

# A Matrix Analytic Solution of a Finite Buffer Queue with PH Distributed Customers' Impatience<sup>\*</sup>

Vicent Pla<sup>1</sup>, Vicente Casares-Giner<sup>1</sup>, Jorge Martínez<sup>1</sup>

Department of Communications. Universidad Politécnica de Valencia (UPV).  
ETSIT Camí de Vera s/n, 46022 Valencia, Spain.  
{vpla,vcasares,jmartinez}@dcom.upv.es

**Abstract.** We study a multiserver finite buffer queue in which customers have a stochastic deadline of phase-type until the beginning of their service. The following service disciplines are considered: *FCFS* (First-Come First-Served), *LCFS* (Last-Come First-Served) and *SIRO* (Service In Random Order) along with a parameterizable probabilistic push-out mechanism. The analysis of the system is performed using a matrix analytic approach and we obtain performance measures such as probabilities of blocking, expulsion and abandonment as well as the sojourn time distribution in different system conditions.

**Keywords:** impatient customers, service discipline, phase-type, matrix-analytic.

## 1 Introduction

Queuing models in which customers abandon the system if their service has not started by a given deadline have many applications in telecommunications as well as in other disciplines (some examples can be found in [1]).

While this topic has attracted the interest of queuing theorists for a few decades, the existing literature is rather limited (see [2] and references therein). To the best of our knowledge, infinite buffer size is assumed in all queuing models considering the impatience phenomenon, except [3] and [4]. Moreover, in [3] only the relatively simple case of exponentially distributed patience time is contemplated.

Intuition seems to indicate that the service order may have an influence on the number of customers that leave the system without being served. Indeed, in [5] the authors give a characterization of the optimal scheduling discipline that minimizes the number of customers that abandon the system before receiving service, and there are some instances where the optimal policy is not the conventional *FCFS*. Zhao and Alfa [6] consider a system in which impatient customers are served on an *LCFS* basis. The analysis in [6] is approximate and the patience time is assumed to be deterministic. Doshi and Heffes [7] study a model quite related to ours, where customers may “turn bad” after some time although they do not abandon the system, i.e. bad customers are served even though

---

<sup>\*</sup> This work has been supported by the Spanish Government (PGE, 30%) and the European Commission (FEDER, 70%) through the projects TIC2003-08272 and TEC2004-06437-C05-01.

they do not count as system goodput. The study in [7] considers *FCFS* and *LCFS* service disciplines as well as various customer rejection schemes including blocking and push-out. Notwithstanding the affinities between the model in [7] and the one of this paper, there are significant differences between them, and the analytical solution of the model in [7] cannot be applied to our model.

In this paper we analyze a queuing model that has the following characteristics: the waiting room is finite; the customers' impatience is modeled by a phase-type (PH) distribution; the service discipline can be *FCFS* (First Come First Served), *LCFS* (Last Come First Served) or *SIRO* (Service In Random Order); when a customer arrives and the buffer is full, the buffer management policy can block the newly arrived customer or push-out the head-of-line (HOL) customer in order to allocate the newly arrived one.

The contribution of our work is twofold. First, we solve the  $M/M/C/K/FCFS + PH$  model using an alternative approach to that of [4] by deploying a *matrix analytic* solution. The model studied in [4] is, in a sense, more general than ours since it considers a general distribution for the patience time. Nevertheless, the family of PH-distributions is highly versatile and can be fitted to a wide range of experimental data and, from a theoretical point of view, the set of PH-distributions is dense in the set of all probability distributions on  $[0, \infty)$  [8]. Secondly, and more importantly, we carry out the analysis for several non-*FCFS* disciplines in the context of customer impatience and a rather general distribution for the patience time, which have not been done before.

The remaining of the paper is structured as follows. In Section 2 the mathematical model of the system is described and its analysis is carried out in Section 3. A numerical example is presented in Section 4. Finally, Section 5 summarizes the paper and draws some conclusions.

## 2 Model Description

We first describe the analytical aspects that are common to all service disciplines and customer rejection schemes, and their specific features are subsequently addressed.

The system has a total of  $C$  identical servers and a waiting line limited to  $N$  positions. Customers arrive according to a Poisson process of rate  $\lambda$  and their service time is exponentially distributed with rate  $\mu$ . Each customer has a stochastic deadline—starting on arrival—until the beginning of service, if service has not begun by his deadline the customer abandons the system and is lost. Deadlines are considered to be *independent and identically distributed* (iid) phase type (PH) random variables (*rv*) with representation  $(\boldsymbol{\beta}, \mathbf{T})$ ; we denote by  $m$  the number of phases in the PH distribution (see [9] for more details on PH distributions).

Throughout this paper the following notation is used:  $x_i$  denotes the  $i$ -th entry of a vector  $\mathbf{x}$ ;  $M_{ij}$  denotes the entry on the  $i$ -th row and the  $j$ -th column of a matrix  $M$ ;  $\otimes$  denotes the Kronecker product of two matrices (example:  $A \otimes B$ ), see [9, p. 17] for further details;  $\text{diag}\{\cdot\}$  is an operator over a vector that yields a diagonal matrix whose diagonal entries are the elements of the vector;  $\mathbf{e}$  is a column vector of ones;  $\mathbf{0}$  is a column vector of zeros and  $\mathbf{I}$  is the identity matrix.

Under the above assumptions the system model fits within the category of non-homogeneous finite Quasi-Birth-and-Death (QBD) processes. Let  $\{X(t) : t > 0\}$  be

the stochastic process for the system with the following two-dimensional state space  $S = \{(l, k) : 0 \leq l \leq N; 0 \leq k \leq m^l - 1\} \cup \{(-1, k) : 0 \leq k \leq C - 1\}$  that can be partitioned into *levels* as  $S = \bigcup_{l=-1}^N L(l)$ , where  $L(l_0) = \{(l, k) : l = l_0; (l_0, k) \in S\}$ . The first coordinate of a state is also referred to as *level* and the second coordinate as *phase*. Level  $(-1)$  groups the states in which not all servers are busy, for this level the state phase indicates the number of customers in the system, i.e. being served. Level  $l \geq 0$  groups the states in which all servers are busy and there are  $l$  customers in the waiting room. In levels  $l \geq 0$  the state phase encodes the current phase of the patience time distribution for each of the  $l$  waiting customers. The mapping between the state phase ( $k$ ) and waiting customers phases is as follows. Let the  $l$ -tuple  $(k_1, \dots, k_l)$ ,  $1 \leq k_i \leq m$ ,  $1 \leq i \leq l$  denote the phases of the Markov process associated to the PH distribution of the waiting customers patience time, being  $k_i$  the phase of the customer at the  $i$ -th position in the waiting room. Then,  $k = \sum_{1 \leq i \leq l} k_i m^{l-i}$ ; in other words, the  $l$ -tuples are numbered in lexicographical order from  $(1, \dots, 1)$  to  $(m, \dots, m)$ .

Let us denote by  $\boldsymbol{\pi}$  the stationary probability vector of the process. In the same way as with states, we partition  $\boldsymbol{\pi}$  by levels into subvectors  $\boldsymbol{\pi}^{(l)}$ ,  $-1 \leq l \leq N$ , where  $\boldsymbol{\pi}^{(-1)}$  has  $C$  components and  $\boldsymbol{\pi}^{(l)}$  ( $l \geq 0$ ) has  $m^l$  components. State transitions are restricted to states in the same level or in two adjacent levels and consequently, the infinitesimal generator  $\mathbf{Q}$  of the process has a block tridiagonal structure,

$$\mathbf{Q} = \left[ \begin{array}{c|ccc} \mathbf{A}_1^{(-1)} & \mathbf{A}_0^{(-1)} & \mathbf{0} & \dots \\ \mathbf{A}_2^{(0)} & & & \\ \mathbf{0} & & \mathbf{Q}_p & \\ \vdots & & & \end{array} \right], \quad \text{where} \quad \mathbf{Q}_p = \left[ \begin{array}{cc} \mathbf{A}_1^{(0)} & \mathbf{A}_0^{(0)} \\ \mathbf{A}_2^{(1)} & \mathbf{A}_1^{(1)} & \mathbf{A}_0^{(1)} \\ & \ddots & \\ & & \mathbf{A}_2^{(N)} & \mathbf{A}_1^{(N)} \end{array} \right].$$

Block matrices which are not in  $\mathbf{Q}_p$  involve the boundary level  $(-1)$  and thus they do not conform to the general construction that will be given for the rest of matrices. Matrix entries for these particular cases, which are independent of the service discipline, are as follows:

$$\mathbf{A}_1^{(-1)} = \left[ \begin{array}{ccc} * & \lambda & \\ \mu & * & \lambda \\ 2\mu & * & \lambda \\ & & \ddots \\ & & & (C-1)\mu & * \end{array} \right], \quad \mathbf{A}_0^{(-1)} = [0 \dots 0 \lambda]^t, \quad \mathbf{A}_2^{(0)} = [0 \dots 0 C\mu]$$

The diagonal entries of  $\mathbf{A}_1^{(-1)}$ , which are represented by asterisks for the ease of display, are such that the corresponding rows of  $\mathbf{Q}$  sum to zero, i.e.  $\mathbf{A}_1^{(-1)} \mathbf{e} + \mathbf{A}_0^{(-1)} \mathbf{e} = \mathbf{0}$ .

## 2.1 FCFS Discipline

The matrices  $\mathbf{A}_0^{(l)}$  correspond to transitions from  $L(l)$  to  $L(l+1)$ ,  $0 \leq l < N$ . These transitions represent the arrival of a customer that will occupy the  $(l+1)$ -th position of the waiting room. The PH distribution for the patience time of the arriving customer will begin at its  $i$ -th phase with probability  $\beta_i$ . It can be easily seen that  $\mathbf{A}_0^{(l)} = \mathbf{I}_{m^l} \otimes \lambda\boldsymbol{\beta}$ .

The matrices  $\mathbf{A}_2^{(l)}$  correspond to transitions from  $L(l)$  to  $L(l-1)$ ,  $0 < l \leq N$ . These transitions represent the departure of a customer from the system which may be due to either a customer abandoning the waiting line (because his deadline has expired) or to a service completion. The former type of transition will be represented by matrix  $\mathbf{U}_1^{(l)}$  and the latter by matrix  $\mathbf{U}_2^{(l)}$ . Then,  $\mathbf{A}_2^{(l)} = \mathbf{U}_1^{(l)} + \mathbf{U}_2^{(l)}$  where

$$\mathbf{U}_1^{(l)} = \begin{cases} \boldsymbol{\tau}, & l = 1 \\ \boldsymbol{\tau} \otimes \mathbf{I}_{m^{l-1}} + \mathbf{I}_m \otimes \mathbf{U}_1^{(l-1)}, & 1 < l \leq N \end{cases}; \quad \mathbf{U}_2^{(l)} = C\mu\mathbf{e}_m \otimes \mathbf{I}_{m^{l-1}}$$

being  $\boldsymbol{\tau} = -\mathbf{T}\mathbf{e}$ .

The matrices  $\mathbf{A}_1^{(l)}$  correspond to transitions between states within  $L(l)$ . These transitions represent phase changes in the PH processes associated to waiting customers. The expression for this matrix is first given ignoring elements on its main diagonal and next they will be computed using the fact that the rows of  $\mathbf{Q}$  must sum to zero. For the sake of clarity we introduce the set of matrices  $\mathbf{D}^{(l)}$  whose entries are equal to the entries of  $\mathbf{A}_1^{(l)}$ , except those on their main diagonal. Now it can be written that

$$\mathbf{D}^{(l)} = \begin{cases} \mathbf{T}, & l = 1 \\ \mathbf{T} \otimes \mathbf{I}_{m^{l-1}} + \mathbf{I}_m \otimes \mathbf{D}^{(l-1)}, & 1 < l \leq N \end{cases} \quad (1)$$

Note that in the expression of  $\mathbf{D}^{(N)}$  it is assumed that customers arriving while the system is at level  $N$ , i.e. when the buffer is full, are blocked. Then  $\mathbf{A}_1^{(l)}$  is given by

$$\mathbf{A}_1^{(l)} = \mathbf{D}^{(l)} - \text{diag} \left\{ \mathbf{A}_2^{(l)}\mathbf{e} + \mathbf{D}^{(l)}\mathbf{e} + \mathbf{A}_0^{(l)}\mathbf{e} \right\}$$

which can be further simplified to

$$\mathbf{A}_1^{(l)} = \begin{cases} \mathbf{D}^{(l)} - (C\mu + \lambda)\mathbf{I}_{m^l}, & l < N \\ \mathbf{D}^{(N)} - C\mu\mathbf{I}_{m^N}, & l = N \end{cases} \quad (2)$$

by virtue of Proposition 1, that will be proved after the following lemma.

**Lemma 1**  $(\mathbf{A} \otimes \mathbf{B})\mathbf{e} = (\mathbf{A}\mathbf{e}) \otimes (\mathbf{B}\mathbf{e})$

**Proposition 1** *The following equalities hold for  $1 \leq l \leq N$*

$$\text{diag} \left\{ \mathbf{A}_0^{(l)}\mathbf{e} \right\} = \lambda\mathbf{I}_{m^l}; \quad \text{diag} \left\{ \mathbf{A}_2^{(l)}\mathbf{e} + \mathbf{D}^{(l)}\mathbf{e} \right\} = C\mu\mathbf{I}_{m^l}$$

*Proof.* The first equality follows immediately by applying the Lemma to  $\mathbf{A}_0^{(l)} = \mathbf{I}_{m^l} \otimes \lambda\beta$  and noting that  $\beta\mathbf{e} = 1$  and  $\text{diag}\{\mathbf{e}\} = \mathbf{I}$ .

To prove the second equality we first observe that  $\mathbf{A}_2^{(l)}\mathbf{e} + \mathbf{D}^{(l)}\mathbf{e} = \mathbf{U}_1^{(l)}\mathbf{e} + \mathbf{U}_2^{(l)}\mathbf{e} + \mathbf{D}^{(l)}\mathbf{e}$  and by applying the Lemma to  $\mathbf{U}_2^{(l)} = C\mu\mathbf{e}_m \otimes \mathbf{I}_{m^{l-1}}$  it is easily seen that

$$\mathbf{U}_2^{(l)}\mathbf{e} = C\mu\mathbf{e}_{m^l}. \quad (3)$$

On the other hand

$$\begin{aligned} \mathbf{U}_1^{(l)}\mathbf{e}_{m^{l-1}} + \mathbf{D}^{(l)}\mathbf{e}_{m^l} &= (\boldsymbol{\tau} \otimes \mathbf{I}_{m^{l-1}} + \mathbf{I}_m \otimes \mathbf{U}_1^{(l-1)})\mathbf{e} + (\mathbf{T} \otimes \mathbf{I}_{m^{l-1}} + \mathbf{I}_m \otimes \mathbf{D}^{(l-1)})\mathbf{e} \\ &= (\boldsymbol{\tau} \otimes \mathbf{e}_{m^{l-1}} + \mathbf{e}_m \otimes (\mathbf{U}_1^{(l-1)}\mathbf{e}_{m^{l-2}})) + (\mathbf{T}\mathbf{e}_m \otimes \mathbf{e}_{m^{l-1}} + \mathbf{e}_m \otimes (\mathbf{D}^{(l-1)}\mathbf{e}_{m^{l-1}})) \\ &= \mathbf{e}_m \otimes (\mathbf{U}_1^{(l-1)}\mathbf{e}_{m^{l-2}} + \mathbf{D}^{(l-1)}\mathbf{e}_{m^{l-1}}) \end{aligned} \quad (4)$$

and by recursive application of (4) we obtain

$$\begin{aligned} \mathbf{U}_1^{(l)}\mathbf{e}_{m^{l-1}} + \mathbf{D}^{(l)}\mathbf{e}_{m^l} &= \mathbf{e}_m \otimes (\mathbf{U}_1^{(l-1)}\mathbf{e}_{m^{l-2}} + \mathbf{D}^{(l-1)}\mathbf{e}_{m^{l-1}}) \\ &\quad \vdots \\ &= \mathbf{e}_{m^{l-1}} \otimes (\boldsymbol{\tau} + \mathbf{T}\mathbf{e}_m) = \mathbf{e}_{m^{l-1}} \otimes \mathbf{0} = \mathbf{0} \end{aligned} \quad (5)$$

Hence, from (3) and (5) it follows that  $\mathbf{A}_2^{(l)}\mathbf{e} + \mathbf{D}^{(l)}\mathbf{e} = C\mu\mathbf{e}_{m^l}$ , and taking the  $\text{diag}\{\cdot\}$  operator on both sides of this equality yields the desired result.

## 2.2 LCFS Discipline

Using the *LCFS* discipline only affects the selection of the queued customer that will start service, and therefore will be removed from the queue, after a service completion. Thus the expression for  $\mathbf{U}_2^{(l)}$  must be modified in the following manner:  $\mathbf{U}_2^{(l)} = \mathbf{I}_{m^{l-1}} \otimes C\mu\mathbf{e}_m$ . The rest remains unchanged since under the new expression for  $\mathbf{U}_2^{(l)}$  it can be proved, in much the same way as before, that (3) holds.

## 2.3 SIRO Discipline

By the same reasoning as in *LCFS* discipline, only  $\mathbf{U}_2^{(l)}$  changes and it can be readily shown to be given by

$$\mathbf{U}_2^{(1)} = C\mu\mathbf{e}_m \quad \text{and} \quad \mathbf{U}_2^{(l)} = \frac{C\mu}{l}\mathbf{e}_m \otimes \mathbf{I}_{m^{l-1}} + \mathbf{I}_m \otimes \frac{(l-1)}{l}\mathbf{U}_2^{(l-1)} \quad 1 < l < N$$

In this case Eq. (3) can be easily proved by induction on  $l$ .

## 2.4 Buffer management scheme

So far we have assumed that arriving customers to a full buffer are blocked. Here we consider a more general buffer management scheme in which a customer arriving when

the buffer is full either is blocked or it pushes-out the HOL waiting customer. The choice between these two options is done randomly and independently for each customer. The random component of the algorithm is tuned by parameter  $q$  which represents the probability that a customer that finds a full queue upon arrival will push-out the HOL customer; therefore, a customer finding a full queue will be blocked with probability  $(1 - q)$ .

The arrival of a customer that finds a full buffer and pushes-out the HOL customer is modeled as a transition within level  $L(N)$  and thus only the values entries of matrix  $\mathbf{D}^{(N)}$  are modified as follows:  $\mathbf{D}^{(N)} \leftarrow \mathbf{D}^{(N)} + \mathbf{e}_m \otimes \mathbf{I}_{m^{N-1}} \otimes \lambda\boldsymbol{\beta}$ . Hence, if the push-out and blocking schemes are probabilistically combined together with probabilities  $q$  and  $(1 - q)$ , respectively, we have that

$$\mathbf{D}^{(l)} = \begin{cases} \mathbf{T}, & l = 1 \\ \mathbf{T} \otimes \mathbf{I}_{m^{l-1}} + \mathbf{I}_m \otimes \mathbf{D}^{(l-1)}, & 1 < l < N \\ \mathbf{T} \otimes \mathbf{I}_{m^{N-1}} + \mathbf{I}_m \otimes \mathbf{D}^{(N-1)} + q(\mathbf{e}_m \otimes \mathbf{I}_{m^{N-1}} \otimes \lambda\boldsymbol{\beta}), & l = N \end{cases} \quad (6)$$

Finally, by the same method we have used to obtain (2) it follows that

$$\mathbf{A}_1^{(l)} = \begin{cases} \mathbf{D}^{(l)} - (C\mu + \lambda)\mathbf{I}_{m^l}, & l < N \\ \mathbf{D}^{(N)} - (C\mu + q\lambda)\mathbf{I}_{m^N}, & l = N \end{cases} \quad (7)$$

Note that when  $q = 0$ , (7) reduces to (2).

### 3 Model Analysis

In this section we describe the method to calculate the stationary state probabilities ( $\boldsymbol{\pi}$ ) of the model. From these, performance evaluation measures are derived.

The stationary probabilities ( $\boldsymbol{\pi}$ ) are obtained as the solution to the set of simultaneous linear equations  $\boldsymbol{\pi}\mathbf{Q} = \mathbf{0}$ ,  $\boldsymbol{\pi}\mathbf{e} = 1$ . Being  $\mathbf{Q}$  a finite dimension matrix the above system could be solved by standard linear algebra methods. However, as the size of  $\mathbf{Q}$  may be very large, it is advisable to use more specific algorithms that take advantage of the structure ( $\mathbf{Q}$  is block-tridiagonal) and the nature of the problem ( $\mathbf{Q}$  is an infinitesimal generator). We used the *Linear Level Reduction* algorithm [9, 10], which can solve level-dependent finite QBDs.

```

1:  $U \leftarrow \mathbf{A}_1^{(N)}$ 
2:  $R^{(N)} \leftarrow \mathbf{A}_0^{(N-1)}(-U)^{-1}$ 
3: for  $l = N - 1, N - 2, \dots, 0, -1$  do
4:    $U \leftarrow \mathbf{A}_1^{(l)} + R^{(l+1)}\mathbf{A}_2^{(l+1)}$ 
5:    $R^{(l)} \leftarrow \mathbf{A}_0^{(l-1)}(-U)^{-1}$ 
6: end for
7: solve  $\boldsymbol{\pi}^{(-1)}$  from  $\{\boldsymbol{\pi}^{(-1)}U = \mathbf{0};$ 
    $\boldsymbol{\pi}^{(-1)}\mathbf{e} = 1\}$ 
8: for  $l = 0, 1, \dots, N$  do
9:    $\boldsymbol{\pi}^{(l)} = \boldsymbol{\pi}^{(l-1)}R^{(l)}$ 
10: end for
```

#### 3.1 Distribution of the Number of Customers

Let  $p_k$  ( $0 \leq k \leq N + C$ ) denote the probability that  $k$  customers are in the system, then we have that

$$p_k = \begin{cases} \pi_k^{-1}, & k = 0, \dots, C - 1 \\ \boldsymbol{\pi}^{(k-C)}\mathbf{e}, & k = C, \dots, C + N \end{cases}$$

### 3.2 Probabilities of Blocking, Expulsion and Reneging

Since arrivals are Poisson, by *PASTA* property [11] we have that the probability that an arriving customer sees the buffer full is  $p_{C+N}$ . Therefore, the blocking probability is given by  $P_b = (1 - q)p_{C+N} = (1 - q)\boldsymbol{\pi}^{(N)}\mathbf{e}$  and the expulsion probability is given by  $P_e = qp_{C+N} = q\boldsymbol{\pi}^{(N)}\mathbf{e}$ . The probability of reneging ( $P_r$ ) is measured by taking the average number of customers who renege divided by the average number of customers that arrived over a sufficiently long period, say  $t_0$ ,  $(\sum_{l=1}^N \boldsymbol{\pi}^{(l)}\mathbf{U}_1^{(l)}t_0\mathbf{e} + o(t_0))/(\lambda t_0 + o(t_0))$  and letting  $t_0 \rightarrow \infty$  we obtain that  $P_r = \sum_{l=1}^N \boldsymbol{\pi}^{(l)}\mathbf{U}_1^{(l)}\mathbf{e}/\lambda$ .

### 3.3 Sojourn Time in Congestion Condition and Blocking Condition

We say that the system is congested if an arriving customer has to wait, i.e. the system is in one of the states in  $\bigcup_{l=0}^N L(l)$ . Let  $T_c$  denote the sojourn time in the congestion condition *rv*. Similarly, let us define the blocking condition as the state in which the number of costumers in the system is at its maximum ( $C + N$ ), so that if a new customer arrives it will be blocked or a waiting customer will be pushed-out. Let  $T_b$  be the sojourn time in the blocking condition *rv*. Next we obtain the distribution of these *rv* and derive their mean values.

A congestion condition period starts when the system enters level  $L(0)$  and lasts until its first visit to level  $L(-1)$ . During this period the system will visit states in  $\bigcup_{l=0}^N L(l)$  whose residence times are all exponential. Thus it may be concluded that the distribution of  $T_c$  is phase-type and it is easy to check that its representation is  $PH(\boldsymbol{\beta}^{(c)}, \mathbf{T}^{(c)})$  where  $\boldsymbol{\beta}^{(c)} = [1 \ 0 \ 0 \ \dots \ 0]$  and  $\mathbf{T}^{(c)} = \mathbf{Q}_p$ .

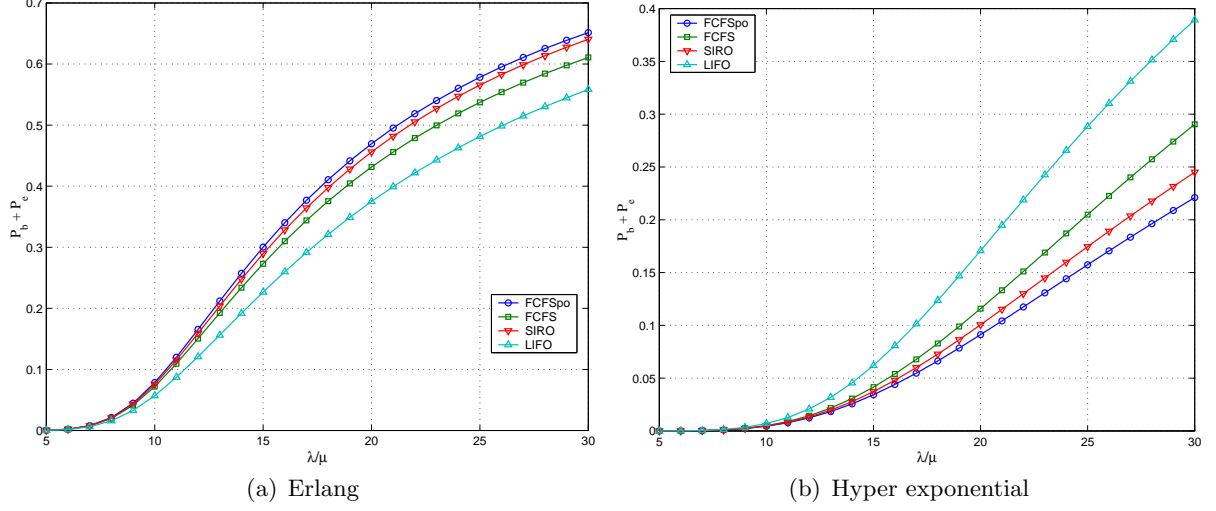
In order to obtain the mean value of  $T_c$  we will use a probabilistic argument instead of using its distribution which would entail inverting matrix  $\mathbf{T}^{(c)}$ . This reasoning is based on the observation that for an infinitely large time period, say  $t_0$ , it holds that, the mean sojourn time per visit to a set of states equals the total sojourn time in that set of states divided by the number of visits. Hence, we can write

$$\bar{T}_c = \lim_{t_0 \rightarrow \infty} \frac{\left(\sum_{k=C}^{C+N} p_k\right) t_0 + o(t_0)}{\lambda p_{C-1} t_0 + o(t_0)} = \frac{1}{\lambda p_{C-1}} \sum_{k=C}^{C+N} p_k.$$

In the same manner we can see that  $T_b$  is phase-type and its representation is  $PH(\boldsymbol{\beta}^{(b)}, \mathbf{T}^{(b)})$  where  $\boldsymbol{\beta}^{(b)} = \boldsymbol{\pi}^{(N-1)}\mathbf{A}_0^{(N-1)}/(\boldsymbol{\pi}^{(N-1)}\mathbf{A}_0^{(N-1)}\mathbf{e})$  and  $\mathbf{T}^{(b)} = \mathbf{A}_1^{(N)}$ . Hence [9, Eq. (2.13)]  $\bar{T}_b = \boldsymbol{\beta}^{(b)} \left(-\mathbf{T}^{(b)}\right)^{-1} \mathbf{e}$ , which noting that  $\mathbf{A}_0^{(N-1)}\mathbf{e} = \lambda\mathbf{e}$  and  $\boldsymbol{\pi}^{(N-1)}\mathbf{A}_0^{(N-1)} + \boldsymbol{\pi}^{(N)}\mathbf{A}_1^{(N)} = \mathbf{0}$  we can rewrite as  $\bar{T}_b = \boldsymbol{\pi}^{(N)}\mathbf{e}/(\lambda\boldsymbol{\pi}^{(N-1)}\mathbf{e})$ .

## 4 Numerical Example

In this section we present a numerical example to illustrate the analysis carried out in the previous sections. In this example the system parameters are:  $C = 10$ ,  $N = 5$ . Four combinations of service discipline and buffer management schemes are considered: *FCFS* (*FCFS* and  $q = 0$ ), *FCFSpo* (*FCFS* and  $q = 1$ ), *SIRO* (*SIRO* and  $q = 0.5$ ) and *LCFS*



**Fig. 1.** Blocking probability plus expulsion probability.

(*LCFS* and  $q = 0$ ). Arrival rate ( $\lambda$ ) and transition rates of the PH distribution ( $\mathbf{T}$ ) are normalized with respect to  $\mu$ . Two different instances of the patience time are examined, one whose hazard rate function<sup>1</sup> is increasing (e.g. Erlang) and the other decreasing (e.g. Hyper exponential). Their PH representations are, respectively,

$$\boldsymbol{\beta} = [1 \ 0 \ 0], \mathbf{T} = \begin{bmatrix} -3 & 3 & 0 \\ 0 & -3 & 3 \\ 0 & 0 & -3 \end{bmatrix}$$

and  $\boldsymbol{\beta} = 1/3 [1 \ 1 \ 1]$ ,  $\mathbf{T} = -56/150 \text{diag} \{ [50 \ 10 \ 1] \}$ .

Figure 1 represents the sum of the probabilities of blocking and expulsion as a function of the offered traffic, and Fig. 2 does the same for the probability of reneging.

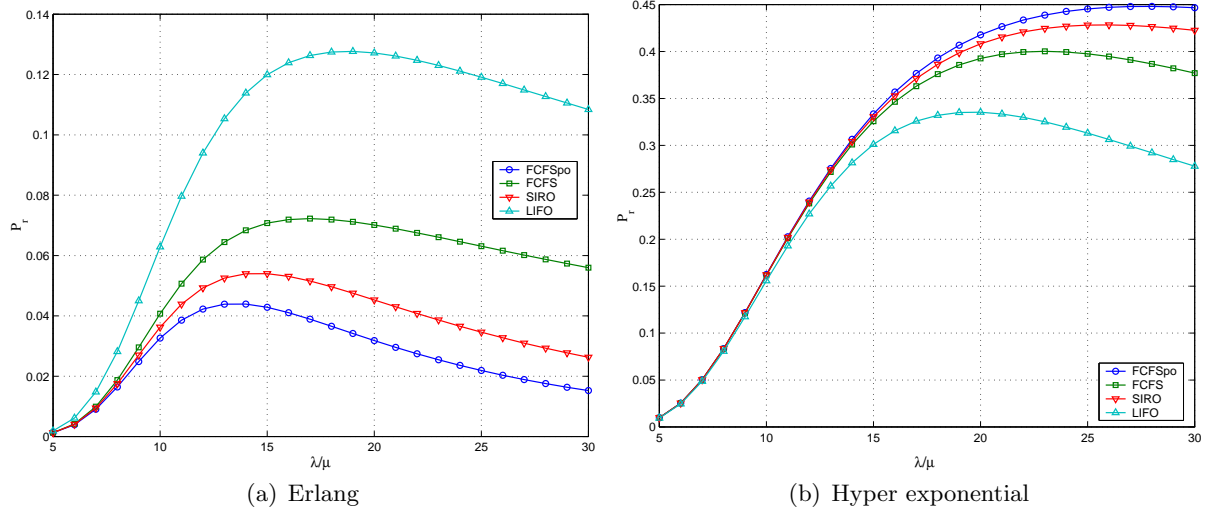
The total computational cost incurred in order to find the state probabilities is dominated by the cost of inverting the matrix  $\mathbf{A}_1^{(N)}$ , which is of size  $m^N$ . Obviously, increasing the buffer size  $N$  and the number of phases  $m$  can make the numerical solution simply unfeasible, which of course is a limitation of our model. Nevertheless, in such situation our model can still provide a better approximation than the simple first order exponential approximation. An example comparing these approximations is shown in figures 3 and 4. In these two examples we consider a 5-phase PH distribution for which the exact performance measures are computed: in Fig. 3 we used a hyper exponential distribution and in Fig. 4 an Erlang distribution. Then, the original PH distribution is approximated (using moment matching) by a 1-phase PH (exponential) and a 2-phase PH. As expected, the 2-phase approximation offers better accuracy than the exponential one.

## 5 Conclusion

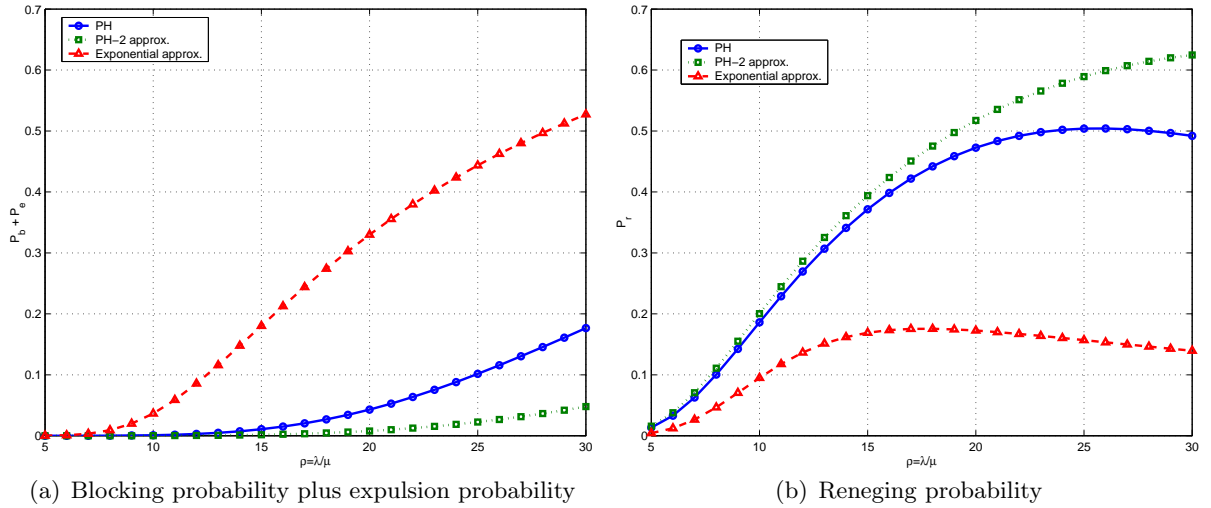
In this paper, we developed a stochastic model for a multiserver finite buffer queue with impatient customers where the patience time is modeled by a *rv* of phase-type. Further-

<sup>1</sup> The hazard rate function is also known as the failure rate function. Let  $f(x)$  and  $F(x)$  be the *pdf* and the *CDF* of a *rv*  $X$ , the hazard rate function  $h(x)$  of  $X$  is defined as  $h(x) = f(x)/(1 - F(x))$ .





**Fig. 2.** Reneging probability.



**Fig. 3.** Approximation of a hyper exponential distribution. LIFO discipline.

more, the model considers different service disciplines (*FCFS*, *LCFS* and *SIRO*) along with a probabilistically weighted buffer management scheme that combines two modes of operation: customers who arrive when system is full are blocked or push-out the HOL customer. The most significant achievement of this paper is to develop a model for the performance analysis of a queue with impatient customers under non-*FCFS* service disciplines. Secondly, our model considers the finite buffer case and a fairly general distribution of the customer patience.

The analytical model is constructed and analyzed using matrix analytic methods and we obtain expressions for the performance evaluation of the system. The complexity of the computations in the analysis of our model is mainly dependent on buffer size and the number of phases of the customer patience distribution. While this complexity can make the numerical analysis unfeasible in some cases, in these cases our model can still

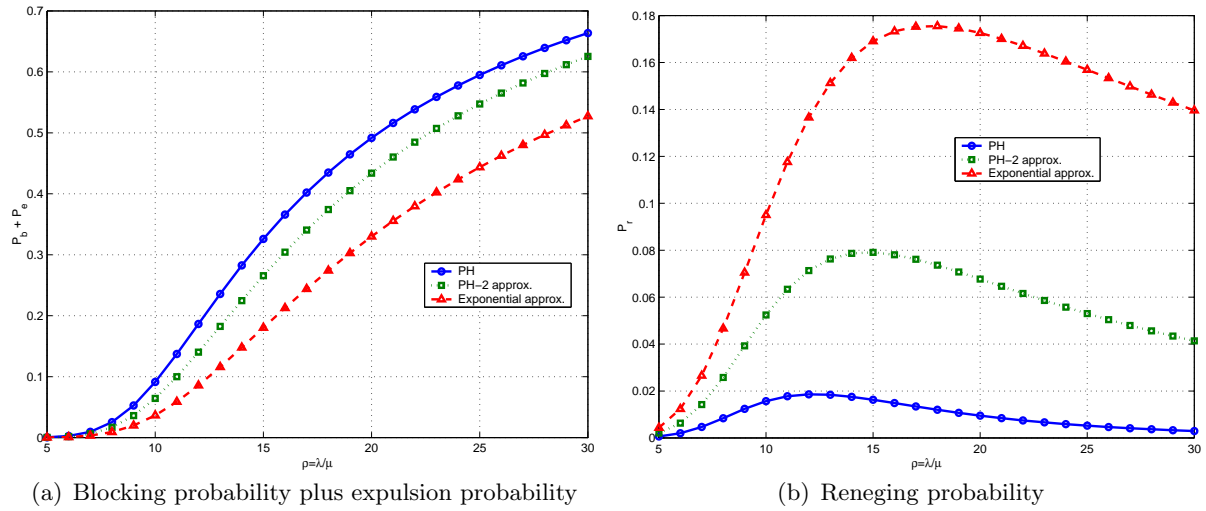


Fig. 4. Approximation of an Erlang distribution. FIFOpo discipline.

provide a better approximation than the one obtained using an exponentially distributed customer patience.

## References

1. Boxma, O., de Waal, P.: Multiserver queues with impatient customers. In: Proceedings of ITC 14, Elsevier Science (1994) 743–756
2. Brandt, A., Brandt, M.: Asymptotic results and a markovian approximation for the  $M(n)/M(n)/s + GI$  system. Queueing Systems **41** (2002) 73–94
3. Ancker, C.J., Gafarian, A.: Some queueing problems with balking and reneging. Operations Research **11** (1963) 88–100
4. Movaghar, A.: On queueing with customer impatience until the beginning of service. Queueing Systems **29** (1998) 337–350.
5. Towsley, D., Panwar, S.: Optimality of the stochastic earliest deadline policy for the  $G/M/c$  queue serving customers with deadlines. In: Proceedins of the 2nd ORSA Telecommunications Conference. (1992).
6. Zhao, Y.Q., Alfa, A.S.: Performance analysis of a telephone system with both patient and impatient customers. Telecommunication Systems **4** (1995) 201–215.
7. Doshi, B.T., Heffes, H.: Overload performance of several processor queueing disciplines for the  $M/M/1$  queue. IEEE Transactions on Communications **COM-34** (1986) 538–546.
8. Neuts, M.: Matrix-geometric Solutions in Stochastic Models: An Algorithmic Approach. The Johns Hopkins University Press (1981)
9. Latouche, G., Ramaswami, V.: Introduction to Matrix Analytic Methods in Stochastic Modeling. ASA-SIAM (1999)
10. Gaver, D., Jacobs, P., Latouche, G.: Finite birth-and-death models in randomly changing environments. Advances in Applied Probability **16** (1984) 715–731
11. Wolff, R.W.: Poisson arrivals see time averages. Operation Research **30** (1982) 223–231