_t/G/\infty$ Queues with Sinusoidal Arrival Rates, *Management Science*, 39 (1993) 241–252.
- [Fede86] A. Federgruen and L. Green, Queueing Systems with Service Interruptions, *Operations Research*, 34 (1986) 752–768.
- [Fede88] A. Federgruen and L. Green, Queueing Systems with Service Interruptions II, *Naval Research Logistics*, 35 (1988) 345–358.
- [Fox88] B. L. Fox and P. W. Glynn, Computing Poisson Probabilities, *Communications of the ACM*, 31 (1988) 440–445.
- [Gris94] S. Grishechkin, GI/G/1 Processor Sharing Queue in Heavy Traffic, *Advances in Applied Probability*, 26 (1994) 539–555.

- [Guil01] F. Guillemin and J. Boyer, Analysis of the M/M/1 Queue with Processor Sharing via Spectral Theory, *Queueing Systems*, 39 (2001) 377–397.
- [He96] Q.-M. He, Queues with Marked Customers, *Advances in Applied Probability*, 28 (1996) 567–587.
- [He98] Q.-M. He and A. S. Alfa, The MMAP[K]/PH[K]/1 Queues with a Last-Come-First-Served Preemptive Service Discipline, *Queueing Systems*, 29 (1998) 269–291.
- [He00a] Q.-M. He, Classification of Markov Processes of M/G/1 Type with a Tree Structure and Its Applications to Queueing Models, *Operations Research Letters*, 26 (2000) 67–80.
- [He00b] Q.-M. He, Classification of Markov Processes of Matrix M/G/1 Type with a Tree Structure and Its Applications to the MMAP[K]/G[K]/1 Queues, *Stochastic Models*, 16 (2000) 407–433.
- [He01] Q.-M. He, The Versatility of MMAP[K] and the MMAP[K]/G[K]/1 Queue, *Queueing Systems*, 38 (2001) 397–418.
- [He03a] Q.-M. He, A Fixed Point Approach to the Classification of Markov Chains with a Tree Structure, *Stochastic Models*, 19 (2003) 75–111.
- [He03b] Q.-M. He and H. Li, Stability Conditions of the MMAP[K]/G[K]/1/LCFS Preemptive Repeat Queue, *Queueing Systems*, 44 (2003) 137–160.
- [Holm82] D. F. Holman, M. L. Chaudhry and B. R. K. Kashyap, On the Number in the System $GI^X/M/\infty$, *Sankhyā : The Indian Journal of Statistics*, 44, Series A, Part 1 (1982) 294–297.
- [Holm83] D. F. Holman, M. L. Chaudhry and B. R. K. Kashyap, On the Service System $M^X/G/\infty$, *European Journal of Operational Research*, 13 (1983) 142–145.
- [Jage91] D. L. Jagerman and B. Sengupta, The GI/M/1 Processor-Sharing Queue and Its Heavy Traffic Analysis, *Stochastic Models*, 7 (1991) 379–395.
- [Lee00] G. Lee and J. Jeon, A New Approach to an N/G/1 Queue, *Queueing Systems*, 35 (2000) 317–322.
- [Liu90] L. Liu, B. R. K. Kashyap and J. G. C. Templeton, On the $GI^X/G/\infty$ System, *Journal of Applied Probability*, 27 (1990) 671–683.
- [Liu91] L. Liu and J. G. C. Templeton, The $GR^{X_n}/G_n/\infty$ System: System Size, *Queueing Systems*, 8 (1991) 323–356.
- [Loyn62] R. M. Loynes, The Stability of a Queue with Non-Independent Inter-Arrival and Service Times, *Proceedings of the Cambridge Philosophical Society*, 58 (1962) 497–520.

- [Luca90] D. M. Lucantoni, K. S. Meier-Hellstern and M. F. Neuts, A Single-Server Queue with Server Vacations and a Class of Non-Renewal Arrival Processes, *Advances in Applied Probability*, 22 (1990) 676–705.
- [Luca91] D. M. Lucantoni, New Results on the Single Server Queue with a Batch Markovian Arrival Process, *Stochastic Models*, 7 (1991) 1–46.
- [Luca93] D. M. Lucantoni, The BMAP/G/1 Queue: A Tutorial, *Models and Techniques for Performance Evaluation of Computer and Communication Systems*, L. Donatiello and R. Nelson, eds., Springer Verlag (1993) 330–358.
- [Mach86] F. Machihara, An Infinitely-Many-Server Queue Having Markov Renewal Arrivals and Hyperexponential Service Times, *Journal of the Operations Research Society of Japan*, 29 (1986) 338–351
- [Mach93] F. Machihara, On the Queue with PH-Markov Renewal Preemptions, *Journal of the Operations Research Society of Japan*, 36 (1993) 13–28.
- [Mach99] F. Machihara, A BMAP/SM/1 Queue with Service Times Depending on the Arrival Process, *Queueing Systems*, 33 (1999) 277–291.
- [Masu02] H. Masuyama and T. Takine, Analysis of an Infinite-Server Queue with Batch Markovian Arrival Streams, *Queueing Systems*, 42 (2002) 269–296.
- [Masu03a] H. Masuyama and T. Takine, Analysis and Computation of the Joint Queue Length Distribution in a FIFO Single-Server Queue with Multiple Batch Markovian Arrival Streams, *Stochastic Models*, 19 (2003) 349–381.
- [Masu03b] H. Masuyama and T. Takine, Sojourn Time Distribution in a MAP/M/1 Processor-Sharing Queue, *Operations Research Letters*, 31 (2003) 406–412.
- [Masu03c] H. Masuyama and T. Takine, Stationary Queue Length in a FIFO Single Server Queue with Service Interruptions and Multiple Batch Markovian Arrival Streams, *Journal of the Operations Research Society of Japan*, 46 (2003) 319–341.
- [Morr85] J. A. Morrison, Response-Time Distribution for a Processor-Sharing System, *SIAM Journal on Applied Mathematics*, 45 (1985) 152–167.
- [Neut79] M. F. Neuts, A Versatile Markovian Point Process, *Journal of Applied Probability*, 16 (1979) 764–779.
- [Neut81] M. F. Neuts, *Matrix-Geometric Solutions in Stochastic Models: An Algorithmic Approach* (Johns Hopkins University Press, Baltimore, 1981).
- [Neut89] M. F. Neuts, *Structured Stochastic Matrices of M/G/1 Type and Their Applications* (Marcel Dekker, New York, 1989).

- [Ott84] T. J. Ott, The Sojourn-Time Distribution in the M/G/1 Queue with Processor Sharing, *Journal of Applied Probability*, 21 (1984) 360–378.
- [Rama78] V. Ramaswami, The N/G/ ∞ Queue, Technical Report, Department of Mathematics, Drexel University, Philadelphia (1978).
- [Rama80] V. Ramaswami and M. F. Neuts, Some Explicit Formulas and Computational Methods for Infinite-Server Queues with Phase-Type Arrivals, *Journal of Applied Probability*, 17 (1980) 498–514.
- [Rama84] V. Ramaswami, The Sojourn Time in the GI/M/1 Queue with Processor Sharing, *Journal of Applied Probability*, 21 (1984) 437–442.
- [Rama88] V. Ramaswami, A Stable Recursion for the Steady State Vector in Markov Chains of M/G/1 Type, *Stochastic Models*, 4 (1988) 183–188.
- [Regt86] G. J. K. Regterschot and J. H. A. de Smit, The Queue M|G|1 with Markov Modulated Arrivals and Services, *Mathematics of Operations Research*, 11 (1986) 465–483.
- [Sche90] H. Schellhaas, On Ramaswami's Algorithm for the Computation of the Steady State Vector in Markov Chains of M/G/1-Type, *Stochastic Models*, 6 (1990) 541–550.
- [Seng85] B. Sengupta and D. L. Jagerman, A Conditional Response Time of the M/M/1 Processor-Sharing Queue, *AT & T Technical Journal*, 64 (1985) 409–421.
- [Seng89] B. Sengupta, Markov Processes Whose Steady State Distribution Is Matrix-Exponential with an Application to the GI/PH/1 Queue, *Advances in Applied Probability*, 21 (1989) 159–180.
- [Seng90] B. Sengupta, A Queue with Service Interruptions in an Alternating Random Environment, *Operations Research*, 38 (1990) 308–318.
- [Seng92] B. Sengupta, An Approximation for the Sojourn-Time Distribution for the GI/G/1 Processor-Sharing Queue, *Stochastic Models*, 8 (1992) 35–57.
- [Shan66] D. N. Shanbhag, On Infinite Server Queues with Batch Arrivals, *Journal of Applied Probability*, 3 (1966) 274–279.
- [Smit72] W. Smith, The Infinitely-Many-Server Queue with Semi-Markovian Arrivals and Customer-Dependent Exponential Service Times, *Operations Research*, 22 (1972) 907–912.
- [Taka62] L. Takács, *Introduction to the Theory of Queues* (Oxford University Press, New York, 1962).
- [Taki94a] T. Takine and T. Hasegawa, The Workload in the MAP/G/1 Queue with State-Dependent Services: Its Application to a Queue with Preemptive Resume Priority, *Stochastic Models*, 10 (1994) 183–204.

- [Taki94b] T. Takine, Y. Matsumoto, T. Suda and T. Hasegawa, Mean Waiting Times in Nonpreemptive Priority Queues with Markovian Arrival and I.I.D. Service Processes, *Performance Evaluation*, 20 (1994) 131–149.
- [Taki95] T. Takine, B. Sengupta and R. W. Yeung, A Generalization of the Matrix M/G/1 Paradigm for Markov Chains with a Tree Structure, *Stochastic Models*, 11 (1995) 411–421.
- [Taki96] T. Takine, A Continuous Version of Matrix-Analytic Methods with the Skip-Free to the Left Property, *Stochastic Models*, 12 (1996) 673–682.
- [Taki97] T. Takine and B. Sengupta, A Single Server Queue with Service Interruptions, *Queueing Systems*, 26 (1997) 285–300.
- [Taki99] T. Takine, The Nonpreemptive Priority MAP/G/1 Queue, *Operations Research*, 47 (1999) 917–927.
- [Taki00] T. Takine, A New Recursion for the Queue Length Distribution in the Stationary BMAP/G/1 Queue, *Stochastic Models*, 16 (2000) 335–341.
- [Taki01a] T. Takine, Distributional Form of Little’s Law for FIFO Queues with Multiple Markovian Arrival Streams and Its Application to Queues with Vacations, *Queueing Systems*, 37 (2001) 31–63.
- [Taki01b] T. Takine, Queue Length Distribution in a FIFO Single-Server Queue with Multiple Arrival Streams Having Different Service Time Distributions, *Queueing Systems*, 39 (2001) 349–375.
- [Taki01c] T. Takine, A Recent Progress in Algorithmic Analysis of FIFO Queues with Markovian Arrival Streams, *Journal of the Korean Mathematical Society*, 38 (2001) 807–842.
- [Taki02] T. Takine, Matrix Product-Form Solution for an LCFS-PR Single-Server Queue with Multiple Arrival Streams Governed by a Markov Chain, *Queueing Systems*, 42 (2002) 131–151.
- [Tijm94] H. C. Tijms, *Stochastic Models: An Algorithmic Approach* (John Wiley & Sons, Chichester, 1994).
- [Yash83] S. F. Yashkov, A Derivation of Response Time Distribution for a M/G/1 Processor-Sharing Queue, *Problems of Control and Information Theory*, 12 (1983) 133–148.
- [Yeun94] R. W. Yeung and B. Sengupta, Matrix Product-Form Solutions for Markov Chains with a Tree Structure, *Advances in Applied Probability*, 26 (1994) 965–987.
- [Zhu91] Y. Zhu and N. U. Prabhu, Markov-Modulated PH/G/1 Queueing Systems, *Queueing Systems*, 9 (1991) 313–322.