

M/G/ ∞ polling systems with random visit times

M. Vlasiou*, U. Yechiali**

* Georgia Institute of Technology,
H. Milton Stewart School of Industrial & Systems Engineering,
765 Ferst Drive, Atlanta GA 30332-0205, USA.

** Department of Statistics and Operations Research,
School of Mathematical Sciences,
Raymond and Beverly Sackler Faculty of Exact Sciences,
Tel Aviv University, Tel Aviv 69978, Israel.

vlasiou@gatech.edu, uriy@post.tau.ac.il

Abstract

We consider a polling system where a group of an infinite number of servers visits sequentially a set of queues. When visited, each queue is attended for a random time. Arrivals at each queue follow a Poisson process, and service time of each individual customer is drawn from a general probability distribution function. Thus, each of the queues comprising the system is, in isolation, an M/G/ ∞ -type queue. A job that is not completed during a visit will have a new service time requirement sampled from the service-time distribution of the corresponding queue. To the best of our knowledge, this paper is the first in which an M/G/ ∞ -type polling system is analysed. For this polling model, we derive the probability generating function and expected value of the queue lengths, and the Laplace-Stieltjes transform and expected value of the sojourn time of a customer. Moreover, we identify the policy that maximises the throughput of the system per cycle and conclude that under the Hamiltonian-tour approach, the optimal visiting order is *independent* of the number of customers present at the various queues at the start of the cycle. In other words (and somewhat surprisingly), additional information regarding the state of the system at the start of a cycle does *not* lead to an improvement of the optimal policy.

1 Introduction

A typical polling system consists of a number of queues, attended by a single server in a cyclic fashion. There is a huge body of literature on polling systems that has developed since the late 1950s, when the papers of Mack et al. [11, 12] concerning a patrolling repairman model for the British cotton industry were published. Rather than giving a partial overview of the literature, we refer the interested reader to the following books, surveys, and papers on polling systems: Takagi [16, 17, 18], Levy and Sidi [10], Yechiali [25], Borst [4], Eliazar and Yechiali [9], Nakdimon and Yechiali [13].

Polling systems have been used as a central model for the analysis of a wide variety of applications in the areas of repair problems [11, 12], telecommunication systems [8], road traffic control [14], computer networks [24], multiple access protocols [3], multiplexing schemes in ISDN [20], satellite systems [1], flexible manufacturing systems [23], and the like.

In many of these applications, as well as in most polling models, it is customary to control the amount of service given to each queue during the server's visit. Common service policies

are the *exhaustive*, *gated*, *globally gated* and *limited* regimes. Under the exhaustive regime, at each visit the server attends the queue until it becomes completely empty, and only then is the server allowed to move on. Under the gated regime, the only customers served during a visit are the ones who are present when the server enters (polls) the queue, while customers arriving when the queue is attended will be served during the next visit. The globally gated regime, introduced by Boxma, Levy and Yechiali [5], is a modification of the gated one: the only customers served during a visit are those who are present at the beginning of a cycle. Finally, under the k -Limited service discipline only a limited number of jobs (at most k) are served at each server's visit to each queue. These service policies imply that the duration of the visit time in a polled queue is a function of the number of customers present there at a given moment (such as the beginning of the cycle or the moment the server enters the queue).

In this paper, we analyse a polling system that differs in two ways from the classical polling model. Rather than considering a *single* server providing service to customers at the various queues, we assume that an *infinite number* of servers is moving as a single group between the queues. Moreover, the service policy we study is *independent* of the queue length. We assume that the group of servers visits each queue for a (possibly random) amount of time that is independent of everything else and which has a distribution that may vary per queue. We further assume that the arrival process of customers to each queue is Poisson and that the service time distribution for customers in each queue is general. To the best of our knowledge, this paper is the first in which an M/G/ ∞ -type polling system is analysed.

The specific application that raised our attention and led us to this model is in the field of road traffic control. Polling models for road traffic are typically along the lines of the classical polling system; namely, they involve a single server rotating around a number of queues. Other assumptions that are typically being made for such models include deterministic service times (i.e. the amount of time that a car needs to pass a traffic light after possibly standing in the queue) and deterministic visit times (i.e. the time the traffic light remains green); see, for example, van der Heijden [21]. Although these models provide fairly good approximations of reality, such assumptions fail to capture the variation both in service times and in visit times. Cars do not need the same amount of time to cross a segment of the road; the ones standing ahead in the queue will inevitably need less time and those that arrive while the queue is empty and the traffic light is still green will not even require the additional time incurred by acceleration. Moreover, recent developments in the technology of traffic lights has led to the design of traffic lights that do not turn green unless there is a queue formed and turn red either when the queue is empty or after a maximum amount of time, which may also vary within a day. As a result, in this paper we provide a framework for studying road traffic control under less restrictive assumptions. We propose an infinite-server polling system, which models the behaviour of traffic: while the traffic light is green all cars present in the queue or approaching the traffic light proceed (receive service) and the time they need to complete service is assumed to be a random variable following a general distribution. Furthermore, we assume that the time the traffic light is green (visit time) is random, although our results are directly applicable in case of deterministic visit times or, more generally, in case the visit times follow a discrete distribution taking positive values.

A common approximation to road traffic is to consider the traffic as *fluid* passing through the road. This approximation is fairly accurate when the traffic is relatively high. Mathematically, high traffic can be modelled by assuming that the arrival rate of customers at each of the queues tends to infinity. The study of such a model provides insights at the queue length (and thus the congestion of a junction) under heavy load. In this paper though we do not study the evolution of the system under heavy load. We assume that the arrival rate at each queue is fixed. This assumption is usually made for the standard polling systems and provides a reasonable

approximation to normal traffic conditions.

The rest of the present paper is organized as follows. Section 2 introduces the model, gives further notation, and describes formally the evolution of the system. In Section 3 we compute recursively the first moment and the probability generating function of the queue length distributions at a polling instant. Later on, in Section 4 we derive the mean and the Laplace-Stieltjes transform of the sojourn time of a customer arriving at queue i , and we show how these expressions simplify in the special case where both the service time and the visit time at queue i are exponentially distributed. Based on the results derived up to that point, in Section 5 we give some numerical results. Specifically, we examine numerically the effect of the first two moments of the visit and service times on the sojourn time of an arbitrary customer. These numerical results indicate that there is an optimal value for the mean visit time to the various queues that minimises the mean sojourn time of an arbitrary customer. In Section 6 we investigate how we can optimise the visit time of the servers at the various queues so that the expected throughput of the system is maximised. It emerges that even when considering *semi-dynamic* control policies, in which the group of servers plans a new route for each cycle, the optimal visiting order that maximises the expected throughput per cycle is *fixed* for all cycles. In other words (and somewhat surprisingly), additional information regarding the system (such as the queue length for all queues at the beginning of the cycle) has no effect on the choice of the optimal strategy.

2 Model description and notation

We consider a polling system with $N \geq 2$ infinite-buffer queues attended by a group of ample number of servers that visits the queues in a fixed cyclic fashion. We index the queues by $i = 1, 2, \dots, N$ in the order of the servers' movement. We shall refer to the polling instant of queue i as the moment when the servers enter that queue. When visiting queue i , the group of servers continues working at this queue for V_i units of time, and acts there as an M/G/ ∞ queue. We assume that the visit times are independent, identically distributed (i.i.d.) random variables.

Customers arrive at all queues according to independent homogeneous Poisson processes with rate λ_i for queue i . After completing their service time, customers leave the system. The service time of each individual customer at queue i is denoted by B_i . It is assumed that all service times in one queue are i.i.d. random variables, which are mutually independent of all service times at any other queue. At the end of a visit to queue i , the group of servers moves to queue $i + 1$, incurring a switch-over time D_i and a realisation of V_{i+1} is drawn. We assume that $\{D_i\}$ is a sequence of independent random variables. The total switch-over time during a full cycle is $D = \sum_{i=1}^N D_i$, and the length of a full cycle is denoted by the random variable C . We assume that all random variables so far are mutually independent.

During the visit time of the group of servers to queue i , a customer present at queue i at the polling instant of that queue will successfully complete his service with probability $p_i(V_i) = \mathbb{P}[B_i \leq V_i \mid V_i]$. We assume that if the service of a customer of queue i is not completed during a single visit, then at the next visit a new service time will be drawn from the service time distribution of B_i for that particular customer.

For a generic random variable Y_i , we denote its first two moments by $\mathbb{E}[Y_i]$ and $\mathbb{E}[Y_i^2]$, respectively. Thus, for example, $\mathbb{E}[V_i]$ is the mean visit time of the servers at queue i . By convention, $\sum_{i \neq j} Y_i = \sum_{\substack{i=1 \\ i \neq j}}^N Y_i$, and similarly for the product operator. All further notation will be introduced when it is first used.

Law of motion

Let X_i^j , $i, j = 1, 2, \dots, N$, denote the number of customers in queue j at the moment when queue i is polled and let $A_j(t)$ denote the number of Poisson arrivals to queue j during a time interval of length t . The law of motion describing the evolution of the system when the server moves from queue i to queue $i + 1$ connects X_{i+1}^j to X_i^j and is given by

$$X_{i+1}^j = \begin{cases} X_i^j + A_{j,1}(V_i) + A_{j,2}(D_i), & j \neq i, \\ \text{Binom}(X_i^i, 1 - p_i(V_i)) + \text{Poisson}(\Lambda_i(V_i)) + A_i(D_i), & j = i, \end{cases} \quad (2.1)$$

where for all k , $A_{j,k}(t)$ is an i.i.d. copy of $A_j(t)$, $\text{Binom}(n, p)$ is a binomial random variable with parameters n and p , and $\text{Poisson}(\Lambda_i(t))$ is a Poisson random variable with rate

$$\Lambda_i(t) = \lambda_i \int_0^t \mathbb{P}[B_i > y] dy.$$

Note that from (2.1) we see that for all j , the random variables X_i^j are independent of V_i and D_i , which is evident, considering that the number of customers in a queue at the beginning of a visit does not depend on the length of the upcoming visit time or switch-over time.

The relation for $j \neq i$ is straightforward. The number of customers at queue j at polling instant of queue $i + 1$ equals the number of customers that were there at polling instant of queue i plus all customers that arrived during the visit time of queue i and the switch-over time from queue i to queue $i + 1$.

For $j = i$, the relation is more involved. When the servers start polling queue i they encounter X_i^i customers. After V_i time units, only a binomial number of customers out of the initial X_i^i is still present. The probability that a single customer does not complete his service after V_i time units is $1 - p_i(V_i) = \mathbb{P}[B_i > V_i | V_i]$. In addition, there is a stream of new arrivals to queue i . The number of customers present at time t in an M/G/ ∞ queue (starting with zero customers at time $t = 0$) is Poisson distributed with rate $\Lambda_i(t)$, as it is given above; see Takács [15]. The last term at the right-hand side of (2.1) incorporates the customers that arrived at the queue during the switch-over time from queue i to queue $i + 1$.

We shall employ this relation to derive the mean queue length and the probability generating function of all queues at a polling instant.

3 Queue lengths at polling instants

One of the main tools used in the analysis of polling systems is the derivation of a set of multi-dimensional probability generating functions of the number of jobs present in the various queues at a polling instant of queue i . The common method is to derive the probability generating function of a given queue at some polling instant in terms of the probability generating function of the same queue at the previous polling instant. Then, from the set of N (implicit) dependent equations of the unknown probability generating functions one can obtain expressions which allow for numerical calculation of the mean queue length at each queue. These equations simplify significantly for several cases of the distribution of the visit times. In this section, we use the law of motion (buffer occupancy), which is given by Equation (2.1) and apply this technique to compute recursively the first moment and the probability generating function of the queue length distributions at a polling instant.

3.1 Mean queue length

From (2.1) we have the following relation for the mean queue length of queue j at two consecutive polling instants.

$$\mathbb{E}[X_{i+1}^j] = \begin{cases} \mathbb{E}[X_i^j] + \lambda_j \mathbb{E}[V_i] + \lambda_j \mathbb{E}[D_i], & j \neq i, \\ (1 - p_i) \mathbb{E}[X_i^j] + \mathbb{E}[\Lambda_i(V_i)] + \lambda_i \mathbb{E}[D_i], & j = i, \end{cases} \quad (3.1)$$

where $p_i = \mathbb{P}[B_i \leq V_i] = \mathbb{E}[p_i(V_i)]$. Summing (3.1) over i we obtain

$$p_j \mathbb{E}[X_j^j] = \lambda_j \sum_{i \neq j} \mathbb{E}[V_i] + \mathbb{E}[\Lambda_j(V_j)] + \lambda_j \mathbb{E}[D]. \quad (3.2)$$

Indeed, in steady state, the mean number of jobs in queue j at a polling instant equals the fraction of jobs $(1 - p_j) \mathbb{E}[X_j^j]$ left behind at the end of the previous visit, plus the mean number of arrivals during the cycle time out of queue j , which is $\lambda_j \left(\sum_{i \neq j} \mathbb{E}[V_i] + \mathbb{E}[D] \right)$, plus the mean number of customers in a M/G/ ∞ queue at time V_j . The mean queue length of queue j at polling instant of queue i is easily derived from (3.1), yielding

$$\mathbb{E}[X_i^j] = \mathbb{E}[X_j^j] (1 - p_j) + \mathbb{E}[\Lambda_j(V_j)] + \lambda_j \sum_{k=j+1}^{i-1} \mathbb{E}[V_k] + \lambda_j \sum_{k=j}^{i-1} \mathbb{E}[D_k]. \quad (3.3)$$

For example, suppose that B_j is exponentially distributed with parameter μ_j . Then,

$$\Lambda_j(V_j) = \lambda_j \int_0^{V_j} e^{-\mu_j y} dy = \frac{\lambda_j}{\mu_j} (1 - e^{-\mu_j V_j}).$$

Thus, $\mathbb{E}[\Lambda_j(V_j)] = \lambda_j (1 - \mathbb{E}[e^{-\mu_j V_j}]) / \mu_j$. So, in particular, if V_j is also exponentially distributed with parameter γ_j , then we have that $\mathbb{E}[\Lambda_j(V_j)] = \lambda_j / (\gamma_j + \mu_j)$, and the mean queue length of each queue can now easily be computed recursively from (3.3).

3.2 Recursive relation for the generating function

Define the generating function of the queue length of all queues at polling instants of queue i as $G_i(\mathbf{z}) = \mathbb{E}[\prod_{j=1}^N z_j^{X_i^j}]$. Then, from (2.1) we have that

$$G_{i+1}(\mathbf{z}) = \mathbb{E}[\prod_{j \neq i} z_j^{X_i^j + A_{j,1}(V_i) + A_{j,2}(D_i)} z_i^{\text{Binom}(X_i^i, 1 - p_i(V_i)) + \text{Poisson}(\Lambda_i(V_i)) + A_i(D_i)}] \quad (3.4)$$

By conditioning on the vector (X_i^1, \dots, X_i^N) , on V_i , and on D_i , Equation (3.4) becomes

$$\begin{aligned} G_{i+1}(\mathbf{z}) &= \mathbb{E}[\prod_{j \neq i} z_j^{X_i^j}] \mathbb{E}[\prod_{j \neq i} z_j^{A_{j,1}(V_i)} \mid V_i] \mathbb{E}[\prod_{j=1}^N z_j^{A_{j,2}(D_i)} \mid D_i] \times \\ &\quad \times \mathbb{E}[z_i^{\text{Binom}(X_i^i, 1 - p_i(V_i))} \mid V_i] \mathbb{E}[z_i^{\text{Poisson}(\Lambda_i(V_i))} \mid V_i]. \end{aligned} \quad (3.5)$$

Since the number of arrivals at any queue during a fixed amount of time is independent of the number of arrivals at any other queue during the same given period, we have that

$$\begin{aligned} \mathbb{E}[\prod_{j=1}^N z_j^{A_j(D_i)} \mid D_i = x] &= \mathbb{E}[\prod_{j=1}^N z_j^{A_j(x)}] = \prod_{j=1}^N \mathbb{E}[z_j^{A_j(x)}] \\ &= \prod_{j=1}^N \sum_{n=0}^{\infty} z_j^n \frac{(\lambda_j x)^n}{n!} e^{-\lambda_j x} = \prod_{j=1}^N e^{-\lambda_j x (1 - z_j)}. \end{aligned}$$

Therefore, we have

$$\mathbb{E}\left[\prod_{j=1}^N z_j^{A_j(D_i)} \mid D_i\right] = e^{-D_i \sum_{j=1}^N \lambda_j (1-z_j)}. \quad (3.6)$$

Likewise, we obtain that

$$\mathbb{E}\left[\prod_{j \neq i} z_j^{A_j(V_i)} \mid V_i\right] = e^{-V_i \sum_{j \neq i} \lambda_j (1-z_j)}. \quad (3.7)$$

Moreover,

$$\begin{aligned} \mathbb{E}[z_i^{\text{Binom}(X_i^i, 1-p_i(V_i))} \mid V_i = x] &= \mathbb{E}[z_i^{\text{Binom}(X_i^i, 1-p_i(x))}] \\ &= \sum_{\ell=0}^{X_i^i} z_i^\ell \binom{X_i^i}{\ell} (1-p_i(x))^\ell p_i(x)^{X_i^i-\ell} = (p_i(x) + z_i[1-p_i(x)])^{X_i^i}, \end{aligned}$$

or in other words,

$$\mathbb{E}[z_i^{\text{Binom}(X_i^i, 1-p_i(V_i))} \mid V_i] = (p_i(V_i) + z_i[1-p_i(V_i)])^{X_i^i}. \quad (3.8)$$

For the last term of the right-hand side of (3.5) we have that

$$\mathbb{E}[z_i^{\text{Poisson}(\Lambda_i(V_i))} \mid V_i = x] = \mathbb{E}[z_i^{\text{Poisson}(\Lambda_i(x))}] = \sum_{n=0}^{\infty} z_i^n \frac{(\Lambda_i(x))^n}{n!} e^{-\Lambda_i(x)} = e^{-\Lambda_i(x)(1-z_i)},$$

which yields that

$$\mathbb{E}[z_i^{\text{Poisson}(\Lambda_i(V_i))} \mid V_i] = e^{-\Lambda_i(V_i)(1-z_i)}. \quad (3.9)$$

Substituting (3.6) – (3.9) into (3.5), we obtain

$$G_{i+1}(\mathbf{z}) = \mathbb{E}\left[\prod_{j \neq i} z_j^{X_j^j} e^{-V_i \sum_{j \neq i} \lambda_j (1-z_j)} e^{-D_i \sum_{j=1}^N \lambda_j (1-z_j)} (p_i(V_i) + z_i[1-p_i(V_i)])^{X_i^i} e^{-\Lambda_i(V_i)(1-z_i)}\right]. \quad (3.10)$$

Recall that for all j , the random variables X_j^j are independent of V_i and D_i . Therefore, Equation (3.10) becomes

$$\begin{aligned} G_{i+1}(\mathbf{z}) &= \mathbb{E}[e^{-D_i \sum_{j=1}^N \lambda_j (1-z_j)}] \times \\ &\quad \times \mathbb{E}\left[\prod_{j \neq i} z_j^{X_j^j} e^{-V_i \sum_{j \neq i} \lambda_j (1-z_j)} (p_i(V_i) + z_i[1-p_i(V_i)])^{X_i^i} e^{-\Lambda_i(V_i)(1-z_i)}\right]. \end{aligned} \quad (3.11)$$

Consequently,

$$\begin{aligned} G_{i+1}(\mathbf{z}) &= \tilde{D}_i\left(\sum_{j=1}^N \lambda_j (1-z_j)\right) \times \\ &\quad \times \mathbb{E}[e^{-V_i \sum_{j \neq i} \lambda_j (1-z_j)} e^{-\Lambda_i(V_i)(1-z_i)} G_i(z_1, z_2, \dots, z_{i-1}, p_i(V_i) + [1-p_i(V_i)]z_i, z_{i+1}, \dots, z_N)], \end{aligned} \quad (3.12)$$

where $\tilde{D}_i(s) = \mathbb{E}[e^{-sD_i}]$ denotes the Laplace-Stieltjes transform of the random variable D_i . Evidently, if V_i follows a discrete distribution, the above expression simplifies significantly. Note that the mean queue length at a polling instant (3.3) can also be obtained by differentiating Equation (3.12).

Remark 1. Applying similar techniques, we can also derive the probability generating function of the number of customers at the end of a visit at queue $i + 1$. If we denote by Y_i^j the number of customers in queue $j = 1, \dots, N$ at the moment when the service at queue $i = 1, \dots, N$ is completed, then the law of motion describing the evolution of the system is given by

$$Y_{i+1}^j = \begin{cases} Y_i^j + A_j(D_i) + A_j(V_{i+1}), & j \neq i + 1, \\ \text{Binom}(Y_i^j + A_j(D_i), 1 - p_j(V_j)) + \text{Poisson}(\Lambda_j(V_j)), & j = i + 1. \end{cases} \quad (3.13)$$

Also note that the expected value of Y_i^j can be easily computed from (3.3) by observing that for all $j \neq i$, $Y_i^j = X_i^j + A_j(V_i)$, while for the i -th queue we have that $Y_i^i = \text{Binom}(X_i^i, 1 - p_i(V_i)) + \text{Poisson}(\Lambda_i(V_i))$.

4 Sojourn time

Let the sojourn time of a customer at queue i be denoted by S_i . We compute its expected value (and thus, by Little's law, also the mean queue length of queue i at an arbitrary moment), and we derive the Laplace-Stieltjes transform of S_i . As stated before, for each queue we assume that if the service of a customer is not completed during a visit, then, for the next visit at that queue, a new service time will be resampled for the same customer from the service time distribution of that queue.

4.1 Mean sojourn time

Recall that the cycle time is given by $C = \sum_{i=1}^N (V_i + D_i)$. In order to derive the mean sojourn time of a customer arriving at queue i , we shall need some further notation. Denote by V_i^{res} the residual visit time of the group of servers at queue i and by $C_{/i}$ the cycle time except the time spent serving queue i , i.e. $C_{/i} = C - V_i$. Similarly, $C_{/i}^{res}$ represents the residual cycle time excluding the visit time of queue i . That is, $C_{/i}^{res}$ measures the length of time from a random moment after leaving queue i until the next polling instant of queue i . Furthermore, let $\{C_m\}$ be a family of i.i.d. random variables distributed like C , and N_i be a (shifted) geometric random variable with success probability $p_i = \mathbb{E}[p_i(V_i)] = \mathbb{P}[B_i \leq V_i]$, i.e. $\mathbb{P}[N_i = n] = (1 - p_i)^n p_i$, for all integer $n \geq 0$. One should observe here that $N_i + 1$ is a stopping time as it is the first time when the service time of a customer in queue i is less than or equal to the visit time at that queue; that is, $N_i + 1 = \inf\{k : B_{i,k} \leq V_{i,k}\}$, where $B_{i,k}$ and $V_{i,k}$ are i.i.d. copies of B_i and V_i respectively. Similarly, a second index is added to a random variable, every time that we explicitly need to indicate that an independent copy is considered. Then the sojourn time of a customer at queue i is given by

$$S_i = \begin{cases} B_{i,0}, & (\text{arrival during } V_i \text{ and } B_{i,0} \leq V_{i,0}^{res}), \\ V_{i,0}^{res} + \sum_{m=1}^{N_i} C_m + C_{/i} + B_{i,N_i+1}, & (\text{arrival during } V_{i,0} \text{ and } B_{i,0} > V_{i,0}^{res}), \\ C_{/i}^{res} + \sum_{m=1}^{N_i} C_m + B_{i,N_i+1}, & (\text{arrival during } C_{/i}). \end{cases} \quad (4.1)$$

Note that the probability of an arrival occurring during the visit time of queue i is $\mathbb{E}[V_i]/\mathbb{E}[C]$, i.e. the expected visit time of queue i over the expected cycle time, and similarly for the other two events. Therefore, from (4.1) we obtain that the expected sojourn time of a customer of

queue i is given by

$$\begin{aligned}\mathbb{E}[S_i] &= \frac{\mathbb{E}[V_i]}{\mathbb{E}[C]} \mathbb{P}[B_i \leq V_i^{res}] \mathbb{E}[B_i \mid B_i \leq V_i^{res}] + \\ &\quad + \frac{\mathbb{E}[V_i]}{\mathbb{E}[C]} \mathbb{P}[B_i > V_i^{res}] \mathbb{E}[V_{i,0}^{res} + \sum_{m=1}^{N_i} C_m + C_{/i} + B_{i,N_i+1} \mid B_{i,0} \geq V_{i,0}^{res}] + \\ &\quad + \frac{\mathbb{E}[C_{/i}]}{\mathbb{E}[C]} \mathbb{E}[C_{/i}^{res} + \sum_{m=1}^{N_i} C_m + B_{i,N_i+1}].\end{aligned}\quad (4.2)$$

In order to compute the second conditional expectation appearing at the right-hand side of the above equation, we think as follows. For N_i cycles, the service of the customer is not completed during that visit because for every visit $B_i > V_i$, while at the $N_i + 1$ st visit the service is completed within that cycle. Therefore, define

$$\bar{C}_m = C_{/i,m} + \min(B_{i,m}, V_{i,m})$$

and observe that

$$\mathbb{E}\left[\sum_{m=1}^{N_i+1} \bar{C}_m\right] = \mathbb{E}\left[\sum_{m=1}^{N_i} C_m + C_{/i} + B_i\right].$$

Thus

$$\begin{aligned}\mathbb{E}[V_{i,0}^{res} + \sum_{m=1}^{N_i} C_m + C_{/i} + B_{i,N_i+1} \mid B_{i,0} \geq V_{i,0}^{res}] &= \mathbb{E}[V_i^{res} \mid B_i \geq V_i^{res}] + \mathbb{E}\left[\sum_{m=1}^{N_i+1} \bar{C}_m\right] \\ &= \mathbb{E}[V_i^{res} \mid B_i \geq V_i^{res}] + \mathbb{E}[N_i + 1] \mathbb{E}[\bar{C}_m] \\ &= \mathbb{E}[V_i^{res} \mid B_i \geq V_i^{res}] + \mathbb{E}[N_i + 1] (\mathbb{E}[C_{/i}] + \mathbb{E}[\min(B_i, V_i)]),\end{aligned}\quad (4.3)$$

where in the second equality we used Wald's equation.

For the third conditional expectation appearing at the right-hand side of (4.2), we have that

$$\begin{aligned}\mathbb{E}[C_{/i}^{res} + \sum_{m=1}^{N_i} C_m + B_{i,N_i+1}] &= \frac{\mathbb{E}[C_{/i}^2]}{2\mathbb{E}[C_{/i}]} + \mathbb{E}\left[\sum_{m=1}^{N_i} (C_{/i,m} + V_{i,m}) + B_{i,N_i+1}\right] \\ &= \frac{\mathbb{E}[C_{/i}^2]}{2\mathbb{E}[C_{/i}]} + \mathbb{E}\left[\sum_{m=1}^{N_i} C_{/i,m}\right] + \mathbb{E}\left[\sum_{m=1}^{N_i+1} \min(B_{i,m}, V_{i,m})\right] \\ &= \frac{\mathbb{E}[C_{/i}^2]}{2\mathbb{E}[C_{/i}]} + \mathbb{E}[N_i] \mathbb{E}[C_{/i}] + \mathbb{E}[N_i + 1] \mathbb{E}[\min(B_i, V_i)].\end{aligned}\quad (4.4)$$

Summarising the above, we have that

$$\begin{aligned}\mathbb{E}[S_i] &= \frac{\mathbb{E}[V_i]}{\mathbb{E}[C]} \mathbb{P}[B_i \leq V_i^{res}] \mathbb{E}[B_i \mid B_i \leq V_i^{res}] + \\ &\quad + \frac{\mathbb{E}[V_i]}{\mathbb{E}[C]} \mathbb{P}[B_i > V_i^{res}] \left(\mathbb{E}[V_i^{res} \mid B_i > V_i^{res}] + \mathbb{E}[N_i + 1] (\mathbb{E}[C_{/i}] + \mathbb{E}[\min(B_i, V_i)]) \right) + \\ &\quad + \frac{\mathbb{E}[C_{/i}]}{\mathbb{E}[C]} \left(\frac{\mathbb{E}[C_{/i}^2]}{2\mathbb{E}[C_{/i}]} + \mathbb{E}[N_i] \mathbb{E}[C_{/i}] + \mathbb{E}[N_i + 1] \mathbb{E}[\min(B_i, V_i)] \right).\end{aligned}\quad (4.5)$$

In Section 5 we shall illustrate through an example the effect of the first two moments of the visit time and the service time on the mean sojourn time of an arbitrary customer.

4.2 The Laplace-Stieltjes transform

We now derive the Laplace-Stieltjes transform of the sojourn time of a customer of queue i . We first rewrite Equation (4.1) in terms of the Laplace-Stieltjes transforms of all variables involved (cf. (4.5)), and thus we get that

$$\begin{aligned} \mathbb{E}[e^{-sS_i}] &= \frac{\mathbb{E}[V_i]}{\mathbb{E}[C]} \mathbb{P}[B_i \leq V_i^{res}] \mathbb{E}[e^{-sB_i} \mid B_i \leq V_i^{res}] + \\ &+ \frac{\mathbb{E}[V_i]}{\mathbb{E}[C]} \mathbb{P}[B_i > V_i^{res}] \mathbb{E}[e^{-sV_i^{res}} \mid B_i > V_i^{res}] \mathbb{E}[e^{-s \sum_{m=1}^{N_i+1} \bar{C}_m}] + \\ &+ \frac{\mathbb{E}[C/i]}{\mathbb{E}[C]} \mathbb{E}[e^{-sC/i}] \mathbb{E}[e^{-s \sum_{m=1}^{N_i} C_{/i,m}}] \mathbb{E}[e^{-s \sum_{m=1}^{N_i+1} \min(B_{i,m}, V_{i,m})}]. \end{aligned} \quad (4.6)$$

We rewrite several of the terms appearing above as follows. The probability density function of V_i^{res} is given by

$$\mathbb{P}[V_i^{res} \leq x] = \frac{1}{\mathbb{E}[V_i]} \int_0^x \mathbb{P}[V_i > y] dy,$$

yielding

$$\mathbb{P}[B_i > V_i^{res}] = \frac{1}{\mathbb{E}[V_i]} \int_0^\infty \mathbb{P}[B_i > x] \mathbb{P}[V_i > x] dx. \quad (4.7)$$

Similarly, we have that

$$\mathbb{P}[C_{/i}^{res} \leq x] = \frac{1}{\mathbb{E}[C/i]} \int_0^x \mathbb{P}[C_{/i} > y] dy,$$

which implies that

$$\mathbb{E}[e^{-sC_{/i}^{res}}] = \frac{1 - \tilde{C}_{/i}(s)}{s\mathbb{E}[C/i]}, \quad (4.8)$$

where $\tilde{C}_{/i}$ denotes the Laplace-Stieltjes transform of the random variable $C_{/i}$. Moreover,

$$\begin{aligned} \mathbb{E}[e^{-s \sum_{m=1}^{N_i+1} \bar{C}_m}] &= \sum_{n=0}^{\infty} \mathbb{E}[e^{-s \sum_{m=1}^{n+1} \bar{C}_m}] (1-p_i)^n p_i = \sum_{n=0}^{\infty} \mathbb{E}[e^{-s\bar{C}}]^{n+1} (1-p_i)^n p_i \\ &= \frac{p_i \mathbb{E}[e^{-s\bar{C}}]}{1 - (1-p_i) \mathbb{E}[e^{-s\bar{C}}]}, \end{aligned} \quad (4.9)$$

where $\bar{C} = C_{/i} + \min(B_i, V_i)$. Likewise, we have that

$$\mathbb{E}[e^{-s \sum_{m=1}^{N_i} C_{/i,m}}] = \frac{p_i}{1 - (1-p_i) \mathbb{E}[e^{-sC_{/i}}]} \quad (4.10)$$

and

$$\mathbb{E}[e^{-s \sum_{m=1}^{N_i+1} \min(B_{i,m}, V_{i,m})}] = \frac{p_i \mathbb{E}[e^{-s \min(B_i, V_i)}]}{1 - (1-p_i) \mathbb{E}[e^{-s \min(B_i, V_i)}]}. \quad (4.11)$$

Substituting (4.7) – (4.11) into (4.6) we have that the Laplace-Stieltjes transform of the sojourn time of a customer of queue i is given by

$$\begin{aligned} \mathbb{E}[e^{-sS_i}] &= \frac{1}{\mathbb{E}[C]} \mathbb{E}[e^{-sB_i} \mid B_i \leq V_i^{res}] \int_0^\infty \mathbb{P}[B_i \leq x] \mathbb{P}[V_i > x] dx + \\ &+ \frac{1}{\mathbb{E}[C]} \mathbb{E}[e^{-sV_i^{res}} \mid B_i > V_i^{res}] \frac{p_i \mathbb{E}[e^{-s\bar{C}}]}{1 - (1-p_i) \mathbb{E}[e^{-s\bar{C}}]} \int_0^\infty \mathbb{P}[B_i > x] \mathbb{P}[V_i > x] dx + \\ &+ \frac{1 - \tilde{C}_{/i}(s)}{s\mathbb{E}[C]} \frac{p_i}{1 - (1-p_i) \tilde{C}_{/i}(s)} \frac{p_i \mathbb{E}[e^{-s \min(B_i, V_i)}]}{1 - (1-p_i) \mathbb{E}[e^{-s \min(B_i, V_i)}]}. \end{aligned} \quad (4.12)$$

Clearly, from the expression above, one can retrieve Equation (4.5) for the mean sojourn time of a customer of queue i .

The transforms appearing in (4.12) may be cumbersome to compute when the service times or the visit times are generally distributed. However, when both B_i and V_i follow a phase-type distribution, all transforms can be computed explicitly since the class of phase-type distributions is closed under finite minima. Phase-type distributions are widely used in computations. The class of phase-type distributions is dense (in the sense of weak convergence) in the class of all distributions on $(0, \infty)$ (cf. [2, Propositions 1 and 2]). As an example, we will derive the Laplace-Stieltjes transform of the sojourn time of a customer of queue i , as well as its mean, in case both the visit time and the service time at queue i are exponentially distributed.

4.3 A special case

Let the service time and the visit time at queue i be exponentially distributed with rates μ_i and γ_i respectively. Then all terms appearing in (4.5) can be easily computed in terms of μ_i and γ_i . For example,

$$\mathbb{E}[B_i \mid B_i \leq V_i^{res}] = \frac{1}{\gamma_i + \mu_i}$$

and $\mathbb{P}[B_i > V_i^{res}] = \gamma_i/(\gamma_i + \mu_i)$. Thus, (4.5) becomes

$$\begin{aligned} \mathbb{E}[S_i] = & \frac{1}{\gamma_i \mathbb{E}[C]} \frac{\mu_i}{\gamma_i + \mu_i} \frac{1}{\gamma_i + \mu_i} + \frac{1}{\gamma_i \mathbb{E}[C]} \frac{\gamma_i}{\gamma_i + \mu_i} \left(\frac{1}{\gamma_i + \mu_i} + \left(\frac{\gamma_i}{\mu_i} + 1 \right) \left(\mathbb{E}[C_{/i}] + \frac{1}{\gamma_i + \mu_i} \right) \right) + \\ & + \frac{\mathbb{E}[C_{/i}]}{\mathbb{E}[C]} \left(\frac{\mathbb{E}[C_{/i}^2]}{2\mathbb{E}[C_{/i}]} + \frac{\gamma_i}{\mu_i} \mathbb{E}[C_{/i}] + \left(\frac{\gamma_i}{\mu_i} + 1 \right) \frac{1}{\gamma_i + \mu_i} \right) \end{aligned}$$

or

$$\mathbb{E}[S_i] = \frac{(\gamma_i \mathbb{E}[C_{/i}] + 1)^2}{\gamma_i \mu_i \mathbb{E}[C]} + \frac{\mathbb{E}[C_{/i}^2]}{2\mathbb{E}[C]}.$$

Similarly, (4.12) reduces to

$$\begin{aligned} \mathbb{E}[e^{-sS_i}] = & \frac{1}{\mathbb{E}[C]} \frac{\gamma_i + \mu_i}{\gamma_i + \mu_i + s} \frac{\mu_i}{\gamma_i(\gamma_i + \mu_i)} + \frac{1}{\mathbb{E}[C]} \frac{\gamma_i + \mu_i}{\gamma_i + \mu_i + s} \frac{\mu_i \mathbb{E}[e^{-s\bar{C}}]}{\gamma_i + \mu_i - \gamma_i \mathbb{E}[e^{-s\bar{C}}]} \frac{1}{\gamma_i + \mu_i} + \\ & + \frac{1 - \tilde{C}_{/i}(s)}{s\mathbb{E}[C]} \frac{\mu_i}{\gamma_i + \mu_i - \gamma_i \tilde{C}_{/i}(s)} \frac{\mu_i \frac{\gamma_i + \mu_i}{\gamma_i + \mu_i + s}}{\gamma_i + \mu_i - \gamma_i \frac{\gamma_i + \mu_i}{\gamma_i + \mu_i + s}}. \end{aligned}$$

Since $\mathbb{E}[e^{-s\bar{C}}] = \tilde{C}_{/i}(s) \mathbb{E}[e^{-s \min(B_i, V_i)}]$ we have that the previous expression reduces to

$$\begin{aligned} \mathbb{E}[e^{-sS_i}] = & \frac{1}{\mathbb{E}[C]} \frac{1}{\gamma_i + \mu_i + s} \frac{\mu_i}{\gamma_i} + \frac{1}{\mathbb{E}[C]} \frac{1}{\gamma_i + \mu_i + s} \frac{\mu_i \tilde{C}_{/i}(s)}{\gamma_i + \mu_i + s - \gamma_i \tilde{C}_{/i}(s)} + \\ & + \frac{1 - \tilde{C}_{/i}(s)}{s\mathbb{E}[C]} \frac{\mu_i}{\gamma_i + \mu_i - \gamma_i \tilde{C}_{/i}(s)} \frac{\mu_i}{\mu_i + s}. \end{aligned}$$

Similar expressions can be easily derived in case both the visit times and the service times follow some phase-type distribution, such as Gamma, hyperexponential, or Coxian distributions.

5 Numerical results

This section is devoted to some numerical results. In particular, we want to examine numerically the effect of the first two moments of the visit and service times on the sojourn time of an arbitrary customer. In all examples, we make the following assumptions. We consider a polling system with two queues. The arrival rate at the first queue is $\lambda_1 = 0.8$ and at the second queue it is $\lambda_2 = 0.5$. The service time and the visit time at the first queue are exponentially distributed with rates $\mu_1 = 1$ and $\gamma_1 = 1$ respectively. Moreover, the total mean switch-over time is taken to be $\mathbb{E}[D] = 0.5$, while its second moment is assumed to be zero. In all figures that follow, we plot the mean sojourn time of an arbitrary customer, which is estimated by $(\lambda_1\mathbb{E}[S_1] + \lambda_2\mathbb{E}[S_2])/(\lambda_1 + \lambda_2)$.

In Figures 1 and 2 we investigate the effect of the first two moments of the service time at the second queue on the mean sojourn time of an arbitrary customer. For these plots, the visit time at the second queue is considered to be exponentially distributed with rate $\gamma_2 = 3/2$. For various values of the squared coefficient of variation of the service time at the second queue, which is denoted by $c_{B_2}^2$, we plot in Figure 1 the mean sojourn time of an arbitrary customer versus the mean service time $\mathbb{E}[B_2]$. The squared coefficient of variation of the service time is chosen to be comparable to the squared coefficient of variation of the (exponentially distributed) visit time, which is equal to 1. In Figure 2, we plot the mean sojourn time of an arbitrary customer versus $c_{B_2}^2$ for three values of $\mathbb{E}[B_2]$, which again are chosen to be comparable to $\mathbb{E}[V_2]$.

For each case of $c_{B_2}^2$, we fit a mixed Erlang or hyperexponential distribution to $\mathbb{E}[B_2]$ and $c_{B_2}^2$, depending on whether the squared coefficient of variation is less or greater than one; see, e.g., Tijms [19]. So, if $1/n \leq c_{B_2}^2 \leq 1/(n-1)$ for some $n = 2, 3, \dots$, then the mean and squared coefficient of variation of the mixed-Erlang distribution

$$G(x) = p \left(1 - e^{-\zeta x} \sum_{j=0}^{n-2} \frac{(\zeta x)^j}{j!} \right) + (1-p) \left(1 - e^{-\zeta x} \sum_{j=0}^{n-1} \frac{(\zeta x)^j}{j!} \right), \quad x \geq 0,$$

matches with $\mathbb{E}[B_2]$ and $c_{B_2}^2$, provided the parameters p and ζ are chosen as

$$p = \frac{1}{1 + c_{B_2}^2} \left(n c_{B_2}^2 - \sqrt{n(1 + c_{B_2}^2) - n^2 c_{B_2}^2} \right), \quad \zeta = \frac{n-p}{\mathbb{E}[B_2]}.$$

On the other hand, if $c_{B_2}^2 > 1$, then the mean and squared coefficient of variation of the hyperexponential distribution

$$G(x) = p(1 - e^{-\zeta_1 x}) + q(1 - e^{-\zeta_2 x}), \quad x \geq 0,$$

match with $\mathbb{E}[B_2]$ and $c_{B_2}^2$, provided the parameters ζ_1 , ζ_2 , p , and q are chosen as

$$\begin{aligned} p &= \frac{1}{2} \left(1 + \sqrt{\frac{c_{B_2}^2 - 1}{c_{B_2}^2 + 1}} \right), & q &= 1 - p, \\ \zeta_1 &= \frac{2p}{\mathbb{E}[B_2]} & \text{and } \zeta_2 &= \frac{2q}{\mathbb{E}[B_2]}. \end{aligned}$$

As is evident from the plot in Figure 1, the expected sojourn time of an arbitrary customer increases as the mean service time at the second queue increases. Moreover, the rate that it increases with is almost linear as $c_{B_2}^2$ grows and the effect of the second moment is less pronounced than the effect of $\mathbb{E}[B_2]$.

In Figure 2, one observes that the mean sojourn time of an arbitrary customer decreases as the squared coefficient of variation of the service time increases, contrary to what is the case for

the M/G/1 queue. This result is due to the fact that the service time of a customer that did not complete his service during one visit time is resampled for the following visit time. Therefore, the larger the variability in the service times, the bigger is the probability that during the next visit time this particular customer will complete his service. Recall that the mean visit time at the second queue is equal to $2/3$ and observe that in case $\mathbb{E}[B_2]$ is less than $\mathbb{E}[V_2]$, the effect of the second moment of the service time on the mean sojourn time of an arbitrary customer is almost negligible.

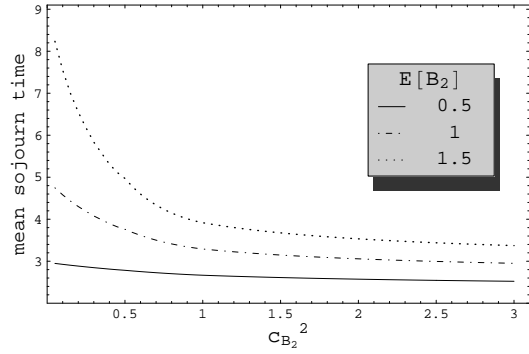
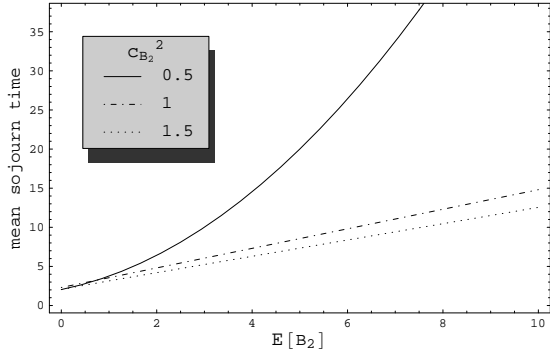


Figure 1: Mean sojourn time of an arbitrary customer against the mean service time $\mathbb{E}[B_2]$. Figure 2: Mean sojourn time of an arbitrary customer against the squared coefficient of variation of the service time B_2 .

In Figures 3 and 4 we now investigate the effects of the first two moments of the visit time at the second queue on the mean sojourn time of an arbitrary customer. For these plots, we now take the service time at the second queue to be exponentially distributed with rate $\mu_2 = 3/2$. For various values of the squared coefficient of variation of the visit time at the second queue, which is denoted by $c_{V_2}^2$, we plot in Figure 3 the mean sojourn time of an arbitrary customer versus the mean visit time $\mathbb{E}[V_2]$. As before, the squared coefficient of variation of the visit time is chosen to be comparable with the squared coefficient of variation of the (exponentially distributed) service time, which is equal to 1. In Figure 4, we plot the mean sojourn time of an arbitrary customer versus $c_{V_2}^2$ for three values of $\mathbb{E}[V_2]$, which again are chosen to be comparable with $\mathbb{E}[B_2]$.

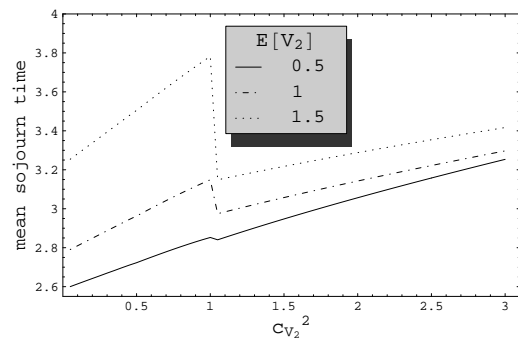
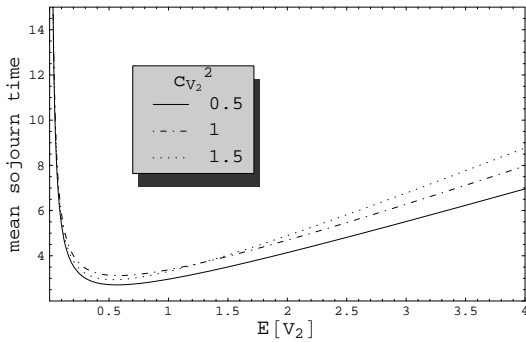


Figure 3: Mean sojourn time of an arbitrary customer against the mean visit time $\mathbb{E}[V_2]$. Figure 4: Mean sojourn time of an arbitrary customer against the squared coefficient of variation of the visit time V_2 .

The plot in Figure 3 is interesting. Evidently, when $\mathbb{E}[V_2]$ is significantly smaller than $\mathbb{E}[B_2]$, only a very small number of customers will be served during a visit. As the mean visit time increases, more customers are served during a visit and the mean sojourn time of an arbitrary

customer is reduced. However, as the mean visit time continues to increase, this trend is reversed after the mean sojourn time of an arbitrary customer reaches a global minimum. In other words, there is an optimal value for the mean visit time to some queue that minimises the mean sojourn time of an arbitrary customer; beyond that value, the mean sojourn time of an arbitrary customer increases at an almost linear rate. This indicates that the polling system under consideration can be optimised in expectation by controlling the visit time to each queue. In the following section, we will develop a policy that minimises the mean sojourn time of an arbitrary customer in the system.

The plot in Figure 4 is also interesting. As is explained above, we fit either a mixture of Erlang distributions or a hyperexponential distribution to each pair of the first two moments of the visit time, depending on the value of the squared coefficient of variation. For every value of $c_{V_2}^2$, we obtain a different visit time distribution. Note that the jump in Figure 4 occurs when the distribution we fit to the first two moments of the visit time shifts from a mixture of Erlang distributions to a hyperexponential distribution. This indicates that the shape of the visit time distribution is important; for example, hyperexponential distributions are always unimodal, which is not the case for mixed Erlang distributions. Consequently, the first two moments cannot capture sufficiently the effect of the visit time distribution on the mean sojourn time of an arbitrary customer; one needs to know the exact distribution.

6 Dynamic control of servers' visits

A basic question that arises when planning efficient polling systems concerns the order of visits performed by the servers. As it is suggested by Figure 3, the polling system we are considering can be optimised in some way so that the mean sojourn time of an arbitrary customer is minimised. Rather than identifying the value of the minimum mean sojourn time for the cyclic processing (visiting) order considered so far, we will investigate whether there exists a fixed static order that the servers visit the various queues so that the mean sojourn time of an arbitrary customer is minimised.

It is immediately clear that a static order will not have the desired effect. The answer to the question whether there exists an optimal fixed order of visits, which will be followed in every cycle, that minimises the sum of the weighted mean sojourn times turns out to be *negative*. As is evident from Equation (4.5), the mean sojourn time of a customer of queue i *does not depend* on the order the queues have been visited, and thus neither does the weighted sum thereof. Since the mean sojourn time of an arbitrary customer remains unaffected when altering the processing order, it does not constitute a practical performance measure of our system.

An appealing approach that leads to a simple and tractable rule is to develop a semi-dynamic control scheme. The idea is to dispatch the group of servers to perform Hamiltonian tours, each tour being possibly different from the previous one, depending on the state of the system at the beginning of the tour, so as to optimise some performance measure. An adequate performance measure is the *throughput* of the system, namely the number of customers served per cycle, as the throughput can be measured per cycle, while the sojourn time of a customer spans over a random number of cycles. The goal is to maximise the throughput of the system for each cycle.

Specifically, suppose that at the beginning of a cycle $n = (n_1, n_2, \dots, n_N)$ is the state of the system, where n_i is the number of jobs waiting in queue i . We shall compute the expected value of the number of customers served per cycle under a specific processing order, and consequently identify the optimal processing order per cycle. The following theorem summarises the result.

Theorem 1. *For the Hamiltonian-tour approach, the optimal visiting order is independent of the number of customers present at the various queues at the start of the cycle and is given by*

the index-type rule

$$\frac{\lambda_i p_i}{\mathbb{E}[V_i] + \mathbb{E}[D_i]}$$

in the sense that the throughput is maximised if and only if the visiting order is arranged according to an increasing sequence of this rule.

Before proceeding with the proof we point out that having more information regarding the system, such as the number of customers present at all queues, has no effect on the optimal strategy, and thus it does not improve the performance of the system. This is a surprising result, since typically more information regarding a system leads to decisions that increase the efficiency of the system.

Proof. The proof follows from an interchange argument. Consider the processing order $\pi_0 = (1, 2, \dots, N)$. Denote by θ_i the throughput of queue i under this processing order, i.e. the number of customers of queue i that are served during a cycle, and denote by θ the total throughput of the system, i.e. the sum of all θ_i . Given n , i.e. the state of the system at the beginning of a cycle, we shall compute the expected value of θ .

The number of customers served after completing a visit at queue i is equal to the portion of customers that were present at polling instant i and successfully completed their service plus the number of customers that arrived during the visit time of queue i and completed their service within that visit. In other words,

$$\theta_i = \text{Binom}\left(n_i + A_i\left(\sum_{k=1}^{i-1} (V_k + D_k)\right), p_i\right) + \left(A_i(V_i) - \text{Poisson}(\Lambda_i(V_i))\right).$$

As a result,

$$\mathbb{E}[\theta_i] = \left(n_i + \lambda_i \sum_{k=1}^{i-1} (\mathbb{E}[V_k] + \mathbb{E}[D_k])\right) p_i + \lambda_i \mathbb{E}[V_i] - \mathbb{E}[\Lambda_i(V_i)],$$

which yields

$$\mathbb{E}[\theta] = c + \sum_{i=1}^N \lambda_i p_i \sum_{k=1}^{i-1} (\mathbb{E}[V_k] + \mathbb{E}[D_k]), \quad (6.1)$$

where

$$c = \sum_{i=1}^N (n_i p_i + \lambda_i \mathbb{E}[V_i] - \mathbb{E}[\Lambda_i(V_i)]).$$

Observe that the constant c that appears in (6.1) does not depend on π_0 , while the second term at the right-hand side of (6.1) does.

Consider now the processing order $\pi_1 = (1, 2, \dots, j-1, j+1, j, j+2, \dots, N)$, where the visit order of queues j and $j+1$ is interchanged and denote by θ'_i and θ' the throughput of queue i and of the whole system under π_1 , respectively. We promptly have that $\mathbb{E}[\theta_i] = \mathbb{E}[\theta'_i]$ for all $i \neq j, j+1$ and that

$$\begin{aligned} \mathbb{E}[\theta'_j] &= \left(n_j + \lambda_j \left(\sum_{k=1}^{j-1} (\mathbb{E}[V_k] + \mathbb{E}[D_k]) + \mathbb{E}[V_{j+1}] + \mathbb{E}[D_{j+1}]\right)\right) p_j + \lambda_j \mathbb{E}[V_j] - \mathbb{E}[\Lambda_j(V_j)], \\ \mathbb{E}[\theta'_{j+1}] &= \left(n_{j+1} + \lambda_{j+1} \sum_{k=1}^{j-1} (\mathbb{E}[V_k] + \mathbb{E}[D_k])\right) p_{j+1} + \lambda_{j+1} \mathbb{E}[V_{j+1}] - \mathbb{E}[\Lambda_{j+1}(V_{j+1})]. \end{aligned}$$

Thus,

$$\begin{aligned} \mathbb{E}[\theta'] = c + \sum_{i \neq j, j+1} \lambda_i p_i \sum_{k=1}^{i-1} (\mathbb{E}[V_k] + \mathbb{E}[D_k]) + \lambda_j p_j \left(\sum_{k=1}^{j-1} (\mathbb{E}[V_k] + \mathbb{E}[D_k]) + \mathbb{E}[V_{j+1}] + \mathbb{E}[D_{j+1}] \right) + \\ + \lambda_{j+1} p_{j+1} \sum_{k=1}^{j-1} (\mathbb{E}[V_k] + \mathbb{E}[D_k]). \end{aligned}$$

Therefore, we have that $\mathbb{E}[\theta] \leq \mathbb{E}[\theta']$ if and only if

$$\begin{aligned} \lambda_j p_j \sum_{k=1}^{j-1} (\mathbb{E}[V_k] + \mathbb{E}[D_k]) + \lambda_{j+1} p_{j+1} \sum_{k=1}^j (\mathbb{E}[V_k] + \mathbb{E}[D_k]) \leq \\ \lambda_j p_j \sum_{k=1}^{j-1} (\mathbb{E}[V_k] + \mathbb{E}[D_k]) + \lambda_j p_j (\mathbb{E}[V_{j+1}] + \mathbb{E}[D_{j+1}]) + \lambda_{j+1} p_{j+1} \sum_{k=1}^{j-1} (\mathbb{E}[V_k] + \mathbb{E}[D_k]), \end{aligned}$$

or

$$\lambda_{j+1} p_{j+1} (\mathbb{E}[V_j] + \mathbb{E}[D_j]) \leq \lambda_j p_j (\mathbb{E}[V_{j+1}] + \mathbb{E}[D_{j+1}]).$$

In other words, we get that the optimal processing order is by visiting the queues according to an *increasing* order of $\lambda_i p_i / (\mathbb{E}[V_i] + \mathbb{E}[D_i])$. \square

Roughly stated, this rule arranges the visit order according to the ratio between new arrivals per unit time that will successfully complete their service, i.e. $\lambda_i p_i$, and the mean duration of a visit there, i.e. $\mathbb{E}[V_i] + \mathbb{E}[D_i]$. It is intuitively clear that if the mean visit and switch time for a queue is relatively long, then one should visit this queue early on. In this way, the number of customers at the other queues during this cycle will also be relatively high, and as a result the throughput will be increased since all customers are served simultaneously by an infinite number of servers.

This is an extremely simple rule, which can be directly implemented. Moreover, suppose that, for one reason or another, the objective is to minimise the throughput of the system for each cycle. Then, the index rule that determines the order of visits to the queues is simply *reversed*; the servers complete a Hamiltonian tour arranged in a decreasing order of $\lambda_i p_i / (\mathbb{E}[V_i] + \mathbb{E}[D_i])$. Observe that under this strategy the servers also visit the queues that are empty at the beginning of the cycle.

One expects that the throughput of the system in the long-run is improved when these queues are not visited within a cycle; namely, it may be more efficient to avoid queues that are empty at the beginning of the cycle in order to allow them to build up.

According to the way the system is designed, even if the servers do not visit a queue that at the beginning of the cycle was empty, the switch time associated with this queue (i.e. the time to switch from this queue to the following one) is still incurred. Therefore, as the number of queues that will not be visited in a cycle grows, the servers spend an increasing amount of time being essentially idle (as they switch between queues).

A possibly more efficient system design is as follows. Rather than envision the group of servers moving from one queue to another, we can think of a central point to which the servers always return after each completion of a visit to a queue. The return time to the central point after visiting queue i is denoted by R_i . The servers depart from that central point and move to the following queue that will be served. The total time from the moment the servers leave the central point until they enter queue i is denoted by E_i . According to this design, the total time to go from queue i to queue j is given by $R_i + E_j$ for any $i \neq j$. The question that arises

is whether there exists a semi-dynamic control of this system. As before, it emerges that a Hamiltonian-tour approach leads to a *static* processing order according to an index rule.

Theorem 2. *For the polling system with a central point, the Hamiltonian-tour approach leads to a fixed optimal visiting order, which is independent of the number of customers present at the various queues at the start of the cycle. The throughput of the system for each cycle is maximised if and only if the visiting order is arranged according to an increasing sequence of the index-type rule*

$$\frac{\lambda_i p_i}{\mathbb{E}[E_i] + \mathbb{E}[V_i] + \mathbb{E}[R_i]}. \quad (6.2)$$

Proof. As before, let $n = (n_1, n_2, \dots, n_N)$ be the state of the system at the start of the tour and denote by L the number of *non-empty* queues, $0 < L \leq N$, at the beginning of the cycle. The throughput of queue i during a Hamiltonian cycle that visits only the non-empty queues according to the order $\pi_0 = (1, 2, \dots, L)$ is given by

$$\theta_i = \text{Binom}\left(n_i + A_i\left(\sum_{k=1}^{i-1} (E_k + V_k + R_k) + E_i\right), p_i\right) + \left(A_i(V_i) - \text{Poisson}(\Lambda_i(V_i))\right).$$

Consequently,

$$\mathbb{E}[\theta_i] = \left(n_i + \lambda_i\left(\sum_{k=1}^{i-1} (\mathbb{E}[E_k] + \mathbb{E}[V_k] + \mathbb{E}[R_k]) + \mathbb{E}[E_i]\right)\right)p_i + \lambda_i\mathbb{E}[V_i] - \mathbb{E}[\Lambda_i(V_i)],$$

which yields

$$\mathbb{E}[\theta] = c' + \sum_{i=1}^N \lambda_i p_i \sum_{k=1}^{i-1} (\mathbb{E}[E_k] + \mathbb{E}[V_k] + \mathbb{E}[R_k]), \quad (6.3)$$

where

$$c' = \sum_{i=1}^N ((n_i + \lambda_i \mathbb{E}[E_i])p_i + \lambda_i \mathbb{E}[V_i] - \mathbb{E}[\Lambda_i(V_i)]).$$

Applying an interchange argument we have that the optimal processing order is constructed by an increasing sequence of the index rule given by (6.2). \square

As before, the optimal tour does not depend on the number of customers present at the beginning of the cycle. This is a direct consequence of the fixed visit times and the underlying M/G/ ∞ process at each queue.

Index rules appear regularly when optimising polling systems. Browne and Yechiali [6, 7] were the first to obtain dynamic control policies for single-server systems under the exhaustive, gated or mixed service regimes. The mechanics of the system are as described here: at the beginning of each cycle the server decides on a new Hamiltonian tour and visits the channels accordingly. The authors showed that if the objective is to optimise the cycle duration under these policies, then an index-type rule applies, which is similar to the one described here. The main difference is that the index rule that is optimal for these policies depends on the state of the system at the beginning of a cycle, *contrary to the results obtained for the fixed-visit-time policy* studied in this paper. The result derived by Browne and Yechiali [6, 7] is a surprising result as the index rule does *not* include the service times at the various channels. It is also surprising that the same index rule holds for both the gated and the exhaustive disciplines although the duration of a cycle starting from the same state is *different* for the two regimes. For a further discussion on other types of index-rule policies see Yechiali [25], van der Wal and Yechiali [22], and references therein.

Acknowledgements

The authors would like to thank Prof. Onno Boxma for several insightful remarks. The authors also acknowledge the hospitality and support of EURANDOM while carrying out this research. The research of the first author was also supported by the Aristotle University of Thessaloniki (full scholarship from the legacy of L. Athanasoula). Part of the second author's work was carried out while visiting EURANDOM in his capacity as Beta Chair.

References

- [1] E. Altman and H. J. Kushner. Control of polling in presence of vacations in heavy traffic with applications to satellite and mobile radio systems. *SIAM Journal on Control and Optimization*, 41(1):217–252, 2002.
- [2] S. Asmussen. Matrix-analytic models and their analysis. *Scandinavian Journal of Statistics. Theory and Applications*, 27(2):193–226, 2000.
- [3] J. M. Bernabéu-Aubàn, M. H. Ammar, and M. Ahamad. Optimizing a generalized polling protocol for resource finding over a multiple access channel. *Computer Networks and ISDN Systems*, 27(10):1429–1445, 1995.
- [4] S. C. Borst. *Polling systems*, volume 115 of *CWI Tract*. Stichting Mathematisch Centrum Centrum voor Wiskunde en Informatica, Amsterdam, 1996.
- [5] O. J. Boxma, H. Levy, and U. Yechiali. Cyclic reservation schemes for efficient operation of multiple-queue single-server systems. *Annals of Operations Research*, 35(3):187–208, 1992.
- [6] S. Browne and U. Yechiali. Dynamic routing in polling systems. In M. Bonatti, editor, *Teletraffic Science for New Cost-Effective Systems, Networks, and Services*, pages 1455–1466, Torino, 1–8 June 1988. Proceedings of the Twelfth International Teletraffic Congress, Elsevier Science Publications.
- [7] S. Browne and U. Yechiali. Dynamic priority rules for cyclic-type queues. *Advances in Applied Probability*, 21(2):432–450, 1989.
- [8] R. B. Cooper and G. Murray. Queues served in cyclic order. *The Bell System Technical Journal*, 48:675–689, 1969.
- [9] I. Eliazar and U. Yechiali. Polling under the randomly timed gated regime. *Communications in Statistics. Stochastic Models*, 14(1-2):79–93, 1998.
- [10] H. Levy and M. Sidi. Polling systems: Applications, modeling, and optimization. *IEEE Transactions on Communications*, 38(10):1750–1760, 1990.
- [11] C. Mack. The efficiency of N machines uni-directionally patrolled by one operative when walking time is constant and repair times are variable. *Journal of the Royal Statistical Society. Series B. Methodological*, 19:173–178, 1957.
- [12] C. Mack, T. Murphy, and N. L. Webb. The efficiency of N machines uni-directionally patrolled by one operative when walking time and repair times are constants. *Journal of the Royal Statistical Society. Series B. Methodological*, 19:166–172, 1957.
- [13] O. Nakdimon and U. Yechiali. Polling systems with breakdowns and repairs. *European Journal of Operational Research*, 149(3):588–613, 2003.

- [14] S. Stidham, Jr. Optimal control of a signalized intersection. Part I: Introduction. structure of intersection models. Part II: Determining the optimal switching policies; Part III: Descriptive stochastic models. Technical Report 94, 95, and 96, Department of Operations Research, Cornell University, Ithaca, New York, 1969.
- [15] L. Takács. *Introduction to the Theory of Queues*. University Texts in the Mathematical Sciences. Oxford University Press, New York, 1962.
- [16] H. Takagi. *Analysis of polling systems*. Series In Research Reports And Notes: Computer Systems Series. MIT Press, Cambridge, MA, USA, 1986.
- [17] H. Takagi. Queueing analysis of polling models: an update. In *Stochastic analysis of computer and communication systems*, pages 267–318. North-Holland, Amsterdam, 1990.
- [18] H. Takagi. Queueing analysis of polling models: progress in 1990–1994. In *Frontiers in queueing*, Probability and Stochastics Series, pages 119–146. CRC, Boca Raton, FL, 1997.
- [19] H. C. Tijms. *A First Course in Stochastic Models*. John Wiley & Sons, Chichester, 2003.
- [20] D.-C. Twu and K.-C. Chen. A novel MAC protocol for broadband communication over CATV-based MANs. *Computer Communications*, 19(11):888–900, 1996.
- [21] M. van der Heijden, A. van Harten, and M. Ebben. Waiting times at periodically switched one-way traffic lanes: a periodic, two-queue polling system with random setup times. *Probability in the Engineering and Informational Sciences*, 15(4):495–517, 2001.
- [22] J. van der Wal and U. Yechiali. Dynamic visit-order rules for batch-service polling. *Probability in the Engineering and Informational Sciences*, 17(3):351–367, 2003.
- [23] M. Van Vuuren and E. M. M. Winands. Iterative approximation of k -limited polling systems. Technical Report 2006-06, Eindhoven University of Technology, May 2006. Available at <http://www.win.tue.nl/math/bs/spor/>.
- [24] Y. T. Wang and R. J. T. Morris. Load sharing in distributed systems. *IEEE Transactions on Computers*, C-34(3):204–217, 1985.
- [25] U. Yechiali. Analysis and control of polling systems. In L. Donatiello and R. Nelson, editors, *Performance Evaluation of Computer and Communication Systems*, volume 729 of *Lecture notes in computer science*, pages 630–650, Berlin, New York, 1993. Springer-Verlag.