



An analytic finite capacity queueing network model capturing blocking, congestion and spillbacks *

C. Osorio M. Bierlaire [†]

June 1st, 2007

Report TRANSP-OR 070604
Transport and Mobility Laboratory
School of Architecture, Civil and Environmental Engineering
Ecole Polytechnique Fédérale de Lausanne
transp-or.epfl.ch

*This research is supported by the Swiss National Science Foundation grant 205321-107838

[†]École Polytechnique Fédérale de Lausanne, Transport and Mobility Laboratory, Station 18, CH-1015 Lausanne, Switzerland. E-mail: {carolina.osoriopizano, michel.bierlaire}@epfl.ch

An analytic finite capacity queueing network model capturing blocking, congestion and spillbacks *

C. Osorio

M. Bierlaire [†]

June 1st, 2007

Abstract

Analytic queueing network models often assume infinite capacity for all queues. For real systems this infinite capacity assumption does not hold, but is often maintained due to the difficulty of grasping the between-queue correlation structure present in finite capacity networks. This correlation structure helps explain bottleneck effects and spillbacks, the latter being of special interest in networks containing loops because they are a source of potential deadlock. We present an analytic queueing network model which acknowledges the finite capacity of the different queues. By explicitly modeling the blocking phase the model yields a description of the congestion effects. The model is adapted for multiple server finite capacity queueing networks with an arbitrary topology and blocking-after-service. A decomposition method allowing the evaluation of the model is described. The method is validated, by comparison to both pre-existing methods and simulation results. A real application to the study of patient flow in a network of operative and post-operative units of the Geneva University Hospital is also presented.

1 Introduction

Modeling complex systems using analytic queueing network models allows us to better understand their behavior, to evaluate and ultimately to improve their performance. The most researched queueing network model is the Jackson network model (Jackson, 1957, 1963) which assumes infinite capacity for all queues. For real systems this infinite capacity assumption does not hold, but is often maintained due to the difficulty of grasping the between-queue correlation structure present in finite capacity networks; e.g. acknowledging the links between the behavior of adjacent queues where chained events can take place. Consider a network of operative and post-operative hospital units where each unit is modeled as a specific queue and where it is the patient flow that is of main interest. For such a network understanding the correlation between the occupation of the different units (e.g. surgical

*This research is supported by the Swiss National Science Foundation grant 205321-107838

[†]École Polytechnique Fédérale de Lausanne, Transport and Mobility Laboratory, Station 18, CH-1015 Lausanne, Switzerland. E-mail: {carolina.osoriopizano, michel.bierlaire}@epfl.ch

intensive care, surgical intermediate care) can help avoid bed blocking and improve a patients recovery procedure. More generally, the between-queue correlation explains bottleneck effects and spillbacks, the latter being of special interest in networks containing loops because they are a source of potential deadlocks (also known as gridlocks) (Daganzo, 1996). This correlation structure is of importance in a variety of networks such as manufacturing networks (Papadopoulos and Heavey, 1996), software architecture networks (Balsamo et al., 2003), circulation systems (e.g. corridors) (Cheah and Smith, 1994) and even prison networks (Korporaal et al., 2000). In order to capture this correlation and to estimate these congestion effects we resort to models with finite capacities. When wanting to model real scale finite capacity networks with arbitrary topologies the main complexity lies in appropriately acknowledging the between-queue correlation while also maintaining a tractable model. We propose a finite capacity queueing network (FCQN) model capturing the between-queue correlation based on a decomposition method which allows the analysis of real scale networks. The intended use of this model is within an optimization framework. We therefore resort to an analytic model which unlike pre-existing methods preserves the network topology and its configuration (number of queues and their capacities) as static parameters. This makes this approach suitable for an optimization framework. We explicitly model the blocking phase within our analytical approach, yielding performance measures such as the probability distribution of the number of blocked jobs in a queue. This paper is structured as follows. We describe the FCQN framework and then review the existing analysis methods. The proposed model and approximation method are then described, followed by their validation versus both pre-existing methods and simulation results. The method is then applied to a real case study.

2 General Framework

In this section we describe the queueing network framework.

A queueing network is composed of a set of linked queues, hereafter called stations. Of interest is the study of the flow of “jobs” throughout the network. A job is the generic name for the units of interest, e.g. a pedestrian, a prisoner, a patient. We consider open queueing networks where jobs are allowed to leave the network and where the external arrivals arise from an infinite population of jobs. We now describe the general process that a job goes through upon arrival to a station. Jobs arriving to a station are either served immediately or queue until a server becomes available. Once a job is served it is routed to its next station, which is chosen according to a probability distribution. If this destination station has finite capacity then it may be full. If it is full then the job will be **blocked** at its current station until a server becomes available at the destination station. Various blocking mechanisms, which are at the heart of spillbacks, have been defined in the literature (Balsamo et al., 2001). They differ either in the moment the job is considered to be blocked (e.g. before or after service) or in the routing mechanism of blocked jobs. The blocking mechanism that we have just described is known as blocking-after-service (BAS). The jobs are unblocked with a First In First Out (FIFO) mechanism. The average arrival rate to station i is denoted λ_i . Station i has c_i parallel servers, each one serves with an average rate μ_i . The total number of jobs allowed in the station is called the capacity of the station, K_i , the buffer size is $K_i - c_i$. The possible routings among stations are given by the transition probability matrix (p_{ij}) , where p_{ij} denotes the probability that a job at station i is routed to station j .

3 Literature review

In this section we describe the existing methods allowing the analysis of FCQN models.

A first survey of FCQN models was made by Perros (1984), who later on also wrote a historical

overview of the research motivations and advances in networks with blocking (Perros, 2003). A detailed introductory book was written by Balsamo et al. (2001). Surveys focusing on specific application fields are given for software architecture performance (Balsamo, De Nitto Persone and Inverardi, 2003), for the production and manufacturing sector (Papadopoulos and Heavey, 1996) and on retrieval queues for the telecommunications sector (Artalejo, 1999).

3.1 Exact methods

The joint stationary distribution of the network, which contains the probability of each possible state of the network, allows us to derive the main network performance measures. Exact analysis of FCQN models, that is exact evaluation of this joint distribution, can be obtained either in analytic closed form or numerically. For an open Jackson network the joint stationary distribution has a product form, thus the stations behave as if they were independent. For a FCQN the between-station correlation suggests a non-product form stationary distribution, thus exact analysis of FCQN models are limited to very small networks.

Closed form analytic expressions for the joint distribution are difficult to obtain and are only available for specific topologies such as single server two or three station tandem topologies (Grassman and Derkic, 2000; Langaris and Conolly, 1984; Latouche and Neuts, 1980; Konheim and Reiser, 1978; Konheim and Reiser, 1976) or two station closed networks (Akyildiz and von Brand, 1994; Balsamo and Donatiello, 1989).

On the other hand exact numerical evaluation of the joint stationary distribution can be obtained by solving the global balance equations. A detailed description of these numerical methods can be found in Stewart (1999). These equations require the construction of the transition rate matrix, i.e. the description of the transition rates between all feasible states of the network. This time consuming task is therefore only conceivable for small networks (i.e. small in the number of stations and their capacity). This approach also lacks flexibility because changes in the network topology require redefining the transition rate matrix. If the networks of interest have a more general topology or an arbitrary size then their analysis is done by approximation methods.

3.2 Approximation methods

Approximation methods can be classified into either analytic approaches or simulation-based methods.

The use of disaggregate models based on simulation is the most popular approach to evaluate the performance of a finite capacity queueing networks. Surveys of simulation models exist for sectors such as transportation (Nagel, 2002; Ben-Akiva et al., 2001), healthcare (Fone et al., 2003; Jun et al., 1999), computer science (Sadoun, 2000; Obaidat, 1990) and the analysis of call centers (Koole and Mandelbaum, 2002; Mandelbaum, 2001). This approach although more realistic and detailed, can be cumbersome to optimize, and its accuracy is strongly dependent on the quality of the calibration data (Korporaal et al., 2000). Analytic models are simpler, less data expensive, more flexible and more suited for an optimization framework (Cochran and Bharti, 2006).

The main motivation of analytic approximation methods is to reduce the dimensionality of the system under study. Decomposition methods achieve this by decomposing the network into subnetworks and analyzing each subnetwork in isolation. The structural parameters of each subnetwork (e.g. average arrival and service rates) depend on the state of other subnetworks and thus capture the correlation with other subnetworks. The main difficulty lies in obtaining good approximations for these parameters so that the stationary distribution of the subnetwork is a good estimate of its marginal stationary

distribution. Given a subnetwork its stationary distribution can be obtained by either establishing a behavioral analogy with a network whose distribution has a closed (and often product) form, or by exact numerical evaluation of the global balance equations which now have a smaller dimension but are often non-linear.

Existing decomposition methods have analyzed small subnetworks consisting of single stations, pairs of stations or triplets. If not stated otherwise the methods concern open finite capacity networks with exponentially distributed service times. The most commonly used decomposition method is single station decomposition, which dates back to the work of Hillier and Boling (1967) who considered tandem single server networks. One of the most used approaches concerns single server feed-forward networks where each station is modeled as an $M/M/1$ station (Takahashi, Miyahara and Hasegawa, 1980). An extension of this method to multiple servers (i.e. $M/M/c$ stations) is given by Koizumi et al. (2005). Here the buffers are considered infinite for each isolated station and their average queue length updates the capacity of the predecessor stations. This approximation holds if the capacity of adjacent predecessor stations can accommodate this average queue length. This constraint is checked only a posteriori. Each station is an $M/M/c$ queue for which closed form expressions of the performance measures exist. A method applicable to networks with an arbitrary topology is given by Korporaal et al. (2000). The individual stations are modeled as $M/M/c/K$ stations for which closed form performance measures are used. As for the method of Koizumi et al. (2005) the capacity of the stations are revised and the validity of these capacity adjustments are verified a posteriori.

The Expansion method (Kerbach and Smith, 1988, 1987), was developed for networks of $M/M/1/K$ stations. Here a network reconfiguration expands all finite capacity stations to artificial infinite capacity holding stations, which register the blocked jobs. This method was later extended to multiple servers and applied to pedestrian traffic flows by Cheah and Smith (1994). Gupta and Kavusturucu (2000) applied this method to production feed-forward systems, where service interruptions are allowed. Singh and Smith (1997) used it to evaluate network performance measures within a buffer allocation problem. A similar transformation where all $GE/GE/c/K$ stations are transformed into $GE/GE/c$ stations, and thus the joint distribution is approximated by a product form joint distribution, was proposed by Tahilramani, Manjunath and Bose (1999). Single server networks with phase-type service distributions have been proposed for tandem (Altiok, 1982) and feed-forward topologies (Altiok and Perros, 1987), with phase-type service distributions. Jun and Perros (1988) have extended this work to an arbitrary topology and have also considered general service times for an open tandem network in Jun and Perros (1990). The use of a phase-type service distribution accounts for all possible blockings but, as stated in Altiok and Perros (1987), it requires the construction of very detailed phase-type service mechanisms, which is a “cumbersome” and CPU time consuming task for large networks. In these methods queue capacity is also augmented in order to allow for storage of all predecessor station capacities.

Few authors have considered subnetworks larger than single stations. Two-station decomposition methods have been proposed for open tandem networks (Alfa and Liu, 2004; Brandwajn and Jow, 1988; Brandwajn and Jow, 1985) and for an arbitrary topology (Lee et al., 1998). Pairwise decomposition was used by van Vuuren, Adan and Resing-Sassen (2005) to study multi-server tandem stations with generally distributed service times. As an extension of the work by Brandwajn and Jow (1988), Schmidt and Jackman (2000) proposed a three-station decomposition method for a single server arbitrary topology network. Subnetworks consisting of more than one station can theoretically provide more accurate results than single station decomposition, but are computationally more intensive (Perros, 1994).

Recent methods, such as those of Koizumi et al. (2005) and Korporaal et al. (2000), have extended the

use of decomposition algorithms mainly to multiple server networks with an arbitrary topology. Nevertheless in order to acknowledge the finite capacity property of these networks the existing methods either revise station capacities or vary the network topologies. The revision of the station capacities renders them dynamic parameters. Moreover, approximations need to be used to ensure their integrality and their positivity is only checked a posteriori. We believe that an optimization-friendly model is one that preserves the network topology and its configuration (number of stations and their capacities) as static parameters. We are also interested in explicitly modeling the blocking phase within our analytical approach, yielding performance measures such as the probability distribution of the number of blocked jobs in a station. Since we have not found methods with these characteristics we have developed the method that we shall now describe.

4 Method

In this section we describe the decomposition method that allows the analysis of a network with finite capacity queues. The model accounts for multiple server queues with an arbitrary topology and blocking-after-service. The method is based on a decomposition of the network into single stations whose structural parameters are approximated so that they can account for the between-station correlation. The general process that a job goes through upon arrival to a station has been described in Section 2. In this paper we are interested in explicitly modeling the blocking phase that a job may go through in a finite capacity network. Thus we now describe in more detail how a job is processed. A job:

1. arrives to a station,
2. waits if all the servers are occupied,
3. is served (this is called the active phase),
4. is blocked if its destination station is full (this is called the blocked phase),
5. leaves the station.

Let $\pi(i)$ denote the stationary distribution of the isolated station i . The main aim of our method is to appropriately approximate $\pi(i)$ so that it is a good estimate of the marginal stationary distribution of station i . $\pi(i)$ can be obtained via the global balance equations along with the use of a normalizing constraint:

$$\begin{cases} \pi(i)Q(i) = 0, \\ \sum_{s \in \mathcal{S}(i)} \pi(i)_s = 1, \end{cases} \quad (1)$$

where $\pi(i)_s$ denotes element number s of $\pi(i)$. The global balance equations involve the state space of station i , $\mathcal{S}(i)$, as well as the transition rate matrix, $Q(i)$. We now define these two elements.

State space, $\mathcal{S}(i)$

The state of station i is described by the number of active jobs A_i , blocked jobs B_i and waiting jobs W_i .

$$\mathcal{S}(i) = \{(A_i, B_i, W_i) \in \mathbb{N}^3, A_i + B_i \leq c_i, A_i + B_i + W_i \leq K_i\}$$

Of interest in the validation runs that will be presented in Section 5 are bufferless stations, ($K_i = c_i$), where the state space reduces to: $\mathcal{S}(i) = \{(A_i, B_i) \in \mathbb{N}^2, A_i + B_i \leq c_i\}$. We denote by $\text{card}(\mathcal{S}(i))$ the cardinal or dimension of the state space.

4.1 Transition rate matrix, $Q(i)$

$Q(i)$ contains the transition rates between all pairs of states in $\mathcal{S}(i)$. Hereafter all rates are average rates. The non diagonal elements, $Q(i)_{sk}$ $s \neq k$, represent the average rate at which the transition between state s and k takes place. The diagonal elements are defined as: $Q(i)_{ss} = -\sum_{k \neq s} Q(i)_{sk}$. Thus $-Q(i)_{ss}$ represents the rate of departure from state s . Each equation of the system of global balance equations can be written as:

$$\sum_{k \in \mathcal{S}(i)} \pi(i)_k Q(i)_{ks} = -\pi(i)_s Q(i)_{ss},$$

it therefore balances the inflow and the outflow for a given state s . We define $Q(i)$ as a function of the following structural parameters:

- the average arrival rate to station i , λ_i ,
- the average service rate of a server at station i , μ_i ,
- the average probability of being blocked at station i , P_i^f .
- the average unblocking rate given that there are b blocked jobs at station i , $\tilde{\mu}(i, b)$,

These four parameters will allow us to describe the transition rates between the different states of station i . We can write $Q(i) = f(\lambda_i, \mu_i, P_i^f, \tilde{\mu}(i, b))$, where μ_i is an endogenous parameter whereas λ_i , $\tilde{\mu}(i, b)$, and P_i^f are exogenous.

To define f let us consider a state s such that $(A_i, B_i, W_i) = (a, b, w)$. The possible transitions with their corresponding rates are displayed in Table 1. The set of possible states to where a transition can take place are tabulated in the second column, the corresponding transition rate is in the third column and the conditions under which such a transition can take place are in the last column. We

initial state s	new state k	rate $Q(i)_{sk}$	condition
(a, b, w)	$(a + 1, b, w)$	λ_i	$a + b + 1 \leq c_i$
(a, b, w)	$(a, b, w + 1)$	λ_i	$a + b == c_i$ & $w + 1 \leq K_i - c_i$
(a, b, w)	$(a - 1, b, w)$	$a\mu_i(1 - P_i^f)$	$w == 0$
(a, b, w)	$(a, b, w - 1)$	$a\mu_i(1 - P_i^f)$	$w \geq 1$
(a, b, w)	$(a - 1, b + 1, w)$	$a\mu_i P_i^f$	always possible
(a, b, w)	$(a, b - 1, w)$	$\tilde{\mu}(i, b)$	$w == 0$
(a, b, w)	$(a + 1, b - 1, w - 1)$	$\tilde{\mu}(i, b)$	$w \geq 1$

Table 1: Transition rates of station i .

now describe the contents of this table. The first two lines of the table distinguish between an arrival that can be served immediately and an arrival that must queue before being served. The next two lines concern the completion of a service (the active phase) that is not followed by a blocking phase, in the first case the freed server remains available whereas in the second case the freed server immediately starts serving a job that was in the queue. The fifth line concerns jobs that have completed their service and become blocked. The last two lines relate to the completion of the blocking phase and differ in whether the server that was blocked stays available or immediately starts serving a queued job.

As emphasized by Korporaal et al. (2000), the main challenge of decomposition methods is to appropriately approximate these structural parameters so that $\pi(i)$ is a good estimate of the marginal stationary distribution of station i . The main complexity lies in appropriately capturing the correlation between the stations via these structural parameters. We now describe how our method revises the structural parameters in order to capture this correlation.

4.1.1 Arrival rate, λ_i

We model each station as a two-dimensional M/M/c/K station (the distributional assumptions will be detailed further on). For these models, known as loss models, all the arrivals that arise while the station is full are considered to be lost. In our model we assume that only external arrivals may be lost, whereas arrivals that arise from within the network are blocked if the destination station is full. We therefore approximate the arrival rates by combining flow conservation with loss model information. We denote by

- λ_i : the total arrival rate to station i (includes potentially lost arrivals),
- λ_i^{eff} : the effective arrival rate to station i (accounts only for the arrivals that are actually processed, i.e. excludes all lost arrivals),
- γ_i : the external arrival rate to station i .

Accounting for the lost arrivals we have:

$$\lambda_i^{\text{eff}} = \lambda_i(1 - P(N_i = K_i)), \quad (2)$$

where N_i denotes the total number of jobs at station i ($N_i = A_i + B_i + W_i$). $P(N_i = K_i)$ is known as the blocking probability.

In most existing decomposition methods the arrival rate is obtained via the flow conservation equations. In the loss model context, the flow conservation laws hold for the effective arrival rates and are approximated as follows:

$$\lambda_i^{\text{eff}} = \gamma_i(1 - P(N_i = K_i)) + \sum_j p_{ji} \lambda_j^{\text{eff}}. \quad (3)$$

Inter-arrival times to station i are assumed to be independent and identically distributed exponential variables with parameter λ_i .

4.1.2 Average probability of being blocked, P_i^f

The average probability of being blocked at station i , P_i^f , helps us describe the rate at which a job gets blocked after service. P_i^f is approximated by the weighted average of the blocking probabilities of all downstream stations:

$$P_i^f = \sum_j p_{ij} P(N_j = K_j). \quad (4)$$

4.1.3 Service and unblocking rates, μ_i and $\tilde{\mu}(i, b)$

The average service rate of a server at station i is μ_i . It accounts for the active phase. It is an exogenous parameter.

We now describe how we approximate $\tilde{\mu}(i, b)$. Suppose that station i is in the state $(A_i, B_i, W_i) = (a, b, w)$. Then the service rate of the station is $a\mu_i$, i.e. the active jobs are being processed by a **parallel** servers. In the state (a, b, w) there are b blocked servers, but they do not all work in parallel, as we now describe. We define:

- $\tilde{\mu}_i^o$: the average unblocking rate of a destination station of station i . (We describe its derivation below.)
- $D(i, b)$: the number of distinct destination stations that are blocking the b jobs at station i .

For each destination station that is blocking a job at station i , we approximate the rate at which it unblocks jobs at station i by $\tilde{\mu}_i^o$. Thus if all b jobs are blocked by the same destination station, then they can be seen as forming a virtual queue in front of the blocking station with a FIFO unblocking mechanism. The average unblocking rate at station i is then $\tilde{\mu}_i^o$. If the jobs are blocked by $D(i, b)$ distinct destination stations then they can be seen as forming $D(i, b)$ virtual **parallel** queues, each with a FIFO unblocking mechanism. The average unblocking rate at station i is then $D(i, b)\tilde{\mu}_i^o$. More specifically we have:

$$\frac{1}{\tilde{\mu}(i, b)} = \sum_{d=1}^{\min(b, \text{card}(\mathcal{I}^+))} P(D(i, b) = d) \frac{1}{d \tilde{\mu}_i^o}, \quad (5)$$

where \mathcal{I}^+ represents the set of destination stations of station i , and $\text{card}(\mathcal{I}^+)$ is its cardinal. Equation (5) holds because we assume that each destination station unblocks at rate $\tilde{\mu}_i^o$. We now describe how we approximate both $\tilde{\mu}_i^o$ and $P(D(i, b) = d)$.

The average unblocking rate of a destination station, $\tilde{\mu}_i^o$

We denote by:

- $\hat{\mu}_i$: the effective service rate of a server at station i (it includes service and blocking). We will describe its approximation further on.
- \tilde{p}_{ij} : the transition probabilities conditional on a job being blocked at station i , i.e.

$$\tilde{p}_{ij} = \frac{p_{ij}P(N_j = K_j)}{\mathcal{P}_i^f}.$$

- r_{ij} : the proportion of arrivals to station j that arise from blocked jobs at station i , i.e.

$$r_{ij} = \frac{\tilde{p}_{ij}\lambda_i^{\text{eff}}}{\lambda_j^{\text{eff}}}.$$

Suppose station j is blocking jobs at predecessor stations. It is therefore full and is serving at rate $\hat{\mu}_j c_j$. It unblocks jobs at station i at the rate $r_{ij}\hat{\mu}_j c_j$. Thus the average time between successive unblockings is:

$$\frac{1}{\tilde{\mu}_i^o} = \sum_j \tilde{p}_{ij} \frac{1}{r_{ij}\hat{\mu}_j c_j},$$

$$\frac{1}{\tilde{\mu}_i^o} = \sum_{j \in \mathcal{I}^+} \frac{\lambda_j^{\text{eff}}}{\lambda_i^{\text{eff}} \tilde{\mu}_j c_j}. \quad (6)$$

Equation (6) is used to approximate $\tilde{\mu}_i^o$.

Probability that d distinct stations are blocking the b blocked jobs, $P(D(i, b) = d)$

We denote by:

- $\delta(i, b, d)$: the random vector containing the b destination stations of the blocked jobs, d of which are distinct, i.e. $\delta(i, b, d)_k$ denotes the destination station of the k^{th} blocked job.
- $\Delta(i, b, d)$: the sample space of $\delta(i, b, d)$.
- \mathbf{d} : a realization of $\delta(i, b, d)$.

In order to approximate $P(D(i, b) = d)$ we sum over all possible realizations of $\delta(i, b, d)$.

$$\begin{aligned} P(D(i, b) = d) &= \sum_{\mathbf{d} \in \Delta(i, b, d)} P(\delta(i, b, d) = \mathbf{d}) \\ &= \sum_{\mathbf{d} \in \Delta(i, b, d)} P(\delta(i, b, d)_1 = \mathbf{d}_1, \delta(i, b, d)_2 = \mathbf{d}_2, \dots, \delta(i, b, d)_b = \mathbf{d}_b) \\ &= \sum_{\mathbf{d} \in \Delta(i, b, d)} \tilde{p}_{i\mathbf{d}_1} \tilde{p}_{i\mathbf{d}_2} \dots \tilde{p}_{i\mathbf{d}_b} \end{aligned}$$

We define $\ell(i, b, d)_j$ as the number of jobs blocked by station j at station i (given that there are a total of b blocked jobs that are blocked by d distinct destination stations). We thus have:

$$P(D(i, b) = d) = \sum_{\mathbf{d} \in \Delta(i, b, d)} \prod_{j \in \mathcal{I}^+} \tilde{p}_{ij}^{\ell(i, b, d)_j}.$$

This last equation shows that for a given realization \mathbf{d} of $\delta(i, b, d)$, what is of interest in determining $P(D(i, b) = d)$ is the occurrence of each destination station (i.e. the vector $\ell(i, b, d)$), the ordering of the destination stations is not important. Thus instead of summing over $\Delta(i, b, d)$, we will sum over the set of $\ell(i, b, d)$ vectors. This reduces the size of the space over which we sum. The set of such vectors is noted $L(i, b, d)$ and is defined by:

$$\ell(i, b, d) \in L(i, b, d) \Leftrightarrow \begin{cases} \sum_{j \in \mathcal{I}^+} \ell(i, b, d)_j = b, \\ \sum_{j \in \mathcal{I}^+} \mathbb{I}(\ell(i, b, d)_j > 0) = d, \\ \ell(i, b, d)_j \geq 0 \quad \forall j \in \mathcal{I}^+, \end{cases} \quad (7)$$

where $\mathbb{I}(x)$ is the indicator function. The first equation of the system of equations (7) means that there are a total of b jobs blocked at station i , and the second means that these jobs are blocked by d different destination stations. For a given vector $\ell(i, b, d)$ that satisfies the system of equations (7) there are $b! / (\prod_{j \in \mathcal{I}^+} \ell(i, b, d)_j!)$ different realizations of $\delta(i, b, d)$ that are associated with it. This corresponds to the number of permutations of a vector of b destination stations, where destination station j is repeated $\ell(i, b, d)_j$ times. Therefore we obtain:

$$P(D(i, b) = d) = \sum_{\mathbf{d} \in \Delta(i, b, d)} P(\delta(i, b, d) = \mathbf{d}) = \sum_{\ell(i, b, d) \in L(i, b, d)} \frac{b!}{\prod_{j \in \mathcal{I}^+} \ell(i, b, d)_j!} \prod_{j \in \mathcal{I}^+} \tilde{p}_{ij}^{\ell(i, b, d)_j}.$$

Coming back to Equation (5) and replacing $P(D(i, b) = d)$ by the approximation that we have just derived we obtain:

$$\frac{1}{\tilde{\mu}(i, b)} = \frac{1}{\tilde{\mu}_i^o} \sum_{d=1}^{\min(b, \text{card}(\mathcal{I}^+))} \frac{1}{d} \sum_{\ell(i, b, d) \in L(i, b, d)} \frac{b!}{\prod_{j \in \mathcal{I}^+} \ell(i, b, d)_j!} \prod_{j \in \mathcal{I}^+} \tilde{p}_{ij}^{\ell(i, b, d)_j}. \quad (8)$$

The size of the space $L(i, b, d)$ is still considerably large therefore when approximating $\tilde{\mu}(i, b)$ we use an exogenous approximation of \tilde{p}_{ij} :

$$\tilde{p}_{ij} = \frac{p_{ij}P(N_j = K_j)}{\mathcal{P}_i^f} = \frac{p_{ij}P(N_j = K_j)}{\sum_k p_{ik}P(N_k = K_k)} \approx \frac{p_{ij}}{\sum_k p_{ik}}.$$

This approximation makes both summations of Equation (8) exogenous. These two summations are therefore evaluated only once when solving the entire system of equations. This approximation is appropriate if the blocking probabilities of the destination stations have the same magnitude, whereas it is inadequate if their magnitudes differ. The only endogenous parameter remaining in Equation (8) is $\tilde{\mu}_i^o$. Thus we have written $\tilde{\mu}(i, b)$ in the form:

$$\tilde{\mu}(i, b) = \tilde{\mu}_i^o \phi(i, b), \quad (9)$$

where $\phi(i, b)$ is estimated exogenously and can be seen as the average number of distinct destination stations that are blocking the b jobs at station i .

When describing the approximation $\tilde{\mu}_i^o$ we came across the effective service rate of a server, $\hat{\mu}_i$. We now describe how we approximate this parameter.

The effective service rate, $\hat{\mu}_i$

The total time spent by a job in front of a server, called the effective service time $1/\hat{\mu}_i$, is composed of the service time (active phase) and for some jobs of the blocked time (blocked phase). We denote by T_i^B the random variable representing the blocked time of a job conditional on it being blocked. For a given station i , all servers serve on average at rate μ_i (active phase). Thus the average time that a job spends in the active phase is $1/\mu_i$. A given job is blocked on average with probability \mathcal{P}_i^f and once he is blocked the average time he spends blocked is $E[T_i^B]$. Accounting for both the service and the possible blocking we obtain the average effective service time $1/\hat{\mu}_i$, which is approximated by:

$$\frac{1}{\hat{\mu}_i} = \frac{1}{\mu_i} + \mathcal{P}_i^f E[T_i^B]. \quad (10)$$

In this equation μ_i is an exogenous parameter, the approximation of \mathcal{P}_i^f was given in Equation (4). We approximate $E[T_i^B]$ by conditioning on the length of the blocked queue:

$$E[T_i^B] = E[E[T_i^B | B_i]] = \sum_{b \geq 0} P(B_i = b | B_i > 0) E[T_i^B | B_i = b] = \sum_{b \geq 1} \frac{P(B_i = b)}{P(B_i > 0)} E[T_i^B | B_i = b].$$

Let $t(i, b)_j$ denote the blocked time of the job that was unblocked in j^{th} position given that there were b blocked jobs. We have:

$$E[T_i^B | B_i = b] = \frac{1}{b} \sum_{j=1}^b E[t(i, b)_j].$$

We know that the average time between successive departures given that there are b blocked jobs at station i is represented by $1/\tilde{\mu}(i, b)$, thus we can approximate the average blocked time of the first job to be unblocked by $1/\tilde{\mu}(i, b)$, that of the second job to be unblocked by $1/\tilde{\mu}(i, b) + 1/\tilde{\mu}(i, b-1)$ and that of the j^{th} by:

$$E[t(i, b)_j] = \sum_{k=b-j+1}^b \frac{1}{\tilde{\mu}(i, k)}.$$

Putting the last two equations together and then interchanging the summations we obtain:

$$\begin{aligned} E[T_i^B \mid B_i = b] &= \frac{1}{b} \sum_{j=1}^b \sum_{k=b-j+1}^b \frac{1}{\tilde{\mu}(i, k)} \\ &= \frac{1}{b} \sum_{k=1}^b \frac{1}{\tilde{\mu}(i, k)} \sum_{j=b-k+1}^b 1 \\ &= \frac{1}{b} \sum_{k=1}^b \frac{k}{\tilde{\mu}(i, k)}. \end{aligned}$$

Therefore our approximation of $E[T_i^B]$ is given by:

$$E[T_i^B] = \sum_{b \geq 1} \frac{P(B_i = b)}{P(B_i > 0)} \sum_{k=1}^b \frac{k}{b} \frac{1}{\tilde{\mu}(i, k)}. \quad (11)$$

Distributional assumptions

Service time and the time between successive unblockings are each assumed to follow an exponential distribution with parameters μ_i and $\tilde{\mu}_i^o$ respectively. For a given station all service times are assumed to be independent and identically distributed, as are all blocked times. By explicitly modeling both of these exponential phases, the number of jobs in front of the servers becomes a two dimensional system (A_i, B_i) composed of the active and the blocked jobs. We are thus in the presence of an M/M/c/K model with a two-dimensional state space. By working in this two-dimensional space we avoid constructing the CPU intensive phase-type service mechanisms defined in some of the pre-existing methods.

4.2 System of equations

The main aim is to obtain the stationary distributions of each station, $\pi(i)$. The main equations consist of the global balance equations which require the definition of the transition rate matrix, for a given station these equations are:

$$\left\{ \begin{array}{l} \pi(i)Q(i) = 0, \\ \sum_{s \in \mathcal{S}(i)} \pi(i)_s = 1, \\ Q(i) = f(\lambda_i, \mu_i, P_i^f, \tilde{\mu}(i, b)). \end{array} \right.$$

The third equation which defines $Q(i)$ is described in Table 1. We have directly implemented these three sets of equations as a single set:

$$\pi(i)g(\lambda_i, \mu_i, \tilde{\mu}(i, b), P_i^f) = 0. \quad (12)$$

The system of nonlinear equations (2-4,6,9-12) is solved simultaneously for all stations. The exogenous parameters are $\{c_i, K_i, p_{ij}, \mu_i, \gamma_i, \phi(i, b)\}$, all other parameters are endogenous. For each station there are seven endogenous parameters: $\lambda_i, \lambda_i^{\text{eff}}, \tilde{\mu}_i^o, \hat{\mu}_i, P_i^f, P(N_i = K_i), P(B_i > 0)$.

For a given station the dimension of its distribution is equal to $\text{card}(\mathcal{S}_i) = (c_i + 1)(K_i + 1 - \frac{c_i}{2})$. Thus the total size of the system of equations is:

$$\sum_i (c_i + 1)(K_i + 1 - \frac{c_i}{2}) + 7,$$

where 7 denotes the seven endogenous parameters for a given station.

Pre-existing methods that require a posteriori validation (e.g. to ensure the integrality of endogenous station capacities) resort to iterative methods. For a given iteration the system of equations for each station is solved sequentially. Since our method requires no a posteriori validation we are able to solve the set of equations associated to all stations simultaneously.

The system is solved by using the Matlab routine *fsolve*, which implements a trust-region dogleg algorithm. The jacobian of the system has been calculated analytically and implemented. In order to ensure the positivity of distributions the system of equations has been implemented in terms of an auxiliary variable $y(i)$ such that $y(i)^2 = \pi(i)$. The initialization of the endogenous parameters are given in Table 2. In this table λ^{FC} corresponds to the arrival rates that satisfy the classical flow conservation laws. π is initialized using a uniform distribution, thus no a priori information concerning the stationary behavior of the stations is required. The other endogenous parameters are deduced from these initializations.

parameter	initialization
$\hat{\mu}$	μ
$\tilde{\mu}^o$	μ
λ	λ^{FC}
λ^{eff}	λ^{FC}
π	\mathcal{U}

Table 2: Parameter initialization

5 Validation

We now present validation results by comparing our method to both pre-existing methods and to simulation results on a set of small networks.

5.1 Validation versus pre-existing methods

Triangular topology

We first compare our method to that of Altioek and Perros (1987) and that of Takahashi, Miyahara and Hasegawa (1980). The latter considered a single server network with triangular topology (depicted in Figure 1) and two cases according to the buffer size of the stations: a null buffer and a buffer of size two. For each case they considered a set of scenarios with increasing service rates for stations two and three. These scenarios are displayed in Table 3. The chosen performance measure was the blocking probability of station one, $P(N_1 = K_1)$. They then compared their estimates to either simulation results or to exact results derived by using the global balance equations of the entire network. The

relative error of the estimates of the different methods are displayed in Figure 2. For both cases all methods yield good estimates, the relative error remaining under 7% for the first case and 4% for the second case. For both cases we yield similar estimates to those of Takahashi, Miyahara and Hasegawa (1980). For the first case Altioik and Perros (1987) yields the most accurate estimates.

$$\forall i \ c_i = 1, \ p_{12} = \frac{1}{2}$$

$$\gamma_1 = 1, \gamma_2 = \gamma_3 = 0$$

scenario	μ_1	μ_2	μ_3
1	1	1.1	1.2
2	1	1.2	1.4
3	1	1.3	1.6
4	1	1.4	1.8
5	1	1.5	2
6	1	1.6	2.2
7	1	1.7	2.4
8	1	1.8	2.6
9	1	1.9	2.8
10	1	2	3

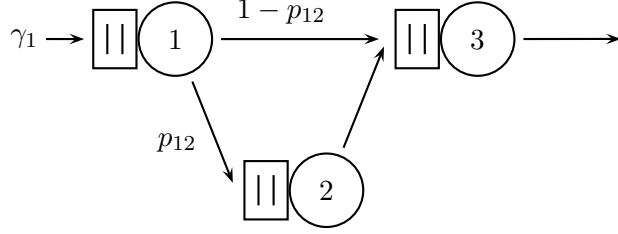


Figure 1: Triangular topology.

Table 3: Increasing service rate scenarios, corresponding to the triangular topology.

Tandem two station topology

Bell (1982) derived a theoretical upper bound on the mean throughput rate of M/M/c/K networks. By considering a tandem two station topology network under a set of scenarios he showed that several decomposition methods “lead to impossible mean throughput rates”. We compare the mean throughput estimates of our method with the methods of Singh and Smith (1997), Kerbache and Smith (1988), Boxma and Konheim (1981), Takahashi, Miyahara and Hasegawa (1980) and Hillier and Boling (1967). The different scenarios and the topology are displayed in Table 4 and the mean throughput estimates of the various methods are depicted in Figure 3. Our mean throughput is estimated by using the effective departure rate at station two, λ_2^{eff} . Figure 3 shows that our mean throughput estimate remains near the upper bound, and is similar to that of the Expansion method of Singh and Smith (1997) and Kerbache and Smith (1988). It slightly violates the bound for the last three scenarios. The relative violations are: 0.3%, 2.2% and 3.8%. Our method therefore yields consistent throughputs unlike the methods of Takahashi, Miyahara and Hasegawa (1980), Hillier and Boling (1967) and Boxma and Konheim (1981).

5.2 Validation versus simulation results

Of main interest in our method are the distributional estimates, which allow us to derive the main performance measures. These could not be compared to pre-existing methods because we know of no method that defines the state space in such a way. We resort to simulation results in order to validate our method on a larger set of scenarios and topologies.

We consider three different topologies. Each network consists of nine stations, all of which are bufferless with three servers. For each network we consider a set of scenarios with increasing external arrival rates. The network configurations and scenario definitions of networks A, B and C are displayed

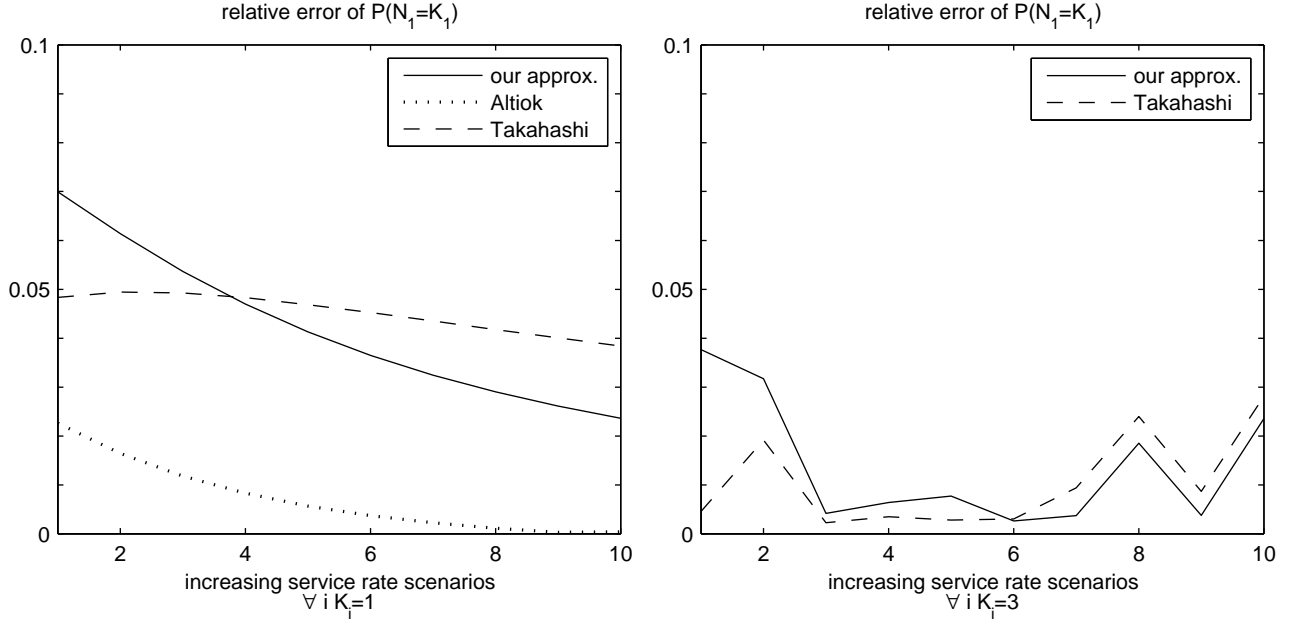


Figure 2: Comparison with the methods of Altioik and Perros (1987) and of Takahashi, Miyahara and Hasegawa (1980) under two capacity configurations.

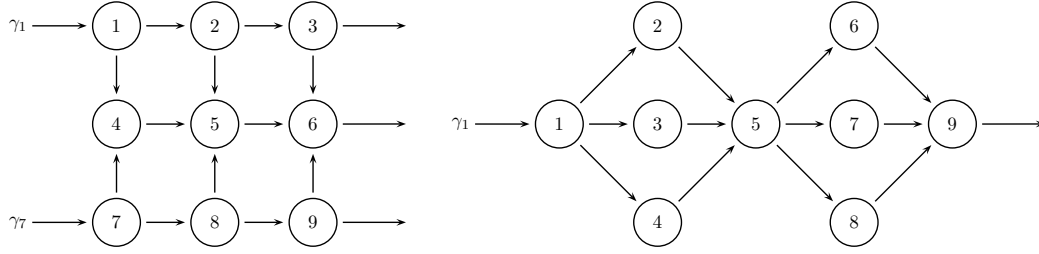


Figure 4: Topologies of networks B and C (left and right hand side respectively).

in Table 5. Network A is a simplified version of the case study network presented in Section 6. Its topology and transition probabilities are the same as that of the case study. The transition probability matrix is in Table 7. The simplifications with regards to the case study concern the number of servers per station and the external arrival rates. The topologies of networks B and C are displayed in Figure 4. For a given station the transition probabilities are uniformly distributed among the possible destinations. In order to validate our results we developed the corresponding simulation models using a discrete event simulator, ProModel version 4.1. Let t_o denote the temporal unit of the transition rates (e.g. minutes, hours). The simulation runs consisted of 20 replications with a warm-up time of 10000 t_o and further run time of 40000 t_o .

Figure 5 displays a histogram of the errors of the distributional estimates for all three networks. For all scenarios, stations and states we consider: $\pi(i)_{(a,b)} - \pi^*(i)_{(a,b)}$, where $\pi(i)_{(a,b)}$ denotes our estimate of the probability that station i is in state (a,b) and π^* is the simulation estimate. There are a total of 1200 estimates. 70% of the absolute errors are smaller than 0.0065, 80% smaller than 0.0125 and 90% smaller than 0.0241. Our method therefore yields good distributional estimates.

In order to illustrate the blocking information derived by our method we consider the scenarios of network C (Table 5). Figure 6 displays the estimates of the distribution of station five given by our

$$\mu_1 = 3, \mu_2 = 1, c_1 = c_2 = 1$$

$$\gamma_1 = 1, \gamma_2 = 0$$

scenario	$K_1 - c_1$	$K_2 - c_2$
1	1	1
2	1	2
3	2	1
4	2	2
5	2	3
6	3	3
7	4	4
8	5	5
9	10	10

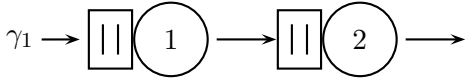


Table 4: Increasing buffer size scenarios that are applied to the tandem two station topology depicted under the table.

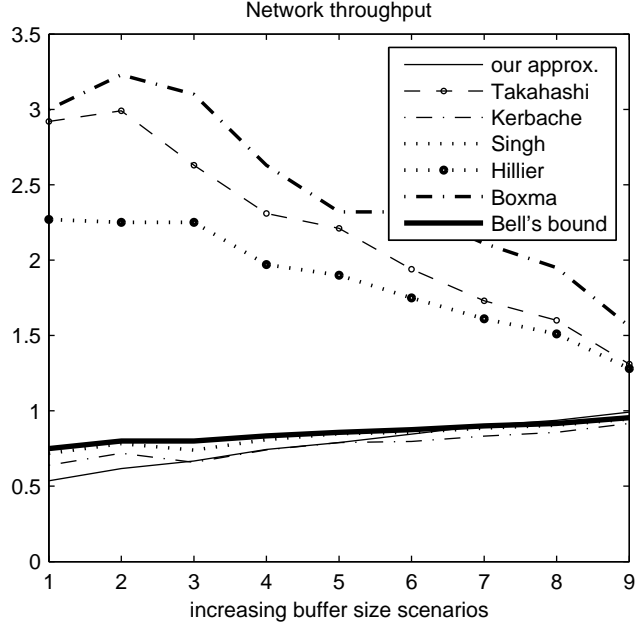


Figure 3: Comparison of the mean throughput estimate of various decomposition methods with the theoretical upper bound derived by Bell (1982).

method and those obtained via simulation. Each plot considers a given state $s = (a, b)$ and plots $\pi(5)_s$ for all scenarios. The scenarios are in a lighter color as the external arrival rate of station one increases. The simulated distribution is depicted as empty squares, whereas our estimates are represented by filled circles. The figure shows that as the external arrival rates increase the states with blocked jobs become more likely, e.g. states (a, b) in $\{(1, 1), (1, 2), (2, 1)\}$. For all states our estimates follow the trend of the simulated probabilities. Overall the estimates are very accurate.

5.3 Convergence of validation runs

A description of the convergence of the algorithm under the different validation runs is tabulated in Table 6. Columns 2-4 summarize the number of iterations, and column 5 gives the average time until convergence across the scenarios. The threshold for the stopping criteria was chosen as 10^{-15} . Convergence was attained when either the first-order optimality condition was smaller than this threshold or when both the relative value and the sum of squares of the system of equations were smaller than the threshold. If after 150 iterations there was no convergence the run was stopped and initialized again. Across all 44 runs 7 required a second initialization to reach convergence.

6 Case study

We now apply our method to a real case study. We consider the patient flow in a network of hospital operative and post-operative units. Clinically, bed blocking may occur for example when a recovered intensive care patient cannot proceed to the intermediate care facility due to unavailable beds, he is said to be blocked until his placement is possible. Studies have acknowledged that bed unavailability

Network A

station index i :	1	2	3	4	5	6	7	8	9	scenario	γ_1
γ_i	-	0.2	0.2	0.0	0.0	0.0	0.0	0.0	0.0	1	0.1
μ_i	0.3	0.3	0.3	0.1	0.01	0.014	0.1	0.4	0.5	2	0.2
										3	0.3
$\forall i \ c_i = K_i = 3, \ card(\mathcal{S}_i) = 10$										4	0.4

Network B

station index i :	1	2	3	4	5	6	7	8	9	scenario	γ_1	γ_7
γ_i	-	0	0	0	0	0	-	0	0	1	0.1	0.1
μ_i	0.3	0.3	0.3	0.6	0.6	0.6	0.3	0.3	0.3	2	0.3	0.3
										3	0.5	0.5
$\forall i \ c_i = K_i = 3, \ card(\mathcal{S}_i) = 10$										4	0.7	0.7
										5	0.9	0.9

Network C

station index i :	1	2	3	4	5	6	7	8	9	scenario	γ_1
γ_i	-	0	0	0	0	0	0	0	0	1	0.1
μ_i	0.3	0.1	0.1	0.1	0.3	0.1	0.1	0.1	0.3	2	0.3
										3	0.5
$\forall i \ c_i = K_i = 3, \ card(\mathcal{S}_i) = 10$										4	0.7
										5	0.9

Table 5: Configuration and scenario definitions for networks A, B and C.

Runs		nb of iterations			average time (sec)	total nb of scenarios
		min	mean	max		
Triangular	case a	7	10.7	20	0.08	10
	case b	7	9	13	0.07	10
Tandem two station		6	6.9	9	0.06	10
Networks A, B and C		11	23.6	39	0.7	14

Table 6: Convergence of validation runs

renders the emergency and surgical admissions procedure less flexible and less responsive (Mackay, 2001). Modeling bed blocking and estimating its effects would bring both patient care and budgetary improvements (Cochran and Bharti, 2006; Koizumi et al., 2005). Thus the importance of modeling the bed blocking phase within a patients recovery procedure. Although few analytic models incorporating blocking have been developed, there is a recently recognized need for them. This is a recent aim defined by Cochran and Bharti (2006): “The next generation of the methodology would include an approximation of the blocking of patients in the queueing model”. The existing analytic models that account for blocking in the healthcare sector have limited their study to feed-forward networks with at most three finite capacity queues (Koizumi et al., 2005; Weiss and McClain, 1987; Hershey, Weiss and Cohen, 1981).

The hospital of interest is the Geneva University Hospital (denoted HUG). The considered units are the emergency operating suite (BO U), elective operating suite (BO OPERA), otorhinolaryngology operating suite (BO ORL), surgical intensive care (IF CHIR), medical intensive care (IF MED), medical intermediate care (IM MED), neuro-surgical intermediate care (IM NEURO), elective recovery (REV OPERA) and otorhinolaryngology recovery (REV ORL). Here the patients are modeled as jobs.

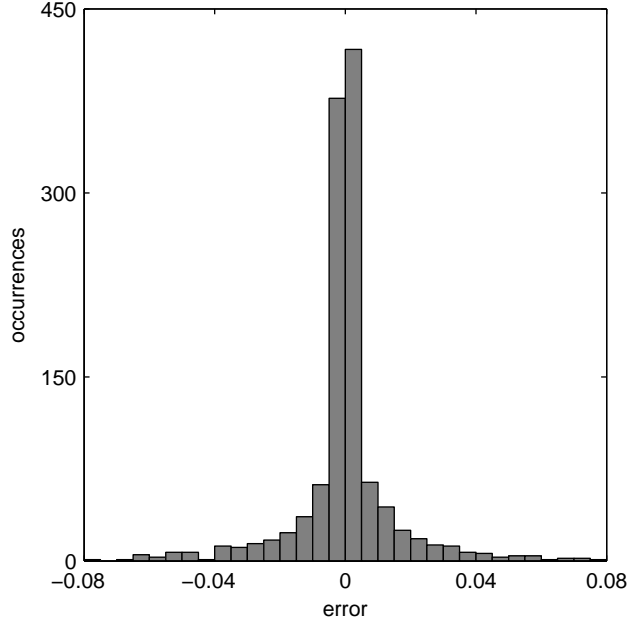


Figure 5: Histogram of the errors of the distributional estimates for all scenarios of networks A, B and C.

Since there is no waiting space each unit is modeled as a bufferless station ($c_i = K_i$). The servers of interest are the beds. The blocking-after-service (BAS) mechanism of our model accurately mimics in-patient bed blocking. The capacities of the different units were estimated according to the evaluations of HUG members. HUG members also extracted patient flow data which we used to estimate the exogenous parameters γ , μ and p_{ij} . Maximum likelihood estimates were used for γ and μ , whereas as the transition probabilities were estimated by the transition frequencies. The data consisted of 25336 patient records ranging over a year. The configuration of the network is presented in Table 7. Note that the sum of the transition probabilities for a given unit (i.e. a given line) may not sum to 1, in this case $1 - \sum_j p_{ij}$ represents the probability of exiting the network given that the job is at station i . The network consists of 9 operative and post-operative units, with 49 possible transitions, containing numerous cycles. This makes the network prone to blocking. We have also carried out this case study using the simulator. This allows us to compare our distributional estimates to those obtained via simulation. The simulation setup is the same as that of Section 5.2. Figure 7 displays the histogram of the errors of the distributional estimates. The 90th, 95th and 99th percentiles of the absolute errors are 0.008, 0.02 and 0.0733 respectively. We have four estimates that have an absolute error larger than 0.1. Figure 8 displays a more detailed error distribution by omitting the four estimates with absolute errors beyond 0.1. These figures show that overall the distributional estimates are very good. The cumulative distribution function for the total number of jobs at each station are depicted in Figure 9. The estimates of our method are represented by filled circles, whereas the simulation estimates are denoted by empty squares. All stations except stations seven and nine have excellent estimates.

Three of the four previously mentioned estimates with large errors concern station seven, the fourth error concerns station nine. Explaining the cause of these large errors is not a straightforward task given the correlation between the endogenous parameters of our system of equations. The detailed distributions of stations seven and nine are displayed in Figure 10. The estimates of our method are represented by filled circles, whereas the simulation estimates are denoted by empty squares. The states (a, b) are ordered by increasing number of active jobs and then increasing number of blocked

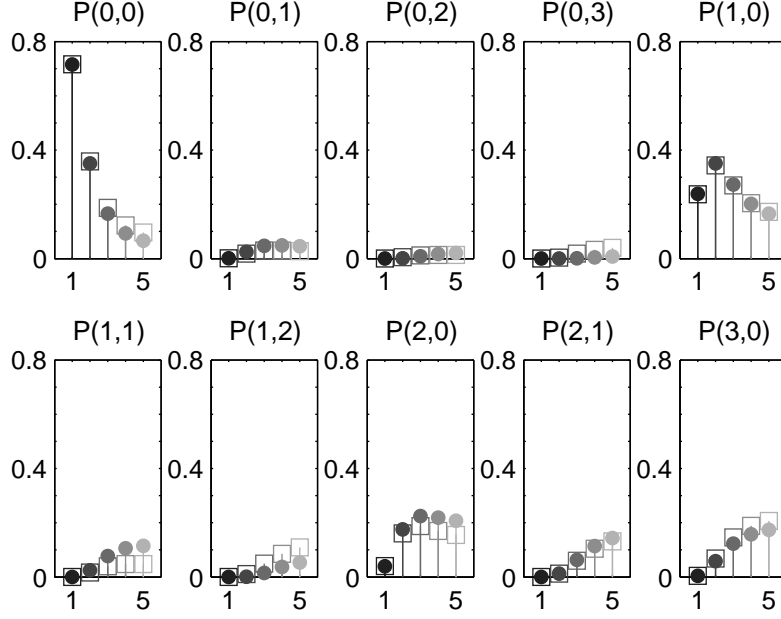


Figure 6: Distribution of station 5 for network C across all scenarios.

jobs. This figure shows that for station seven the state (4,0) is underestimated and for station nine it is the blocked states (0,1) and (0,2) that are underestimated. These misestimations may be correlated since $\tilde{p}_{97} = 0.82$ (displayed in Table 8 and discussed later on), i.e. given that a job is blocked at station nine the probability that it has been blocked by station seven is 0.82. Thus the underestimation of the occupation of station seven may lead to an underestimation of the blocking at station nine.

The outputs of our model can help us quantify the blocking, and also investigate its sources. The transition probabilities conditional on a patient being blocked, \tilde{p}_{ij} , displayed in Table 8, can help us determine the source of blocking. The probabilities have been rounded to 10^{-2} , those smaller than 0.01 are denoted by a dashed line. For a given unit (i.e. a given line in the table) we can identify the destination units that are more likely to block patients. This table helps us detect three main sources of blocking. IF MED and IM MED mutually block each others patients. The same holds for IF CHIR and IM NEURO. This first type of blocking (mutual blocking) may be irrelevant in practice given that the swapping of patients can be identified and carried out easily. The second source of blocking which may be more difficult to solve is the blocking at operating units (BO U, BO OPERA or BO ORL) due to IF CHIR. Moreover, the performance of BO U is strongly linked to its responsiveness, which will be deteriorated by blocking. The third source of blocking occurs at the recovery units (REV OPERA and REV ORL) and is due to IM NEURO.

Other performance measures of the different units are depicted in Table 9. It is important to notice that although P_i^f quantifies the occurrence of blocking at a given unit, it does not capture the impact that a given blocking event may have on the unit or the patient which is blocked. Take for example the ORL recovery unit where $P_i^f = 0.03$, that is on average the probability of a patient getting blocked at that unit is 0.03. In this unit the average service time is 1.9 hours ($1/\mu_9$) and blocking is mainly due to IM NEURO ($\tilde{p}_{97} = 0.82$) where the average service time is 66.67 hours ($1/\mu_7$). Thus the average blocked time at the ORL recovery due to IM NEURO will have a strong impact on the ORL recovery

	BO U	BO OPERA	BO ORL	IF CHIR	IF MED	IM MED	IM NEURO	REV OPERA	REV ORL
c_i	4	8	5	18	18	4	4	10	6
γ_i	0.392	0.502	0.246	0.059	0.176	0.025	0.013	0.155	0
μ_i	0.317	0.255	0.335	0.013	0.015	0.014	0.015	0.22	0.518
$card(\mathcal{S}_i)$	15	45	21	190	190	15	15	66	28
$\forall i K_i = c_i$									

$$(p_{ij}) = \begin{pmatrix} 0 & 0 & 0 & 0.16 & 0.02 & 0 & 0 & 0.71 & 0 \\ 0 & 0 & 0 & 0.07 & 0 & 0 & 0 & 0.84 & 0 \\ 0 & 0 & 0 & 0.03 & 0.01 & 0 & 0 & 0 & 0.95 \\ 0.18 & 0.01 & 0.03 & 0 & 0.03 & 0.01 & 0.11 & 0.03 & 0 \\ 0.05 & 0.01 & 0.01 & 0.01 & 0 & 0.07 & 0 & 0 & 0 \\ 0.02 & 0 & 0 & 0.01 & 0.1 & 0 & 0 & 0 & 0 \\ 0.05 & 0 & 0.05 & 0.04 & 0 & 0 & 0 & 0.01 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0.01 & 0 & 0 \\ 0 & 0 & 0 & 0.05 & 0 & 0 & 0.05 & 0.02 & 0 \end{pmatrix}$$

Table 7: Configuration of the Geneva University Hospital network of operative and post-operative units.

unit. This can also be seen when comparing $E[B_i]/E[N_i]$ and P_i^f . The fact that $E[B_i]/E[N_i]$ is much larger than P_i^f also indicates that although blocking may be rare the impact that it may have on the unit or on the job is not to be ignored. In the case of the ORL recovery unit these performance measures are 0.11 and 0.03 respectively.

The threshold for the stopping criteria of the algorithm was chosen as 10^{-9} . Over a set of 20 runs the average convergence time was 20.5 min, and the average number of iterations required was 2200.4. The jacobian at the solution is ill-conditioned. The 2-norm condition number is 1.3^{10} . The application of preconditioning methods is a source of further improvement.

7 Conclusions and future work

We have presented a method allowing the analysis of network flows via the use of analytic queueing networks that acknowledge the finite capacity property of the real system. The model is adapted for multiple server finite capacity queueing networks with an arbitrary topology and blocking-after-service. The analysis method is based on a decomposition of the network into single queues whose structural parameters are approximated so that they can account for the between-queue correlation. Unlike pre-existing methods the network topology and its configuration are preserved throughout the analysis thus no constraints need to be checked a posteriori. This renders the method suitable for use within an optimization framework. The originality of this method also lies in its capacity to explicitly model the blocking phase that jobs may go through under congested traffic conditions. Performance measures have been validated by comparison with both pre-existing methods and with a theoretical upper bound on the average throughput, on networks with varying buffer size or service rates. The distributional approximations have been compared to those obtained via simulation on a set of networks under a set of scenarios with varying arrival rates, namely under high intensity traffic. This has allowed us to validate distributional information concerning blocked jobs, which will be used in the description of congestion effects. In both types of validations the results are very encouraging.

Pre-existing methods that allow for feedback topologies have assumed that no deadlock occurs or that it is solved instantaneously (e.g. by swapping). The latter approach, although more realistic, violates the FIFO service mechanism assumption. Such as other methods our method does not detect nor solve

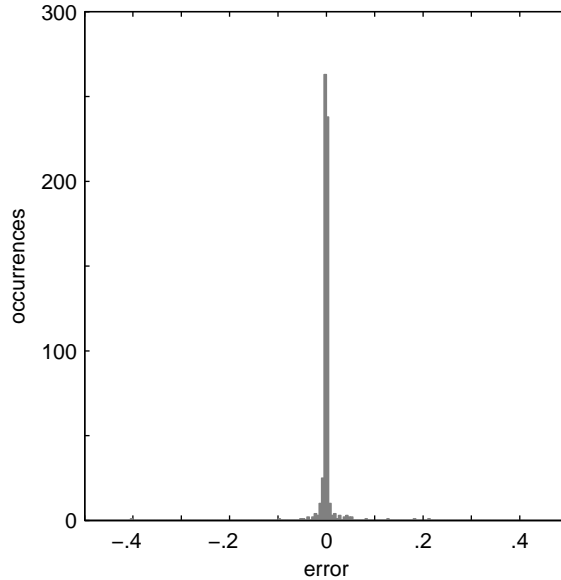


Figure 7: Histogram of errors for the distribution estimates of the HUG network.

deadlock occurrence. Nevertheless we believe that it is of interest to investigate analytic deadlock detection methods.

Further improvements of the method will focus on two main topics. Firstly we will consider aggregating the state space for stations with large capacities, this will considerably reduce the size of the system of equations. Secondly we wish to improve the approximation of the transition probabilities conditional on the job being blocked. This model will be combined with a simulation model within an optimization framework, while ensuring consistency between the two models. The aim of this framework is to allow us to benefit from an optimization friendly analytic model, while accounting for fine details that can be reproduced by the simulation tool.

Acknowledgments

The authors thank Philippe Garnerin and Pau Perez from the Division of Anaesthesiology at the Geneva University Hospitals.

References

- Akyildiz, I. F. and von Brand, H. (1994). Exact solutions to networks of queues with blocking-after-service, *Theoret. Comput. Sci.* **125**(1): 111–130.
- Alfa, A. S. and Liu, B. (2004). Performance analysis of a mobile communication network: the tandem case, *Comp. Comm.* **27**(3): 208–221.
- Altioek, T. (1982). Approximate analysis of exponential tandem queues with blocking, *Europ. J. Operational Res.* **11**(4): 390–398.
- Altioek, T. and Perros, H. G. (1987). Approximate analysis of arbitrary configurations of open queueing networks with blocking, *Ann. Oper. Res.* **9**(1): 481–509.

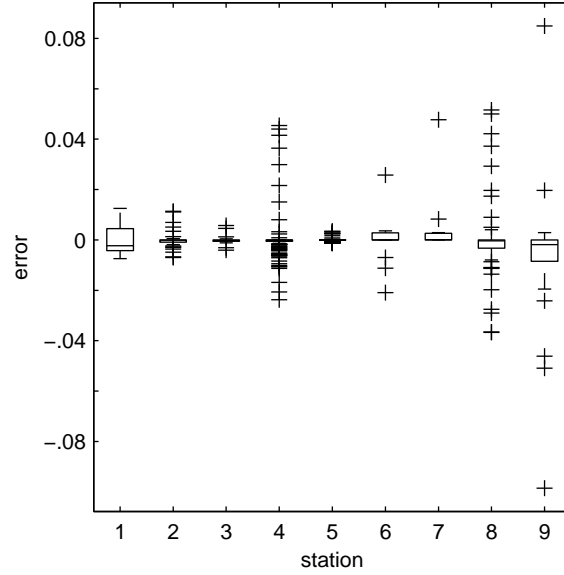


Figure 8: Boxplot of the errors for the distribution estimates, omitting the four errors that are larger in magnitude than 0.1

- Artalejo, J. R. (1999). Accesible bibliography on retrial queues, *Math. Comput. Modelling* **30**: 1–6.
- Balsamo, S., De Nitto Persone, V. and Inverardi, P. (2003). A review on queueing network models with finite capacity queues for software architectures performance prediction, *Perf. Evaluation* **51**: 269–288.
- Balsamo, S., De Nitto Persone, V. and Onvural, R. (2001). *Analysis of Queueing Networks with Blocking*, Vol. 31 of *International Series in Operations Res. and Management Sci.*, Kluwer Academic Publishers.
- Balsamo, S. and Donatiello, L. (1989). On the cycle time distribution in a two-stage cyclic network with blocking, *IEEE Trans. Software Eng.* **15**(10): 1206–1216.
- Bell, P. C. (1982). Use of decomposition techniques for the analysis of open restricted queueing networks, *Operations Res. Letters* **1**(6): 230–235.
- Ben-Akiva, M., Bierlaire, M., Burton, M., Koutsopoulos, H. and Mishalani, R. (2001). Network state estimation and prediction for real-time transportation management applications, *Networks and Spatial Economics* **1**(3-4): 293–318.
- Boxma, O. J. and Konheim, A. J. (1981). Approximate analysis of exponential queueing systems with blocking, *Acta Inform.* **15**(1): 19–66.
- Brandwajn, A. and Jow, Y. (1985). Tandem exponential queues with finite buffers, *Comp. Networking and Perf. Evaluation* pp. 245–258.
- Brandwajn, A. and Jow, Y. (1988). An approximation method for tandem queues with blocking, *Operations Res. Letters* **36**(1): 73–83.
- Cheah, J. Y. and Smith, M. G. J. (1994). Generalized M/G/C/C state dependent queueing models and pedestrian traffic flows, *Queueing Syst.* **15**: 365–386.
- Cochran, J. and Bharti, A. (2006). Stochastic bed balancing of an obstetrics hospital, *Health Care Management Sci.* **9**(1): 31–45.

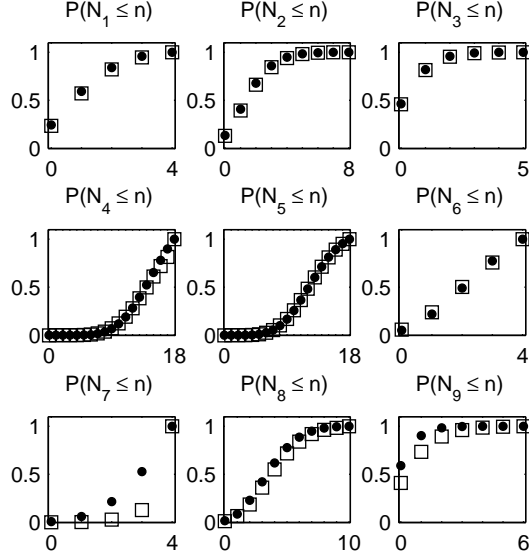


Figure 9: Comparison of the cumulative distribution function, $P(N_i \leq n)$ for all stations.

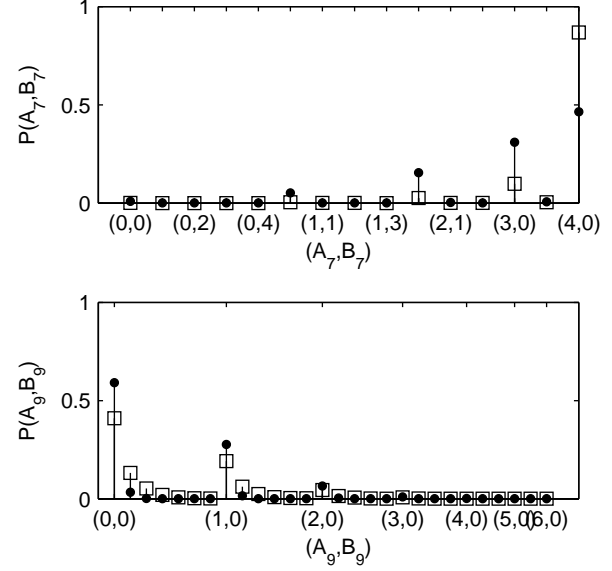


Figure 10: Distributions of stations seven and nine.

unit id	1	2	3	4	5	6	7	8	9
unit	BO U	BO OPERA	BO ORL	IF CHIR	IF MED	IM MED	IM NEURO	REV OPERA	REV ORL
BO U	-	-	-	0.76	0.04	-	-	0.19	-
BO OPERA	-	-	-	0.59	-	-	-	0.41	-
BO ORL	-	-	-	0.87	0.13	-	-	-	0.01
IF CHIR	0.12	-	-	-	0.02	0.04	0.82	-	-
IF MED	0.11	-	-	0.05	-	0.83	-	-	-
IM MED	0.13	-	-	0.16	0.71	-	-	-	-
IM NEURO	0.34	-	0.01	0.65	-	-	-	0.01	-
REV OPERA	-	-	-	-	-	-	1.00	-	-
REV ORL	-	-	-	0.18	-	-	0.82	-	-

Table 8: Transition probabilities conditional on a patient being blocked, \tilde{p}_{ij} .

- Daganzo, C. F. (1996). The nature of freeway gridlock and how to prevent it, *Proceedings of the 13th International Symposium on Transportation and Traffic Theory* pp. 629–646.
- Fone, D., Hollinghurst, S., Temple, M., Round, A., Lester, N., Weightman, A., Roberts, K., Coyle, E., Bevan, G. and Palmer, S. (2003). Systematic review of the use and value of computer simulation modelling in population health and health care delivery, *Journal of Public Health Medicine* **25**(4): 325–335.
- Grassman, W. and Derkic, S. (2000). An analytical solution for a tandem queue with blocking, *Queueing Syst.* **36**: 221–235.
- Gupta, S. M. and Kavusturucu, A. (2000). Production systems with interruptions, arbitrary topology and finite buffers, *Ann. Oper. Res.* **93**: 145–176.
- Hershey, J. C., Weiss, E. N. and Cohen, M. A. (1981). A stochastic service network model with application to hospital facilities, *Operations Res.* **29**(1): 1–22.
- Hillier, F. S. and Boling, R. W. (1967). Finite queues in series with exponential or Erlang service times—a numerical approach, *Operations Res.* **15**(2): 286–303.
- Jackson, J. R. (1957). Networks of waiting lines, *Operations Res.* **5**(4): 518–521.

unit id	1	2	3	4	5	6	7	8	9
unit	BO U	BO OPERA	BO ORL	IF CHIR	IF MED	IM MED	IM NEURO	REV OPERA	REV ORL
K_i	4	8	5	18	18	4	4	10	6
P_i^f	0.02	0.01	0.00	0.06	0.02	0.01	0.01	0.00	0.03
$E[B_i]$	0.04	0.01	0.01	0.22	0.04	0.01	0.01	0.00	0.06
$E[N_i]$	1.37	2.00	0.77	14.03	12.56	2.46	3.19	4.04	0.53
$\frac{1}{\mu_i}$	3.15	3.92	2.99	76.92	66.67	71.43	66.67	4.55	1.93

Table 9: Performance measures for the HUG network. We recall the units capacities, K_i , and average service time, $\frac{1}{\mu_i}$, which are exogenous parameters.

- Jackson, J. R. (1963). Jobshop-like queuing systems, *Management Sci.* **10**: 131–142.
- Jun, J. B., Jacobson, S. H. and Swisher, J. R. (1999). Application of discrete-event simulation in health care clinics: A survey, *J. Oper. Res. Soc.* **50**: 109–123.
- Jun, K. P. and Perros, H. G. (1988). Approximate analysis of arbitrary configurations of queueing networks with blocking and deadlock, *Queueing Networks with Blocking: Proceedings of the First international workshop*.
- Jun, K. P. and Perros, H. G. (1990). An approximate analysis of open tandem queueing networks with blocking and general service times, *Europ. J. Operational Res.* **46**(1): 123–135.
- Kerbach, L. and Smith, M. G. J. (1987). The generalized expansion method for open finite queueing networks, *Europ. J. Operational Res.* **32**(3): 448–461.
- Kerbach, L. and Smith, M. G. J. (1988). Asymptotic behaviour of the expansion method for open finite queueing networks, *Comp. and Operations Res.* **15**(2): 157–169.
- Koizumi, N., Kuno, E. and Smith, T. E. (2005). Modeling patient flows using a queueing network with blocking, *Health Care Management Sci.* **8**(1): 49–60.
- Konheim, A. G. and Reiser, M. (1976). Queueing model with finite waiting room and blocking, *J. Assoc. Comput. Mach.* **23**: 328–341.
- Konheim, A. G. and Reiser, M. (1978). Finite capacity queueing systems with applications in computer modeling, *SIAM J. Computing* **7**: 210–229.
- Koole, G. and Mandelbaum, A. (2002). Queueing models of call centers: An introduction, *Ann. Oper. Res.* **113**(1-4): 41–59.
- Korporaal, R., Ridder, A., Klopogge, P. and Dekker, R. (2000). An analytic model for capacity planning of prisons in the Netherlands, *J. Oper. Res. Soc.* **51**(11): 1228–1237.
- Langaris, C. and Conolly, B. (1984). On the waiting time of a two-stage queueing system with blocking, *J. Appl. Probab.* **21**(3): 628–638.
- Latouche, G. and Neuts, M. F. (1980). Efficient algorithmic solutions to exponential tandem queues with blocking, *SIAM. J. Alg. Disc. Meth* **1**: 93–106.
- Lee, H. S., Bouhchouch, A., Dallery, Y. and Frein, Y. (1998). Performance evaluation of open queueing networks with arbitrary configuration and finite buffers, *Ann. Oper. Res.* **79**: 181–206.
- Mackay, M. (2001). Practical experience with bed occupancy management and planning systems: an Australian view, *Health Care Management Sci.* **4**: 47–56.
- Mandelbaum, A. (2001). Call centers (centres): Research bibliography with abstracts, *Electronically available: <http://ie.technion.ac.il/~serveng/References/ccbib.pdf>*.
- Nagel, K. (2002). Traffic networks, in S. Bornholdt and H. G. Schuster (eds), *Handbook of Graphs and Networks*, Wiley VCH.
- Obaidat, M. S. (1990). Simulation of queueing models in computer systems, *Queueing Theory and Applications*, Taylor & Francis/Hemisphere, pp. 111–151.
- Papadopoulos, H. T. and Heavey, C. (1996). Queueing theory in manufacturing systems analysis

- and design: A classification of models for production and transfer lines, *Europ. J. Operational Res.* **92**(1): 1–27.
- Perros, H. (1984). Queueing networks with blocking: A bibliography, *Perf. Evaluation Review, ACM SIGMETRICS* **12**: 8–12.
- Perros, H. (1994). *Queueing networks with blocking: Exact and Approximate Solutions*, Oxford Press.
- Perros, H. (2003). *Open queueing networks with blocking a personal log. Performance Evaluation: Stories and Perspectives*, Austrian Computer Society.
- Sadoun, B. (2000). Applied system simulation: a review study, *Information Sciences* **124**(1-4): 173–192.
- Schmidt, L. C. and Jackman, J. (2000). Modeling recirculating conveyors with blocking, *Europ. J. Operational Res.* **124**(2): 422–436.
- Singh, A. and Smith, M. G. J. (1997). Buffer allocation for an integer nonlinear network design problem, *Comp. and Operations Res.* **24**(5): 453–472.
- Stewart, W. J. (1999). *Numerical methods for computing stationary distribution of finite irreducible Markov chains*, Advances in Computational Prob., Kluwer Academic Publishers.
- Tahilramani, H., Manjunath, D. and Bose, S. K. (1999). Approximate analysis of open network of GE/GE/m/N queues with transfer blocking, *Seventh IEEE International Symposium on Modeling; Analysis, and Simulation of Computer and Telecommunications Systems (MASCOTS'99)*.
- Takahashi, Y., Miyahara, H. and Hasegawa, T. (1980). An approximation method for open restricted queueing networks, *Operations Res.* **28**(3): 594–602.
- van Vuuren, M., Adan, I. J. B. F. and Resing-Sassen, S. A. E. (2005). Performance analysis of multi-server tandem queues with finite buffers and blocking, *OR Spectrum* **27**: 315–338.
- Weiss, E. N. and McClain, J. O. (1987). Administrative days in acute care facilities: A queueing-analytic approach, *Operations Res.* **35**(1): 35–44.