# Centrum voor Wiskunde en Informatica

Centre for Mathematics and Computer Science

P.R. de Waal

Performance analysis and optimal control of
an M/M/1/k queueing system with impatient customers

CWI

1987

P.R. de Waal

Performance analysis and optimal control of
an M/M/1/k queueing system with impatient customers

# Performance Analysis and Optimal Control of
# an M/M/1/k Queueing System with Impatient Customers

Peter R. de Waal

*Centre for Mathematics and Computer Science*
*P.O. Box 4079, 1009 AB Amsterdam, The Netherlands*

A simple M/M/1/k queue with impatient customers is presented as a model for communication systems operating under overload conditions. The performance analysis and optimal control problem for this model are discussed. An efficient algorithm for computing the optimal control is presented along with numerical results.

## 1. INTRODUCTION

In this paper we discuss the performance analysis and optimal control of an M/M/1/k queueing system with impatient customers. The motivation of this investigation is the problem of overload for communication systems.

The most familiar communication system is the telephone network. The operational units in this network are the so-called telephone-switches or -exchanges. During the last few years sophisticated exchanges of the *Stored Program Controlled (SPC)* type have been developed and installed. Most SPC-exchanges are composed of several modules: the *Central Module (CM)*, the *Switching Module (SM)* and one or more *Peripheral Modules (PM)* (for example see figure 1.1). The Central Module takes care of the overall control functions of the system, the Peripheral Module is the interface to the various types of digital and analog user equipment and the Switching Module connects the modules to each other.

The operations in the Central Module are carried out by one or more processors according to a stored program. Besides typical system operations, such as error handling, maintenance and I/O functions, the main part of the workload of the processor is initiated through subscribers' call requests. Examples of such operations are generating a dial tone, receiving the requested phone number digits, checking the validity of the received digits and allocating the available hardware resources.

An SPC-exchange in operation is a typical example of a queueing system where customers compete for a number of limited resources. The limitations stem from both the finite processor capacity and the limited number of hardware resources in the exchange. The performance of an SPC-exchange may therefore degrade significantly during periods in which the demands for service exceed the design capacity (cf. [6]). The response time of tasks scheduled for the processor may become relatively long and this may cause impatient customers to abandon their call request prematurely. Another type of impatient behaviour arises from the use of time-out mechanisms, for instance in the search for free allocatable hardware resources such as senders and receivers. When the time limit that was set for

2

such an operation, expires, a call request may be abandoned before the connection is established. In both cases processor capacity and memory space are wasted on tasks that do not lead to a successful connection. The aim of this investigation is to develop a control mechanism that regulates admission to the exchange to maximize the successful throughput under conditions of overload.
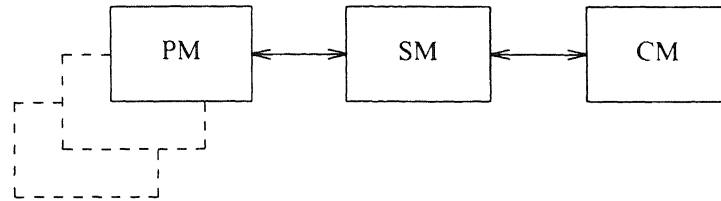


FIGURE 1.1. A Stored Program Controlled Exchange.

In [12] an M/G/1 queue with batch arrivals and service time discretization is presented as a model for switching systems. The successful throughput is expressed in the steady-state probabilities and a suggestion for an overload control is given. In [11] an approximate analysis is presented of an M/G/c queue where customers' call requests become successful according to a probability distribution that is dependent on the waiting time of a customer. A GI/G/1 queue is introduced in [1] with limitations on the waiting and sojourn times of customers. Functional equations for the distribution functions of waiting times and stability conditions are established. The model we present differs from the models in [1], since we allow impatient customers to offer some workload to the server even if they abandon the queue. In [2,3] models for call request processing and the stochastic control problem for these models are presented.

In this paper we present a simple queueing system model of an SPC-exchange with impatient customers. In Section 2 we give a description of the model. It consists of an M/M/1/k queueing system where customers are served according to the First-Come-First-Served discipline. A newly arriving customer is admitted only if the waiting space of the queue is not fully occupied upon arrival, otherwise he is lost. Once a customer has been accepted, his sojourn time limit is generated according to a general probability distribution, and the customer joins the queue if the server is busy or is served immediately otherwise. Subsequently the customer stays in the queueing system until his service is completed. The service completion of a customer is defined to be *successful* if the actual sojourn time of the customer has not exceeded his sojourn time limit. Our *objective* is to maximize the call completion rate which is defined as the mean number of successful service completions in equilibrium per time unit. In section 3 this performance measure is expressed as a function of the steady-state probabilities of an M/M/1/k queueing system and its dependence of the parameters is discussed. In section 4 we introduce the optimal control problem for an approximating queueing system and derive the Hamilton-Jacobi equations for this problem. It is shown that the optimal control is *bang-bang*, i.e. arriving customers are either accepted or rejected without randomization. Furthermore sufficient conditions are given to ensure that the optimal control is of the form where newly arriving customers are accepted if and only if the current number of customers does not exceed a certain level. In section 5 an efficient algorithm for computing this optimal queue size is given and some numerical results are presented. We conclude this paper with some remarks and suggestions for further study.

## 2. DESCRIPTION OF THE MODEL.

In this section we propose a model for analysing the influence of impatient behaviour of customers in a simple queueing system. The service center consists of one server and $k-1$ waiting places. Customers are served in order of arrival and their service demands are exponentially distributed with mean $1/\mu$. Arrivals are generated by a Poisson process with rate $\lambda$. An arriving customer is admitted to the queueing system if the waiting space is not fully occupied, otherwise he is rejected. Rejected customers are assumed not to return.

We now extend this ordinary M/M/1/k queueing system with a concept of impatience. When a customer is admitted to the queueing system, he generates a random sojourn time limit with distribution function $F_S$ (to be discussed later in this section) on the set of positive real numbers with mean $1/\sigma$. All random variables are assumed to be independent. Subsequently he joins the queue, if the server is busy, or is served immediately otherwise. At the moment of his service completion his actual *sojourn time* - i.e. waiting time plus service time - is compared with the sojourn time limit that was set at the time of his arrival. If the actual sojourn time has not exceeded this limit, then the service completion is called *successful*.

We can view this M/M/1/k queueing system with impatient customers alternatively as a queueing system consisting of two parallel queues (see figure 2.1). Queue 1 is an M/G/∞ queue used to model the sojourn of customers and queue 2 is an M/M/1/k queue for service demands of customers. Both queues have identical arrival processes. Upon completion of a task queue 1 is examined to check if the customer who generated the task is still present. If so then the service completion is successful.



FIGURE 2.1. An M/M/1/k queueing system with impatient customers.

Note that, although the sojourn time limit may already have been exceeded while waiting in the queue, the customer still remains in the queueing system to have his service demand fullfilled. This approach has been chosen, since in telephone exchanges even abandoned call requests offer some load to the processor. Furthermore, note that the time limit may be exceeded while the customer is being served. This is to account for the fact that a calling subscriber is not aware of the moment his service begins (e.g. when he is waiting for a dial tone, the exchange is executing tasks, although the subscriber is not aware of this).

We conclude this description with a remark about the probability distribution function $F_S$ of the sojourn time limit. We consider only *Erlang* and *deterministic* distributions, mainly because the Erlang-3 distribution seems a reasonable choice for describing customer impatience in telephone exchanges [1] and the deterministic distribution because this is convenient for describing time-out mechanisms.

## 3. Performance Analysis.

In this section we give a performance analysis of the queueing system we introduced in the preceding section. The performance measures we are interested in are the call completion rate and the rejection probability.

The *rejection probability* $R$ is simply the probability that an arbitrary arriving customer will be rejected, or equivalently finds the queue fully occupied. This measure will become of interest when considering so-called reattempts, i.e. attempts by customers to reenter the communication system after having been refused access. In this report we will however not consider reattempts.

The *call completion rate* $\Lambda$ is defined as the throughput of successful service completions, i.e. the mean number of successful service completions in equilibrium per time unit. We will express both measures in the equilibrium probabilities of an M/M/1/k queue.

In this section we will take the viewpoint of the queueing system as in figure 2, i.e. customers have a *sojourn* time probability distribution function $F_S$ and tasks have a *service* time with an exponential distribution with mean $1/\mu$. Let $p(n)$, $n=0,...,k$, denote the probability that $n$ tasks are present in the queue in equilibrium. It is well known that $p(n)$ is given by

$$p(n) = \frac{1-\rho}{1-\rho^{k+1}}\rho^n, \qquad\qquad n=0,\cdots,k, \qquad (3.1)$$

where $\rho = \lambda/\mu$.

Due to *PASTA (Poisson Arrivals See Time Averages)* [13] the rejection probability $R$ is equal to $p(k)$, i.e. the probability that $k$ tasks are present.

When considering the call completion rate it is more convenient to look at arriving customers (or tasks) rather than departing tasks. Suppose a customer arrives and is accepted in the queue, where he finds $n,n=0,...,k-1$ tasks in front of him. The sojourn time of his task is the sum of $n+1$ independent exponentially distributed random variables, each with mean $1/\mu$, or equivalently an Erlang-$(n+1)$ distributed random variable with mean $(n+1)/\mu$. Denote this variable by $S_{n+1}$. Let $S$ denote the random variable corresponding to the sojourn time of the customer, i.e. a random variable with probability distribution function $F_S$ and mean $1/\sigma$. Furthermore let $q(n)$ denote the probability that the sojourn time of the customer exceeds the sojourn time of his task. It is clear that $q(n) = P\{S>S_{n+1}\}$. Since in equilibrium the mean number of successful service completions per time unit equals the mean number of arriving customers per time unit whose sojourn times exceed the sojourn times of their corresponding tasks, we have

$$\Lambda = \lambda\sum_{n=0}^{k-1} p(n)q(n). \qquad (3.2)$$

The following lemma gives expressions for the probabilities $q(n)$.

**Lemma 2.1.** *If $S$ has an Erlang-m distribution then*

$$q(n) = 1 - \left[\frac{m\sigma}{m\sigma+\mu}\right]^m \sum_{j=0}^{n} \binom{m+j-1}{j} \left[\frac{\mu}{m\sigma+\mu}\right]^j, \qquad n=0,\cdots,k-1. \qquad (3.3)$$

*If $S$ has a deterministic distribution then*

$$q(n) = 1 - e^{-\mu/\sigma} \sum_{j=0}^{n} \frac{\mu^j}{\sigma^j j!}, \qquad\qquad n=0,\cdots,k-1. \qquad (3.4)$$

**Proof**

Let $F_S$ be an Erlang-m distribution function with mean $1/\sigma$, i.e.

$$F_S(t) = \frac{1}{(m-1)!}\int_0^t (m\sigma)^m x^{m-1} e^{-m\sigma x}\, dx, \qquad t\geq 0 \qquad (3.5)$$

and $F_{S_{n+1}}$ an Erlang-$(n+1)$ distribution function with mean $(n+1)/\mu$, i.e.

$$F_{S_{n+1}} = \frac{1}{n!} \int_0^t \mu^{n+1} x^n e^{-\mu x} \, dx \,, \qquad\qquad t \geqslant 0 \qquad\qquad (3.6)$$

We will prove by induction that

$$q(n) = \begin{cases} 1 - \left[ \dfrac{m\sigma}{m\sigma+\mu} \right]^m \,, & n = 0 \\[2ex] q(n-1) - \left[ \begin{matrix} n+m-1 \\ n \end{matrix} \right] \left[ \dfrac{m\sigma}{m\sigma+\mu} \right]^m \left[ \dfrac{\mu}{m\sigma+\mu} \right]^n \,, & n > 0 \end{cases} \qquad (3.7)$$

For $n = 0$ we have

$$q(0) = \int_0^\infty \int_0^t dF_{S_1}(x) \, dF_S(t)$$

$$= \frac{1}{(m-1)!} \int_0^\infty \int_0^t \mu e^{-\mu x} (m\sigma)^m t^{m-1} e^{-m\sigma t} \, dx \, dt$$

$$= 1 - \frac{1}{(m-1)!} \int_0^\infty (m\sigma)^m t^{m-1} e^{-(m\sigma+\mu)t} \, dt$$

$$= 1 - \left[ \frac{m\sigma}{m\sigma+\mu} \right]^m .$$

Let $n > 0$, then

$$q(n) = \int_0^\infty \int_0^t dF_{S_{n+1}}(x) \, dF_S(t)$$

$$= \frac{1}{n!(m-1)!} \int_0^\infty \int_0^t \mu^{n+1} x^n e^{-\mu x} (m\sigma)^m t^{m-1} e^{-m\sigma t} \, dx \, dt$$

$$= q(n-1) - \frac{(m\sigma)^m \mu^n}{n!(m-1)!} \int_0^\infty t^{n+m-1} e^{-(m\sigma+\mu)t} \, dt$$

$$= q(n-1) - \frac{(n+m-1)!}{n!(m-1)!} \left[ \frac{m\sigma}{m\sigma+\mu} \right]^m \left[ \frac{\mu}{m\sigma+\mu} \right]^n .$$

From (7) one can easily derive (3).

Let $S$ be equal to $1/\sigma$ (deterministic), then for $n = 0$ we have

$$q(0) = \int_0^{1/\sigma} \mu e^{-\mu x} \, dx$$

$$= 1 - e^{-\mu/\sigma} .$$

For $n > 0$ we have

$$q(n) = \frac{1}{n!} \int_0^{1/\sigma} \mu^{n+1} x^n e^{-\mu x} \, dx$$

$$= q(n-1) - e^{-\mu/\sigma} \frac{\mu^n}{n! \sigma^n} . \qquad\qquad \blacksquare$$

With the expressions (1)-(4) we can easily evaluate the influence of the parameters of the model on $\Lambda$ and $R$. Some numerical examples are presented in section 5.

## 4. Optimal Control

In this section we discuss the optimal control of an $M/M/1$ queueing system with impatient customers. We first introduce some definitions and notations of counting processes and subsequently discuss the optimal control of an $M/M/1$ queue with a general reward structure. A queueing system with impatient customers is then treated as a special example.

### 4.1. Counting Processes.

We let $N$ denote the set of natural numbers $\{0,1,2,...\}$, $(\Omega,\mathcal{F},P)$ be some probability space and $(\mathcal{F}_t, t \geq 0)$ a family of increasing $\sigma$-algebra's on $\mathcal{F}$, i.e. for all $0 \leq s \leq t$ $\mathcal{F}_s \subset \mathcal{F}_t \subset \mathcal{F}$. A stochastic process $(X_t, t \geq 0)$ is called *adapted* to $\mathcal{F}_t$ if for all $t > 0$ $\mathcal{F}_t^X = \sigma(X_s, s \in [0,t]) \subset \mathcal{F}_t$. In the sequel all stochastic processes are assumed to be adapted to $\mathcal{F}_t$.

A *counting process* $(n_t, t \geq 0)$ is a process taking values in $N$ that has unit jumps. An example of a counting process is a *Poisson process*, where the intervals between successive jumps are independent and exponentially distributed with a constant parameter.

In most applications counting processes can be represented as

$$n_t = \int_0^t \lambda_s \, ds + m_t, \qquad\qquad t \geq 0$$

or

$$dn_t = \lambda_t \, dt + dm_t, \qquad\qquad t \geq 0$$

where $(\lambda_t, t \geq 0)$ is a nonnegative process adapted to $\mathcal{F}_t$ and $(m_t, t \geq 0)$ a $(P,\mathcal{F}_t)$-*martingale*. The process $(\lambda_t, t \geq 0)$ is called the *rate* or *intensity* process of $(n_t, t \geq 0)$. For example a $(P,\mathcal{F}_t)$-Poisson process is a counting process with a deterministic rate process $(\lambda(t), t \geq 0)$. Finally a counting process $(n_t, t \geq 0)$ is called *non-explosive* if

$$n_t < \infty, \qquad\qquad t \geq 0, \ P-a.s.$$

A thorough treatment of counting processes can be found in [5].

### 4.2. The optimal control of an $M/M/1$-queueing system.

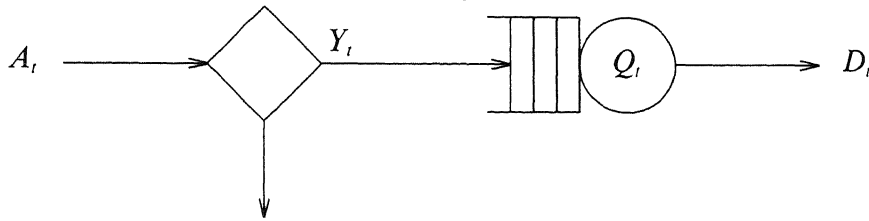Consider the $M/M/1$-queueing system as shown in figure 4.1.



FIGURE 4.1. An $M/M/1$-queueing system with regulated arrivals.

The number of customers present in the queue at time $t$ is denoted as $Q_t$. This queueing process $(Q_t, t \geq 0)$ is an $N$-valued process defined on $(\Omega,\mathcal{F},P)$ and is of the form

$$Q_t = Q_0 + Y_t - D_t, \qquad\qquad t \geq 0 \qquad\qquad (4.1)$$

where $Y_t$ and $D_t$ are non-explosive counting processes *without common jumps*. The above definition implies that $P-a.s.$, $D_t \leqslant Q_0 + Y_t$, $t \geqslant 0$ since $Q_t$ is nonnegative. $Q_t$ is called the *state process* and $Q_0$ is the *initial state*. For each $t \geqslant 0$ $Y_t$ is the number of admitted arrivals in $(0, t]$ and $D_t$ the number of departures in $(0, t]$. These processes are constructed as follows.

Customers arrive according to a point process $(A_t, t \geqslant 0)$ which admits a constant deterministic intensity $\lambda > 0$. Let $\{t_n\}_{n \in N}$ denote the sequence of stopping times that correspond to arrival times, then

$$0 < t_1 < t_2 < \cdots < \infty$$

Upon arrival customers may be admitted to the queue or refused access. Let $\{X_n\}_{n \in N}$ be a sequence of $\{0,1\}$-valued random variables, where if $X_n = 1$, then the customer arriving at time $t_n$ is admitted to the queueing system; otherwise he is refused access. The number of admitted customers over the interval $(0, t]$ is represented by the counting process $(Y_t, t \geqslant 0)$ where

$$Y_t = \sum_{n \geqslant 1} X_n I_{(t_n \leqslant t)}. \tag{4.2}$$

The service times $\{s_n\}_{n \in N}$ of customers at the queue are assumed to be a sequence of i.i.d. random variables which are independent of $A_t$ (and hence of $Y_t$), but which may depend on $Q_t$. If we let the service time of a customer, when $k$ customers are present, be exponentially distributed with mean $1/\mu_k > 0$, then the departure process $D_t$ is a counting process that admits the intensity $\mu_{Q_{t-}}$, where $\mu_0 = 0$.

For the information patterns of these counting processes let $\mathscr{F}_s^A = \sigma(A_s, 0 \leqslant s \leqslant t)$, $\mathscr{F}_s^Q = \sigma(Q_s, 0 \leqslant s \leqslant t)$ and $\mathscr{F}_t = \mathscr{F}_t^A \vee \mathscr{F}_t^Q$. Denote $G_n = \mathscr{F}_{t_n-}^Q \vee \mathscr{F}_{t_n}^A$, $n \geqslant 1$, where

$$\mathscr{F}_{t_n-}^Q = \sigma(A_s \cap \{s < t_n\}; A_s \in \mathscr{F}_s^Q, s \geqslant 0) \tag{4.3}$$

and

$$\mathscr{F}_{t_n}^A = \{A \in \mathscr{F} | A \cap \{t_n \leqslant t\} \in \mathscr{F}_t^A\}. \tag{4.4}$$

We can interpret $G_n$ as the cumulative information of $(A_t, t \geqslant 0)$ up to the instant $t_n$ as well as that of $(Q_t, t \geqslant 0)$ up to the instant right before $t_n$. Based on the information pattern $G_n$, $n \geqslant 1$, the admissible admission policies can be defined as follows.

DEFINITION 4.1. *An $\mathscr{F}_t$-predictable process $(u_t, t \geqslant 0)$ with $u_t \in [0, 1]$ is said to be an admissible admission policy if there exists a probability measure $P^u$ on $(\Omega, \mathscr{F})$ such that $(A_t, t \geqslant 0)$, $(D_t, t \geqslant 0)$ and $(Y_t, t \geqslant 0)$ are counting processes with $(P^u, \mathscr{F}_t)$-intensities $\lambda$, $\mu_{Q_{t-}}$ and $\lambda u_t$ respectively.*

The class of all admissible (random) admission policies will be denoted by $U$. With this definition and $t_n$ an arrival time $u_{t_n}$ can be interpreted as the probability of admitting the $n$th arriving customer given the information $G_n$. Moreover definition 4.1 ensures that any admissible policy $u$ induces a probability measure $P^u$ under which the original distributions of $(A_t, t \geqslant 0)$ and $(D_t, t \geqslant 0)$ are not altered. For details on this approach see [4,5 VII].

Consider now the following optimal admission problem.

PROBLEM P. *Given $Q_0 = x \in N$, $\alpha > 0$, and a function $g: N \to [0, \infty)$, find a $u^* \in U$ such that*

$$J_x^\alpha(u^*) = \sup_{u \in U} J_x^\alpha(u) \tag{4.5}$$

*where*

$$J_x^\alpha(u) = \lim_{t \to \infty} E^u [\int_0^t e^{-\alpha s} g(Q_{s-}) dD_s] \tag{4.6}$$

8

*and $E^u$ denotes the expectation with respect to $P^u$.*

Note that problem P is formulated as an optimization problem for a discounted reward, although the call completion rate is formulated as a performance measure of the queueing system in equilibrium, i.e. a time-average reward. We study discounted rewards, however, mainly because the optimal control for this problem can easily be solved and because the average reward can be treated as the limit of the discounted reward as $\alpha \downarrow 0$.

The function $g$ corresponds with a reward associated with the departure of customers. One can view $g(k)$ as the probability that a service completion (or departure) is successful, given there are $k$ customers in the service center immediately before the departure. In principle this probability can be determined exactly, but it will be a function of the past of the state process. Therefore it is approximated by

$$g(k) = I_{(k > 0)} \left[ \frac{\mu_{k-1}}{\mu_{k-1} + \sigma} \right]^{k-1}, \qquad\qquad k \in N \qquad (4.7)$$

where $1 / \sigma$ is the mean sojourn time of a customer as defined in section 3. Equation (4.7) may be a reasonable approximation, because in equilibrium a departing customer on the average leaves behind as many customers as he faced on arrival. Considering this we can approximate $g(k)$ by the probability that a customer's service will become successful, given he finds $k - 1$ customers in the queue at his arrival epoch, i.e. $q(k - 1)$ as defined in section 3. Since $g(Q_{t-})$ is a non-negative $\mathcal{F}_t$-predictable process, (4.6) can also be written as

$$J_x^\alpha(u) = \lim_{t \to \infty} E^u \left[ \int_0^t e^{-\alpha s} g(Q_{s-}) \mu_{Q_{s-}} \, ds \right]$$

$$= \lim_{t \to \infty} E^u \left[ \int_0^t e^{-\alpha s} g(Q_s) \mu_{Q_s} \, ds \right]$$

since $D_t$ admits the $(P^u, \mathcal{F}_t)$-intensity $\mu_{Q_{t-}}$.

Note that in the original description of the problem, the admission policy was of the impulsive control type, meaning that a decision $X_n$ had to be made at time $t_n$ for each $n \geq 1$. With an appropriate transformation (see e.g. [5] VII.3) the optimal admission problem is equivalent to the intensity control problem formulated above.

### 4.3. The dynamic programming equation.

In this subsection sufficient conditions for an optimal admission policy for problem P are given. These conditions are expressed in terms of a dynamic programming equation. With this equation we show that the optimal control is of *bang-bang* type, i.e. new customers are either accepted or rejected without randomisation. Furthermore we give sufficient conditions for the optimal control being of the type where new customers are admitted if and only if the present number of customers does not exceed a certain number (like window flow-control in communication networks, cf [8]).

First we define the following two transition maps $A, D : N \to N$ by $A k = k + 1$ and $D k = \max(0, k - 1)$ respectively. Admission of a customer at time $t$ then corresponds to a transition $Q_{t-} \to A Q_{t-}$ and a departure corresponds to a transition $Q_{t-} \to D Q_{t-}$.

The next theorem, which is presented without proof, gives sufficient conditions for the optimal admission policy.

THEOREM 4.2. *(Dynamic programming equation) If the function $V : N \to [0, \infty)$ solves the following equation*

$$0 = -\alpha V(k) + g(k)$$

$$+ \lambda \max_{u \in [0,1]} \{ u [V(Ak) - V(k)] \} + \mu_k [V(Dk) - V(k)], \qquad (4.8)$$

*then* $J_x^\alpha(u^*) = V(x)$ *and the optimal control* $u^*$ *is given by*

$$u_t^* = \begin{cases} 1 & \text{if } V(AQ_{t-}) - V(Q_{t-}) > 0 \\ 0 & \text{if } V(AQ_{t-}) - V(Q_{t-}) \leq 0 \end{cases} \tag{4.9}$$

The solution $V$ is called the *value function* of problem P. The proof of the theorem is analogous to that of ([4] Lemma 3), and can be found by standard dynamic programming techniques. It can be shown that the solution to equations (4.8) and (4.9) exists and is unique (cf. [10]). Furthermore one can see that the optimal control is bang-bang, since the term in (4.8) that has to be maximized is linear in $u$. From (4.9) one can also see that the optimal control value depends only on the number of customers present. We can therefore also represent the optimal control as a control law $u^* : N \to \{ 0, 1 \}$, which uses only the state of the queueing process. From now on we won't distinguish between the control process and the control law.

Although theorem 4.2 gives us a sufficient condition for the optimal control, solving the equations is a formidable computational task. Firstly one has to impose a sufficiently large maximum queue size, say $l$, to avoid dealing with an infinite set of equations. Secondly, finding the optimal policy amounts to proposing a possible control $\bar{u}$ (i.e. a function $\bar{u} : N_l \to \{ 0, 1 \}$, a total of $2^{l+1}$ possibilities), solving equation (4.8) with $\max_{u \in [0,1]} u [ V(Ak) - V(k)]$ replaced by $\bar{u}_k [ V(Ak) - V(k)]$ and checking the solution $V(k)$ whether it is consistent with the control $\bar{u}$, i.e. $\bar{u}_k = 1$ if and only if $V(Ak) - V(k) > 0$. Without any knowledge about the optimal control, finding the optimal control in the worst case amounts to solving $2^{l+1}$ sets of $l + 1$ linear equations.

From these considerations we may deduce that any preliminary knowledge about the structure or form of the optimal control might be very useful in finding the optimal control. For instance, if we know that the optimal control is of the type $u_t^* = I_{(Q_{t-} < l^*)}$ for some $l^*$, $0 \leq l^* \leq l$, then we only have to solve $l + 1$ sets of $l + 1$ equations, a significant reduction in computational effort.

The following lemma gives sufficient conditions for the optimal control being of this type.

LEMMA 4.3. *Let* $u^* : N \to \{ 0, 1 \}$ *be the optimal control for problem P.*

(i)    *If there exists an* $l_1 \in N$ *such that for all* $k$, $1 \leq k \leq l_1$, $g(k) \geq g(k-1)$, *then for all* $k$, $0 \leq k < l_1$, *we have* $u_k^* = 1$.

(ii)   *If there exists an* $l_2 \in N$ *such that* $u_{l_2}^* = 0$ *and for all* $k$, $k \geq l_2$, $g(k) \geq g(k+1)$, *then for all* $k$, $k \geq l_2$, *we have* $u_k^* = 0$.

PROOF. Both proofs are given by complete induction. Let $V^* : N \to R$ be the solution of (4.8) and (4.9) corresponding to the optimal control $u^*$. Define $x^* : N \to R$ by

$$x^*(k) = \begin{cases} V^*(0) & , k = 0 \\ V^*(k) - V^*(k-1) & , k > 0. \end{cases} \tag{4.10}$$

Equations (4.8) and (4.9) can now be written as

$$0 = -\alpha \sum_{i=0}^{k} x^*(i) + g(k) + \lambda u_k^* x^*(k+1) - \mu_k x^*(k), \quad k \in N \tag{4.11}$$

and

$$u_k^* = \begin{cases} 1 & \text{if } x^*(k+1) > 0 \\ 0 & \text{if } x^*(k+1) \leq 0. \end{cases} \tag{4.12}$$

(i) Let $l_1 \in N$ be such that for all $k$, $1 \leq k \leq l_1$, $g(k) > g(k-1)$. We have to prove that $x^*(k) > 0$ for all $k$, $1 \leq k \leq l_1$. Equation (4.11) for $k = 0$ and $k = 1$ reads

$$0 = -\alpha x^*(0) + g(0) + \lambda u_0^* x^*(1) \tag{4.13}$$

$$0 = -\alpha(x^*(0) + x^*(1)) + g(1) + \lambda u_1^* x^*(2) - \mu_1 x^*(1) \tag{4.14}$$

respectively. Subtracting (4.13) from (4.14) gives

$$\lambda(u_0^* x^*(1) - u_1^* x^*(2)) = -\alpha x^*(1) + (g(1) - g(0)) - \mu_1 x^*(1). \tag{4.15}$$

Suppose now that $u_0^* = 0$, or, equivalently, $x^*(1) \leq 0$. Since $u_1^* x^*(2) \geq 0$, we then have the left-hand side of (4.15) smaller than or equal to zero and the right-hand side strictly greater than zero. From this contradiction we may deduce that $u_0^* = 1$. The proof of $u_k^* = 1$ for $k$, $1 \leq k < l_1$ proceeds in a similar way.

(ii) Let $l_2 \in N$ be such that $u_{l_2}^* = 0$, or, equivalently, $x^*(l_2 + 1) \leq 0$, and $g(k) \geq g(k + 1)$ for all $k$, $k \geq l_2$. Define $x : N \to R$ by

$$x(k) = \begin{cases} x^*(k) & 0 \leq k \leq l_2 \\ [g(k) - \alpha \sum_{i=0}^{k-1} x(i)][\alpha + \mu_k]^{-1} & k > l_2 \end{cases} \tag{4.16}$$

We first prove that $x(k) \leq 0$ for $k > l_2$.

(ii.a) Let $k = l_2 + 1$. From the definition of $x$ and $x^*$ we have

$$x(l_2 + 1)(\alpha + \mu_{l_2+1})$$

$$= g(l_2 + 1) - \alpha \sum_{i=0}^{l_2} x(i)$$

$$= g(l_2 + 1) - \alpha \sum_{i=0}^{l_2} x^*(i)$$

$$= \alpha x^*(l_2 + 1) - \lambda u_{l_2+1}^* x^*(l_2 + 2) + \mu_{l_2+1} x^*(l_2 + 1)$$

$$\leq 0$$

by the induction assumption $x^*(l_2 + 1) \leq 0$ and by $u_k^* x^*(k + 1) \geq 0$ for all $k$.

(ii.b) Suppose that for $k \in N$, $k > l_2$, $x(k) \leq 0$. It will be shown that then $x(k + 1) \leq 0$.

$$x(k + 1)(\alpha + \mu_{k+1})$$

$$= g(k + 1) - g(k) - \alpha x(k) + g(k) - \alpha \sum_{i=0}^{k-1} x(i)$$

$$= g(k + 1) - g(k) - \alpha x(k) + (\alpha + \mu_k) x(k)$$

$$= g(k + 1) - g(k) + \mu_k x(k)$$

$$\leq 0$$

by the induction assumption and $g(k + 1) - g(k) \leq 0$, hence $x(k + 1) \leq 0$. Now the definition of $x$ and the fact that $x(k) \leq 0$ for $k > l_2$ imply that $x$ is the solution of the system of equations

$$0 = -\alpha \sum_{i=0}^{k} x(i) + g(k) + \lambda \max_{u_k \in [0,1]} u_k x(k + 1) - \mu_k x(k)$$

Because this system of equations has a unique solution, it follows that $x(k) = x^*(k)$ for all $k \in N$ so $u_t^* = 0$ if $Q_{t-} \geqslant l_2$.      ▯

COROLLARY 4.4. *If there exists an* $l^* \in N$ *such that*

(i)     $g(k-1) < g(k)$     *for all* $k$, $1 \leqslant k \leqslant l^*$

(ii)    $g(k-1) \geqslant g(k)$     *for all* $k$, $k > l^*$               (4.17)

*then there exists a* $k^* \geqslant l^*$ *such that the optimal control* $u^*$ *of problem* $P$ *is of the form*

$$u_t^* = \begin{cases} 1 & \text{if } Q_{t-} < k^* \\ 0 & \text{if } Q_{t-} \geqslant k^* \end{cases}$$

PROOF. Part (i) of lemma 4.6 ensures that $u_t^* = 1$ if $Q_{t-} < l^*$ and part (ii) ensures that $u_t^* = 0$ if $Q_{t-} \geqslant k^*$ with $k^* = \inf_{k \in N}\{k \,|\, u^*(k) = 0\}$.

For instance the example of $g$ given in equation (4.7) satisfies (4.17) if $\mu_k \leqslant \mu_{k-1}$, for all $k$, $k \geqslant 1$. This can easily be seen, since $g(0) = 0$ and for $k \geqslant 1$ we have

$$g(k+1) = \left[ \frac{\mu_k}{\mu_k + \sigma} \right]^k$$

$$= \left[ \frac{1}{1 + \sigma/\mu_k} \right]^k$$

$$\leqslant \left[ \frac{1}{1 + \sigma/\mu_{k-1}} \right]^k$$

$$\leqslant \left[ \frac{1}{1 + \sigma/\mu_{k-1}} \right]^{k-1}$$

$$= g(k)$$

The condition $\mu_k \leqslant \mu_{k-1}$ is not an unrealistic one for one-server queues with queue-dependent service rates, since the service rate is likely to decrease with the number of customers in the queue, mainly because of the increasing amount of overhead. In the next subsection we present an algorithm for computing $k^*$ if $g$ satisfies (4.17).

As we stated in subsection 4.2, the average reward can be treated as the limit of the discounted reward as $\alpha \downarrow 0$. To do this we first have to extend the notation of the value function to $V_\alpha : N \to R$ to denote its dependency on the discount factor $\alpha$. The result is stated in the next theorem.

THEOREM. 4.5. *If* $g$ *is bounded and if there exists an* $M < \infty$ *such that*

$$|V_\alpha(i) - V_\alpha(0)| < M \tag{4.18}$$

*for all* $\alpha > 0$ *and* $i \in N$, *then the optimal average reward is equal to* $\lim_{\alpha \downarrow 0} \alpha V_\alpha(0)$.

The proof of the theorem can be found in [9], Theorem 7.7. We have not been able to prove whether (4.18) holds, but we have observed this bound in numerical examples. Furthermore, in most cases, as presented in section 5, the optimal policies for all $\alpha \leqslant 10^{-10}$ are equal.

*4.4. An algorithm for determining the optimal policy.*

Suppose we have a reward rate $g : N \to [0, \infty)$ that satisfies $g(0) = 0$ and $0 < g(k+1) \leqslant g(k)$ for $k \geqslant 1$. These conditions are satisfied for example if $g(k) = I_{\{k > 0\}} q(k-1)$ with $q(k)$ defined as

in (3.3) or (3.4). From corollary 4.4 we know that the optimal control is of the form $u_t^* = I_{\{Q_t < k^*\}}$ for some $k^* \geq 1$.

Let us suppose that we know that $k^* \geq k$ for some $k$. In order to check whether $k^* = k$ we have to solve (4.8) with the corresponding $u$. This can be rewritten as the matrix equation

$$g_k = A_k \, V_k$$

with $g_k$, $V_k : N_{k+1} \to [0, \infty)$ defined by

$$g_k = [g(0), \dots, g(k+1)]^T$$

and

$$V_k = [V_k(0), \dots, V_k(k+1)]^T$$

and $A_k : N_{k+1} \times N_{k+1} \to R$ given by

$$\begin{bmatrix} \alpha+\lambda & -\lambda & & & & \\ -\mu_1 & \alpha+\mu_1+\lambda & -\lambda & & & \\ & -\mu_2 & \alpha+\mu_2+\lambda & -\lambda & & \\ & & & \cdot & \cdot & \cdot \\ & & & & \cdot & \cdot & \cdot \\ & & & & \cdot & \cdot & \cdot \\ & & & & -\mu_{k-1} & \alpha+\mu_{k-1}+\lambda & -\lambda \\ & & & & & -\mu_k & \alpha+\mu_k \\ & & & & & & -\mu_{k+1} & \alpha+\mu_{k+1} \end{bmatrix} \qquad (4.19)$$

and check the solution whether

$$V_k(l) > V_k(l-1), \qquad\qquad\qquad 1 \leq l \leq k$$

and

$$V_k(k) \geq V_k(k+1).$$

Solving (4.18) can be done by standard *LU-decomposition* [7], using the diagonal entries of $A_k$ as pivots and starting with column zero.

Let $A_k^{(n)} = M_n \cdots M_0 A_k$, $0 \leq n \leq k+1$, denote the matrix that is obtained after the $n$th step of the LU-decomposition, i.e. the matrix where all the entries under the diagonal in the columns 0 to $n$ are eliminated. $A_k^{(k-1)}$ has the following form

$$\begin{bmatrix} A_k^{(k-1)}[0,0] & -\lambda & & & & \\ & A_k^{(k-1)}[1,1] & -\lambda & & & \\ & & \cdot & \cdot & & \\ & & & \cdot & \cdot & \\ & & & A_k^{(k-1)}[k-1,k-1] & -\lambda & \\ & & & & A_k^{(k-1)}[k,k] & \\ & & & & -\mu_{k+1} & \alpha+\mu_{k+1} \end{bmatrix} \qquad (4.20)$$

Let $g_k^{(n)} = M_n \cdots M_0 g_k$. Then from (4.20) one can easily check whether the solution $V_k$ satisfies $V_k(k+1) \leq V_k(k)$ since $V_k$ must also satisfy $A_k^{(k-1)} V_k = g_k^{(k-1)}$, and consequently

$$V_k[k] = g_k^{(k-1)}[k] / A_k^{(k-1)}[k,k]$$

and

$$V_k[k+1] = [g(k+1) + \mu_{k+1} V_k(k)] / \{\alpha + \mu_{k+1}\}$$

If this solution would not satisfy the inequality, then our next guess for $k^*$ would be $k + 1$, in which case we have to solve the matrix equation $A_{k+1} V_{k+1} = g_{k+1}$. One can easily check that $A_{k+1}^{(k-1)}$ is almost equal to $A_k^{(k-1)}$, namely

$$\begin{bmatrix} A_k^{(k-1)}[0,0] & -\lambda & & & & & \\ & A_k^{(k-1)}[1,1] & -\lambda & & & & \\ & & \cdot & \cdot & & & \\ & & & \cdot & \cdot & & \\ & & & A_k^{(k-1)}[k-1,k-1] & -\lambda & & \\ & & & & A_k^{(k-1)}[k,k]+\lambda & -\lambda & \\ & & & & -\mu_{k+1} & \alpha+\mu_{k+1} & -\lambda \\ & & & & & -\mu_{k+2} & \alpha+\mu_{k+2} \end{bmatrix}$$

$A_{k+1}^{(k)}$ can thus easily be computed from $A_k^{(k-1)}$, enabling us to solve the dynamic programming equation through one set of linear equations by incorporating the inequality check into the LU-decomposition of the solution. From the above discussion one can easily see that we can compute the optimal value $k^*$ by the following algorithm.

ALGORITHM 4.6.

```
{declarations}
HUGE    : {arbitrary large integer}
i,k     : integer;
g,diag,V : array[0..HUGE] of real;
found   : boolean;

{initialisation}
found: = false;
k: = 0;
for i: = 0 to HUGE
do
  g[i]: = gᵢ;
  diag[i]: = α + μᵢ
od;

{LU decomposition}
while ((k<HUGE) and not found)
do
  if (g[k+1]+μ_{k+1}g[k]/diag[k])/diag[k+1]≤g[k]/diag[k]
  then
    found: = true;
    g[k+1]: = g[k+1]+μ_{k+1}g[k]/diag[k]
  else
    diag[k]: = diag[k]+λ;
    diag[k+1]: = diag[k+1]-λμ_{k+1}/diag[k];
    g[k+1]: = g[k+1]+μ_{k+1}g[k]/diag[k]
  fi;
  if not found then k: = k+1 fi;
od;
```

14

```
{computation of V(i)}
V[k+1]:=g[k+1]/diag[k+1]
V[k]:=g[k]/diag[k]
for i:=k-1 downto 0
do
   V[i]:=(g[i]+λV[i+1])/diag[i];
od
```

## 5. NUMERICAL RESULTS.

In this section we present some numerical examples of the performance analysis and computation of the optimal control law.

The first system is an M/M/1/k queue with mean service time equal to 1.0 and the sojourn times of customers Erlang-3 distributed with mean 20.0. The optimal value for $k$ as derived from the dynamic programming equation is equal to 8, where we have used in the equations an arrival rate $\lambda = 0.8$, discount factor $\alpha = 10^{-10}$ and reward rate $g(k)$ as defined in equation (4.7).

The call completion rates and blocking probabilities of this queue for various values of $k$ and $\rho$ (or equivalently $\lambda$ since $\mu = 1.0$) are shown in figures 5.1 and 5.2 respectively. From figure 5.1 we may conclude that, although $k = 8$ is optimal for $\rho = 0.8$, the M/M/1/8 queue behaves well also under overload, i.e. for values of $\rho$ ranging from 1.0 to 2.0.

The second system we present is equal to the first one, with the exception that the mean sojourn time is now equal to 50.0. The optimal value for $k$, using the same $\lambda$, $\alpha$ and $g(k)$, in this case is equal to 15. The call completion rate is given in figure 5.3. The graph of the blocking probability is the same as of the first system, since the equilibrium distribution of the number of tasks in both queues does not depend on the sojourn time distribution of the customers.
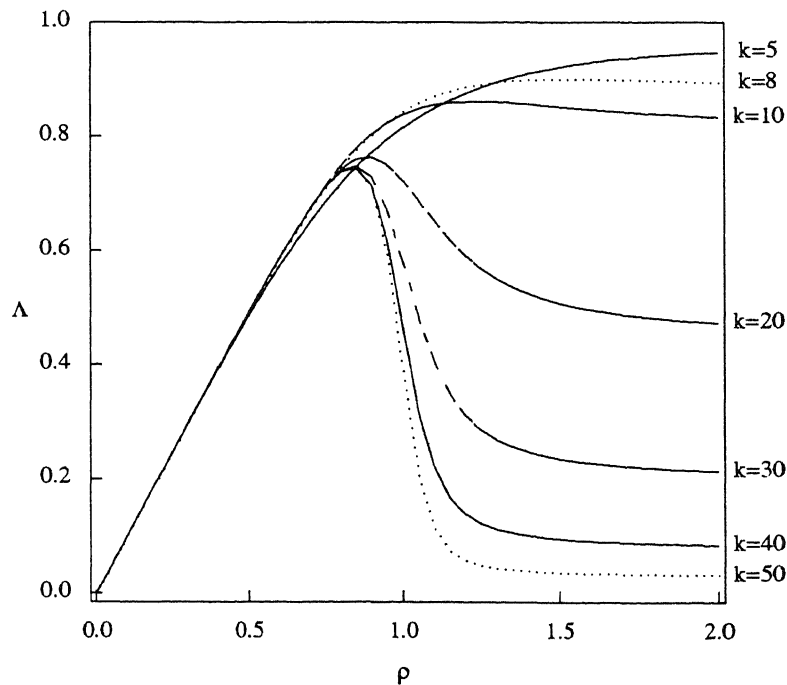


FIGURE 5.1. Call completion rate for an M/M/1/k queue with Erlang-3 distributed sojourn times with mean 20.0.
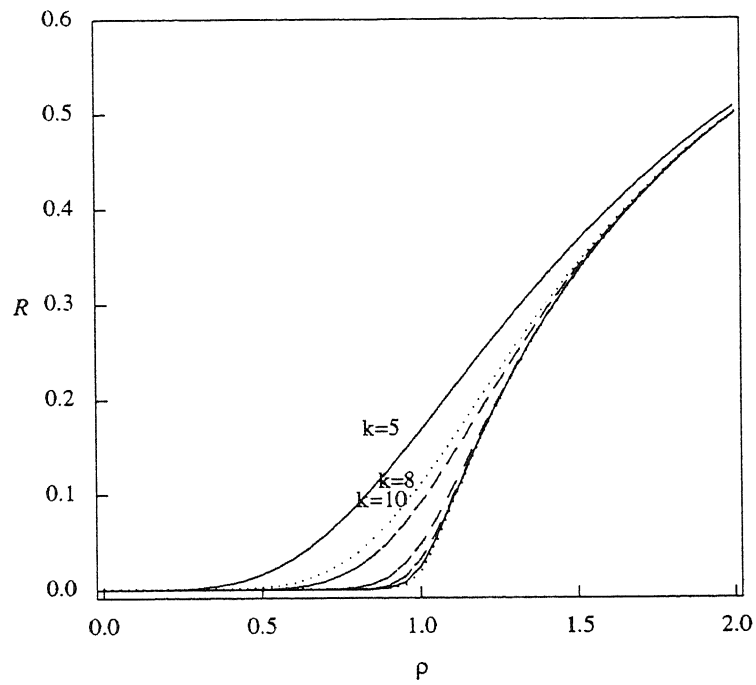
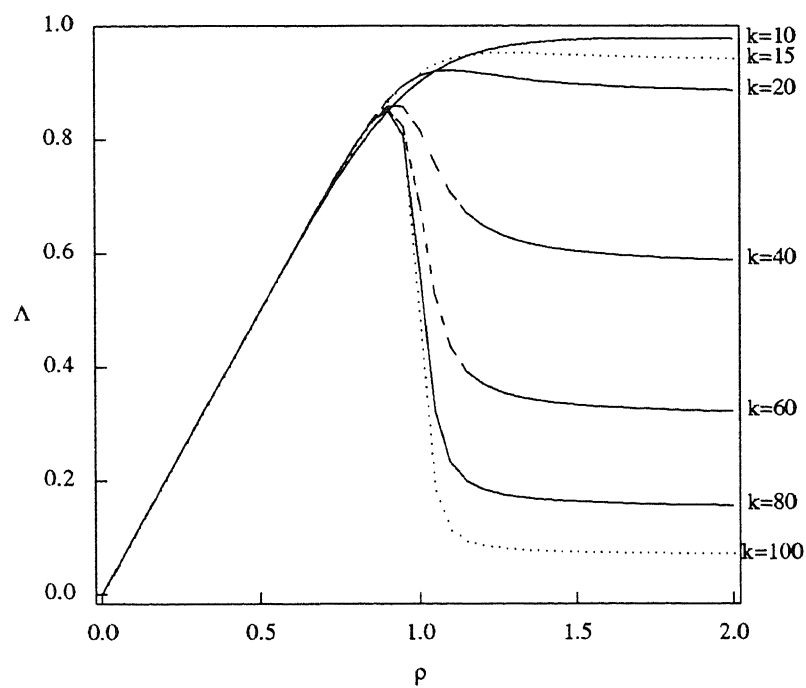FIGURE 5.2. Blocking probabilities for an M/M/1/k queue.



FIGURE 5.3. Call completion rate for an M/M/1/k queue with Erlang-3 distributed sojourn times with mean 50.0.

16

REFERENCES

1. F. BACCELLI (1983). *Modèles probabilistes de systèmes informatiques distribués*, Thèse d'état, INRIA, Paris.

2. R.K. BOEL and J.H. VAN SCHUPPEN (1985). Overload control for SPC telephone exchanges - refined models and stochastic control, in *Proceedings of the 3rd Bad Honnef Conference*, 100-110, ed. N. Christopeit, K. Helmes, M. Kohlmann, Springer-Verlag, Berlin.

3. R.K. BOEL and J.H. VAN SCHUPPEN (1986). Overload control for switches of communication systems - A two-phase model for call request processing, in *Proceedings of the International Seminar on Teletraffic Analysis and Computer Performance Evaluation*, 209-224, ed. O.J. Boxma, J.W. Cohen, H.C. Tijms, Elsevier, Amsterdam.

4. P. BRÉMAUD (1979). Optimal thinning of a point process, *SIAM J.Control and Optimization*, 17, 222-230.

5. P. BRÉMAUD (1981). *Point processes and queues - martingale dynamics*, Springer-Verlag, Berlin.

6. L.J. FORYS (1982). Performance analysis of a new overload strategy, in *10th International Teletraffic Congres*.

7. G.H. GOLUB and C.F. VAN LOAN (1983). *Matrix computations*, The Johns Hopkins University Press, Baltimore.

8. T.G. ROBERTAZZI and A.A. LAZAR (1985). On the modeling and optimal flow control of the Jacksonian network, *Perf.Eval.*, 5, 29-43.

9. S.M. ROSS (1970). *Applied probability models with optimization applications*, Holden-Day, San Francisco.

10. C. STRIEBEL (1975). *Optimal control of discrete time stochastic systems*, Springer-Verlag, Berlin.

11. D.Y. SZE (1986). A queueing model for overload analysis, in *Computer Networking and Performance Evaluation*, 413-422, ed. T. Hasegawa, H. Takagi, Y. Takahashi, Elsevier (North-Holland).

12. PHUOC TRAN-GIA and M.H. VAN HOORN (1986). Dependency of service time on waiting time in switching systems - A queueing analysis with aspects of overload control, *IEEE Trans. Comm.*, 34, 357-364.

13. R. WOLFF (1982). Poisson arrivals see time averages, *Oper.Res.*, 30, 223-231.