



An analytic model for capacity planning of prisons in the Netherlands

R Korporaal¹, A Ridder², P Kloprogge³ and R Dekker^{1*}

¹Erasmus Universiteit, Rotterdam, ²Vrije Universiteit, Amsterdam and ³Point Logic Systems, Rotterdam, The Netherlands

In this paper we describe a decision support system developed to help in assessing the need for various types of prison cells. In particular we predict the probability that a criminal has to be sent home because of a shortage of cells. The problem is modelled through a queuing network with blocking after service. The main objective of our study is to describe our analytical method and an approximate algorithm to solve this network. Through simulation studies we evaluate our method. Both the analytic and the simulation tool are elements of the decision support system.

Keywords: capacity planning; prisons; queuing network; approximate algorithm

Introduction

Crime puts high costs on society, however, confining criminals in prisons, costs the society a great deal as well! Calculations executed by the Dutch Ministry of Justice¹ in 1997 indicated that the building of a standard prison cell requires about 125 000 Euro while keeping someone imprisoned is about 36 000 Euro per year. Accordingly, it is quite important to make a proper estimate for the need of prison cells.

According to Dutch law, no two severely punished criminals are put in the same cell. If no cells is available, then offenders of small crimes are sent home where they wait for their sentence. They have to report at a later time to the authorities at so-called self-reporting institutions. Due to this shortage of cells the total number of home-sendings in 1994 amounted to more than 5000, at that time there were about 8700 detainees. The huge number of home-sendings caused very negative publicity. However, the problem is not that easy since the demand for cells varies from day to day and there are several types of prisons.

In this paper we describe a decision support system which has been built for the Dutch Penitentiary Agency to support them in their assessments of the need for prison cells. In this system the problem is modelled as a queuing network with blocking after service. The two tools in the decision support system are a numerical program which calculates performance measures of the queuing network, and a simulation program of the penitentiary system. Because the type of blocking in the queuing network is a new one, we have developed a new analytical method to

solve the network. Therefore we mainly focus on the first tool of the decision support system.

There does exist some literature in operational research applied to crime and justice (see the overview given in Maltz²) but very few papers describe the assessment of the need for prison cells. Cuvelier³ presented a simulation model for use in the United States but does not model the possibility of sending criminals home because of cell shortage, which is the essential aspect in the Dutch case. Although the Dutch situation may be very specific, we think the problem is interesting both from a methodological and a problem oriented point of view.

The structure of the paper is as follows. In the following section we describe the penitentiary system in the Netherlands, then we describe the decision support system. Thereafter we give the details of the queuing model and the solution procedure. Then we sketch the evaluation of the numerical and simulation tools of the decision support system, and finally we give an example of usage.

The Dutch penitentiary system

The penitentiary system in the Netherlands consists of several prisons or institutions of different types. These prisons may be viewed as connected to each other and to society through flows of prisoners.

A person who is arrested on the suspicion of having committed a crime is first put into a police cell. There he (we will use the male form as 95% of the criminals are male) is allowed to stay for at most 72 hours. In some cases this may be extended once with again 72 hours. Thereafter he has to be transferred to a remand centre where he waits his trial. If no cell in such an institution is available, or if his crime is not serious enough, he is sent home. This is

*Correspondence: Dr R Dekker, Econometrisch Instituut, Erasmus Universiteit Rotterdam, PO Box 1738, 3000 DR Rotterdam, the Netherlands. E-mail: rdekker@few.eur.nl

called a *home-sending*. A suspect who is sent home, is adjudged later, and when he gets sentenced he reports to a self-reporting institution. He will spend the first term of his sentence in such a prison. If he does not report, he is arrested again and brought to an institution for arrested people where he will spend the first term of his sentence.

The other suspects wait in the remand centres for their trial. After trial and sentence a criminal is sent to a short-term or a long-term detention prison, depending on the length of his sentence. If he behaves well, he may go after some time to either a half or a fully open detention. For very short sentences the criminals go to the self-reporting institutions as mentioned before. Finally, there are specially protected prisons and prisons for women and youngsters. As the latter groups are rather small we do not include them in this paper. Figure 1 illustrates the possible flows and transfers of detainees in this described penitentiary system. The dotted line is the flow of suspects who are sent home because of cell shortage. The small arrows pointing outside represent the return to society because of completion of detention.

For instance, a person may follow the *route* society \rightarrow PC \rightarrow RC \rightarrow LT \rightarrow society. This means that a person was arrested and put in a police cell. Then he waited in the remand centre for his trial (apparently there was place for him). After conviction he is transferred to the long-term prison where he completes his detention. The latter is not preferred by the Ministry of Justice: one feels the need to let the prisoner adjust to society before the final release, this is only possible in prison types HO and FO. However, when there is no place available in these prison types, the prisoner stays in his current cell.

The moments of transfer between institutions are laid down in official rules. Yet if no place in the next institution is available the prisoner stays in his present cell. We say that he waits for transfer, and we call him a *waiting person*, and the time lag his *waiting time*. It may even happen that he

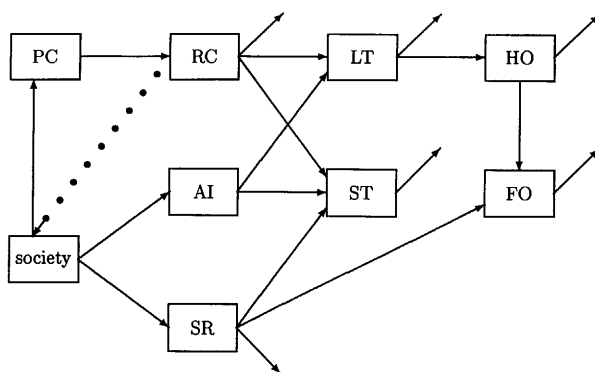


Figure 1 The prison types and the (most common) prisoner flows in the Dutch penitentiary system. PC: police cell; RC: remand centres; AI: institution of arrested people; SR: self-reporting institution; ST: short-term institution; LT: long-term institution; HO: half-open institution; FO: fully-open institution.

skips an institution, for example, go over from a prison for long sentences to directly an open institution, instead of staying for some time in a half-open institution.

The decision support system

To make optimal use of all cell capacity the Dutch Penitentiary Agency developed a nation wide information system indicating which cells in which institutions were available at which times. The information system only collects data of all events in the prisons. The next step was to build a more advanced system which could determine the need for each type of prison cell given estimated inflow of prisoners and their sentences, and given a required performance. The *performance* of the penitentiary system is expressed most importantly in

- the number of home-sendings,
- the waiting times of prisoners for transfers,
- and the cell occupancy rate.

The purpose of this system was to support the political decision to which amount to invest in building of new cells. In an ideal prison system the home-sendings and waiting times are low, the occupancy rate is high.

This decision support system has been developed by the consulting company Point Logic Systems jointly with people of the Econometrisch Instituut of the Erasmus University in Rotterdam at the end of 1995 and the beginning of 1996. The system consists of two major components, CellSim and CellNet. Each component is a PC-based, stand-alone program. CellSim is a simulation program of the penitentiary system, CellNet is a numerical program of an analytic model of the penitentiary system.

In order to run the simulation program CellSim, one has to specify its input: (i) the number of different prison cells; (ii) the rate of incoming detainees; (iii) the distribution of sentences; and (iv) the rules which control the processes of the prisoners. The output of the simulation gives estimates of the performance of the system. The major disadvantage of the program is that it requires long runs, basically because in a realistic system the prisons are almost constantly full. Hence the so-called warming-up period (before the simulation reaches equilibrium) takes a long time. Furthermore, in the program each individual must be monitored throughout his presence in the system. In a realistic simulation study this requires monitoring some eight to nine thousand people. This causes a heavy use of memory and processor capacity, because of the frequent sorting of events. Therefore, we developed CellNet, which calculates the performance of an approximate queuing model of the penitentiary system. In this model individuals are grouped into classes or types. The calculations are based on equilibrium properties of queuing models resulting in an iterative algorithm. In the following sections we

describe the queuing model and the iterative algorithm in more detail.

An indication of different running times might be that a typical run of the CellSim program takes two hours and that it takes about half a minute to run CellNet (on a Pentium 75 MHz PC), the programs are used in practice in serial order. Through CellNet many scenarios of capacity extensions are considered and their performances are calculated (approximately). Those scenarios which seem to result in required performance, are simulated in CellSim in order to get more accurate answers. As an example, one of the striking results was that the scenario of extending the long-term prisons resulted in a much smaller number of people waiting than the scenario of extending the remand centres. In the last paragraph of the paper we come back to this example.

Finally we remark that it took five months to develop the CellSim program. This was much longer than initially planned. The reason was that the program should simulate the penitentiary system with all its complex rules as realistic as possible. (We shall report on the validation in the last section.) As a consequence, the program is somewhat rigid since any change in rules needs reprogramming. Contrary, we think that CellNet is much more robust because it is based on an approximation of the system. And it is interesting to remark that this development took only three months.

The mathematical model in CellNet

We model the penitentiary system as a network of J finite capacity service facilities (*queues*). Each queue represents one of the prisons which we described earlier. The detainees form the *customers* of the network. The *servers* at the queues represent the cells of detention.

In reality each detainee follows a specific route through the penitentiary system and knows in advance his minimal duration in each of the prisons and his maximal sojourn time in the system. This may be modelled in the queuing network by introducing different customers classes each with his own route and service times.⁴ However, we are interested in average long-run performances of the system, therefore we observe the detainees only as a singular flow through the system.² In other words, in our queuing model we aggregate all customers into one class and we average over all possible routes and service times.

In this way we apply the classical modelling concepts of a queuing network⁵ which we will summarise in the following paragraph along with the performance measures of interest.

Queuing network

Customers from outside the network arrive at queue j according to a Poisson process with rate γ_j ,

$j = 1, 2, \dots, J$. The service times at the j th queue form a sequence of i.i.d. random variables which are exponentially distributed with mean μ_j^{-1} . We call these service times the *scheduled* service times, all arrival and service processes are independent. After service completion at queue i , the customer routes to queue j with probability r_{ij} or leaves the network with probability r_{i0} , where

$$\sum_{j=0}^J r_{ij} = 1, \quad i = 1, 2, \dots, J.$$

Each queue is a bufferless queue with finitely many servers. That means that each customer currently present at the queue has his own server to whom he has been assigned immediately upon entering the queue. Let c_j be the number of servers at the j th queue. A blocking protocol is required which regulates the evolution of the network when all c_j servers are occupied. An external arrival at queue j finding all c_j servers occupied, is blocked and lost. Now suppose that a customer completes his service at queue i and chooses queue j to be the next queue on his route. When, at that moment, all c_j servers are occupied, the customer is blocked and remains in queue i until a server becomes available in queue j . During his blocking time, the customer keeps holding his server of queue i . This server is blocked and not available for other customers who wish to enter queue i . As soon as a space becomes free at queue j , the blocked customer moves to it and frees his server at queue i . The time during which a customer blocks a server at queue i due to waiting time for a free server at queue j is denoted as the *blocking* time at queue i and the *transfer* time to queue j . The implemented blocking protocol is known as *blocking after service*.⁶ When more customers wait for entering queue j , the protocol assigns them according to the *first-blocked-first-enter* rule. In an arbitrary network configuration it is possible that a deadlock occurs. In that case the first-blocked-first-enter rule is violated because the deadlock is resolved by exchanging blocked customers.⁷

There is one more specific property of the penitentiary system that we have to model. It deals with the fact that the scheduled term of imprisonment in a prison is diminished with the time duration of waiting for a transfer when the prison turns out to be full at the planned time of the transfer. As soon as a customer in our queuing model enters a free server at queue j , he subtracts his transfer time from his scheduled service time. Hence, the *effective* service time of a customer at a queue is the resultant of the scheduled service time (at queue j) minus his transfer time (to queue j) plus his blocking time (at queue j). Of course, a detainee is released as soon as his sentence has been reached. Translated in the queuing model: as soon as the total effective service times (summed over all visited queues) exceeds the total scheduled service times (summed over all queues) the customer leaves the network.

Model performance

The queuing network given above is an example of a stochastic discrete event system. Events occur every time a new customer arrives or a service is completed. A generalised semi-Markov process (GSMP) is suited to model such a system.⁸ The states of a GSMP reflect both physical states and remaining durations in these physical states (so-called *clock times*). Let $\mathbf{X} = \{X(t): t \geq 0\}$ be such a GSMP that describes the evolution of the queuing network. A state of the process gives information on

- the current number of occupied and blocked servers at all queues;
- the remaining service time of the occupied servers (scheduled minus transfer);
- the destination of all blocked customers, that is the queue for which they wait for a server to become free;
- the waiting time of blocked customers;
- the order of waiting customers who wish to enter a full queue.

Clearly, this is a nontrivial state description of the network. However, modelled as a GSMP, one can formulate easily performance measures of the system, and one has a framework for executing a simulation of the system.⁸ In our queuing network the performance measures of interest are: (i) the mean numbers of occupied and blocked servers; (ii) the mean sojourn and blocking times; and (iii) the loss probabilities. These are denoted by (for $j = 1, 2, \dots, J$)

- L_j := mean number of occupied (including blocked) servers at queue j ;
- B_j := mean number of blocked servers at queue j ;
- S_j := mean sojourn time at queue j ;
- W_j := mean blocking time at queue j ;
- p_j^{loss} := loss probability at queue j ;
- U_j := utilisation at queue j (fraction of time that a server is busy).

The throughput of the j th queue is the mean number of departures per unit time, denoted by Λ_j . Then by Little's formula the mean sojourn time and the mean number at queue j are related via

$$L_j = \Lambda_j S_j, \quad j = 1, 2, \dots, J.$$

Since only external arriving customers are lost, the throughputs $\Lambda_j, j = 1, 2, \dots, J$ satisfy the following sets of linear (*traffic*) equations,

$$\Lambda_j = \gamma_j(1 - p_j^{\text{loss}}) + \sum_{i=1}^J \Lambda_i r_{ij}, \quad j = 1, 2, \dots, J. \quad (1)$$

Translation to penitentiary system

The queuing network is a mathematical model of the real-world penitentiary system, hence, it makes some simplifying assumptions. The assumption of Poisson arrivals and

exponential service times is based on the Poisson-exponential model of criminality proposed by Avi-Itzhak and Sinar,⁹ that is, (i) a criminal commits crimes throughout his career (while not in prison) under a Poisson process with constant rate; (ii) after committing a crime he is imprisoned with a fixed probability; and (iii) the length of his career is exponentially distributed.

The performance measures of the mathematical queuing network are translated to the performances of interest in the penitentiary system through the following obvious relations. B_j , the mean number of blocked servers, corresponds to the number of detainees waiting for transfer in prison j , W_j , the mean blocking time, corresponds to the waiting time of a detainee in prison j ; p_j^{loss} , the loss probability, corresponds to the percentage of home-sendings at prison j ; and U_j , the utilisation, corresponds to the occupancy rate of prison j .

Model solution

When we would know the stationary distribution of the states of the GSMP \mathbf{X} , we could calculate the performance measures of interest. Moreover, open networks without blocking exhibit in many cases a so-called product form solution of the stationary distribution,^{4,5} making the calculation a lot easier. However, due to the blocking phenomena, there is no analytical expression of the stationary distribution in the model we studied. Furthermore, the state space of the model is too complicated to consider a numerical approach.

In this section we propose an approximation algorithm. The idea is to decompose the network model in isolated independent queues, in other words, we approximate the model by a product form network of *simple* queues. We made an extensive literature search of studies on queuing networks with blocking and found that the general idea of approximating is a decomposition of the network in smaller subsystems. Because our type of blocking was not reported before, we adapted this idea and developed a dedicated algorithm which we shall describe in the following paragraphs. In the appendix to this paper we give an overview of the literature on approximating queuing networks with blocking.

An algorithm for tandem configurations

In this section we describe our algorithm for the model of a tandem series of J queues. Recall that the j th queue has c_j servers, no buffers, and exponential scheduled service time distributions with rate μ_j . Upon service completion the customer moves downstream to queue $j+1$ while it blocks his server as long as there is no server available at this queue. Customers arrive at the network at queue 1 according to a Poisson process with rate γ_1 and leave the network after the J th queue.

Our approximate algorithm decomposes the network in J isolated (independent) exponential queues of type $M(\lambda)/M(\mu)/c/N$. This is a queuing model with Poisson (λ) arrivals, exponential (μ) service times, c servers, and N buffer or waiting places. In other words, we approximate the process of arrival epochs at the network queues by Poisson processes in the isolated queues. The effective service time distributions are approximated by exponential distributions in the isolated queues. The mean number of customers waiting in the j th isolated $M(\lambda)/M(\mu)/c/N$ queue approximates the number of blocked servers in the $j-1$ st network queue. These servers are not usable and, therefore, we let the number of servers in the $j-1$ st isolated queue to be equal to the original c_{j-1} minus the mean length of the waiting line of the j th isolated queue. Finally, since all c_{j-1} servers in the original network may produce a blocked customer waiting to enter queue j , we let the buffer size of the j th isolated queue to be c_{j-1} . The main effort lies in determining the arrival and service rates of the isolated queues.

Equilibrium performance measures of an $M(\lambda)/M(\mu)/c/N$ queue are easily obtained (see Section 2.4 of Reference 10). Of interest are

- $L(\lambda, \mu, c, N)$ for the mean number of customers in the system;
- $Q(\lambda, \mu, c, N)$ for the mean number of waiting customers;
- $S(\lambda, \mu, c, N)$ for the mean sojourn time;
- $W(\lambda, \mu, c, N)$ for the mean waiting time;
- $B(\lambda, \mu, c, N)$ for the blocking (or loss) probability (full system);
- $U(\lambda, \mu, c, N)$ for the utilisation.

Because we deal with a tandems series, the throughput of all queues in the network is the same, $\Lambda_j = \Lambda$, where

$$\Lambda := \gamma_1(1 - p_1^{\text{loss}}).$$

Suppose that we know this number, then the j th isolated queue becomes an

$$M(\lambda_j)/M(\varepsilon_j)/s_j/N_j$$

queue. According to our heuristic reasoning about the service and buffer size we let

$$s_j = c_j - Q(\lambda_{j+1}, \varepsilon_{j+1}, s_{j+1}, N_{j+1}), \quad N_j = c_{j-1}, \quad (2)$$

with $s_j = c_j$ and $N_1 = 0$. The first isolated queue should reflect the (only) entrance queue of the network, with the possibility of loss. Once in the network no more customers are lost, that is the arrival rates of the other queues are equal to the network throughput. Hence,

$$\lambda_1 = \gamma_1, \quad \lambda_j = \Lambda \quad (j = 2, 3, \dots, J). \quad (3)$$

The service rates ε_j in the isolated queues approximate the effective service rates in the network queues. Recall that the effective service time is the scheduled service time minus the transfer time plus the blocking time. By decreasing

the number of servers in (2) the blocking time is already taken care of. Hence,

$$\varepsilon_1 = \mu_1, \quad \frac{1}{\varepsilon_j} = \frac{1}{\mu_j} - W(\lambda_j, \varepsilon_j, s_j, N_j) \quad (j = 2, 3, \dots, J). \quad (4)$$

Note that the mean sojourn time in the j th isolated queue equals the mean waiting time plus the mean service time. Therefore, by (4), $S(\lambda_j, \varepsilon_j, s_j, N_j) = 1/\mu_j$.

The algorithm is an iterative method, it is initialised by estimating the loss probability p_1^{loss} by $p^{(0)} = 0$. Let $p^{(n)}$ be the estimation of the loss probability in the n th iteration. The network throughput is estimated by $\gamma_1(1 - p^{(n)})$. The isolated queues are analysed recursively $j = J, J - 1, \dots, 1$ by implementing (2) and (3), and by solving (4). The blocking probability of the first isolated queue becomes

$$p_1 = B(\gamma_1, \mu_1, s_1, 0),$$

which should be compared with the estimated loss probability $p^{(n)}$ at the beginning of the iteration. If the difference is relatively large, a new iteration starts with estimated loss probability $p^{(n+1)} := (p^{(n)} + p_1)/2$, otherwise the iteration is stopped, and we approximate the performance measures of the system by

$$\begin{aligned} L_j &\approx L(\lambda_j, \varepsilon_j, s_j, N_j) - Q(\lambda_j, \varepsilon_j, s_j, N_j) \\ &\quad + Q(\lambda_{j+1}, \varepsilon_{j+1}, s_{j+1}, N_{j+1}), \\ B_j &\approx Q(\lambda_{j+1}, \varepsilon_{j+1}, s_{j+1}, N_{j+1}), \\ S_j &\approx \frac{1}{\mu_j} - W(\lambda_j, \varepsilon_j, s_j, N_j) + W(\lambda_{j+1}, \varepsilon_{j+1}, s_{j+1}, N_{j+1}) \\ &= \frac{1}{\varepsilon_j} + W(\lambda_{j+1}, \varepsilon_{j+1}, s_{j+1}, N_{j+1}), \\ W_j &\approx W(\lambda_{j+1}, \varepsilon_{j+1}, s_{j+1}, N_{j+1}), \\ U_j &\approx U(\lambda_j, \varepsilon_j, s_j, N_j), \\ p_1^{\text{loss}} &\approx B(\gamma_1, \mu_1, s_1, 0). \end{aligned}$$

Remark. We could not prove the convergence of our algorithm, a fact that happens more often in such approximations.¹¹⁻¹⁵ However, in all our experiments we attained rather quickly (less than 25 iterations) the stopping criterion. In these experiments we applied linear interpolation (of performance measures) to deal with non-integer service capacities s_j in (2).

Remark. We have considered a minor adaptation in reasoning for the arrival rates λ_j of the isolated queues. Because the throughput of the j th isolated queue should be equal to the network throughput Λ and because the j th isolated queue loses blocked customers, we should better impose

$$\Lambda = \lambda_j(1 - B(\lambda_j, \varepsilon_j, s_j, N_j)), \quad j = 2, 3, \dots, J.$$

We have tried to solve the system of two equations originated by this equation to (4) for pairs $(\lambda_j, \varepsilon_j)$. Although solutions could be determined for separate queues, in many cases we found that the iteration algorithm was not robust, so we decided to use (3) and (4).

An algorithm for arbitrary configurations

The principles of our algorithm for arbitrary network configurations are similar to those described above. The approximate network consists of J independent (isolated) queues of type $M(\gamma) + M(\lambda)/M(\mu)/c/N$. This queue differs from the previous $M(\lambda)/M(\mu)/c/N$ model in the additional Poisson (γ) arrival process. Customers arriving via the γ stream are called X -type, and customers arriving via the λ stream are called I -type. When all c servers are occupied, arriving I -type customers join the waiting list, arriving X -type customers are lost. An arriving I -type customer is lost when he finds all server and all N waiting places occupied. The performance measures are now labelled with 5-tuples, $PF(\gamma, \lambda, \mu, c, N)$, and are determined easily by noticing that the $M(\gamma) + M(\lambda)/M(\mu)/c/N$ queue is a finite birth-death process. Moreover, we distinguish two loss probabilities, $B_X(\gamma, \lambda, \mu, c, N)$ of the X -type customers, and $B_I(\gamma, \lambda, \mu, c, N)$ of the I -type customers.

Suppose that the network loss probabilities p_j^{loss} are known, then the throughputs in the network are given as the (unique) solution of (1). In our algorithm we approximate the j th network queue by the queue

$$M(\gamma_j) + M(\lambda_j)/M(\varepsilon_j)/s_j/N_j.$$

The X -type customers of this queue represent the external arrivals to network queue j , whereas the I -type customers represent arrivals from other queues in the network. The parameters $\lambda_j, \varepsilon_j, s_j, N_j$ are determined by the same line of heuristic reasoning as in the tandem configuration.

$$\begin{aligned} s_j &= c_j - \sum_{i=1}^J \frac{\Lambda_i r_{ji}}{\Lambda_i} Q(\gamma_i, \lambda_i, \varepsilon_i, s_i, N_i), \\ N_j &= \sum_{i=1}^J c_i r_{ij}, \\ \lambda_j &= \sum_{i=1}^J \Lambda_i r_{ij}, \\ \frac{1}{\varepsilon_j} &= \frac{1}{\mu_j} - W(\gamma_j, \lambda_j, \varepsilon_j, s_j, N_j). \end{aligned} \tag{5}$$

The algorithm is an iterative method. The n th iteration starts with estimates $p_j^{(n)}$ of the loss probabilities p_j^{loss} ($j = 1, 2, \dots, J$). After solving the traffic (1) for getting approximate throughputs, the set of (5) is solved for the parameters of the isolated queues. For the mean numbers of waiting customers $Q(\gamma_i, \lambda_i, \varepsilon_i, s_i, N_i)$ in the equation of the server size s_j we used the most recently updated. Therefore, some are determined in the n th iteration,

others previously in the $n - 1$ st. This depends on the routing probabilities and on the order of considering the queues. Notice that the performance measures $Q(\cdot)$ and $W(\cdot)$ in these equations refer to the waiting queue, hence to the internal customers. The n th iteration ends with calculating the loss probabilities of X -type customers of the isolated queues:

$$p_j = B_X(\gamma_j, \lambda_j, \varepsilon_j, s_j, N_j), \quad j = 1, 2, \dots, J.$$

We stop the iteration when

$$\max_{j=1,2,\dots,J} (|p_j - p_j^{(n)}|) < \varepsilon,$$

for some small given ε , otherwise we continue the iteration with estimated loss probabilities $p_j^{(n+1)} = (p_j + p_j^{(n)})/2$.

System evaluation

In this section we describe the evaluation of the simulation program CellSim and the numerical program CellNet of the decision support system.

CellSim evaluation

We verified the simulation component of the system during the implementation by testing the outcomes of modules of the program. Furthermore, when the sojourn times of people in prisons are not adapted for waiting for transfer, the model becomes a traditional network which has been studied. We checked that the outcomes in that case agree. Finally, we did tracing of individuals in simulation runs.

Validation took place by simulation of the penitentiary system with input (arrival processes and sojourn times) estimated from actual inflow and sojourn times data. The data were collected in 1995 at a daily basis by the department of police information of the Dutch Ministry of Justice. They compromised the three different types of inflow of prisoners: (1) into the remand centres; (2) into the self-reporting institutions; and (3) into the institutions for arrested people. These inflow data were tested for (daily) stationarity and Poisson distributions. The tests showed that the arrival streams may well be described by Poisson processes, although we could not conclude the stationarity to hold.

The other data concerned the terms of imprisonment providing the distributions of the durations scheduled to spend in prisons (the *scheduled* service or sojourn times). Unfortunately, these data were not kept in the information system, but instead the realised sojourn times were available (which correspond to the effective service times in the queuing model). Although we ran a number of tests, we could only estimate the means of the scheduled sojourn times from the given realised sojourn times. We could not fit theoretical distributions for these sojourn times, and therefore we took simply exponential distributions.

The simulation results were compared with the actual performance measures. The outcomes were satisfactory: using 95% confidence intervals we obtain satisfying agreements between the simulation outcomes and the given data for the sojourn times. As an example we show in Table 1 the realised and simulated sojourn time in prisons, although, one (important) performance measure did not come out too well. The percentage of home-sendings at a full remand centre in the simulation ([8.8%:9.9%] 95%-confidence interval) is much lower than 16.4% realised. The two reasons for this disagreement are (i) that we assumed stationarity of inflow where it is actually not; and (ii) that we assumed exponentially distributed scheduled sojourn times where they are actually not known, and where even no averages are known.

CellNet evaluation

Two series of experiments were executed for testing the quality of our approximate algorithm.

- (a) Simulation of the penitentiary system (CellSim) with Poisson-exponential arrivals and service times.
- (b) The approximate algorithm of the queuing model (CellNet).

The means of the assumed Poisson-exponential distributions in these experiments were estimated by the data averages. Furthermore, the routing probabilities (r_{ij}) in the queuing network model are estimated by the average transfers between the penal institutions. Also these data were supplied by the Ministry of Justice.

Both the experiments (a) and (b) assume the same stochastic modelling. All these tests do assure that the algorithm is indeed a good approximation, in the sense that the majority of the calculated performance measures agree with the simulation outcomes. To illustrate our findings, we shall give the comparison in case the input parameters of arrivals (γ_i), terms of imprisonment (μ_i), routing (r_{ij}), and capacities (c_i) are given by the actual figures of 1995, which are summarised in Table 2.

The results of both the algorithm and the simulations are listed in Table 3. The simulations were executed with the batch means method and a fixed sample size: 40 batches,

Table 1 Average times spent in prisons (in days): realised and 95% confidence intervals of estimated by simulations

	<i>Realised</i>	<i>Simulated</i>
RC	116.5	[112.7: 113.6]
AI	84.9	[82.0: 83.8]
SR	48.8	[49.2: 50.0]
LT	235.1	[231.1: 236.0]
ST	83.0	[81.4: 84.1]
HO	162.2	[157.7: 162.4]
FO	100.5	[100.4: 104.5]

Table 2 Input data. Empty entries are 0 (are not possible)

<i>i</i>	c_i	γ_i (daily)	$1/\mu_i$ (days)	r_{ij}							
				RC	AI	SR	LT	ST	HO	FO	free
RC	5044	45.6	102				0.12	0.06	0.05	0.02	0.75
AI	592	9.0	70	0.09			0.08	0.18	0.02		0.63
SR	394	12.0	48	0.05				0.01		0.03	0.91
LT	1061		228	0.14				0.02	0.09	0.06	0.69
ST	182		79	0.01							0.99
HO	348		150	0.07						0.34	0.59
FO	171		102	0.09							0.91

each batch realises 20 000 events due to routing (transfers from one prison to another or back to society). The table gives the performances after translation to the terms of the penitentiary system: the mean realised sojourn times S_i , the mean waiting times W_i , the mean number of waiting people B_i , the percentage of home-sendings p_i^{loss} , and the times of occupation U_i . As we can see, the latter ones are underestimated by the approximation algorithm. A reason may be that the algorithm gives full capacities (100% utilisation) only if the inflow is exceptionally large.

The results of the simulations are reported only through the 95% confidence intervals of the estimates.

Application

Using CellNet, performance of difference scenarios of capacity extensions can be calculated quickly. As an example, consider the situation given above in Table 2. Suppose that the total capacity of all prisons can be extended by 1000 cells and suppose that the cost of a cell is the same in all the different kinds of prisons, the question is which prison to extend to what amount. The conditions are that certain performance criteria should be met. Suppose that the allocation scenarios of Table 4 are considered.

We do not present the detailed performances of each of the scenarios as in Table 3. We summarise these by showing in Table 5 the cumulative total of the waiting times, that is, $\sum_i W_i$, the cumulative total of the numbers of waiting people, $\sum_i B_i$, and the percentage of home-sendings at the remand centres. (The cumulative total of the sojourn times, $\sum_i S_i$, remains the same in all scenarios).

It is now up to the decision maker whether it is more important to have a small percentage of home-sendings, meaning that most arrestants can be imprisoned immediately, but resulting in longer waiting times (scenario 1). Or that the waiting times are low, meaning that most scheduled transfers can be realised immediately, but resulting in more home-sendings (scenario 2). Scenario 3 seems to take both criteria into account and its performances are averages of the other two scenarios. (The average occupation times in scenarios 1 and 3 is more than 99%, in scenario 2 almost

Table 3 Results of CellNet and CellSim. Empty entries are 0 (do not occur)

i	S_i (days)	W_i (days)	B_i	U_i (%)	p_i^{loss} (%)
RC	113.5 [112.7 : 113.5]	11.4 [11.2 : 11.8]	506.1 [501.3 : 523.7]	99.7 [99.7 : 99.8]	7.8 [7.0 : 8.0]
AI	81.6 [81.6 : 82.5]	11.6 [11.3 : 12.0]	83.4 [81.1 : 85.9]	99.3 [99.3 : 99.4]	
SR	49.4 [48.7 : 49.6]	1.4 [1.2 : 1.4]	10.7 [10.0 : 11.4]	99.5 [99.5 : 99.5]	
LT	179.9 [176.4 : 182.1]	6.4 [4.8 : 5.3]	37.9 [28.5 : 31.5]	98.8 [100.0 : 100.0]	
ST	43.6 [43.4 : 44.9]			99.5 [100.0 : 100.0]	
HO	119.8 [123.0 : 127.3]	11.3 [8.0 : 9.4]	32.6 [22.2 : 26.4]	99.5 [100.0 : 100.0]	
FO	68.9 [77.5 : 81.3]			99.2 [100.0 : 100.0]	

Table 4 Allocation scenarios of 1000 extra cells

	RC	AI	SR	LT	ST	HO	FO
Scenario 1	1000	0	0	0	0	0	0
Scenario 2	0	0	0	500	200	100	100
Scenario 3	300	0	0	500	100	100	0

Table 5 Performances of 3 extension scenarios compared to the original situation (in Table 3)

	$\sum_i W_i$	$\sum_i B_i$	p^{loss}
Original	42	670	7.8
Scenario 1	48	983	1.7
Scenario 2	1	14	5.8
Scenario 3	22	289	3.4

98%). Also we see that the number of waiting persons decreases dramatically in scenario 2.

We like to stress that apart from providing support for a specific decision concerning the extension of the penitentiary system with a specific prison the CellSim and CellNet programs were also used to create awareness among a large group of policy makers and politicians of the relationships between the various performance measures of the problem. In the second case the programs serve an educational purpose, where it is not that necessary to model all details of the problem exactly, but to represent the relations in the correct way instead. Note that an analytical program, such as CellNet, usually performs much better than a simulation program in this respect. In the first type of application one needs to realise that the outputs of CellSim and CellNet are only parts of the input into the decision making. As a cost analysis of options is likely to be problem specific we did not include it in the programs. The added value of a more accurate program as CellSim is, can only be realised if the input data is equally accurate. If the latter is not the case one may also rely on less accurate but better understandable output from CellNet.

Conclusions

In this paper we considered the problem of assessing the required number of prison cells in various detention centres to meet several performance criteria, like fraction of criminals sent home, waiting time for transfers, etc. The problem is highly relevant because of the high costs associated to prison cells and the public sensitivity of sending home criminals. Proper decision making should rely on proper information and tools to assess the effects of various possible actions. Generating proper information is costly and a difficult task because of the many different interpretations of the concepts used. Yet, even while not enough information is available, decisions have to be taken to deal with the actual problem. Quoting Mr van Gemmert from the Ministry of Justice on the value of the CellNet and CellSim programs: ‘The programs have shown the policy makers the discrepancies between the occupancy and home sending objectives in the sense that it is not possible to achieve a zero home-sending objective with a 100% occupancy of the cells. The programs have been used to support our large investment decisions in the recent years. Furthermore, they have increased the appreciation of quantitative approaches in this area’.

It is well-known that simulation is a very flexible tool, nevertheless, its usefulness in this case was hindered by a lack of proper information about the actual process and by long development and computation times. Although the approximate method underlying the CellNet program uses rough approximations of the actual process, it produced reasonable outcomes within relatively short computation and development times. Furthermore, we think that the awareness and educational function of CellNet is also very important, especially in showing various relationships to policy makers. For instance, it can evaluate quickly the performance when the numbers of inflow of detainees change, or when the terms of sentences change. Finally we like to remark that since the time of development of

CellNet and our involvement with testing, the government of the Netherlands decided to extend the number of cells considerably: from 8567 in 1995 to 12 491 in 1998. At the time of writing there are about 11 700 people imprisoned, while there were only 400 home-sendings during 1998. We must add here that the occupation was only 94%. The government expects, though, that after year 2000 shortage of cells will happen again. Since 1995 one sees a trend in longer sentences and thus longer sojourn times in prisons, causing the need for more cells.

Appendix

An overview of approximating queuing networks with blocking

Many approximate algorithms for queuing networks with blocking have been reported in the literature.¹⁶ The general idea is to decompose the network in smaller subsystems or in individual queues which are analysed in isolation.¹⁷ The main effort of each method lies in finding the revised arrival rates and the revised service time distributions so that the stationary distributions of the individual queues approximate the marginal stationary distributions of the network.

For tandem configurations of queues with a single exponential server, finite capacity and blocking after service, Hillier and Boling¹⁸ implemented exponential service times of the isolated queues. The capacity of an isolated queue is augmented by one in order to accommodate the blocked customer of the upstream queue. The algorithm of Altioik¹¹ approximates the effective service time distributions by Coxian-2 distribution by looking only one queue ahead which may be blocked. The effective service time distributions are approximated by general Coxian distributions in the algorithm of Perros and Altioik¹² by taking into account that blocking is backlogged over a number of successive queues. In Altioik¹⁹ the original service times are allowed to be of a general phase-type.

An alternative approach decomposes the tandem network in successive pairs of adjacent queues. This so-called two-node decomposition algorithm has been developed by Brandwajn and Jow¹³ in case of exponential interarrival and service time distributions with state dependent rates.

Single-node decomposition algorithms for arbitrary network configurations of queues with a single exponential server, finite capacity, blocking after service, and first-blocked-first-enter rule, have been developed by Takahashi *et al*²⁰ (without deadlock detection). Again, the approximate effective service times follow an exponential distribution. Altioik and Perros¹⁴ use phase-type distributions to approximate the effective service time distributions in feed-forward networks with either external or internal arrivals at any queue. This is generalised in Jun and Perros⁷ for arbitrary network models with deadlock detection and resolution. The

phase-type mechanism becomes quite complex because many different blocking configurations have to be incorporated. Lee and Pollock¹⁵ proposed an algorithm for feed-forward networks which approximates the marginal occupancy probabilities, the queues in isolation are analysed as finite birth-death processes.

Also other approximation methods have been developed to include more general service times. The isolation method²¹ is based on diffusion approximations of queues with repetitive-service blocking. The maximum entropy method²² decomposes the network in isolated queues and maximises constrained entropy functionals. Again, the repetitive-service blocking is assumed. The expansion method²³ is based on (i) a reconfiguration of the network by adding holding queues with appropriate routings; and (ii) parameters adjustment. The method applies to feed-forward models.

Much less work has been done on approximate algorithms of queuing networks with multiple servers and blocking. The maximum entropy method²⁴ applies to models with repetitive-service blocking. The method in Han and Smith²⁵ applies to two-layered feed-forward networks. The effective service time distributions in the first layer are approximated by Coxian distributions. They provide approximate blocking probabilities and system throughput. Closed networks with multiple exponential servers and blocking are analysed in Akyildiz.²⁶ The approximate algorithm modifies the network into a nonblocking network by a state space transformation.

Networks of bufferless queues are models applied to specific production lines in manufacturing systems.²⁷ However, the analysis is restricted to single server models and concerned with the throughput of the system.^{28,29} To our best knowledge, no algorithms have been reported on networks of bufferless multiple server queues with the blocking-after-service rule. Furthermore, a complicating factor in our model of the penitentiary system is the subtraction of the transfer time from the regular service time. Consequently, the effective service time in our model has a different interpretation and implementation than in the literature we have found so far.

Acknowledgements—The authors would like to thank Arnold van Gemmert of the Department of Policy Information of the Ministry of Justice for his comments and suggestions, and anonymous referees for valuable remarks.

References

- 1 Ministry of Justice (1996). *Calculation of Nominal Prices of Prisons 1997*. Ministry of Justice of the Netherlands, Department of policy information, Den Haag.
- 2 Maltz MD (1994). Operations research in studying crime and justice: its history and accomplishments. In: Pollock SM, Rothkopf MH and Barnett A (eds). *Handbooks in Operations Research and Management Science*, Vol. 6. *Operations Research in the Public Sector*. North-Holland: Amsterdam, pp 201–262.

- 3 Cuvelier SJC (1991). Transforming projections into policy, the evolution of computer simulation in correctional research. *International Seminar on Prison Population Projections*. Manchester, England. (See www.shsu.edu/cjcenter/vitas/cuvelier.html)
- 4 Kelly FP (1976). Networks of queues. *Advan Appl Probab* **8**: 416–432.
- 5 Walrand J (1988). *An Introduction to Queueing Networks*. Prentice Hall: Englewood Cliffs.
- 6 Balsamo S and de Nitto Personè V (1994). A survey of product form queueing networks with blocking and their equivalences. *Annals Opns Res* **48**: 31–61.
- 7 Jun KP and Perros HG (1989). Approximate analysis of arbitrary configurations of queueing networks with blocking and deadlock. In: Perros HG and Altiok T (eds). *Proceedings First International Conference on Queueing Networks with Blocking*. North-Holland: Amsterdam, pp 259–279.
- 8 Glasserman P (1991). *Gradient Estimation via Perturbation Analysis*. Kluwer: Norwell.
- 9 Avi-Itzhak B and Shinnar R (1973). Quantitative models in crime control. *J of Criminal Justice* **1**: 185–217.
- 10 Gross D and Harris CM (1985). *Fundamentals of Queueing Theory*, 2nd edn. Wiley: New York.
- 11 Altiok T (1982). Approximate analysis of exponential tandem queues with blocking. *Eur J Opl Res* **11**: 390–398.
- 12 Perros HG and Altiok T (1986). Approximate analysis of open networks of queues with blocking: tandem configurations. *IEEE Trans on Software Engng* **12**: 450–461.
- 13 Brandwajn A and Jow YLL (1988). An approximation method for tandem queues with blocking. *Opns Res* **36**: 73–83.
- 14 Altiok T and Perros HG (1987). Approximate analysis of arbitrary configurations of open queueing networks with blocking. *Annals Opns Res* **9**: 481–509.
- 15 Lee HS and Pollock SM (1990). Approximation analysis of open acyclic exponential queueing networks with blocking. *Opns Res* **38**: 1123–1134.
- 16 Perros HG (1994). *Queueing Networks with Blocking*. Oxford University Press: Oxford.
- 17 Dallery Y and Frein Y (1993). On decomposition methods for tandem queueing networks with blocking. *Opns Res* **41**: 386–399.
- 18 Hillier FS and Boling RW (1967). Finite queues in series with exponential or Erland service times—a numerical approach. *Opns Res* **15**: 286–303.
- 19 Altiok T (1989). Approximate analysis of queues in series with phase-type service times and blocking. *Opns Res* **37**: 601–610.
- 20 Takahashi Y, Miyahara H and Hasegawa T (1980). An approximation method for open restricted queueing networks. *Opns Res* **28**: 594–602.
- 21 Labetoulle J and Pujolle G (1980). Isolation method in a network of queues. *IEEE Trans on Software Engng* **6**: 373–381.
- 22 Kouvatso DD and Xenios NP (1989). Maximum entropy analysis of general queueing networks with blocking. In: Perros HG and Altiok T (eds). *Proceedings First International Conference on Queueing Networks with Blocking*. North-Holland: Amsterdam, pp 281–309.
- 23 Kerbache L and Smith JM (1987). The generalized expansion method for open finite queueing networks. *Eur J Opl Res* **32**: 448–461.
- 24 Kouvatso DD and Xenios NP (1989). MEM for arbitrary queueing networks with multiple general servers and repetitive-service blocking. *Perform Eval* **10**: 169–195.
- 25 Han Y and Smith JM (1993). Approximate analysis of open M/M/C/K queueing networks. In: Onvural and Akyildiz IF (eds). *Proceedings of the Second International Conference on Queueing Networks with Finite Capacity*. North-Holland: Amsterdam, pp 113–126.
- 26 Akyildiz IF (1989). Product form approximations for queueing networks with multiple servers and blocking. *IEEE Trans on Comput* **38**: 99–115.
- 27 Papadopoulos HT and Heavey C (1996). Queueing theory in manufacturing systems analysis and design: a classification of models for production and transfer lines. *Eur J Opl Res* **92**: 1–27.
- 28 Muth EJ (1984). Stochastic processes and their network representations associated with a production line model. *Eur J Opl Res* **15**: 63–83.
- 29 Papadopoulos HT (1995). The throughput of multistation production lines with no intermediate buffers. *Opns Res* **43**: 712–717.

Received March 1999;
accepted April 2000 after one revision