

# Queues with correlated service and inter-arrival times and its application to optical buffers

J. Lambert, B. Van Houdt, C. Blondia

University of Antwerp, Department of Mathematics and Computer Science,  
Performance Analysis of Telecommunication Systems Research Group,  
Middelheimlaan, 1, B-2020 Antwerp - Belgium,  
{joke.lambert,benny.vanhoudt,chris.blondia}@ua.ac.be

August 24, 2005

## Abstract

This paper presents an algorithmic procedure to calculate the queue length and delay distribution of customers in a discrete time D-MAP/PH/1 queue, where the service time distribution of a customer depends on the inter-arrival time between himself and his predecessor. Setting up a Markov chain that keeps track of the contents of such a queue will result in a state space explosion as the inter-arrival times of all customers present in the system must be remembered. We avoid these difficulties by making use of the age process. This process keeps track of the “age” of the customer in the service facility. From this process, which we solve by means of matrix analytic methods, we compute the queue length and sojourn time distribution by means of a simple formula and obtain an expression for the stability of the system. We also demonstrate that the D-MAP arrival process can be easily replaced by the more general semi-Markovian arrival process, without any additional computational costs. Queueing systems of this type arise in the domain of synchronous optical buffers. Based on the numerical analysis of such a queueing system, some guidelines for the design of optical buffers are presented. We also show the impact on the numerical results when the cross-correlation that exists between the service and inter-arrival times is neglected.

**Index Terms:** Matrix analytic methods, D-MAP arrival process, semi-Markovian arrivals, phase-type services, correlated service and inter-arrival times, delay distribution, queue length distribution, optical buffer, fibre delay line.

## 1 Introduction

This paper introduces an algorithmic procedure to calculate the delay and the queue length distribution of a discrete-time first-come-first-serve queue with correlated arrivals (D-MAP and SM, see Sections 2 and 8), and phase-type (PH) service times that depend on the inter-arrival time between a customer and his predecessor. Moreover, the fact that a customer finds the server idle upon arrival may also affect his service time. Earlier work on queueing systems with correlated inter-arrival and service times can be found in [17, 6] and the references therein. The continuous time analog of the queueing system discussed in this paper nearly forms a subclass of the semi-Markovian queues introduced in [17] (see Remark 2, Section 3), where most of the results were derived from the theory developed in [16]. It is not a true subclass as in our case customers who experience no waiting time have different service requirements (therefore, we have a more general boundary condition). We restrict ourselves to this set of queues as it suffices for our main purpose, that is, the analysis of the synchronous optical buffer. Although there also exists an interest in analyzing asynchronous optical buffers, the framework in [17, 16] cannot be applied directly as PH distributions have only a finite (or infinite but countable) number of phases.

The basic model in [6] considers a system where the service time  $s_k$  of customer  $k$  is proportional to the inter-arrival time  $\tau_k$ ; thus, the service time of the  $k$ -th customer  $s_k = \xi\tau_k$  ( $\xi < 1$ ). This model was extended in two different ways: (i) by adding an independent, generally distributed, nonnegative random variable to the service time, (ii) by allowing the proportionality constant  $\xi$  to be itself a random variable, that takes a value  $\xi_1 > 0$  with probability  $g_1$  and  $\xi_2 = 0$  with probability  $1 - g_1$ . The arrival process in [6] is Poisson, meaning that  $\tau_k$  is exponentially distributed (actually, by setting  $\xi_2 = 0$  the arrival process behaves like an ON-OFF source). Hwang and Sohraby [9] also considered proportional service times with  $\xi = 1$  and a more general ON-OFF arrival process.

In this paper the service times do not need to be proportional to the inter-arrival time, but are PH distributed random variables characterized by  $(\alpha_{\tau_k}, T)$  (that is, the initial phase is determined by the inter-arrival time  $\tau_k$ ). In a synchronous optical buffer the service times are far from proportional; therefore, neither [6, 9] applies. Also, the arrival process considered in this paper is Markovian (D-MAP) and therefore has correlated inter-arrival times. Furthermore, we demonstrate that introducing semi-Markovian arrivals within this framework is straightforward. Queues with Markovian arrivals and correlated service and inter-arrival times have also been considered in [1]. In this paper recursive equations for the calculation of the moments of the waiting time and queue length are derived. Finally, in many studies the Laplace-Stieltjes transform (LST) of the waiting time and/or queue length distribution is determined, whereas we use matrix analytic methods that allow us to compute the distributions of interest directly (as in [17]).

Both [6, 9] were motivated by the correlation that the finite speed of communication links introduces between inter-arrival and service times. Our motivation lies in understanding the behavior of buffers present in optical communication networks. More specifically, as opposed to classical electronic buffers, optical buffers create voids on the outgoing channel. When the optical packets (called bursts) have a fixed size, one can show that the void between two packets is a function of the buffer granularity  $D$  and their inter-arrival time. By considering these voids as additional service, one obtains a queue with correlated service and inter-arrival times (see Sections 9-10).

The paper is structured as follows. Section 2 introduces the queueing system of interest.

The construction of the GI/M/1 type Markov chain is presented in Section 3. Section 4 is devoted to calculating the steady-state probabilities of such a system, while Section 5 focuses on the ergodicity of the GI/M/1 type Markov chain of interest. The delay density function is determined in Section 6, whereas in Section 7 we develop a simple formula to compute the queue length distribution. Some comments on how to incorporate semi-Markovian arrivals are provided in Section 8. Queueing systems of this type occur in the domain of optical buffers. Section 9 explains why a Fibre Delay Line - Non Void Filling system fits within this framework, while Section 10 contains some numerical examples. A comparison with a queueing model that neglects the cross-correlation between the inter-arrival and service times is also included.

## 2 The D-MAP/PH/1 queue with correlated service and inter-arrival times

The arrival process of the queueing system of interest is a discrete time Markov arrival process, commonly known as the D-MAP process [3, 4], that does not allow batch arrivals; therefore, it is a subclass of the D-BMAP arrival process, which allows batch arrivals. Formally, a D-MAP is characterized — similar to its continuous time variant the MAP process [12] — by two  $m \times m$  matrices  $\mathbf{D}_0$  and  $\mathbf{D}_1$ , where  $m$  is a positive integer. The  $(j, j')^{th}$  entry of the matrix  $\mathbf{D}_1$  represents the probability that a customer arrives and the underlying Markov chain makes a transition from state  $j$  to state  $j'$ . The matrix  $\mathbf{D}_0$  covers the case when there is no arrival. The matrix  $\mathbf{D}$ , defined as

$$\mathbf{D} = \mathbf{D}_0 + \mathbf{D}_1,$$

represents the stochastic  $m \times m$  transition matrix of the underlying Markov chain of the arrival process. Let  $\boldsymbol{\theta}$  be the stationary probability vector of  $\mathbf{D}$ , that is,  $\boldsymbol{\theta}\mathbf{D} = \boldsymbol{\theta}$  and  $\boldsymbol{\theta}\mathbf{e} = 1$ , where  $\mathbf{e}$  is a column vector with all entries equal to one. The stationary arrival rate is given by  $\lambda = \boldsymbol{\theta}\mathbf{D}_1\mathbf{e}$ . It is easy to see that

$$\boldsymbol{\theta} = \boldsymbol{\theta}\mathbf{D}_1(\mathbf{I} - \mathbf{D}_0)^{-1}, \tag{1}$$

where the inverse of the matrix  $\mathbf{I} - \mathbf{D}_0$  exists as  $\mathbf{D}_0$  is substochastic.

The service time of a customer depends upon the inter-arrival time between himself and the previous customer. If the server is empty when a new customer, with inter-arrival time  $k$  time units, arrives, the service time of this customer has a common phase-type distribution function [14] with a matrix representation  $(m_{ser}, \boldsymbol{\alpha}_{k,0}, \mathbf{T})$ , where  $m_{ser}$  is a positive integer,  $\boldsymbol{\alpha}_{k,0}$  is an  $1 \times m_{ser}$  nonnegative stochastic vector and  $\mathbf{T}$  is an  $m_{ser} \times m_{ser}$  substochastic matrix. On the other hand, if the server is busy when a new customer, with inter-arrival time  $k$ , arrives, the service process can be described by a phase-type distribution with matrix representation  $(m_{ser}, \boldsymbol{\alpha}_k, \mathbf{T})$ , where  $\boldsymbol{\alpha}_k$  is an  $1 \times m_{ser}$  nonnegative stochastic vector. The  $i^{th}$  component of the vector  $\boldsymbol{\alpha}_{k,0}$  is the probability that a customer, with an inter-arrival time of  $k$  time units, starts his service in phase  $i$  given that the server was empty when this customer arrived. Analogous we can say that the  $i^{th}$  component of the vector  $\boldsymbol{\alpha}_k$  is the probability that such a customer starts his service in phase  $i$  given that he saw the server busy upon arrival. If  $\mathbf{t} = \mathbf{e} - \mathbf{T}\mathbf{e}$ , then the  $i^{th}$  entry of the vector  $\mathbf{t}$  denotes the probability that a customer completes his service provided that he is in the  $i^{th}$  phase at the current time instant. The  $(i_1, i_2)^{th}$  entry of  $\mathbf{T}$ , on the other hand, is the probability that a customer continues his service in phase  $i_2$  at the next time instant provided that he is in phase  $i_1$  at the current time

instant. Notice, the minimum service time of a customer is one time unit (because the vectors  $\boldsymbol{\alpha}_k$  and  $\boldsymbol{\alpha}_{k,0}$  are assumed to be stochastic). The mean service time of a customer with an inter-arrival time of  $k$  time units given that the server was busy when he arrived, is given by  $1/\mu_k = \boldsymbol{\alpha}_k(\mathbf{I} - \mathbf{T})^{-1}\mathbf{e}$ . Thus, assuming that all customers see the server busy upon arrival, the mean service time of a customer equals

$$E[S_b] = \sum_k \omega_k / \mu_k, \text{ with } \omega_k = \frac{\boldsymbol{\theta} \mathbf{D}_1 \mathbf{D}_0^{k-1} \mathbf{D}_1}{\lambda} \mathbf{e}. \quad (2)$$

This expression is used to formulate the stability condition.

Finally, in the case of a simultaneous arrival and departure we assume that the departure occurs first. For further use we define  $m_{tot} = m_{ser}m$ . While constructing the GI/M/1 type Markov chain in the next section, we will always observe the system just prior to possible phase changes, arrivals or departures. Thus, if we refer to the system state at time  $n$ , such events happening at time  $n$  are not yet taken into account by the system state.

**Remark 1:** The system described above is, from a theoretical point of view, equivalent to the following system. A new customer, with inter-arrival time  $k$ , finding the server busy (resp. idle) has a service process that can be described by a phase-type distribution with matrix representation  $(m_k, \boldsymbol{\alpha}'_k, \mathbf{T}_k)$  (resp.  $(m_k, \boldsymbol{\alpha}'_{k,0}, \mathbf{T}_k)$ ). By defining the phase-type distribution function with matrix representation  $(m_{ser}, \boldsymbol{\alpha}_k, \mathbf{T})$ , with  $m_{ser} = \sum_{k=1}^{\infty} m_k$ ,

$$\boldsymbol{\alpha}_k = (\mathbf{0}_{k,f}, \boldsymbol{\alpha}'_k, 0, 0, \dots),$$

$$\mathbf{T} = \begin{bmatrix} \mathbf{T}_1 & \mathbf{0} & \mathbf{0} & \dots \\ \mathbf{0} & \mathbf{T}_2 & \mathbf{0} & \ddots \\ \mathbf{0} & \mathbf{0} & \mathbf{T}_3 & \ddots \\ \vdots & \ddots & \ddots & \ddots \end{bmatrix},$$

where  $\mathbf{0}_{k,f}$  is a  $1 \times \sum_{i=1}^{k-1} m_i$  vector filled with zeros, we see that both cases are equivalent. An analogous set of vectors  $\boldsymbol{\alpha}_{k,0}$  for the case when the server is idle upon arrival is found easily. Notice, for practical purposes the dimension of the matrix  $\mathbf{T}$  has to be finite. Often the matrix  $\mathbf{T}$  can easily be reduced to a finite, moderate size matrix (as many matrix representations of the same PH distribution function exist). For instance, when all the service time distributions have a general distribution with a finite support  $a$ , the dimension of  $\mathbf{T}$  should not exceed  $a$ . The study of the synchronous optical buffer presented in Sections 9 and 10 is of this type. Also, if some of the matrices  $\mathbf{T}_{k'}$  are identical to  $\mathbf{T}_k$ , for  $k \neq k'$ , it suffices to have a single copy of  $\mathbf{T}_k$  in  $\mathbf{T}$ . E.g., in a system where a customer who arrives more than 100 time units after his predecessor, needs some additional service, the matrix  $\mathbf{T}$  contains only two blocks.

### 3 Constructing the GI/M/1 Type Markov chain

In this section, we consider the D-MAP/PH/1 queue with service times depending on the inter-arrival times. Instead of observing the system at each time instant, including the instants when the server is idle, we create a GI/M/1 type Markov chain (MC) by observing the system state only when the server is occupied.

Consider an MC with an infinite number of states labeled  $1, 2, \dots$ . The set of states  $\{(i-1)m_{tot} + 1, \dots, im_{tot}\}$  is referred to as level  $i$  of the MC, for  $i > 0$ . The states of each

level are labeled as  $(s, j)$ , where  $1 \leq s \leq m_{ser}$  and  $1 \leq j \leq m$ . Let state  $(s, j)$  of level  $i$  of the MC correspond to the situation in which there is a customer in service, who arrived  $i$  time units ago, while the service process is currently in phase  $s$  and the D-MAP arrival process was in state  $j$  at time  $n - i + 1$ , where  $n$  is the current time instant. Recall, we observe the system just prior to possible phase changes, arrivals or departures.

The level of the Markov chain can never increase by more than one during a transition between time instant  $n$  and the next time instant  $n + x$  where the server is busy (because the customers are served in a FCFS order). As a result, the system can be described by a transition matrix  $\mathbf{P}$  with the following structure:

$$\mathbf{P} = \begin{bmatrix} \mathbf{C}_0 & \mathbf{A}_0 & \mathbf{0} & \mathbf{0} & \mathbf{0} & \dots \\ \mathbf{C}_1 & \mathbf{A}_1 & \mathbf{A}_0 & \mathbf{0} & \mathbf{0} & \dots \\ \mathbf{C}_2 & \mathbf{A}_2 & \mathbf{A}_1 & \mathbf{A}_0 & \mathbf{0} & \dots \\ \mathbf{C}_3 & \mathbf{A}_3 & \mathbf{A}_2 & \mathbf{A}_1 & \mathbf{A}_0 & \dots \\ \vdots & \vdots & \vdots & \ddots & \ddots & \ddots \end{bmatrix}, \quad (3)$$

where  $\mathbf{A}_i$  and  $\mathbf{C}_i$  are  $m_{tot} \times m_{tot}$  matrices.

Next, let us derive an expression for each of the matrices  $\mathbf{A}_i$  and  $\mathbf{C}_i$ . Suppose that the MC is in state  $(s, j)$  of level  $i$  at time  $n$ . Then, a transition to state  $(s', j)$  of level  $i + 1$  occurs if the customer remains in the service facility (with probability  $(\mathbf{T})_{s,s'}$ ). Notice, in this case the state of the D-MAP remains the same, therefore,

$$\mathbf{A}_0 = \mathbf{T} \otimes \mathbf{I}_m, \quad (4)$$

where  $\otimes$  denotes the Kronecker product between matrices and  $\mathbf{I}_m$  denotes the  $m \times m$  unity matrix. A transition from level  $i$  to  $i + 1$  cannot occur if there is a service completion at time  $n$ , because there are no batch arrivals (implying that the next customer has an age of at most  $i$  at time  $n + 1$ ).

A transition to level  $i - l$ , with  $0 \leq l < i - 1$ , occurs if the customer in service completes his service (with probability  $\mathbf{t}_s$ ) and there is no arrival until time  $n + 1 - (i - l)$ , that is, there is no arrival at time  $n - i + 1, \dots, n - i + l$  and at time  $n + 1 - i + l$  we have an arrival<sup>1</sup>. The inter-arrival time between the customer completing service and the next one is  $l + 1$ . Because the new customer saw the server busy upon arrival, his service requirements are described by  $(m_{ser}, \boldsymbol{\alpha}_{l+1}, \mathbf{T})$ . Transitions from level  $i$  to  $i - l$  are governed by the  $\mathbf{A}_{l+1}$ ; hence, for  $i > 1$  and  $0 \leq l < i - 1$ ,

$$\mathbf{A}_{l+1} = \mathbf{t}\boldsymbol{\alpha}_{l+1} \otimes (\mathbf{D}_0^l \mathbf{D}_1), \quad (5)$$

for hereon, any matrix to the power 0 is taken to be the identity matrix of appropriate dimension. Finally, a transition to level 1 occurs if the customer in service completes his service (with probability  $\mathbf{t}_s$ ) and a new customer, who arrived at or after time  $n$ , starts his service. Such a customer sees the server empty upon arrival, so his service time is determined by  $(m_{ser}, \boldsymbol{\alpha}_{k,0}, \mathbf{T})$ . Because we observe the system only at time instants when the server is busy, the inter-arrival time  $k$  can take on any value larger than or equal to  $i$ . Hence,

$$\mathbf{C}_{i-1} = \sum_{k \geq i-1} \mathbf{t}\boldsymbol{\alpha}_{k+1,0} \otimes (\mathbf{D}_0^k \mathbf{D}_1). \quad (6)$$

This concludes the description of the transition matrices.

<sup>1</sup>Indeed, we are at level  $i - l$  at time  $n + 1$  if the customer in service has an age  $i - l$  at time  $n + 1$ , thus, he arrived at time  $n + 1 - i + l$ .

**Remark 2:** As indicated in the introduction, if  $\alpha_{k,0} = \alpha_k$  for all  $k$ , this system is the discrete time analog of (a special case) of the semi-Markovian queue considered in [17]. This can be seen as follows.  $\mathbf{A}_0$  plays the same role as  $\mathbf{S}$  in [17], the difference being that  $\mathbf{A}_0$  corresponds to the matrix of a discrete PH variable. Notice, by definition of  $\mathbf{A}_0$ , the state of the D-MAP is also part of the phase of  $\mathbf{A}_0$ ; therefore, we refer to these phases as *meta* phases  $(j, v)$  where  $j$  is the phase of the customer in the service facility and  $v$  the state of the D-MAP immediately after his arrival. Let  $P(t)_{(j,v),(j',v')}$  be the probability that the inter-arrival time between customer  $k$  and  $k + 1$  is less than or equal to  $t$  time units and that customer  $k + 1$  starts service in meta phase  $(j', v')$ , provided that customer  $k$  ended this service in meta phase  $(j, v)$ . Then, the matrix  $P(t)$  holding these probabilities as its entries can be written as:

$$\mathbf{P}(t) = \sum_{i=1}^t \mathbf{e} \alpha_i \otimes (\mathbf{D}_0^{i-1} \mathbf{D}_1). \quad (7)$$

The matrix  $P(t)$  now plays the same role as in [17]. The same remark can be made with respect to the system with semi-Markovian arrivals considered in Section 8 (if we replace  $\mathbf{D}_0^{i-1} \mathbf{D}_1$  by  $\mathbf{L}_i$ ).

## 4 Calculating the Steady-State Probabilities

In this section, we consider the D-MAP/PH/1 queue with service times depending on the inter-arrival times and indicate how to calculate the steady-state probabilities of the Markov chain characterized by the transition matrix  $\mathbf{P}$ . From equation (3), we see that the Markov chain is a Markov chain of the GI/M/1 type [14]. Define the  $1 \times m_{tot}$  vectors  $\boldsymbol{\pi}_i = (\pi_i(1, 1), \pi_i(1, 2), \dots, \pi_i(1, m), \pi_i(2, 1), \dots, \pi_i(m_{ser}, m))$ , for  $i \geq 1$ . The steady state vector  $\boldsymbol{\pi} = (\boldsymbol{\pi}_1, \boldsymbol{\pi}_2, \dots)$  of  $\mathbf{P}$  exists if and only if  $\rho_b < 1$ , where  $\rho_b = \lambda \sum_k \omega_k / \mu_k = \lambda E[S_b]$  and  $\lambda, \omega_k$  and  $\mu_k$  were defined in Section 2. A proof of this theorem based on Neuts' stability condition [13] is provided in Section 5. Moreover,  $\boldsymbol{\pi}_i = \boldsymbol{\pi}_{i-1} \mathbf{R}$ , for  $i > 1$ , where  $\mathbf{R}$  is an  $m_{tot} \times m_{tot}$  matrix that is the smallest nonnegative solution of the following equation:

$$\mathbf{R} = \sum_{i=0}^{\infty} \mathbf{R}^i \mathbf{A}_i. \quad (8)$$

There are several techniques to calculate  $\mathbf{R}$ , a brief overview is given in the Appendix. In order to obtain  $\boldsymbol{\pi}_1$  we solve the boundary equation:

$$\boldsymbol{\pi}_1 = \boldsymbol{\pi}_1 \sum_{i=0}^{\infty} \mathbf{R}^i \mathbf{C}_i, \quad (9)$$

where the vector  $\boldsymbol{\pi}_1$  is normalized as

$$\boldsymbol{\pi}_1 (\mathbf{I} - \mathbf{R})^{-1} \mathbf{e} = 1. \quad (10)$$

## 5 Algebraic proof for the ergodicity of a GI/M/1 type Markov chain

In this section we will prove that the GI/M/1 type Markov chain introduced in Section 3 is ergodic if and only if  $\rho_b = \lambda E[S_b] = \lambda \sum_k \frac{\omega_k}{\mu_k} < 1$ . We start by defining the  $1 \times m_{ser}$  stochastic

vectors  $\beta_k$  for  $k \geq 1$ :

$$\beta_k = \beta_k(\mathbf{T} + \mathbf{t}\alpha_k). \quad (11)$$

Some elementary manipulations shown that they obey the following equality:

$$\beta_k \mathbf{t} = \mu_k. \quad (12)$$

**Lemma 1** *The vector  $\Pi_g$ , defined as*

$$\frac{1}{\rho_b} \sum_{k=1}^{\infty} \left( \frac{\beta_k}{\mu_k} \otimes (\theta \mathbf{D}_1 \mathbf{D}_0^{k-1} \mathbf{D}_1) \right), \quad (13)$$

*is an invariant vector of  $\sum_{l=0}^{\infty} \mathbf{A}_l$  and  $\Pi_g$  is stochastic.*

*Proof.* The sum  $\sum_{l=0}^{\infty} \mathbf{A}_l$  can be written as  $\mathbf{T} \otimes \mathbf{I}_m + \sum_{l=1}^{\infty} \mathbf{t}\alpha_l \otimes \mathbf{D}_0^{l-1} \mathbf{D}_1$ . The product  $\Pi_g(\mathbf{T} \otimes \mathbf{I}_m)$  equals

$$\frac{1}{\rho_b} \sum_{k=1}^{\infty} \left( \frac{\beta_k}{\mu_k} \mathbf{T} \otimes (\theta \mathbf{D}_1 \mathbf{D}_0^{k-1} \mathbf{D}_1) \right) = \frac{1}{\rho_b} \sum_{k=1}^{\infty} \left( \frac{\beta_k - \mu_k \alpha_k}{\mu_k} \otimes (\theta \mathbf{D}_1 \mathbf{D}_0^{k-1} \mathbf{D}_1) \right),$$

by means of Eqns. (11) and (12). Second, we calculate  $\Pi_g(\mathbf{t}\alpha_l \otimes \mathbf{D}_0^{l-1} \mathbf{D}_1)$ . This equals

$$\frac{1}{\rho_b} \sum_{k=1}^{\infty} \left( \frac{\beta_k}{\mu_k} \mathbf{t}\alpha_l \otimes (\theta \mathbf{D}_1 \mathbf{D}_0^{k-1} \mathbf{D}_1)(\mathbf{D}_0^{l-1} \mathbf{D}_1) \right).$$

Using Eqns. (12) and (1) we can further simplify this to

$$\begin{aligned} & \frac{1}{\rho_b} \sum_{k=1}^{\infty} \left( \alpha_l \otimes (\theta \mathbf{D}_1 \mathbf{D}_0^{k-1} \mathbf{D}_1)(\mathbf{D}_0^{l-1} \mathbf{D}_1) \right) \\ &= \frac{1}{\rho_b} \left( \alpha_l \otimes (\theta \mathbf{D}_1 (\mathbf{I} - \mathbf{D}_0)^{-1} \mathbf{D}_1)(\mathbf{D}_0^{l-1} \mathbf{D}_1) \right) \\ &= \frac{1}{\rho_b} \alpha_l \otimes (\theta \mathbf{D}_1 \mathbf{D}_0^{l-1} \mathbf{D}_1). \end{aligned} \quad (14)$$

From hereon it is straightforward to prove that  $\Pi_g$  is an invariant vector of  $\sum_{l=0}^{\infty} \mathbf{A}_l$ :

$$\begin{aligned} \Pi_g \sum_{l=0}^{\infty} \mathbf{A}_l &= \frac{1}{\rho_b} \sum_{k=1}^{\infty} \left( \frac{\beta_k - \mu_k \alpha_k}{\mu_k} \otimes (\theta \mathbf{D}_1 \mathbf{D}_0^{k-1} \mathbf{D}_1) \right) + \sum_{l=1}^{\infty} \frac{1}{\rho_b} \alpha_l \otimes (\theta \mathbf{D}_1 \mathbf{D}_0^{l-1} \mathbf{D}_1) \\ &= \frac{1}{\rho_b} \sum_{k=1}^{\infty} \left( \frac{\beta_k}{\mu_k} \otimes (\theta \mathbf{D}_1 \mathbf{D}_0^{k-1} \mathbf{D}_1) \right) \\ &= \Pi_g. \end{aligned}$$

$\Pi_g$  is clearly a stochastic vector as  $\beta_k$  is stochastic and  $\theta \mathbf{D}_1 \mathbf{D}_0^{k-1} \mathbf{D}_1 \mathbf{e}$  equals  $\lambda \omega_k$ .  $\square$

The vector  $\Pi_g$  holds the probabilities that, when observing the system at an arbitrary time instant, the phase of the server equals  $s$  (for  $s = 1, \dots, m_{ser}$ ) and the state of the D-MAP was  $j$  (for  $j = 1, \dots, m$ ) immediately after the arrival of the customer occupying the service facility, provided that the server is busy during this time instant and under the assumption that  $\alpha_{k,0} = \alpha_k$  for all  $k$ .

**Lemma 2**  $\Pi_g \left( \sum_{l=1}^{\infty} l \mathbf{A}_l \right) \mathbf{e} = 1/\rho_b$ .

*Proof.* We can rewrite  $\left( \sum_{l=1}^{\infty} l \mathbf{A}_l \right) \mathbf{e}$  using the stochastic nature of the vectors  $\alpha_k$  and the equality  $\mathbf{D}_1 \mathbf{e} = (\mathbf{I} - \mathbf{D}_0) \mathbf{e}$  as

$$\begin{aligned} \left( \sum_{l=1}^{\infty} l \mathbf{A}_l \right) \mathbf{e} &= \sum_{l=1}^{\infty} l (\mathbf{t} \alpha_l \otimes \mathbf{D}_0^{l-1} \mathbf{D}_1) \mathbf{e} = \sum_{l=1}^{\infty} \left( \mathbf{t} \alpha_l \mathbf{e}_{m_{ser}} \otimes l \mathbf{D}_0^{l-1} \mathbf{D}_1 \mathbf{e}_m \right) \\ &= \mathbf{t} \otimes \sum_{l=1}^{\infty} l \mathbf{D}_0^{l-1} \mathbf{D}_1 \mathbf{e}_m = \mathbf{t} \otimes (\mathbf{I} - \mathbf{D}_0)^{-1} \mathbf{e}_m. \end{aligned} \quad (15)$$

Thus, applying Eqns. (12) and (1), we have

$$\begin{aligned} \Pi_g \left( \sum_{l=1}^{\infty} l \mathbf{A}_l \right) \mathbf{e} &= \frac{1}{\rho_b} \sum_{k=1}^{\infty} \left( \frac{\beta_k}{\mu_k} \mathbf{t} \otimes (\boldsymbol{\theta} \mathbf{D}_1 \mathbf{D}_0^{k-1} \mathbf{D}_1) ((\mathbf{I} - \mathbf{D}_0)^{-1} \mathbf{e}_m) \right) \\ &= \frac{1}{\rho_b} \sum_{k=1}^{\infty} \left( \boldsymbol{\theta} \mathbf{D}_1 \mathbf{D}_0^{k-1} \mathbf{D}_1 ((\mathbf{I} - \mathbf{D}_0)^{-1} \mathbf{e}_m) \right) \\ &= \frac{1}{\rho_b} (\boldsymbol{\theta} \mathbf{D}_1 (\mathbf{I} - \mathbf{D}_0)^{-1} \mathbf{D}_1) ((\mathbf{I} - \mathbf{D}_0)^{-1} \mathbf{e}_m) = \frac{1}{\rho_b}. \end{aligned} \quad (16)$$

□

**Theorem 1** *The Markov chain characterized by the transition matrix  $P$  (see Eqn. (3)) is ergodic if and only if  $\rho_b < 1$ .*

*Proof.* Neuts [13] has shown that a GI/M/1 type MC is ergodic if and only if the product of the stochastic invariant vector of  $\sum_{l=0}^{\infty} \mathbf{A}_l$  with the vector  $\left( \sum_{l=0}^{\infty} l \mathbf{A}_l \right) \mathbf{e}$  is larger than one. Therefore, Lemma 1 and Lemma 2 suffice to proof the theorem. □

## 6 Calculating the Delay Density Function

As opposed to the general approach in many queueing systems, we calculate the delay distribution without obtaining the steady state probabilities of the queue length. The term delay is used here as a synonym for the sojourn time, thus it encompasses both the waiting and the service time. Let  $X$  be the random variable that denotes the delay suffered by a customer.

**Theorem 2** *If  $\boldsymbol{\pi}$ , the steady state vector of the Markov chain characterized by  $\mathbf{P}$ , exists, the delay distribution  $X$  of an arbitrary customer is given by*

$$P[X = i] = \frac{c}{\lambda} \sum_{s=1}^{m_{ser}} \left( \sum_{j=1}^m \boldsymbol{\pi}_i(s, j) \right) (\mathbf{t})_s. \quad (17)$$

where  $c$  is a normalization constant that equals the probability that the system is busy. If the vectors  $\alpha_{k,0}$  equal  $\alpha_k$  for all  $k$ , then  $c = \lambda E[S_b]$ .

*Proof.* The probability that a customer has a delay of  $i$  time units, denoted by  $P[X = i]$ , equals the expected number of customers with an “age” of  $i$  time units that complete their service at an arbitrary time instant divided by the expected number of customers that complete their service during an arbitrary time instance — that is,  $\lambda$  for a stable queue. This explains



the expression given by (17). If the vectors  $\alpha_{k,0}$  equal  $\alpha_k$  for every  $k$ , we use the identity  $\sum_i \pi_i = \mathbf{\Pi}_g$  as follows, where  $\mathbf{\Pi}_g$  is the invariant vector of  $\sum_{l=0}^{\infty} \mathbf{A}_l$ :

$$\begin{aligned}
1 &= \sum_i P[X = i] = \frac{c}{\lambda} \sum_i \sum_{s=1}^{m_{ser}} \left( \sum_{j=1}^m \pi_{i(s,j)} \right) \mathbf{t}_s = \frac{c}{\lambda} \sum_{s=1}^{m_{ser}} \left( \sum_{j=1}^m \mathbf{\Pi}_g(s,j) \right) \mathbf{t}_s \\
&= \frac{c}{\lambda} \sum_{s=1}^{m_{ser}} \left( \sum_{j=1}^m \frac{1}{\rho_b} \sum_k (\boldsymbol{\theta} \mathbf{D}_1 \mathbf{D}_0^{k-1} \mathbf{D}_1)_j \frac{(\boldsymbol{\beta}_k)_s}{\mu_k} \right) \mathbf{t}_s \\
&= \frac{c}{\lambda} \sum_{s=1}^{m_{ser}} \left( \sum_k \frac{1}{E[S_b]} \omega_k \frac{(\boldsymbol{\beta}_k)_s}{\mu_k} \right) \mathbf{t}_s = \frac{c}{\lambda E[S_b]} \sum_k \omega_k = \frac{c}{\lambda E[S_b]}.
\end{aligned}$$

□

**Theorem 3** *If the queue is stable and  $\alpha_{k,0} = \alpha_k$  for all  $k$ , the delay distribution is a phase-type distribution with representation  $(\gamma, \mathbf{Q})$ , where  $\gamma = E[S_b] \mathbf{\Pi}_g (\mathbf{I} - \mathbf{R}) \boldsymbol{\Delta}$ ,  $\mathbf{Q} = \boldsymbol{\Delta}^{-1} \mathbf{R} \boldsymbol{\Delta}$ ,  $\boldsymbol{\Delta} = \text{diag}(\boldsymbol{\delta})$  and  $\boldsymbol{\delta} = (\mathbf{I} - \mathbf{R})^{-1} (\mathbf{t} \otimes \mathbf{e})$ .*

*Proof.* Using the equality  $\sum_i \pi_i = \mathbf{\Pi}_g$ , we find  $\pi_1 = \mathbf{\Pi}_g (\mathbf{I} - \mathbf{R})$ . The remainder of the proof follows from Theorem 2 and is similar to [7, Theorem 4.4]. □

## 7 Calculating the Queue Length Distribution

Let  $Q$  be the random variable that denotes the queue length, i.e., the number of customers present in the waiting room. Define the  $m \times 1$  vectors  $\mathbf{d}_{q,i}$ , where the  $j^{\text{th}}$  component of the vector  $\mathbf{d}_{q,i}$  is the probability that  $q$  arrivals occur in an interval of length  $i$  that started in state  $j$ . These vectors can be computed by means of the following recursion:

$$\begin{aligned}
\mathbf{d}_{q,0} &= 1[q = 0] \mathbf{e}, \\
\mathbf{d}_{0,i} &= (\mathbf{D}_0)^i \mathbf{e}, \\
\mathbf{d}_{q,i} &= \mathbf{D}_1 \mathbf{d}_{q-1,i-1} + \mathbf{D}_0 \mathbf{d}_{q,i-1}.
\end{aligned}$$

By means of these vectors we can compute the queue length distribution as follows:

$$P[Q = q] = c \sum_{i \geq q} \left( \sum_{s=1}^{m_{ser}} \sum_{j=1}^m \pi_{i(s,j)} (\mathbf{d}_{q,i})_j \right), \quad (18)$$

for  $q > 0$  and

$$P[Q = 0] = 1 - c + c \sum_{i \geq 1} \left( \sum_{s=1}^{m_{ser}} \sum_{j=1}^m \pi_{i(s,j)} (\mathbf{d}_{0,i})_j \right), \quad (19)$$

with  $c$  equaling the probability that the server is busy, see Section 6. Recall, in case the vectors  $\alpha_{k,0}$  equal  $\alpha_k$  for every  $k$ , the factor  $c$  equals  $\rho_b$ .

## 8 Semi-Markovian Arrivals

Consider a semi-Markov chain  $(\xi_n, \tau_n)$  with  $m$  phases. The variable  $\xi_n$  represents the phase of the semi-Markov chain immediately after the  $n$ -th transition. The variable  $\tau_n$  denotes the number of time slots between the  $(n-1)$ -th and the  $n$ -th transition (inter-transition time). The semi-Markovian (SM) arrival process is constructed as follows. Define

$$P[\xi_n = j', \tau_n = t | \xi_{n-1} = j] = l_{j,j'}(t), \quad (20)$$

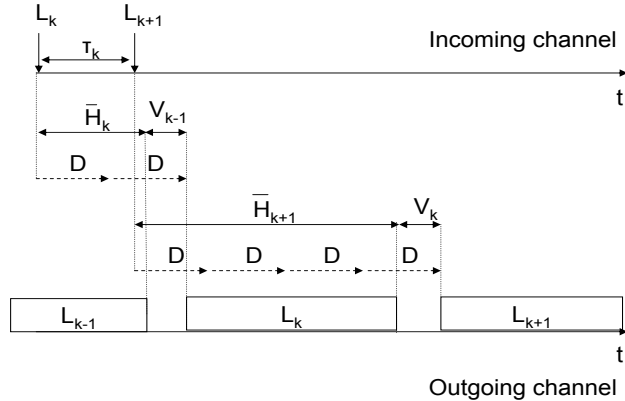
for  $j, j' \in \{1, \dots, m\}$  and  $t \geq 1$ . The variable  $l_{j,j'}(t)$  equals the probability that a customer arrives  $t$  time units after his predecessor, while the underlying semi-Markov chain makes a transition from state  $j$  to  $j'$ . Let  $\mathbf{L}_t$  be the  $m \times m$  matrix whose  $(j, j')$ -th element equals  $l_{j,j'}(t)$ . The matrix  $\mathbf{L} = \sum_t \mathbf{L}_t$  is the transition matrix of the underlying Markov chain. Let  $\boldsymbol{\xi}$  be its invariant probability vector, then the mean inter-arrival time  $1/\lambda = \boldsymbol{\xi} \sum_t t \mathbf{L}_t \mathbf{e}$ . In the special case where two  $m$  dimensional square matrices exist  $\mathbf{D}_0, \mathbf{D}_1 \geq 0$  such that  $\mathbf{L}_k = \mathbf{D}_0^{k-1} \mathbf{D}_1$ , the SM arrival process is a D-MAP. Notice, in this particular case  $\mathbf{L} = (\mathbf{I} - \mathbf{D}_0)^{-1} \mathbf{D}_1$ ; therefore,  $\boldsymbol{\xi} = \boldsymbol{\theta} \mathbf{D}_1$ , where  $\boldsymbol{\theta}$  is the invariant probability vector of  $\mathbf{D} = \mathbf{D}_0 + \mathbf{D}_1$ .

In order to generalize the results presented in Sections 3 and 4, one simply needs to replace the products  $\mathbf{D}_0^k \mathbf{D}_1$  by  $\mathbf{L}_{k+1}$  in Eqns. (5) and (6). This provides us with the probabilities  $\pi_i(s, j)$  that the server holds an age  $i$  customer at an arbitrary point in time, the phase of the customer in service equals  $s$  and the state of the SM arrival process immediately after his arrival is  $j$ . The delay distribution  $X$  can be computed from these vectors by means of Eqn. (17), the constant  $c$  is found by normalizing  $X$ . However, for the queue length distribution we need to make some minor changes to Eqns. (18) and (19). Define the following recursion:

$$\begin{aligned} \bar{\mathbf{F}}_{q,0} &= 1[q=0] \mathbf{I}, \\ \bar{\mathbf{F}}_{q,i} &= \sum_{i'=q-1}^{i-1} \bar{\mathbf{F}}_{q-1,i'} \mathbf{L}_{i-i'}, \\ \bar{\mathbf{d}}_{q,0} &= 1[q=0] \mathbf{e}, \\ \bar{\mathbf{d}}_{0,i} &= \sum_{i'>i} \mathbf{L}_{i'} \mathbf{e}, \\ \bar{\mathbf{d}}_{q,i} &= \sum_{i'=q}^i \bar{\mathbf{F}}_{q,i'} \bar{\mathbf{d}}_{0,i-i'} \mathbf{e}. \end{aligned}$$

for  $i > 0$ . Notice, the  $(j, j')$ -th element of  $\bar{\mathbf{F}}_{q,i}$  represents the probability of having  $q$  arrivals in an interval of length  $i$ , where the last arrival occurred at the end of the interval, while the state of the semi-Markov chain at the start, resp. end, is  $j$ , resp.  $j'$ . The  $j$ -th element of  $\bar{\mathbf{d}}_{q,i}$  equals the probability that  $q$  arrivals occur in an interval of length  $i$  that started in state  $j$ . Replacing the vectors  $\mathbf{d}_{q,i}$  by  $\bar{\mathbf{d}}_{q,i}$  in Eqns. (18) and (19) gives us the queue length distribution  $Q$ .

**Remark 3:** Assume that  $\boldsymbol{\alpha}_{k,0} = \boldsymbol{\alpha}_k$  for all  $k$ . The class of queues obtained by further restricting ourselves to the  $GI$  arrival process (i.e., the case where the matrices  $\mathbf{L}_i$  are scalars), is the discrete time analog of the subclass of the semi-Markovian queues analyzed in [17] for which the inter-arrival time and initial phase of customer  $k+1$  are unaffected by the end phase of customer  $k$ . In this case,  $\mathbf{T}$  plays the role of  $\mathbf{S}$ , while  $\mathbf{J}(t) = \mathbf{e} \sum_{i=1}^t \boldsymbol{\alpha}_i \mathbf{L}_i$ .



**Figure 1:** Evolution of the scheduling horizon  $\bar{H}$  from one arrival to the next.  $L_k$  is the length of the  $k$ -th OB and  $\tau_k$  the burst inter-arrival time

## 9 Fibre Delay Lines

In this section, we study a single Wavelength Division Multiplexing (WDM) channel and assume contention for it is resolved by means of a Fibre Delay Line (FDL) buffer, which can delay, if necessary, data packets, called optical bursts (OBs), until the channel becomes available again. Unlike conventional electronic buffers, however, it cannot delay bursts for an arbitrary period of time, but only for multiples of a basic unit  $D$ , called the granularity of the FDL buffer. In [10, 5] it is shown that this leads to voids on the outgoing channel. We do not attempt to fill these voids (as this would alter the order of the bursts), hence, the term Non Void Filling system.

Define the scheduling horizon as the earliest time by which all previously arrived OBs will have left the system and denote it by  $\bar{H}$ . When the  $k$ -th burst sees a scheduling horizon  $\bar{H}_k$  upon arrival, it will have to be delayed by at least  $\bar{H}_k$ . An FDL buffer can delay bursts only for multiples of  $D$ , so  $\bar{H}_k$  cannot always be realized exactly and this leads to voids on the outgoing channel, as illustrated in Figure 1. In this figure the horizon seen by the  $k$ -th burst lies somewhere between  $D$  and  $2D$ , meaning that the  $k$ -th burst is delayed for  $2D$  time units. Similarly, the  $k + 1$ -th burst is delayed by  $4D$ . It is not difficult to show (see [8], lemma 1) that the length of a void between OB  $k$  and  $k + 1$  is distributed according to  $(\tau_k - L_k) \bmod D$ , where  $\tau_k$  represents the inter-arrival time between these two OBs and  $L_k$  is the length of the  $k$ -th OB. Unless an OB burst finds the FDL buffer empty upon arrival, this void can be regarded as additional service required for OB  $k + 1$ . As such, the system evolves like of a FIFO queue with the same D-MAP arrival process, but with a modified service time for customers who find the server busy:  $L'_k = L_k + V_k$ , where  $V_k$  is the void between the  $k$  and  $k + 1$ -th OB. Given that all bursts have the same length  $L$ , the void  $V_k$  depends solely on the inter-arrival time and the constants  $L$  and  $D$ . Thus, the algorithms developed in Sections 4, 6 and 7 can be used to analyze FDL buffers.

## 10 Numerical example

In this section we present some numerical examples of an FDL buffer with Markovian arrivals. The results are compared with those obtained by the classic D-MAP/PH/1 queue that does

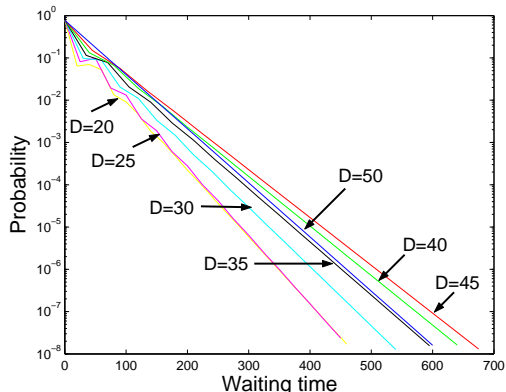


Figure 2: *Waiting time distribution*

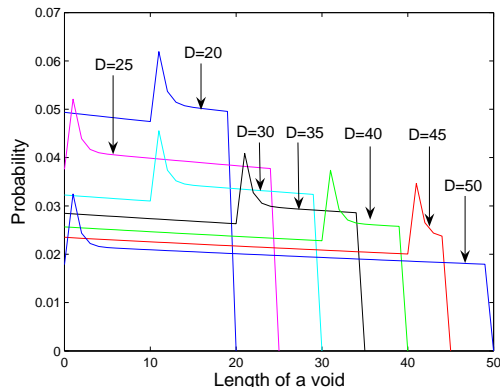


Figure 3: *Void distribution*

not take the cross-correlation between the inter-arrival and service times into account. A similar comparison for the M/M/1 queue was presented in [8]. The OBs are assumed to have a fixed length  $L$ . New incoming OBs follow an interrupted Bernoulli process (IBP). An IBP process consists of two states: new OBs arrive with probability  $p$  in state two, whereas in state one the IBP generates no traffic. The sojourn time in state  $i = 1, 2$  is geometrically distributed with a mean  $s_i$  and we fix  $s_1 = 10s_2$ . This process is characterized by two  $m \times m$  matrices  $\mathbf{D}_0$  and  $\mathbf{D}_1$ :

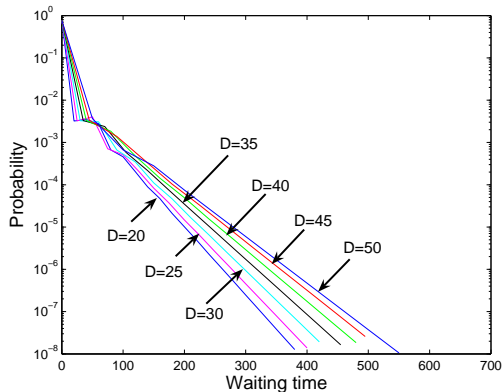
$$\mathbf{D}_0 = \begin{bmatrix} 1 - \frac{1}{s_1} & \frac{1}{s_1} \\ (1-p)\frac{1}{s_2} & (1-p)\left(1 - \frac{1}{s_2}\right) \end{bmatrix}, \quad \mathbf{D}_1 = \begin{bmatrix} 0 & 0 \\ \frac{p}{s_2} & p\left(1 - \frac{1}{s_2}\right) \end{bmatrix}. \quad (21)$$

In this example  $L_k$  equals  $L$  for each  $k$ , so the service time for OBs who find the server busy equals  $L'_k = L + V_k = L + [(\tau_k - L) \bmod D]$  and thus depends on the inter-arrival times. Customers who find the server idle have a deterministic service time equal to  $L$ . A matrix representation of the modified PH distributions is given next. Define  $\mathbf{T}$  as the following  $(L + D - 1) \times (L + D - 1)$  matrix

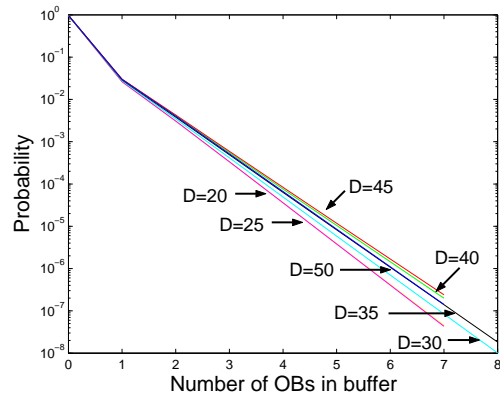
$$\mathbf{T} = \begin{bmatrix} 0 & 1 & 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & 1 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & & & \ddots & & \\ 0 & 0 & 0 & 0 & \dots & 0 & 1 \\ 0 & 0 & 0 & 0 & \dots & 0 & 0 \end{bmatrix}.$$

Notice,  $L + D - 1$  is the maximum service time of an OB (as the length of a void is between zero and  $D - 1$ ). Let  $\alpha_i$  be an  $1 \times (L + D - 1)$  nonnegative stochastic vector, the entries of which all equal zero except for the  $(D - [(i - L) \bmod D])$ -st which equals one. The vectors  $\alpha_{i,0}$  are independent of the value of  $i$ , that is,  $\alpha_{i,0} = \alpha_{L+1}$ . Therefore, the service requirements of an OB who finds the server busy, resp. idle, upon arrival has a common phase-type distribution function with matrix representation  $(m_{\text{ser}}, \alpha_i, \mathbf{T})$ , resp.  $(m_{\text{ser}}, \alpha_{i,0}, \mathbf{T})$ .

Figure 2 shows the waiting time distribution for  $L = 50$  (slots), the mean on period  $s_2$  equals 1.5,  $s_1 = 10s_2$  and  $p = 0.2\frac{11}{L}$ ; therefore  $\lambda = 0.2/L = 0.004$ . The granularity parameter



**Figure 4:** *Approximated waiting time distribution*



**Figure 5:** *Queue length distribution*

$D$  varies between 20 and 50. Looking at Figure 2 one sees that the waiting time has a tendency to grow as a function of the granularity  $D$ . This is logical as an infinite FDL buffer with a granularity  $D$  can only realize a subset of the delays that an FDL with granularity  $D/n$  can introduce (provided that  $D/n \in \mathbf{N}$ ). The interest in FDLs with a large granularity  $D$  stems from that fact that real systems can only support a finite number of fibres. The curves for  $D = 25$  and 50 deviate from the general tendency mentioned above. This can be understood by looking at the void length distribution, which shows that setting the granularity  $D$  equal to 25 or 50 assures that the most probable void length equals one (see Figure 3, this figure was generated using Lemma 1 of [8]). For this particular example, the density function of the inter-arrival time distribution is decreasing. Therefore, it is best to choose  $D \approx (L - 1)/n$ , for some  $n \in \mathbf{N}$  as this will guarantee that bursts who have a small inter-arrival time, cause small voids. This observation was also made in [11, 18], where a different approach was taken to study FDL buffers.

Figure 4 depicts the approximated waiting time distribution, by relying on the classic D-MAP/PH/1 queue (with a modified boundary behavior). To generate these results we can make use of Section 4 and 6 by setting  $\alpha_k = \alpha$  for all  $k$  (and leaving all the other parameters unaltered). The last  $D - 1$  entries of  $\alpha$  equal zero, while the  $i$ -th equals the probability that a void has length  $D - i$  (for  $i = 1, \dots, D$ ). Comparing Figure 2 and 4 clearly shows that neglecting the cross-correlation makes a substantial difference. More importantly, the beneficial properties of setting  $D = 25$  or 50 do not appear in Figure 4. Thus, the influence of the correlation is critical when selecting the granularity  $D$ .

Finally, Figure 5 depicts the queue length distribution for various granularity values  $D$ . Again, the granularities  $D = 25$  and 50 are somewhat special. Thus, the shorter voids caused by this particular choice of  $D$  does not only reduce the waiting time, but also causes fewer bursts to be present in the buffer (as bursts followed by a short void are more rapidly served). This observation seems to suggest that we can also expect smaller loss rates when choosing  $D \approx (L - 1)/n$  (with  $n \in \mathbf{N}$ ). This idea was confirmed in [10, 18]. In both these papers a Markov chain was formed by keeping track of the scheduling horizon  $\bar{H}$  (see Section 9). Although the approach taken here can only be used for infinite buffers with fixed length optical bursts, it allows us to study the number of bursts in the FDL. This distribution cannot be obtained from the steady state of the scheduling horizon. Also, using this approach one can

easily incorporate semi-Markovian arrivals at no additional computational costs (see Section 8).

**A note on the implementation:** For large values of  $L$  and  $D$  one requires a lot of memory to store all the  $\mathbf{A}_i$  and  $\mathbf{C}_i$  matrices simultaneously. In this case the following method can be used to reduce the memory complexity. Store a matrix  $\mathbf{X} = [\mathbf{I} \ \mathbf{D}_0 \ \mathbf{D}_0^2 \ \mathbf{D}_0^3 \ \dots]$  and use this matrix to determine the necessary  $\mathbf{A}_i$  matrices during each iteration while computing  $\mathbf{R}$  (Eqn. (8)). The disadvantage being that the computation of the steady-state probabilities is somewhat slower by repeatedly computing the  $\mathbf{A}_i$  matrices. Hence, it is a trade-off between the time and memory consumption of the algorithm. Roughly speaking, the overall time and memory complexity per iteration is  $O(m^3(L+D)^3)$  and  $O(m^2(L+D)^2)$ . The number of iterations required depends on various system parameters—e.g., the load, correlation structure, etc.—as well as on the algorithm used (see Appendix).

## Acknowledgment

B. Van Houdt is a post-doctoral fellow of the FWO-Flanders. This work was partly funded by the IWT-GBOU project 010058 “Optical Networking and Node Architectures”.

## References

- [1] I.J.B.F. Adan and V.G. Kulkarni. Single-server queue with markov-dependent inter-arrival and service times. *Queueing Systems and its Applications*, 45:113–134, 2003.
- [2] A.S. Alfa, B. Sengupta, T. Takine, and J. Xue. A new algorithm for computing the rate matrix of GI/M/1 type Markov chains. In *Proc. of the 4th Int. Conf. on Matrix Analytic Methods*, pages 1–16, Adelaide, Australia, 2002.
- [3] C. Blondia. A discrete-time batch markovian arrival process as B-ISDN traffic model. *Belgian Journal of Operations Research, Statistics and Computer Science*, 32(3,4):3–23, 1993.
- [4] C. Blondia and O. Casals. Statistical multiplexing of VBR sources: A matrix-analytical approach. *Performance Evaluation*, 16:5–20, 1992.
- [5] F. Callegati. Optical buffers for variable length packet switching. *IEEE Communications Letters*, 4:292–294, 2002.
- [6] I. Cidon, R. Gurin, A. Khamisy, and M. Sidi. Analysis of a correlated queue in a communication system. *IEEE Trans. on Information Theory*, IT-39(2):456–465, March 1993.
- [7] Q. HE. Age process, workload process, sojourn times, and waiting times in a discrete-time SM[K]/PH[K]/1/FCFS queue. *Queueing Systems*, 49:363–403, 2005.
- [8] D. Hong, F. Poppe, J. Reynier, F. Bacelli, and G. Petit. The impact of burstification on TCP throughput in optical burst switching networks. In *Proc. of the 18th International Teletraffic Congress (ITC)*, Berlin (Germany), Sept. 2003.

- [9] G.U. Hwang and K. Sohraby. Performance of correlated queues: impact of correlated service and inter-arrival times. *Performance Evaluation*, 55:129–145, 2004.
- [10] K. Laevens and H. Bruneel. Analysis of a single wavelength optical buffer. In *Proceedings of Infocom*, San Francisco, April 2003.
- [11] K. Laevens, B. Van Houdt, H. Bruneel, and C. Blondia. On the sustainable load of fibre delay line buffers. *IEE Electronics Letters*, 40(2):137–138, 2004.
- [12] 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*, 22:676–705, 1990.
- [13] M.F. Neuts. Markov chains with applications in queueing theory, which have a matrix geometric invariant probability vector. *Adv. Appl. Prob.*, 10:185–212, 1978.
- [14] M.F. Neuts. *Matrix-Geometric Solutions in Stochastic Models, An Algorithmic Approach*. John Hopkins University Press, 1981.
- [15] V. Ramaswami. Nonlinear matrix equations in applied probability - solution techniques and open problems. *SIAM review*, 30(2):256–263, June 1988.
- [16] B. Sengupta. Markov processes whose steady state distribution is matrix exponential with an application to the GI/PH/1 queue. *Adv. Appl. Prob.*, 21:159–180, 1989.
- [17] B. Sengupta. The semi-Markovian queue: theory and applications. *Stochastic Models*, 6(3):383–413, 1990.
- [18] B. Van Houdt, K. Laevens, J. Lambert, C. Blondia, and H. Bruneel. Channel utilization and loss rate in a single-wavelength fibre delay line (FDL) buffer. In *Proc. of Globecom 2004*, paper OC05-7, Dallas (US), Nov 2004.

## Appendix

In this section we present a brief overview of two techniques, discussed in [13, 15], to calculate the smallest nonnegative solution of Eqn. (8). Various other algorithms exist which we shall not discuss (see [2] and the references therein). The direct iteration scheme computes  $\mathbf{R}$  as

$$\mathbf{R}_0 = \mathbf{A}(\tau) = \sum_{l=0}^{\infty} \tau^l \mathbf{A}_l, \quad \mathbf{R}_{k+1} = \sum_{\nu=0}^{\infty} \mathbf{R}_k^\nu \mathbf{A}_\nu, \quad k \geq 0. \quad (22)$$

The number of iterations required by this scheme may grow considerably, especially when the Perron Frobenius (PF) eigenvalue  $\tau$  of  $\mathbf{R}$  is close to one. The choice for  $\mathbf{R}_0$  guarantees that the PF eigenvalue of all the matrices  $\mathbf{R}_k$  equals  $\tau$ . Let  $\xi(\mathbf{A}(z))$  denote the PF eigenvalue of  $\mathbf{A}(z) = \sum_i A_i z^i$ . Neuts [14] has shown that  $\tau$  is the unique solution in  $(0, 1)$  of the equation

$$x = \xi(\mathbf{A}(x)), \quad (23)$$

which can be solved using a bisection algorithm. If the PF eigenvalue is close to one, it is often more efficient to use the following variant of the Newton-Kantorovich scheme [15]. Define the

function  $F(\mathbf{X})$  as

$$F(\mathbf{X}) = \mathbf{X} - \sum_{\nu=0}^{\infty} \mathbf{X}^{\nu} \mathbf{A}_{\nu}. \quad (24)$$

Then the following scheme can be used:

$$\mathbf{R}_0 = 0, \quad \mathbf{R}_{k+1} = \mathbf{R}_k + \mathbf{Y}_k, \quad (25)$$

with

$$\mathbf{Y}_k = -F(\mathbf{R}_k) + \mathbf{Z}_k \mathbf{A}_1 + (\mathbf{R}_k \mathbf{Z}_k + \mathbf{Z}_k \mathbf{R}_k) \mathbf{A}_2, \quad (26)$$

and

$$\mathbf{Z}_k = -F(\mathbf{R}_k)(\mathbf{I} - \mathbf{A}_1)^{-1}. \quad (27)$$

The most efficient method to compute  $\mathbf{R}$  in terms of the time complexity is probably found by taking the time inverse of the GI/M/1 process and to compute the  $\mathbf{G}$  matrix of the resulting M/G/1 type MC by means of a quadratically converging algorithm. We did not explore this possibility as it generally requires more memory, which was, to some extent, the limiting factor in Sections 9 and 10.