



Response time in a tandem queue with blocking, Markovian arrivals and phase-type services

B. Van Houdt^{a,*}, Attahiru Sule Alfa^b

^a*Department of Mathematics and Computer Science, University of Antwerp, Performance Analysis of Telecommunication Systems Research Group, Middelheimlaan 1, B-2020 Antwerp, Belgium*

^b*Department of Electrical and Computer Engineering, University of Manitoba, Winnipeg, Manitoba, R3T 5V6, Canada*

Received 22 January 2004; accepted 20 August 2004

Available online 6 October 2004

Abstract

A novel approach for obtaining the response time in a discrete-time tandem-queue with blocking is presented. The approach constructs a Markov chain based on the age of the leading customer in the first queue. We also provide a stability condition and carry out several numerical examples.

© 2004 Elsevier B.V. All rights reserved.

Keywords: Tandem queue; Blocking; Markovian arrivals; Phase-type services; Response time; Matrix analytic method

1. Introduction

In most communication networks we are interested in obtaining the response time of jobs as well as the number of jobs in the system. It has been common practice to first obtain the queue length and then use that to obtain the response time. The two part procedure could be cumbersome in many situations,

especially when dealing with tandem queues. This is because in order to obtain the queue length we have to set up the associated Markov chain, obtain its stationary vector and then use that in a process that involves a considerable amount of effort to obtain the response time. In this paper we present a different and novel approach in which we set up the Markov chain, based on the age of the leading job in the first queue and other auxiliary variables. The stationary distribution of this Markov chain will lead us to the response time easily and the queue length can also be computed from it using a simple procedure. We believe that the effort required to compute the two quantities is less using this age process approach we are presenting, when compared to the traditional approach of using the Markov chain of queue length. More importantly, since the

* Corresponding author. Tel.: +32 3 265 3891; fax: +32 3 265 3777.

E-mail address: benny.vanhoudt@ua.ac.be (B. Van Houdt).

¹ B. Van Houdt is a post-doctoral Fellow of the FWO-Flanders. This work was performed while Houdt was visiting the University of Manitoba and TR Labs in Winnipeg, Manitoba. We would like to thank the University and TR Labs for their hospitality.

response time is sometimes the only key measure of interest, our approach is very favorable in such cases.

Tandem queues with blocking have received considerable attention in the queueing, communications and manufacturing literature because of their pervasiveness and significance in real life. A number of survey papers have been published during the last two decades [7,20,10]. Some of the earlier works on this include those of Hunt [12], who first studied the blocking effects in a sequence of waiting lines. Later Avi-Itzhak [1] studied the system with arbitrary input and regular service times. There has been a plethora of studies on this subject and an additional literature survey can be found in [21]. A continuous-time tandem queue with blocking, Markovian arrivals (MAP) and no intermediate buffer was considered by Gómez-Corral [6]. Phase-type service was assumed at the first infinite queue whereas the other queue is assumed to have a general service time. Results for the joint queue length distribution were provided and the stability issues were partially addressed. Gómez-Corral used matrix analytical methods (MAMs), as we do, to obtain his results, the difference in methodology being that we propose a novel approach that keeps track of the age of the leading customer in the first queue as opposed to the number of customers, as this more easily leads to the response-time distribution. MAMs have been used on a variety of occasions when studying tandem queues with blocking, going back to Latouche and Neuts [15].

The bulk of the work done in this area of tandem queues focussed on continuous-time models. As pointed out by Daduna [4], the introduction of the asynchronous transfer mode as a multiplexing technique for broadband integrated services digital networks has increased the studies in the areas of discrete-time queueing models. Even though carrying out discrete-time analysis of tandem queues may be done, to some extent, in a similar manner as their continuous-time counterparts their analysis often introduce some additional challenges. Setting up their transition matrices is more complex, and obtaining the response times of such a system after that requires a more considerable effort. Gün and Makowsky [8,9] considered a discrete-time tandem queue with blocking (and failures), Bernoulli arrivals and phase-type services. They used the MAM approach also and assumed both waiting rooms to be finite. Daduna [4]

considered the case with an infinite waiting room for the first queue, but restricted himself to Bernoulli service processes. Desert and Daduna [5] focused on discrete-time tandem queues with state-dependent Bernoulli service rates and a state-dependent Bernoulli arrival stream at the first node. They obtained the joint sojourn time distribution for a customer traversing the tandem system under consideration. In our current paper we consider an infinite waiting room in front of the first server and allow Markovian arrivals (D-MAP), enabling us to model arrival processes that have some elements of correlations, which is more common in the telecommunication field where arrivals are usually bursty. Moreover, while Gün and Makowsky focus on the joint queue length distribution, we provide an algorithm for the total response time of a customer and address the stability issues raised by the infinite queue.

Other related works, in the sense that they focus on the response time as opposed to the joint queue length, are those by van der Mei et al. [22] and Knessl and Tier [13], who studied the first two moments of the response time in an open two-node queueing network with feedback for the case with an exponential processor sharing (PS) node and a FIFO node (while the arrivals at the PS node are Poisson). Chao and Pinedo [3] considered the case of two tandem queues with batch Poisson arrivals and no buffer space in the second queue. They allowed the service times to be general and obtained the expected time in system.

We start with a description of the model under consideration in Section 2, while the GI/M/1-type Markov chain constructed to obtain the performance measures is given in Section 3. An efficient method to compute the response time and the joint queue length distribution from the steady-state vector is presented in Section 4, whereas Section 5 addresses the stability issues surrounding our model. We end by demonstrating the strength of our model through a variety of numerical examples.

2. Model description

Consider two queues in tandem, where the first queue has an infinite waiting line and the second has a finite one with capacity B . Customers arrive (to the first queue) according to a discrete-time Markovian arrival process (D-MAP), characterized by the

$l \times l$ sub-stochastic matrices D_0 and D_1 . The matrix $D = D_0 + D_1$ is the stochastic matrix of the underlying Markov chain that governs the arrival process. The element $(D_k)_{i,j}$, $1 \leq i, j \leq l$, $k = 0, 1$, represents the probability of making a transition from state i to j with k arrivals. Let γ be the stationary distribution associated with D , then the arrival rate is given by $\lambda = \gamma D_1 \mathbf{1}_l$, where $\mathbf{1}_l$ is a $l \times 1$ column vector of ones. For more details on MAP see [2,16].

The service required by a customer in the i th queue is phase-type (PH) distributed with matrix representation (m_i, α_i, T_i) , for $i = 1, 2$. It is well known that PH distributions are very good for representing most of the types of services encountered in communication systems [14]. The mean service time of a PH is given as $\mu_i^{-1} = \alpha_i (I_{m_i} - T_i)^{-1} \mathbf{1}_{m_i}$, where I_x is an x -dimensional identity matrix. The matrix T_i is sub-stochastic and is of order m_i , while t_i is defined as $\mathbf{1}_{m_i} - T_i \mathbf{1}_{m_i}$. The elements $(\alpha_i)_s$ of the stochastic vector α_i represent the probability that a customer starts his service in phase s . Let $r_i^{(k)}$ be the probability that the service time at node i lasts for k or more units of time, then $r_i^{(k)} = \alpha_i T_i^{k-1} \mathbf{1}_{m_i}$, $k \geq 1$. Notice, the minimum service time at node i is at least 1 and the probability that the service time equals exactly k is found as $\alpha_i T_i^{k-1} t_i$. For more details on the phase-type distribution see [19].

Whenever a customer finishes service in the first queue it advances to the second queue (at no switching cost), unless the waiting line of the second queue is already fully occupied. In this case, the customer remains within the service facility of the first queue until there is a service completion in the second queue. Thereby, preventing any other customers waiting in the waiting line of queue 1 from entering the server (meaning, we adopt the *blocking-after-service* mechanism, see [20, p. 6]). Both queues serve their customers in a FCFS order. All events such as arrivals, transfers from a waiting line to the server and service completions are assumed to occur at instants immediately after the discrete time epochs. This implies, amongst others, that the age of a customer in service at some time epoch n is at least 1.

3. The GI/M/1-type Markov Chain

A Markov chain (MC) that allows us to efficiently obtain the response-time distribution of an arbitrary customer, is constructed next. The state space of this

MC will be subdivided into an infinite number of groups, called levels. Level zero will contain all the states that correspond to a situation in which the first server is idle. Whereas level i , for $i > 0$, reflects the fact that the first server is occupied by a customer of age i (either because he is being served or because he is blocked by the second queue). To be more specific, the states of level 0 are divided into 2 sets:

- $BI = \{j \mid 1 \leq j \leq l\}$: The MC is said to be in state j of the set BI at time n , if both servers are idle and the arrival process is in state j at time n .
- $FI = \{(b, s_2, j) \mid 0 \leq b \leq B, 1 \leq s_2 \leq m_2, 1 \leq j \leq l\}$: If the first server is idle, the second server is occupied by a customer whose service is in phase s_2 and b customers are waiting to be served by the second server, while the state of the D-MAP is j at time n , then the MC is said to be in state (b, s_2, j) of the set FI at time n .

The states of level i , for $i > 0$, are further subdivided into three sets:

- $SI = \{(s_1, j) \mid 1 \leq s_1 \leq m_1, 1 \leq j \leq l\}$: The MC is said to be in state (s_1, j) of the set SI , at time n , in case the second server is idle, an age i customer is served by server 1, the phase of his service equaling s_1 , while the arrival process is, at time $n - i + 1$, in state j .
- $BS = \{(b, s_1, s_2, j) \mid 0 \leq b \leq B, 1 \leq s_v \leq m_v, v = 1, 2; 1 \leq j \leq l\}$: The situation in which, at time n , a customer of age i is being served by server 1, there are b customers waiting for service in the second waiting line, the phase of service in the v th server equals s_v , for $v = 1, 2$, and the state of the D-MAP at time $n - i + 1$ equals j , will correspond to the state (b, s_1, s_2, j) of level i .
- $BL = \{(s_2, j) \mid 1 \leq s_2 \leq m_2, 1 \leq j \leq l\}$: The scenario where, at time n , there is an age i customer blocked in server 1, a customer is served by server 2, whose current phase is s_2 , and the state of the D-MAP at time $n - i + 1$ equals j is represented by state (s_2, j) of the set BL .

Let $|S|$ denote the number of elements in a set S . Define d_t and d_b as $|SI| + |BS| + |BL| = (B+1)m_1m_2l + (m_1 + m_2)l$ and $|BI| + |FI| = (B+1)m_2l + l$, respectively. As we shall explain later on, the transition matrix P

of this MC has the following form:

$$P = \begin{bmatrix} B_1 & B_0 & 0 & 0 & 0 & \dots \\ B_2 & A_1 & A_0 & 0 & 0 & \dots \\ B_3 & A_2 & A_1 & A_0 & 0 & \dots \\ B_4 & A_3 & A_2 & A_1 & A_0 & \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}, \quad (3.1)$$

In case (i), the number of customers in the second queue will either remain the same or decrease by one, depending on whether there is a service completion in server 2. In cases (ii) and (iii), the number of customers in the second waiting line has to remain equal to B , implying that there can be no service completion. As a result, we find $A_0 = K_0 \otimes I_l$, where

$$K_0 = \begin{bmatrix} T_1 & 0 & 0 & \dots & 0 & 0 & 0 \\ T_1 \otimes t_2 & T_1 \otimes T_2 & 0 & \dots & 0 & 0 & 0 \\ 0 & T_1 \otimes t_2 \alpha_2 & T_1 \otimes T_2 & \ddots & 0 & 0 & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \ddots & T_1 \otimes T_2 & 0 & 0 \\ 0 & 0 & 0 & \dots & T_1 \otimes t_2 \alpha_2 & T_1 \otimes T_2 & t_1 \otimes T_2 \\ 0 & 0 & 0 & \dots & 0 & 0 & T_2 \end{bmatrix}, \quad (3.2)$$

where the matrices A_k are $d_t \times d_t$ matrices, B_k , for $k \geq 2$, is a $d_t \times d_b$ matrix, B_0 a $d_b \times d_t$ and B_1 a square matrix of dimension d_b . Notice, the matrix A_k represents the transition probabilities of going from level $i > 0$ to level $i - k + 1$, for $k \leq i$, whereas the matrices B_k are related to transitions from and/or to level 0.

Let us now discuss the matrices A_k and B_k in detail, the structure of P will be apparent from this discussion. Assume the MC is in some state of level $i > 0$ at time n , meaning that an age i customer, referred to as customer c , is occupying server 1. In order to get a transition to level $i + 1$, customer c has to remain in server 1. This is because customers are served in a FCFS order and there are no batch arrivals; hence, the age of the very next customer who arrives after c cannot be larger than i at time $n + 1$. There are three scenarios that would cause customer c to remain in server 1: (i) his service (in server 1) did not finish at time n , (ii) his service finished, but he is blocked by the second queue, or (iii) customer c remains blocked.

I_k represents the identity matrix of dimension k , $t_i = \mathbf{1}_{m_i} - T_i \mathbf{1}_{m_i}$, for $i = 1, 2$; and $\mathbf{1}_k$ is a $1 \times k$ vector with each entry set to 1. Notice, this matrix does not depend on the age i of customer c .

Next, we consider the transitions from level i to $i - k + 1$, for $i \geq k \geq 1$. Thus, as before we have a customer, called c , occupying server 1 at time n . To get a transition to level $1 \leq i - k + 1 \leq i$, customer c has to leave server 1. Moreover, a new customer, whose age should equal $i - k + 1$ at time $n + 1$, has to enter the server at time n , call him customer c' . Meaning, the interarrival time between c and c' has to equal k . The fact that customer c leaves server 1 implies that the MC cannot make a transition to one of the states $SI \cup BL$ of level $i - k + 1$. Also, given that the waiting line of server 2 was fully occupied at time n , there should have been a service completion in server 2 (otherwise c would become/remain blocked). Finally, if there still was a vacancy in the waiting line of the second queue at time n , the number of waiting customers there either increases by one or remains the same, depending on whether there is a service completion in server 2. Hence, transitions from level i to $i - k + 1$ are governed by the matrix $A_k = K_1 \otimes D_0^{k-1} D_1$, where

$$K_1 = \begin{bmatrix} 0 & t_1 \alpha_1 \otimes \alpha_2 & 0 & \dots & 0 & 0 & 0 \\ 0 & t_1 \alpha_1 \otimes t_2 \alpha_2 & t_1 \alpha_1 \otimes T_2 & \dots & 0 & 0 & 0 \\ 0 & 0 & t_1 \alpha_1 \otimes t_2 \alpha_2 & \ddots & 0 & 0 & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \ddots & t_1 \alpha_1 \otimes t_2 \alpha_2 & t_1 \alpha_1 \otimes T_2 & 0 \\ 0 & 0 & 0 & \dots & 0 & t_1 \alpha_1 \otimes t_2 \alpha_2 & 0 \\ 0 & 0 & 0 & \dots & 0 & \alpha_1 \otimes t_2 \alpha_2 & 0 \end{bmatrix}. \quad (3.3)$$

The matrices B_k , for $k \geq 0$, describing the transitions to and from level 0 can be obtained using similar arguments (see [11] for details). B_k , for $k > 1$, equals $K_1^c(\alpha_1 \leftarrow 1) \otimes D_0^{k-1}$, where $K_1^c(\alpha_1 \leftarrow 1)$ is obtained from K_1 by replacing all vectors α_1 by the scalar 1 and removing the last m_2 columns. The matrix B_1 can be written as $K_0^r(T_1 \leftarrow \alpha_1, t_1 \leftarrow 0) \otimes D_1$. To obtain $K_0^r(T_1 \leftarrow \alpha_1, t_1 \leftarrow 0)$ we replace all matrices T_1 appearing in the expression for K_0 into the vector α_1 , the vector t_1 into the scalar 0 and remove the last m_2 rows of K_0 . Finally, the matrix B_1 obeys the following equation:

$$B_1 = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 & 0 \\ t_2 & T_2 & 0 & \dots & 0 & 0 \\ 0 & t_2\alpha_2 & T_2 & \ddots & 0 & 0 \\ \vdots & \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & 0 & 0 & \ddots & T_2 & 0 \\ 0 & 0 & 0 & \dots & t_2\alpha_2 & T_2 \end{bmatrix} \otimes D_0. \tag{3.4}$$

Let $\pi = [\pi_0, \pi_1, \pi_2, \dots]$ be the steady-state vector of P , where π_0 is a $1 \times d_b$ vector and π_i , for $i > 0$, a $1 \times d_t$ vector. Since, the MC characterized by P is a GI/M/1-type MC, π exists if and only if $\theta\beta > 1$, where $\theta = \sum_{k \geq 0} A_k = \theta$, $\theta \mathbf{1}_{d_t} = 1$ and $\beta = \sum_{k \geq 1} k A_k \mathbf{1}_{d_t}$. The stability condition $\theta\beta > 1$ is discussed in Section 5. The steady-state vector of a GI/M/1-type MC can be found by iteratively solving for the minimal non-negative R in the non-linear equation $R = \sum_{k \geq 0} R^k A_k$ (see [19]). Having solved this equation we find π as

$$\begin{aligned} & [\pi_0, \pi_1] \\ &= [\pi_0, \pi_1] \\ & \times \begin{bmatrix} B_1 & B_0 \\ \sum_{k \geq 1} R^{k-1} B_{k+1} & \sum_{k \geq 1} R^{k-1} A_k \end{bmatrix}, \tag{3.5} \end{aligned}$$

$$\pi_i = \pi_{i-1} R, \tag{3.6}$$

where $i > 1$, π_0 and π_1 are normalized as $\pi_0 \mathbf{1}_{d_b} + \pi_1 (I - R)^{-1} \mathbf{1}_{d_t} = 1$. However, the computation of π can be done much more efficiently (both in terms of the time and memory complexity) by constructing a quasi-Birth–Death (QBD) Markov chain and applying the cyclic reduction algorithm [17] to solve this QBD. We refer to [11] for details on the QBD reduction.

4. Performance measures

In this section we demonstrate how to get the response-time distribution from the steady-state vector π . We start by introducing the following set of random variables:

- T_{W_i} : The amount of time a tagged customer has to wait in the i th waiting line, for $i = 1, 2$.
- T_{S_i} : The service time duration of a tagged customer in server i , for $i = 1, 2$.
- T_B : The time that elapses while a tagged customer is blocked in server 1.

Having defined these variables the total response time T_R of a tagged customer is defined as $T_{W_1} + T_{S_1} + T_B + T_{W_2} + T_{S_2}$, while the response time in queue 1, denoted as T_{R_1} , equals $T_{W_1} + T_{S_1} + T_B$. We need two more variables before we can proceed:

- F_N : The number of customers still requiring full service by server 2 before a tagged customer who just left server 1 can start his service.
- F_T : The remaining service time of the customer occupying server 2 when a tagged customer leaves server 1.

Write π_i as $[\pi_i^{SI}, \pi_i^{BS}, \pi_i^{BL}]$ in accordance with the three sets of states of level i , then

$$\begin{aligned} & P[T_{R_1} = r, F_N = b, F_T = h] \\ &= \sum_{s_1} \frac{(t_1)_{s_1}}{\lambda} \left\{ 1_{\{b=0 \& h=0\}} \left(\sum_j \pi_r^{SI}(s_1, j) \right) \right. \\ & \quad \left. + 1_{\{b < B \mid h=0\}} \right. \\ & \quad \left. \times \left(\sum_{s_2, j} \pi_r^{BS}(b, s_1, s_2, j) (T_2^h t_2)_{s_2} \right) \right\} \\ & \quad + \frac{1_{\{b=B \& h=0\}}}{\lambda} \left(\sum_{s_2, j} \pi_r^{BL}(s_2, j) (t_2)_{s_2} \right), \end{aligned}$$

where $(x)_i$ denotes the i th component of the vector x . From this it is straightforward to compute the probabilities $P[T_{R_1} + F_T = r, F_N = b]$. The total

response-time distribution is then found by

$$P[T_R = i] = \sum_b \sum_{r \leq i-b} P[T_{R_1} + F_T = r, F_N = b] \times P[S_2^{(*b+1)} = i - r], \tag{4.7}$$

where $S_2^{(*x)}$ is the x -fold convolution of the PH service time distribution of server 2 characterized by (m_2, α_2, T_2) . The blocking probability p_{BL} can be computed as

$$p_{BL} = \frac{1}{\lambda} \sum_{r>0} \sum_{s_1, s_2, j} \pi_r^{BS}(B, s_1, s_2, j)(t_1)_{s_1}(T_2)_{s_2}, \tag{4.8}$$

and the blocking time distribution as

$$P[T_B = t] = \frac{1}{\lambda} \sum_{r>0} \sum_{s_1, s_2, j} \pi_r^{BS}(B, s_1, s_2, j)(t_1)_{s_1} \times (T_2^t)_{s_2}, \tag{4.9}$$

for $t > 0$ and $P[T_B = 0] = 1 - p_{BL}$. Finally, notice that the probability of having a joint queue contents equal to (q_1, q_2) , equals the probability of having a customer of age a in the first service center, q_2 customers in the intermediate queue and having q_1 arrivals during a time interval of length a (that starts immediately after the customer in the first service center arrived). Hence, a simple procedure can be devised to obtain the joint queue contents distribution from the steady state probability vector π (see [11]).

5. Stability condition

The GI/M/1-type Markov chain introduced in Section 3 is ergodic if and only if $\theta\beta > 1$, where $\theta \sum_{k \geq 0} A_k = \theta$, $\theta \mathbf{1}_{d_1} = 1$ and $\beta = \sum_{k \geq 1} k A_k \mathbf{1}_{d_1}$. The following theorem summarizes the stability results (see [11] for details).

Theorem 5.1. *The MC developed in Section 3 is stable if and only if*

- (1) $\lambda/k < 1$, where λ is the D-MAP mean arrival rate and

$$k = \mu_1(1 - \kappa_{BL} \mathbf{1}_{m_2}) = \mu_2(1 - \kappa_{SI} \mathbf{1}_{m_1}), \tag{5.10}$$

with μ_i the service rate of the i th service center and $\kappa = [\kappa_{SI}, \kappa_0, \kappa_1, \dots, \kappa_B, \kappa_{BL}]$ is the invariant probability vector of $K = K_0 + K_1$. Thus, the arrival process only affects the stability through its mean arrival rate λ .

- (2) the queue contents process of the infinite waiting line of queue 1 has a steady state.

Remark 1. Eq. (5.10) suffices to prove that interchanging both service-time distributions does not affect the stability of the system. Intuitively, when studying the stability, we may assume that there are always customers ready to be served in the infinite waiting line. Now, using an argument by Melamed [18], we can regard the empty places (holes) in the intermediate buffer as dual customers. For each regular customer that moves through the system, a dual customer receiving identical service moves in the opposite direction. Hence, the maximum stable throughput of the regular customers and the dual customers must be identical. Clearly, the dual system is identical to the interchanged system. This type of interchangeability result was already established long ago for exponential servers (and Poisson arrivals) [19, Section 5.2].

Remark 2. In the special case where the size of the intermediate waiting line B equals 0, one can prove that $1/k$ is nothing but the mean of the maximum of both PH service time distributions.

6. Numerical results

In this section we present some numerical examples that provide insight on the system behavior.

6.1. Influence of the capacity B and the correlation of the D-MAP arrival process

Consider an interrupted Bernoulli process (IBP), where the mean sojourn time in both states equals x_c and an arrival occurs in the on-state with probability x_p . Thus,

$$D_0 = \begin{bmatrix} 1 - 1/x_c & 1/x_c \\ (1 - x_p)/x_c & (1 - x_p)(1 - 1/x_c) \end{bmatrix}, \tag{6.11}$$

$$D_1 = \begin{bmatrix} 0 & 0 \\ x_p/x_c & x_p(1 - 1/x_c) \end{bmatrix}.$$

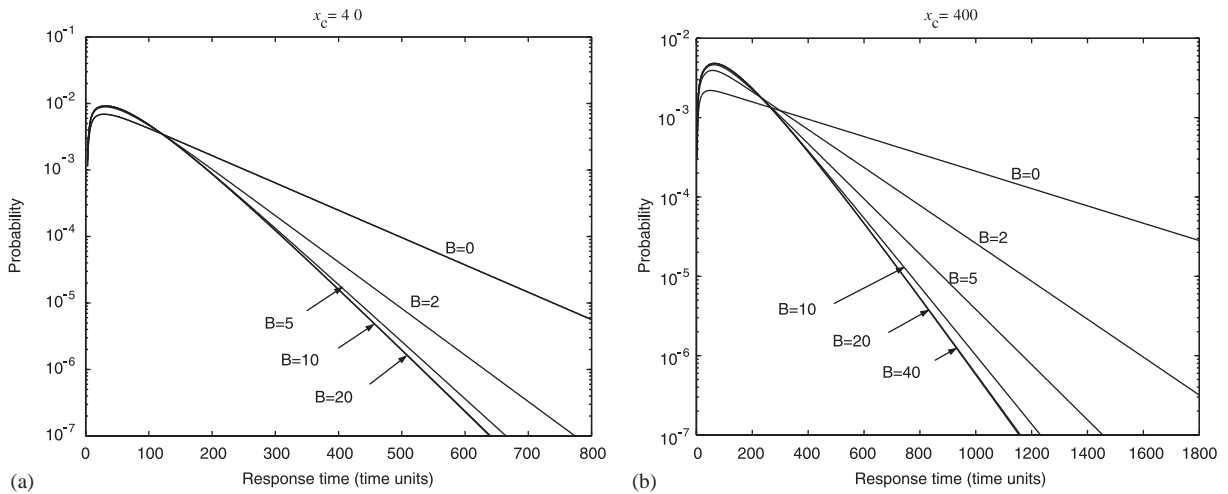


Fig. 1. Response-time distribution for various capacities B , $x_p = \frac{1}{16}$ and (a) $x_c = 40$, (b) $x_c = 400$.

The service time distribution in server 1 is hypergeometric with parameters:

$$\alpha_1 = [\frac{1}{3}, \frac{2}{3}], \quad T_1 = \begin{bmatrix} \frac{4}{5} & 0 \\ 0 & \frac{19}{20} \end{bmatrix}. \quad (6.12)$$

Thus, with probability $\frac{1}{3}$ and $\frac{2}{3}$ the service time is geometrically distributed with a mean of 5 and 20 time units, respectively. The mean of this distribution is 15 time units. The service-time distribution of the second server is characterized by

$$\alpha_2 = [\frac{1}{16}, \frac{1}{8}, \frac{13}{16}], \quad T_2 = \begin{bmatrix} \frac{3}{4} & \frac{1}{8} & 0 \\ \frac{1}{10} & \frac{4}{5} & 0 \\ 0 & 0 & x_l \end{bmatrix}, \quad (6.13)$$

where $x_l = 0.93887$ is chosen such that the mean service time equals 15. We have chosen the mean service time identical in both servers as this ought to create a strong coupling between both queues.

Fig. 1 depicts the response-time distribution for various capacities B , for $x_p = \frac{1}{16}$ (meaning that the arrival rate $\lambda = \frac{1}{32}$), and x_c either 40 or 400. Clearly, the larger x_c the more correlated the arrival process. Obviously, stronger correlated arrivals give rise to slower response times, while adding more capacity B between both servers reduces the response time. However, at some point there is little use in further augmenting the capacity as the response time seems to converge for B large. This is easily understood as the blocking prob-

ability tends to decrease to zero while increasing B . The figure further illustrates that the rate of convergence is affected by the correlation of the arrival process: stronger correlated arrival processes more easily justify increasing the capacity B .

6.2. The maximum arrival rate λ and the variation of the service times

We consider the same IBP process as in the previous section. The service-time distribution is either geometric (Geo), Erlang-5 (Er5)¹ or hypergeometric (HypGeo) characterized by

$$\alpha_1 = [\frac{9}{10}, \frac{1}{10}], \quad T_1 = \begin{bmatrix} \frac{4}{5} & 0 \\ 0 & \frac{104}{105} \end{bmatrix}. \quad (6.14)$$

The mean of each of these service-time distributions is 15 time units. The HypGeo distribution is the most variable of the three, followed by Geo and Er5.

Fig. 2 shows the maximum stable load ρ_{\max} , defined as $E[S_1]\lambda_{\max} = 15\lambda_{\max}$, as a function of B for different server configurations. λ_{\max} is the maximum arrival rate λ for which the system is stable. Recall, $\lambda_{\max} = \mu_1(1 - \kappa_{BL}\mathbf{1}_{m_2}) = \mu_2(1 - \kappa_{SI}\mathbf{1}_{m_1})$, see Section 5. The notation (X, Y) indicates that the service time in servers 1 and 2 is distributed as X and Y , respectively.

¹ That is, the sum of 5 independent and identically distributed geometric random variables.

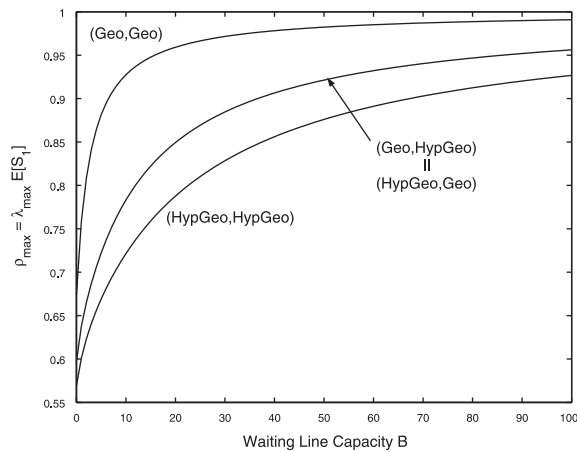


Fig. 2. Maximum stable load ρ_{\max} as a function of the capacity B for different server configurations.

Fig. 2 clearly demonstrates that more variable service times, that is, HypGeo, give rise to a lower maximum stable input rate λ_{\max} . This result stems from the fact that a more variable service time, whether in the first or second server, causes a higher degree of blocking in comparison with a more deterministic service time distribution. As proved in Section 5, interchanging both service time distributions does not alter the system stability. Notice, the actual nature of the D-MAP arrival process is irrelevant as the stability is only affected by the arrival process through its mean.

The response-time distribution for different server configurations is depicted in Fig. 3. We assume that arrivals occur according to a IBP with $x_c = 40$ and $x_p = \frac{1}{16}$, see Section 6.1. The capacity of the intermediate waiting line B is assumed to be 10. Fig. 3 confirms that the response of the system slows down as the service times become more variable. It further demonstrates that interchanging the service times generally causes a (limited) change in the response time (as opposed to the stability). Placing the more variable server first seems to result in a somewhat slower response. This might be explained by noticing that the output process of server 1 is more bursty in such case, causing a higher blocking probability. On the other hand, less variability in server 2 decreases the blocking probability, so when we interchange both service time distributions both these effects influence the degree of blocking in the system. Various numerical experi-

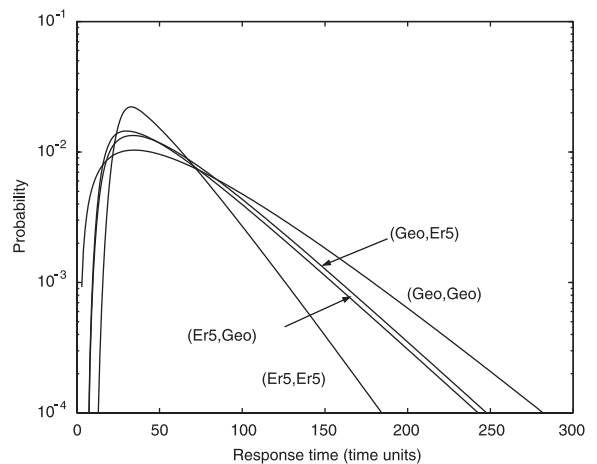


Fig. 3. Response-time distribution for different server configurations for $B = 10$ and IBP arrivals ($x_c = 40$, $x_p = \frac{1}{16}$).

ments, including Fig. 3, seem to indicate that reducing the variation of server 1 should be slightly favored.

References

- [1] B. Avi-Itzhak, A sequence of service stations with arbitrary input and regular service times, *Manage. Sci.* 11 (5) (1965) 565–571.
- [2] C. Blondia, A discrete-time batch markovian arrival process as B-ISDN traffic model, *Belg. J. Oper. Res. Statist. Comput. Sci.* 32 (3,4) (1993).
- [3] X. Chao, M. Pinedo, Batch arrivals to a tandem queue without an intermediate buffer, *Stochastic Models* 6 (4) (1990) 735–748.
- [4] H. Daduna, *Queueing Networks with Discrete Time Scale*, Springer, Berlin, 2001.
- [5] B. Desert, H. Daduna, Discrete time tandem networks of queues: effects of different regulation schemes for simultaneous events, *Performance Evaluation* 47 (2002) 73–104.
- [6] A. Gómez-Corral, A tandem queue with blocking and markovian arrival process, *Queueing Systems* 41 (2002) 343–370.
- [7] L. Gün, Annotated bibliography of blocking systems. Technical Report, Institute for Systems Research ISR 1987-187, 1987.
- [8] L. Gün, M. Makowski, Matrix geometric solution for finite capacity queues with phase-type distributions, in: *Proceedings of Performance 87*, Brussels, 1987, (pp. 269–282).
- [9] L. Gün, M. Makowski, Matrix geometric solution for two node tandem queueing systems with phase-type servers subject to blocking and failures. Technical Report, Institute for Systems Research, ISR 87-210, 1987.

- [10] N.G. Hall, C. Sriskandarajah, A survey of machine scheduling problems with blocking and no-wait in process, *Operations Research* 44 (1996) 510–525.
- [11] B. van Houdt, A.S. Alfa, Response time in a tandem queue with blocking, markovian arrivals and phase-type services: extended version. Technical Report, TR04/01, University of Antwerp, 2004.
- [12] G.C. Hunt, Sequential arrays of waiting lines, *Operations Research* 4 (6) (1956) 674–683.
- [13] C. Knessl, C. Tier, Approximation to the moments of the sojourn time in a tandem queue with overtaking, *Stochastic Models* 6 (3) (1990) 499–524.
- [14] A. Lang, J.L. Arthur, Parameter approximation for phase-type distributions, in: S.R. Chakravarty, A.S. Alfa (Eds.), *Matrix-Analytic Methods in Stochastic Models*, Marcel-Dekker, New York, 1996, pp. 151–206.
- [15] G. Latouche, M.F. Neuts, Efficient algorithmic solutions to exponential tandem queues with blocking, *SIAM J. Algebraic and Discrete Meth.* 1 (1) (1980) 93–106.
- [16] D.M. Lucantoni, New results on the single server queue with a batch markovian arrival process, *Stochastic Models* 7 (1) (1991) 1–46.
- [17] B. Meini, Solving QBD problems: the cyclic reduction algorithm versus the invariant subspace method, *Advances in Performance Analysis* 1 (1998) 215–225.
- [18] B. Melamed, A note on the reversibility and duality of some tandem blocking queueing systems, *Management Science* 32 (12) (1986) 1648–1650.
- [19] M.F. Neuts, *Matrix-Geometric Solutions in Stochastic Models, An Algorithmic Approach*, John Hopkins University Press, MD, 1981.
- [20] H. Perros, *Queueing Networks with Blocking*, Oxford University Press, New York, 1994.
- [21] M. Schamber, Decomposition methods for finite queue networks with a non-renewal arrival process in discrete time, thesis, University of Manitoba, 1997.
- [22] R.D. van der Mei, B.M.M. Gijsen, N. in't Veld, J.L. van den Berg, Response times in a two-node queueing network with feedback, *Performance Evaluation* 49 (2002) 99–110.