

# Bounding the Mean Response Time of the Minimum Expected Delay Routing Policy: An Algorithmic Approach

John C.S. Lui, *Member, IEEE*, Richard R. Muntz, *Fellow, IEEE*, and Don Towsley, *Member, IEEE*

**Abstract**—Balancing loads in a multi-server system can have a significant impact on performance. In this paper, we model such a system as a heterogeneous multi-server queueing system. We study the behavior of such a system operating under the minimum expected delay (MED) routing policy, i.e., an arriving customer is assigned to the queue which has the minimal expected value of unfinished work. This routing discipline can be viewed as a generalization of the join-the-shortest queue (SQ) discipline for homogeneous servers. There is no closed-form solution for this class of queueing problem. In this paper, we provide a methodology to compute upper and lower bounds on the mean response time of the system. This methodology allows one to tradeoff the tightness of the bounds and computational cost. Applications and numerical examples are presented which show how to use this methodology for deriving performance measures and also illustrating that the excellent accuracy of the computational algorithm which is achievable with modest computational cost.

**Index Terms**—Load balancing, parallel systems, scheduling, queueing models, shortest delay routing.

## I. INTRODUCTION

WITH the advent of multiprocessors and multicomputer systems, there has been considerable interest in the problem of balancing the load among processors or computers. In this paper, we model such systems as heterogeneous multi-server queueing systems. Using such models, we investigate the performance of such systems operating under the minimum expected delay (MED) routing policy. Although not optimal, this policy can provide excellent performance in these systems. Some major difficulties in analyzing this kind of a routing policy, even under Markovian assumptions, are

- 1) each queue in the system is correlated because the arrival process to each server depends on the state of the entire system and,
- 2) since each queue has infinite capacity, the state space of the system is multi-dimensional in nature and is infinite in each of the dimensions.

In its general form, there is no known closed-form solution, and it is impossible to exactly solve the problem numerically due to

the infinite state space. One way to approach this problem is to construct a modified model which provably bounds the performance of the original policy and for which the performance measures of the modified model can be easily computed.

The goal of this paper is to analyze the minimum expected delay (MED) routing algorithm (a natural generalization of the join-the-shortest-queue (SQ) policy for homogeneous servers) in a multiprocessor or multicomputer system, which is modeled as a heterogeneous multi-server queueing system. Let  $K$  be the number of servers, where  $K \geq 2$ . Each server has an infinite capacity queue, and service rates are exponentially distributed with rates  $\mu_i$ ,  $i = 1, 2, \dots, K$ . Without loss of generality, we assume  $\mu_1 \geq \mu_2 \geq \dots \geq \mu_K$ . The customer arrival process is a Poisson process with a mean rate of  $\lambda$ . We propose a methodology which provides upper and lower bounds on the mean number of customers (and thereby the mean response time) in the system and which can be used to trade off the tightness of the bounds with the computational cost. By virtue of providing bounds, rather than simply an approximation, our results are distinguished from previous work on this problem.

We begin with a brief review of the published literature on the join-the-shortest-queue routing problem. The optimality of the SQ policy for homogeneous multi-server systems has been established in numerous papers [10], [27], [29]. In its most general form, SQ has been shown to minimize the queue length vector in the sense of Schur-convex ordering [27]. Of more interest to us is the literature dealing with the performance evaluation of the SQ and MED policies. In [4], MED policy was mentioned in the user search sequences. In the case of the SQ policy for two identical servers, numerous authors have provided exact, though not necessarily computable solutions [14], [11], [31], [7], [1]. Several authors have provided similar solutions to the heterogeneous server problem, for example, [15], [2]. The last paper (also [1]) is interesting because it can generate a sequence of increasingly more accurate approximations with error bounds that decrease exponentially. Recently, Adan, et al. [3] have provided an error bound for a homogeneous server system.

Numerous authors have proposed approximations for the SQ and MED policies. These include Conolly [6], Rao and Posner [24], and Towsley and Chen [26] in the case of the SQ policy. The first of these treats both queues as having bounded capacity whereas the last two treat only one queue as having bounded capacity. The last two papers produce solutions that can be expressed in a matrix-geometric form [23]. The last paper, [26], is also noteworthy in that it provides upper and

Manuscript received July 16, 1993; revised May 22, 1994.

J.C.S. Lui is with the Department of Computer Science, The Chinese University of Hong Kong.

R.R. Muntz is with the Computer Science Department, University of California Los Angeles.

D. Towsley is with the Computer Science Department, University of Massachusetts.

To order reprints of this article, e-mail: transactions@computer.org, and reference IEEECS Log Number C95143.

lower bounds on various performance statistics that are established using less sophisticated sample path techniques than are used in this paper. Grassmann [13] studied the same problem with  $K = 2$  and solved for transient and steady state behavior. Halfin [12] studied the two servers problem and used a linear programming technique to compute bounds on the mean number of customers in the system. Blanc [5] studied the SQ routing policy with an arbitrary number of heterogeneous servers. He proposed an approximation method which was based on power series expansions and a recursion which required a substantial computational effort. Various approximations for computing the mean response time of  $K$  homogeneous servers have been proposed by Lin and Raghavendra [18], Nelson and Philips [21], [22], and Wang and Morris [28]. Zhao and Grassmann [30] studied the shortest queue model with jockeying. This problem has the matrix-geometric form and an explicit solution can be obtained. None of the work cited above treated more than two heterogeneous servers and simultaneously provided error bounds. Lui and Muntz [19] were the first to propose a methodology to bound the mean response time of a minimum expected delay routing system. This paper differs from [19] in several ways. First, we derive improved bounds for the homogeneous servers case, and secondly, we use sample path analysis to prove the bounds, yielding more elegant and intuitive proofs.

This work distinguishes itself from previous published results in that it simultaneously

- 1) allows more than  $K \geq 2$  servers,
- 2) allows heterogeneous servers,
- 3) includes a scheduling policy based on queue lengths and service rates (thus, we treat a generalization of the join-the-shortest queue for homogeneous systems) and
- 4) provides error bounds on the mean number of customers (and thereby mean response time) in the system.

The bounding methodology has the desirable property that it allows one to tradeoff accuracy and computational cost, as will be demonstrated.

The organization of the paper is as follows. In Section II we formally define the queueing model. Sections III and IV present the modified models and prove that they do provide bounds. In Section V, we provide a methodology for obtaining tighter bounds in the special case of homogeneous servers. In Section VI, we present an applications with a numerical example which shows the excellent accuracy of this methodology. Conclusions are given in Section VII.

## II. MODEL

We model a multiprocessor or multicomputer system as a queueing system, as depicted in Fig. 1, with  $K$  heterogeneous servers with associated queues being fed by a Poisson process<sup>1</sup> with mean rate  $\lambda$ . The service times at servers form mutually independent sequences of exponential random variables that are also independent of arrival times with rates  $\mu_1 \geq \mu_2 \geq \dots \geq \mu_K$ .

1. Although in the computational algorithm we assume the arrival process to be Poisson, the proofs for the bounds can accommodate a general arrival process.

Let  $Y$  denote a policy that routes an arriving customer to a server on the basis of the server queue lengths and mean service time and service rate of the servers. Let  $l_Y^*(N)$  denote the identity of the queue to which the customer is routed under policy  $Y$  when the queue length vector is  $N$ . When we are interested in the joint queue length at time  $t$  under a specific policy, we will denote it as  $N^Y(t)$ . We assume that  $p$  is stationary.

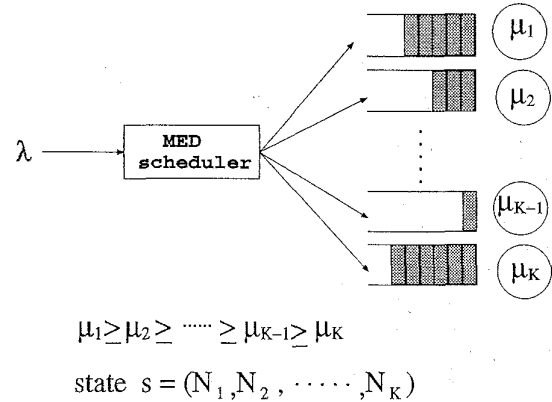


Fig. 1. Minimum expected delay routing policy queueing model.

We define the minimum expected delay (MED) routing policy as follow. Let  $N_i(t)$  be the number of customers at server  $i$  (on that server or in the server's queue) at time  $t$ . We define  $u_i(t) = (1 + N_i(t))/\mu_i$ , which is the mean unfinished work at the  $i$ th server if a customer arrives at time  $t$  and is assigned to the  $i$ th server. Let us define  $u^*(t) = \min\{u_i(t), i = 1, \dots, K\}$ . Upon arrival of a customer at time  $t$ , the customer joins a server  $j$  where  $u_j(t) = u^*(t)$ . If a tie occurs, the customer chooses the server with the lowest index. When all service rates are equal, this MED routing policy reduces to the classic join-the-shortest queue (SQ) routing algorithm.

Assume the system is stable, that is  $\lambda < \sum_{i=1}^K \mu_i$ . Then

$\lim_{t \rightarrow \infty} N_i(t) = N_i$ . We can construct a Markov model,  $M$ , for this queueing system with state space:

$$\{(N_1, N_2, \dots, N_K) \mid N_i \geq 0, i = 1, \dots, K\}$$

The unique steady state probability vector for this continuous-time Markov model satisfies the following system of linear equations:

$$\bar{\pi}G = \bar{0} \quad \text{and} \quad \bar{\pi}e = 1 \quad (1)$$

where  $\bar{\pi}$  is the  $K$ -dimensional steady state probability vector,  $e$  denotes an appropriately dimensioned column vector of 1s, and  $G$  is the transition rate matrix having the following structure:

$$\begin{aligned} (N_1, \dots, N_i, \dots, N_K) &\rightarrow (N_1, \dots, N_i + 1, \dots, N_K) \quad 1\{i = \min\{k \mid u_k = u^*\}\} \lambda \\ (N_1, \dots, N_i, \dots, N_K) &\rightarrow (N_1, \dots, N_i - 1, \dots, N_K) \quad 1\{N_i > 0\} \mu_i \end{aligned}$$

and the balance equation can be expressed as:

$$\begin{aligned} \lambda + \sum_{i=1}^K 1\{N_i > 0\} \mu_i \pi(N_1, \dots, N_K) &= \sum_{i=1}^K \pi(N_1, \dots, N_i + 1, \dots, N_K) \\ + \lambda \sum_{i=1}^K 1\{i = l_{MED}^*((N_1, \dots, N_i - 1, \dots, N_K))\} &\pi(N_1, \dots, N_i - 1, \dots, N_K) \end{aligned}$$

The above model does not possess a known closed form solution, and it is not possible to solve the problem numerically due to its infinite state space cardinality. Since the Markov process lacks the appropriate special structure, techniques such as matrix-geometric methods do not apply. One natural way to approach this problem is to *construct* other models that closely bound the performance of the original problem and which, at the same time, have either known closed form solutions or at least can be efficiently evaluated by numerical methods.

It is intuitive that the stationary state probabilities for the model  $M$  are highly *skewed* or, in other words, the probability mass of the system is concentrated in some relatively small subset of the state space rather than distributed nearly uniformly over the entire state space. For example, consider a system of four homogeneous servers. The purpose of using the routing policy discussed above is to balance the load of the system as much as possible; therefore it is reasonable to assume that a highly unbalanced state (e.g., (8, 4, 3, 1)) has a much smaller probability mass than a balanced state (e.g., (4, 4, 4, 4)). This crucial insight provides the rationale for constructing two modified versions of the original model which can be shown to bound the mean response time of the original system. In both cases we represent the exact behavior (transition rates) for the most "popular" states (where most of the probability mass resides). The number of states in the most popular subset is a function of the accuracy demanded and the computational cost one is willing to pay. When the system leaves this subset we modify the behavior of the system in such a way that

- 1) the modified system has an efficient solution and
- 2) the modified model's behavior can be shown to bound the behavior of the original model from above or from below.

Therefore, one modified model provides an upper bound on the mean response time while another provides a lower bound on the mean response time. In the next section, we discuss the upper bound model and then, in the following section, we cover the lower bound model.

### III. UPPER BOUND

In this section, we present a modified model,  $M^u$ , which provides an upper bound for the mean response time and the mean number of customers in the system for the original model,  $M$ . The upper bound model has the same system configuration, namely that the customer arrival process is a Poisson process, and  $K$  servers with service rates  $\mu_i$ ,  $i = 1, 2, \dots, K$ , where  $\mu_1 \geq \mu_2 \geq \dots \geq \mu_K$ .<sup>2</sup> The upper bound model  $M^u$  has two additional parameters. The first parameter we term the *artificial capacity vector*  $\bar{C} = (C_1, \dots, C_K)$ . The second parameter is a threshold setting  $d$ , which is the maximum allowable difference between the longest queue and the shortest queue in

2. We require that the service rates be rational numbers such that they can be expressed as integers after normalization, i.e., the service rates are mutually commensurable.

$M^u$ . We first give the formal definition of these two parameters and, in the following paragraph, we give the intuitive idea of how these two parameters can be used to construct the upper bound model  $M^u$ .

**DEFINITION 1.** Let  $\bar{C}^* = (C_1^*, \dots, C_K^*)$  be a vector where  $C_i^* = \frac{\mu_i}{\mu_1} C$  and  $C$  is chosen to be the minimum positive integer such that  $C_i^*$  is a positive integer for  $i = 1, \dots, K$ . Then the artificial capacity vector  $\bar{C}$  is an integer multiple of  $\bar{C}^*$ , i.e.,  $\bar{C} = j\bar{C}^*$  for some  $j \geq 1$ .

For example, if  $\mu_1 = \frac{3}{2}$  and  $\mu_2 = 1$ , then  $C = 3$  and  $\bar{C}^* = (3, 2)$ . So the artificial capacity vector can be  $\bar{C} = j\bar{C}^*$  for any  $j \geq 1$ .

**DEFINITION 2.** Let a state of the model be  $s = (N_1, \dots, N_K)$ . Let  $N_i^*$  be the number of active customers<sup>3</sup> in the  $i$ th server. Define  $q(s)$  to be the degree of imbalance for state  $s$ , as:

$$q(s) = \max\{N_i^* - N_j^* \mid \text{where } i, j \in \{1, \dots, K\}\}$$

**DEFINITION 3.** Let  $d$  be the threshold setting in the modified model. We require that  $q(s) \leq d$  for each state  $s$  in model  $M^u$ .

We first give an intuitive idea of the construction of model  $M^u$ . In  $M^u$ , the degree of imbalance is required to be less than or equal to the parameter  $d$ . A customer may depart from the system only if its departure does not violate the maximum degree of imbalance permitted. If the customer departure would violate the threshold setting, the customer restarts its service within the same server. Intuitively, this mechanism forces a customer to stay in the system at least as long as in the original model and thereby increases the number of customers in the system. Note that, due to the routing policy, an arrival never causes the degree of imbalance to exceed  $d$ . The rationale behind the threshold parameter is to generate a model with a state space which is a subset of the state space of the original model.

The second parameter is the *artificial capacity*,  $C_i$ ,  $i = 1, 2, \dots, K$  for each server. In model  $M^u$ , there are two classes of customers, active customers and suspended customers. At any point in time, there are never more than  $C_i$  active customers in queue  $i$ ; all of the remaining customers are suspended. Whenever a customer arrives to the system and finds that each server  $i$ ,  $i = 1, \dots, K$ , has *exactly* an integer multiple of  $C_i$  customers, all active customers in the system (except for the arriving customer) are put into a *suspended mode* and a new "busy cycle" is started. This busy cycle ends when all servers complete all active customers.  $C_i$  suspended customers are then released from queue  $i = 1, \dots, K$  and can be served. Note that the definition here is recursive; during the busy period following suspension of a set of customers, the capacities  $C_i$  can *again* be exceeded, causing *another* set of customers to be suspended. When a busy cycle ends, only the set of customers suspended at the initiation of that busy cycle is released for service. The purpose of the  $C_i$ ,  $1 \leq i \leq K$ , is to create a matrix with a *repeti-*

3. Definition of active customers will be defined in a later paragraph.

itive structure; based on that structure, we will be able to derive an efficient numerical solution algorithm. The computation algorithm is based on partitioning the state space of  $M^u$  into  $\bigcup_{i=0}^{\infty} S_i \dots$  where:

$$S_0 = \{(N_1, \dots, N_K) \mid 0 \leq N_j \leq C_j \text{ for } j = 1, \dots, K\}$$

$$S_i = \{(N_1, \dots, N_K) \mid iC_j \leq N_j \leq (i+1)C_j \text{ for } j = 1, \dots, K\}$$

$$- \{(iC_1, \dots, iC_K)\} \quad i \geq 1$$

Due to the routing of arrivals and the constraint on departures, we can show that all transitions from  $S_i$  to  $S_{i+1}$  occur through one state in  $S_i$  and the transitions from  $S_{i+1}$  to  $S_i$  can only go to one state in  $S_i$ . As will be shown later, this property allows us to efficiently solve the model via *exact* decomposition based on the partition  $\{S_0 \cup S_1 \cup \dots\}$ . Intuitively, this second modification to the model should also increase the mean number of customers in the system compared to the original model since additional server idle time is introduced and service of a suspended customers can only be resumed when all active customers depart from the system.

As an example, assume that we have a system with four homogeneous servers, and we let  $C_i = 10$ , for  $i = 1, 2, 3, 4$ . It is easy to see that  $S_0$  includes all states for which each queue has between 0 to 10 customers;  $S_1$  consists of all states for which each queue has 10 suspended customers, and has between 0 to 10 active customers and at least one queue has an active customer. Observe that the only transition from  $S_0$  to  $S_1$  is from state (10, 10, 10, 10). This is due to the shortest expected delay routing of arrivals. The only non-zero transitions from  $S_1$  to  $S_0$  are from states (11, 10, 10, 10), (10, 11, 10, 10), (10, 10, 11, 10) and (10, 10, 10, 11) to state (10, 10, 10, 10). This is due to the rule introduced in  $M^u$  to the effect that suspended customers are only served when the busy period (corresponding to states in  $S_1$ ) has completed. An important point is that the parameters  $d$  and  $C_i$ , for  $i = 1, \dots, K$ , can be chosen to control the extent to which  $M^u$  behaves like the original model  $M$ , i.e., the larger  $d$  and the  $C_i$ s are, the larger the portion of the state space that has behavior identical to the original model.

### A. Proof of Upper Bound

In this section, we prove that the model  $M^u$  provides an upper bound on the number of customers in the system at any point in time. In the case that the model exhibits stationary behavior, Little's result can be invoked to show that  $M^u$  provides an upper bound for the mean response time. We therefore concentrate on the mean number in the system in the remainder of this section. It is important to point out that the proofs can accommodate a general arrival process. We start by defining an auxiliary concept that will be useful in the proof.

**DEFINITION 4.** A policy  $Y$  is a proper policy if  $N \leq N'$  (here " $\leq$ " is taken to mean componentwise) implies that  $N + \underline{e}_{k(N)}^T \leq N' + \underline{e}_{k(N')}^T$  where  $\underline{e}_k$  is the column vector of all 0s except for a 1 in position  $k$ .

It is easy to see that the minimum expected delay routing policy is a proper routing policy.

In establishing an upper bound, it is useful to look at the times when events such as arrivals and departures occur. In the latter case, it is useful to think of each server as continuously serving customers. If the queue is empty, then the server serves a *fictitious* customer. Hence *service events* at server  $k$  occur as a Poisson process with parameter  $\mu_k$ . (Note that a service event is a departure event only when there is a customer in the queue.) Furthermore, if a customer is routed to an empty queue, then it is assigned the remaining service time of the fictitious customer on the server. The exponential assumption guarantees that the time to the next service event is an exponential random variable with the same parameter. It follows that, under this interpretation, the service times are still i.i.d. exponential with the same mean.

Consider the  $i$ th event. Let  $N_i = (N_{i,1}, \dots, N_{i,K})$  be the joint queue lengths immediately after the  $i$ th event. Let  $N_0$  denote the initial queue lengths. We have the following evolution equations. If the  $(i+1)$ st event corresponds to an arrival,

$$N_{i+1,k} = N_{i,k} + 1 \{I_p^*(N_i) = k\}, \quad 1 \leq k \leq K \quad (2)$$

If the  $(i+1)$ st event corresponds to a service event at server  $j$ ,

$$N_{i+1,k} = \begin{cases} n_{i,k}, & k \neq j, \\ (N_{i,j} - 1)^+, & k = j. \end{cases} \quad (3)$$

Now suppose that we have a modified system for which we define a new binary valued random variable  $Y_i$  that takes on the value 0 if no customer is allowed to depart and the value 1 if a customer is allowed to depart at the  $i$ th event (provided that it is a service event). In the original model  $M$ , the random variable  $Y_i$  is always equal to 1. On the other hand, in the upper bound model  $M^u$  presented above,  $Y_i$  can be 0 or 1 depending on the model state. Let  $N^u(t)$  be the joint queue lengths for the model  $M^u$ . We have the following evolution equations at the time of arrival and service events. If the  $(i+1)$ st event corresponds to an arrival,

$$N_{i+1,k}^u = N_{i,k}^u + 1 \{I_p^*(N_i^u) = k\}, \quad 1 \leq k \leq K \quad (4)$$

If it is a service event at server  $j$ ,

$$N_{i+1,k}^u = \begin{cases} N_{i,k}^u, & k \neq j, \\ (N_{i,j}^u - Y_{i+1})^+, & k = j. \end{cases} \quad (5)$$

**LEMMA 1.** If  $N(0) \leq_{st} N^u(0)$  and  $p$  is a proper routing policy, then  $N(t) \leq_{st} N^u(t)$  for  $t \geq 0$ .

**PROOF.** Couple the initial queue lengths so that  $N(0) \leq N^u(0)$ . Condition on the initial queue lengths, arrival times, and service event times. The proof is by induction on the event times to establish the deterministic relation  $N_i \leq N_i^u$  for  $i \geq 0$ .

For  $i = 0$ ,  $N(0) \leq N^u(0)$ . For the induction step, assume  $N_i \leq N_i^u$  holds for  $i = k$ . For  $i = k+1$ , if the  $i$ th event is an arrival event, then by the definition of a proper policy the relationship holds. If the  $i$ th event is a service event, then due to  $Y_{k+1} \leq 1$ , the relationship holds. Therefore, the upper bound model  $M^u$  satisfies the assumptions described above,

and we have  $N(t) \leq N^u(t)$ . By removing the conditions on initial queue lengths, arrival times, and service event times, we have  $N(t) \leq_{st} N^u(t)$  for  $t \geq 0$ .  $\square$

Let  $N_i = \lim_{t \rightarrow \infty} N_i(t)$  when it exists,  $1 \leq i \leq K$  and  $N = \sum_{i=1}^K N_i$ . Based on this lemma, we have  $E[N] \leq E[N_u]$ . If  $R$  and  $R_u$  denote the stationary customer response times, when they exist, then by Little's result, we have  $E[R] \leq E[R_u]$ .

**B. Computational Algorithm for Solving the Model  $M^u$**

In this section, we provide an algorithm for computing the mean response time of the upper bound model when the arrival process is Poisson with mean rate  $\lambda$ . We partition the state space of  $M^u$ ,  $S^u = \bigcup_{i=0}^{\infty} S_i$  and  $S_i \cap S_j = \emptyset, \forall i \neq j$ , where:

$$S_0 = \{(N_1, \dots, N_K) \mid 0 \leq N_j \leq C_j \text{ for } j = 1, \dots, K\}$$

$$S_i = \{(N_1, \dots, N_K) \mid iC_j \leq N_j \leq (i+1)C_j \text{ for } j = 1, \dots, K\}$$

$$- \{(iC_1, \dots, iC_K)\}$$

$Q_{S_i, S_j}$  = transition rate matrix from states in  $S_i$  to states in  $S_j$ .

The transition rate matrix  $Q^u$  has the form depicted in Fig. 2 when the states are ordered in the natural way.

$$Q^u = \begin{bmatrix} Q_{S_0, S_0} & Q_{S_0, S_1} & 0 & 0 & 0 & \dots \\ Q_{S_1, S_0} & Q_{S_1, S_1} & Q_{S_1, S_2} & 0 & 0 & \dots \\ 0 & Q_{S_2, S_1} & Q_{S_2, S_2} & Q_{S_2, S_3} & 0 & \dots \\ 0 & 0 & Q_{S_3, S_2} & Q_{S_3, S_3} & Q_{S_3, S_4} & \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix}$$

Fig. 2. Transition rate matrix for the upper bound model.

This is a block tridiagonal transition rate matrix and therefore represents a quasi-birth-death process. By aggregating each partition  $S_i$ , a birth-death process is formed. First, we show how to obtain the exact conditional state probability vector, given that the system is in partition  $S_i$ . Once we have this information, it follows easily that we can obtain the exact aggregate transition rates. We can then obtain the exact stationary state probabilities for the aggregate model. The aggregate state probabilities and the conditional state probabilities together are a complete solution for the stationary state probabilities for the upper bound model  $M^u$ .

There are several important features of the upper bound model,  $M^u$ . First, there is only a single state in  $S_i$  that has a non-zero transition rate into any state in  $S_{i+1}, i \geq 0$ . Let us call this state  $s_i(C_0)$ . State  $s_i(C_0)$  is:

$$s_i(C_0) = (N_1, N_2, \dots, N_K) \in S_i \text{ where}$$

$$N_j = (i+1)C_j, \forall j = 1, 2, \dots, K$$

This follows from the rule used to assign an arriving customer to a server. Also, there are  $K$  states from  $S_i$  that have nonzero transition rates to a state in  $S_{i-1}$  where  $i \geq 1$ . Each corresponds to a state in which a server is the last to complete its "active" (nonsuspended) customer. Let us call these states  $s_i(l), 1 \leq l \leq K, i \geq 1$ . These states are:

$$s_i(l) = (N_1, N_2, \dots, N_K) \in S_i \quad l = 1, \dots, K$$

where

$$N_l = iC_l + 1 \text{ and } N_j = iC_j \text{ for } j \neq l \text{ and } j = 1, 2, \dots, K$$

This follows from the restrictions on departures in the upper bound model. The following are easily seen to be the transition rates between  $s_i(C_0)$  and  $s_{i+1}(l), l = 1, 2, \dots, K$ :

$$s_i(C_0) \rightarrow s_{i+1}(l) \left( t_p^*(s_i(C_0)) \right) \quad \lambda$$

$$s_{i+1}(l) \rightarrow s_i(C_0) \quad \mu_l \text{ for } l = 1, 2, \dots, K$$

The second important observation is that the submatrices  $Q_{S_i, S_i}$ , for  $i \geq 1$ , are all identical. The conditional state probabilities  $P\{s \in S_i \mid S_i\}$  can now be computed exactly using the following lemma from [9]:

LEMMA 2. Given an irreducible Markov process with state space  $S = A \cup B$  and transition rate matrix:

$$\begin{bmatrix} Q_{A,A} & Q_{A,B} \\ Q_{B,A} & Q_{B,B} \end{bmatrix}$$

where  $Q_{i,j}$  is the transition rate submatrix from partition  $i$  to partition  $j$ . If  $Q_{B,A}$  has all zero entries except for some non-zero entries in the  $i$ th column, the conditional steady state probability vector, given that the system is in partition  $A$ , is the solution to the following system of linear equations:

$$\bar{\pi}_{1A} [Q_{A,A} + Q_{A,B} \underline{e} \underline{e}^T] = \bar{0} \quad ; \quad \bar{\pi}_{1A} \underline{e} = 1$$

where  $\underline{e}^T$  is a row vector with a 0 in each component, except for the  $i$ th component which has value 1.

We are now in a position to compute the conditional state probabilities for each partition  $S_i$  of  $M^u$  exactly. Without loss of generality, let us consider  $S_i$ , for some  $i \geq 1$ .

THEOREM 1. Let  $\tilde{Q}_{S_i, S_i}$  be the transition rate matrix which is equal to  $Q_{S_i, S_i}$  except for the following modifications:

$$\tilde{q}_{s_i(C_0), s_i(C_0)} = q_{s_i(C_0), s_i(C_0)} + \lambda \tag{6}$$

$$\tilde{q}_{s_i(l), s_i(1)} = q_{s_i(l), s_i(1)} + \mu_l \quad \text{where } 1 \leq l \leq K \tag{7}$$

The solution to the following system of linear equations:

$$\tilde{\pi} \tilde{Q}_{S_i, S_i} = \bar{0} \quad \text{and} \quad \tilde{\pi} \underline{e} = 1$$

is the conditional steady state probability vector for states in  $S_i$ , that is:

$$\tilde{\pi}(s) = \frac{\pi(s)}{\sum_{s \in S_i} \pi(s)} \quad \forall s \in S_i$$

PROOF. Let us partition the state space  $S^u = \{S'_i \cup S''_i\}$  where  $S'_i = \bigcup_{j=0}^{i-1} S_j$  and  $S''_i = \{S^u - S'_i\}$ . There is only a single return state in  $S'_i$ , which is  $s_i(C_0)$ , from the states in  $S''_i$ . Based on Lemma 2, the modification of (6) provides the conditional steady state probability, given the system is in  $S'_i$ . Now partition the state space  $S'_i = \{S'_i \cup S_i\}$  where

$S_i^1 = \bigcup_{j=0}^{i-1} S_j$ . Based on (7) and the definition of the MED routing policy, we obtain the conditional state probability vector, given the system is in state  $S_i$ .  $\square$

Since we can compute the conditional state probabilities for each partition  $S_i$  exactly, we can exactly aggregate each  $S_i$  into a single state  $s_i$ ,  $i \geq 0$ . The aggregated process is depicted in Fig. 3 where,  $\lambda_0$ ,  $\lambda_{agg}$ , and  $\mu_{agg}$  are:

$$\begin{aligned}\lambda_0 &= \tilde{\pi}(s_0(C_0))\lambda \\ \lambda_{agg} &= \tilde{\pi}(s_i(C_0))\lambda \\ \mu_{agg} &= \sum_{l=1}^K \tilde{\pi}(s_i(l))\mu_l\end{aligned}$$

Solving this chain, we have:

$$\pi^*(s_0) = \left[ 1 + \frac{\lambda_0}{\mu_{agg} - \lambda_{agg}} \right]^{-1} \quad (8)$$

$$\pi^*(s_i) = \left[ 1 + \frac{\lambda_0}{\mu_{agg} - \lambda_{agg}} \right]^{-1} \left( \frac{\lambda_0}{\mu_{agg}} \right) \left( \frac{\lambda_{agg}}{\mu_{agg}} \right)^{i-1} \quad \text{for } i = 1, 2, \dots \quad (9)$$

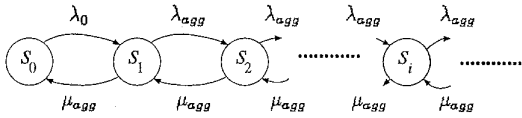


Fig. 3. Aggregate process for the upper bound model.

To obtain the mean number of customers,  $N_u$ , in the upper bound model, let us define the following:

$$\begin{aligned}r(s) &= \sum_{i=1}^K N_i \quad \text{for state } s \in S^u \\ C_0 &= \sum_{i=1}^K C_i \\ \tilde{r}(s) &= r(s) - iC_0 \quad s \in S_i \\ \tilde{N}(s_i) &= \sum_{s \in S_i} \tilde{r}(s)\tilde{\pi}(s)\end{aligned}$$

where  $\tilde{\pi}(s)$  is the solution of the following Markov chain:

$$\tilde{\pi} \tilde{Q}_{S_i, S_i} = \tilde{0} \quad \text{and} \quad \tilde{\pi} \underline{e} = 1$$

Then we have:

$$N_u = \tilde{N}(s_0)\pi^*(s_0) + \sum_{i=1}^{\infty} [\tilde{N}(s_i) + iC_0]\pi^*(s_i) \quad (10)$$

Since  $\tilde{N}(s_i) = \tilde{N}(s_j)$  for  $i \neq j$  where  $i, j \geq 1$ , we can simplify the expression above for  $N_u$  to:

$$\begin{aligned}N_u &= \tilde{N}(s_0)\pi^*(s_0) + \tilde{N}(s_i)(1 - \pi^*(s_0)) \\ &\quad + C_0\lambda_0 \frac{\mu_{agg}}{(\mu_{agg} - \lambda_{agg})^2} \pi^*(s_0)\end{aligned} \quad (11)$$

From Little's result [17], the upper bound mean system response time  $R_u$  is:

$$R_u = \frac{1}{\lambda} \left[ \tilde{N}(s_i)\pi^*(s_0) + \tilde{N}(s_i)(1 - \pi^*(s_0)) + C_0\lambda_0 \frac{\mu_{agg}}{(\mu_{agg} - \lambda_{agg})^2} \pi^*(s_0) \right] \quad (12)$$

It is important to note that the upper bound model  $M^u$  has a different stability condition compared to the original model  $M$ . The original model is stable if:

$$\rho = \frac{\lambda}{\sum_{i=1}^K \mu_i} < 1$$

but the stability condition of the upper bound model is:

$$\rho^u = \frac{\lambda_{agg}}{\mu_{agg}} < 1$$

In general,  $\rho^u > \rho$  but as we increase  $d$  and  $\bar{C}$ , we have  $\rho^u \rightarrow \rho$  from below.

Lastly, to comment about the computational complexity of the upper bound model  $M^u$ . First, we have to obtain the conditional state probabilities for rate matrices  $S_0$  and  $S_1$ . This can be accomplished by using numerical methods, such as the power iteration method, as suggested in [25]. Although the theoretical complexity is  $O(n^3)$  where  $n$  is the dimension of the rate matrix. In practice, the number of operations is much less. After we obtain the conditional state probabilities, we can use (10), (11), and (12) to obtain the expected response time of the upper bound model.

#### IV. LOWER BOUND

In this section we present a model  $M^l$ , which provides a lower bound on the mean response time of the original model. As before, the arrival process is Poisson, the service times are exponentially distributed and that there are  $K$  servers with service rates  $\mu_i$ ,  $i = 1, 2, \dots, K$ , where  $\mu_1 \geq \mu_2 \geq \dots \geq \mu_K$ . For the lower bound model, in addition to the two parameters introduced for the upper bound model  $M^u$ , we define  $C_f = \sum_{i=1}^K C_i$ .

We first give an intuitive idea of the construction of the lower bound model  $M^l$ . The modified system alternates between two phases. The *normal service phase* begins when the system is empty and continues until either the maximum degree of imbalance  $d$  is exceeded or until the total number of customers exceeds  $C_f$ . Once either event occurs, the system enters a *full service phase* where it behaves as a heterogeneous M/M/K system in which, if there are  $j$  customers, where  $j \leq K$ , these  $j$  customers are executed on the  $j$  fastest servers (i.e., customers are moved to the faster servers instantaneously). The system operates in this mode until the system becomes idle. Once the system empties, it returns to the normal service mode. Intuitively, these modifications yield a lower bound on the mean response time since the modifications are an idealization in which either the model behaves exactly as the original model or the best possible service rate is delivered. While

the result is intuitive, we will also formally prove that the modified model  $M^l$  yields a lower bound on the mean response time. Of course, it is intended that  $d$  and  $C_i$ ,  $i = 1, 2, \dots, K$ , be chosen large enough so that most of the time  $M^l$  behaves like the original model. On the other hand, to be able to solve the model efficiently, we would like to keep these parameters small. Numerical examples are given later to illustrate the tradeoffs between the size of a model solved and the spread in the bounds obtained.

#### A. Proof that $M^l$ Provides a Lower Bound for $M$

In this section, we prove that  $M^l$  provides a lower bound on the mean response time of the system for all proper routing policies  $p$ . Again, we point out that the proof can accommodate a general arrival process. The proof is based on the following two lemmas. The first is straightforward and requires little explanation. It will be used to establish the bound during the normal service phase.

**LEMMA 3.** *If  $N(0) \leq_{st} N'(0)$  and  $p$  is a proper routing policy, then*

$$N(t) \leq_{st} N'(t), \quad 0 \leq t.$$

**PROOF.** Without loss of generality, we can couple the systems so that  $N(0) \leq N'(0)$ . Condition on the arrival times and on the service event times at the different servers during the time interval  $[0, t]$ . A simple induction argument using that  $p$  is a proper policy suffices to establish that  $N(t) \leq N'(t)$ . Removal of the conditioning yields the desired result.  $\square$

Consider the system operating solely in the full service mode of operation and let  $N^f(t)$  denote the *total number of customers* in the system. Let  $N^p(t) = \sum_{i=1}^K N_i^p(t)$  denote the total number of customers in the original system under policy  $p$  (henceforth referred to as the normal service system).

**LEMMA 4.** *If  $N^f(0) \leq_{st} N^p(0)$  and  $p$  is proper routing policy, then*

$$N^f(t) \leq_{st} N^p(t), \quad 0 \leq t.$$

**PROOF.** As before, we couple the initial queue lengths so that  $N^f(0) \leq N^p(0)$  and condition on the arrival and departure times. Let  $\{t_n\}$  be a sequence of times where each  $t_i$  corresponds to an arrival or service event. Let  $M^p(t_n)$  denote the number of busy servers at time  $t_n$  in the system under policy  $p$ . Define  $\{1, \dots, K\} \rightarrow \{1, \dots, K\}$  to be a mapping such that  $\gamma_n^p(k)$  is the index of the  $k$ th fastest busy server in the system, provided  $k \leq M^p(t_n)$ . In the case that  $K \geq k > M^p(t_n)$ ,  $\gamma_n^p(k)$  is the index of the  $(k - M^p(t_n))$ th fastest idle server. (Actually, the idle servers can be mapped in an arbitrary manner.) We introduce the following sequences of random variables (r.v.),

- $\{A_n\}$  is a sequence of r.v. such that  $A_n = 1$  if the  $n$ th event is an arrival and 0 if it is a service event.
- $\{I_n\}$  is an independent and identically distributed sequence of r.v. taking values from  $\{1, \dots, K\}$  such that  $\Pr\{I_n = k\} = 1/K$ ,  $k = 1, 2, \dots, K$ , and 0 otherwise.

- $\{B_n\}$  is an i.i.d. sequence of uniformly distributed r.v. in the interval  $[0, 1]$ .

The evolution of the two systems is described as follows. Let  $N_n^p$  denote the joint queue lengths under  $p$  immediately after the  $n$ th event and let  $N_n^f$  denote the total number of customers in the full service system immediately after the  $n$ th event. Let  $N_{n,k}^p$  be the  $k$ th component of  $N_n^p$ . We have:

$$\begin{aligned} N_{n,k}^p = & \\ & (N_{n-1,k}^p - 1 \{ (A_n = 0) \wedge (I_n = l) \wedge (\gamma_{n-1}^p(l) = k) \wedge (B_n < \mu_k / \mu) \})^+ \\ & + A_n 1 \{ I_n^*(N_{n-1}^p) = k \} \end{aligned} \quad (13)$$

$$N_n^f = (N_{n-1}^f - 1 \{ (A_n = 0) \wedge (I_n = l) \wedge (B_n < \mu_l / \mu) \})^+ + A_n \quad (14)$$

It remains to establish that  $N_n^f$  is less than  $N_n^p$  (the total number of customers under policy  $p$ ) immediately after the  $n$ th event for  $n = 1, 2, \dots$ . This is easily done by induction.

**Basis step.** For  $t_0 = 0$ , the result follows from the coupling of the initial queue lengths.

**Inductive step.** Assume that the hypothesis holds for the first  $n-1$  events. We must distinguish between arrivals and service events. If an arrival occurs at time  $t_n$  ( $A_n = 1$ ), then the result follows immediately from the above evolution equations. In the case of a service event, we distinguish between four cases depending on whether  $I_n$  corresponds to a busy or idle server in each system.

**Case (1).** In both systems the server in the chosen position is idle. Then there is no departure from either system and the full service system model continues to have a lower total number of customers, i.e.,  $N_n^f = N_{n-1}^f \leq N_{n-1}^p = N_n^p$ .

**Case (2).** In the normal service system, the chosen position corresponds to a busy server, but in the full service model it corresponds to an idle server. In the normal service system there can be customers waiting in queues while some servers are idle. This does not occur with the full service system. It follows that the total number of customers in the full service system is strictly less than the total number of customers in the normal service model in the interval  $t_{n-1} \leq t < t_n$ , i.e.,  $N_{n-1}^f < N_{n-1}^p$ . Hence,  $N_n^f \leq N_n^p$  since the normal service system only "catches up" by 1.

**Case (3).** The server is busy in the full service system, but it is not busy in the normal service system. Clearly

$$N_n^f \leq N_{n-1}^f \leq N_{n-1}^p = N_n^p.$$

**Case (4).** The servers are busy in both systems. In this case let  $j$  be the label of the server in the full service system and let  $k$  be the index of the server in the normal service system. Since, in the full service system, the fastest servers are always being utilized it follows that  $j \leq k$ , i.e., the chosen server in the full service system is at least as fast as the chosen

server in the normal mode system. Therefore, if  $B_n \leq \mu_d/\mu$ , then  $B_n \leq \mu_j/\mu$ . Hence it follows from the evolution equations, if there is a departure from the normal service system, then there is also a departure from the full service system. So we conclude that  $N_n^f \leq N_n^p$ .

This completes the inductive step. Removal of the conditioning on the initial queue lengths, the arrival times, and the service events completes the proof.  $\square$

Lastly, let  $N^l(t)$  denote the total number of customers in the lower bound system at time  $t$ . We have the following result.

**THEOREM 2.** *If  $N^l(0) \leq_{st} N^p(0)$ , then  $N^l(t) \leq_{st} N^p(t)$  for  $t \geq 0$ , for any proper routing policy  $p$ .*

**PROOF.** This follows directly from the above two lemmas by noting that  $M^l$  goes through alternating intervals in which it operates in normal mode and full service mode. When the transition is made from the full service phase to the normal phase,  $N^l(t) = 0$  which implies that  $N^l(t) \leq N^p(t)$  and so the first lemma can be applied during each normal service mode interval. Similarly, when there is a transition from the normal service phase to the full service phase,  $N^l(t) \leq N^p(t)$  which implies that  $N^l(t) \leq N^p(t)$  and so the second lemma is applicable during every full service mode interval.  $\square$

It is important to note that the stability conditions for the lower bound model  $M^l$  and the original model  $M$  are the same.

**B. Computational Algorithm for Solving the Model  $M^l$**

In this section, we describe an algorithm for computing the mean response time of the lower bound model  $M^l$ . Let us define the following notation:

$S_0$  = set of states with  $0 \leq N_j \leq C_j, j = 1, 2, \dots, K$  such that the threshold  $d$  is satisfied.

$G_1 = \{S_0 - (0, 0, \dots, 0)\}$ .

$a_i$  = a state not in the set  $S_0$ , in which the system contains exactly  $i$  customers.

$Q_{G_1, a_i}$  = transition rate matrix between  $G_1$  and state  $a_i$ .

$g_{a_i, a_j}^*$  = transition rate from state  $a_i$  to state  $a_j$ .

The transition rate matrix of the model  $M^l$  is depicted in Fig. 4. (Note that some of the  $Q_{G_1, a_i} = 0$  but this will not effect the development that follows.)

$Q_{a_0, a_0}$	$Q_{a_0, a_1}$	0	0	0	0	...
$Q_{a_1, a_0}$	$Q_{a_1, a_1}$	$Q_{a_1, a_2}$	$Q_{a_1, a_3}$	$Q_{a_1, a_4}$	...	
$g_{a_1, a_0}^*$	0	$g_{a_1, a_1}^*$	$\lambda$	0	0	...
0	0	$g_{a_2, a_1}^*$	$g_{a_2, a_2}^*$	$\lambda$	0	...
0	0	0	$g_{a_3, a_2}^*$	$g_{a_3, a_3}^*$	$\lambda$	...
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$

Fig. 4. Transition rate matrix for lower bound model.

Since  $S_0$  represents all possible states during the normal mode and states  $a_i, i \geq 1$ , represent all possible states during the full service mode, it is easy to see that the transition rate  $g_{a_i, a_{i-1}}^*$  is:

$$g_{a_i, a_{i-1}}^* = \begin{cases} \sum_{j=1}^i \mu_j & 1 \leq i \leq K \\ \sum_{j=1}^K \mu_j & i > K \end{cases} \quad (15)$$

Observe that if we know the conditional state probabilities for the states in  $S_0$  (where  $S_0 = \{a_0 \cup G_1\}$ ), then we can aggregate  $S_0$  as a single state,  $s_0$ , and we will have a simple aggregated process from which the mean number of customers in the system can be easily derived. Note that there is only a single entry to  $S_0$  from all states outside  $S_0$  because the system must be idle to switch from full service mode to the normal mode. Based on Lemma 2, the state probabilities conditioned on the system being in  $S_0$  can be obtained by solving the following system of linear equations:

$$\begin{aligned} \tilde{\pi}(S_0) \left[ Q_{S_0, S_0} + \sum_{i=1}^{C_f+1} Q_{S_0, a_i} e_i^T \right] &= \vec{0} \\ \tilde{\pi}(S_0) e &= 1 \end{aligned}$$

where  $\tilde{\pi}(S_0)$  is the steady state probability vector, given that the system is in  $S_0$ . We can now apply exact aggregation; the aggregated process is depicted in Fig. 5.

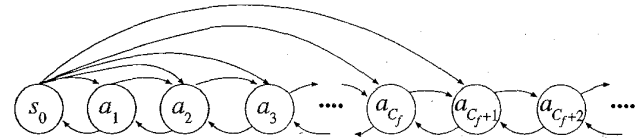


Fig. 5. Aggregate process for the lower bound model.

The transition rates for the aggregated chain are:

$$\begin{aligned} g_{s_0, a_i}^* &= \tilde{\pi}(S_0) Q_{S_0, a_i} & i = 1, \dots, C_f + 1 \\ g_{a_i, a_{i+1}}^* &= \lambda & i \geq 1 \\ g_{a_1, s_0}^* &= \mu_1 \\ g_{a_i, a_{i-1}}^* &= \begin{cases} \sum_{j=1}^i \mu_j & i = 2, 3, \dots, K \\ \mu^* & \text{otherwise} \end{cases} \end{aligned}$$

where  $\mu^* = \sum_{i=1}^K \mu_i$ .

Solving the chain, we have:

$$\begin{aligned} \pi^*(s_0) &= \left[ 1 + \sum_{i=1}^{C_f+1} \sum_{j=1}^i \lambda^{i-j} \left( \sum_{k=j}^{C_f+1} g_{s_0, a_k}^* \right) \left( \prod_{k=j}^i g_{a_k, a_{k-1}}^* \right) \right]^{-1} + \\ & \left[ \frac{\lambda}{\mu^* - \lambda} \sum_{j=1}^{C_f+1} \lambda^{C_f+1-j} \left( \sum_{k=j}^{C_f+1} g_{s_0, a_k}^* \right) \left( \prod_{k=j}^{C_f+1} g_{a_k, a_{k-1}}^* \right) \right]^{-1} \end{aligned} \quad (16)$$



$$\pi^*(a_i) = \pi(s_0) \sum_{j=1}^i \left[ \lambda^{i-j} \left( \sum_{k=j}^{C_f+1} g_{s_0, a_j}^* \left( \prod_{k=j}^i g_{a_k, a_{k-1}}^* \right) \right)^{-1} \right] \quad i = 1, \dots, C_f + 1 \quad (17)$$

$$\pi^*(a_i) = \pi(s_0) \left( \frac{\lambda}{\mu^*} \right)^{i-C_f-1} \sum_{j=1}^{C_f+1} \left[ \lambda^{C_f+1-j} \left( \sum_{k=j}^{C_f+1} g_{s_0, a_j}^* \left( \prod_{k=j}^{C_f+1} g_{a_k, a_{k-1}}^* \right) \right)^{-1} \right] \quad i = C_f + 2, \dots \quad (18)$$

To obtain the mean number of customers in the system,  $N_i$ , and the mean response time,  $R_i$ , let

$$\tilde{N}(S_0) = \sum_{s \in S_0} r(s) \tilde{\pi}(s)$$

where  $r(s) = \sum N_i$ , then we have:

$$\begin{aligned} N_i &= \tilde{N}(S_0) \pi^*(s_0) + \sum_{i=1}^{C_f} i \pi^*(a_i) + \sum_{i=C_f+1}^{\infty} i \pi^*(a_i) \\ &= \tilde{N}(S_0) \pi^*(s_0) + \sum_{i=1}^{C_f} i \pi^*(a_i) + \sum_{i=C_f+1}^{\infty} i \pi^*(a_{C_f+1}) \left( \frac{\lambda}{\mu^*} \right)^{i-C_f-1} \end{aligned}$$

After simplifying, the mean number of customers  $N_i$  is:

$$\begin{aligned} N_i &= \tilde{N}(S_0) \pi^*(s_0) + \sum_{i=1}^{C_f} i \pi^*(a_i) \\ &\quad + \frac{C_f \pi^*(a_{C_f+1}) \mu^*}{\mu^* - \lambda} + \frac{\pi^*(a_{C_f+1}) \mu^*}{(\mu^* - \lambda)^2} \end{aligned} \quad (19)$$

From Little's result [17], the lower bound mean response time is:

$$\begin{aligned} R_i &= \\ \frac{1}{\lambda} &\left[ \tilde{N}(S_0) \pi^*(s_0) + \sum_{i=1}^{C_f} i \pi^*(a_i) + \frac{C_f \pi^*(a_{C_f+1}) \mu^*}{\mu^* - \lambda} + \frac{\pi^*(a_{C_f+1}) \mu^*}{(\mu^* - \lambda)^2} \right] \end{aligned} \quad (20)$$

Lastly, to comment about the computational complexity of the lower bound model  $M^l$ . First, we have to obtain the conditional state probabilities for rate matrix  $S_0$ . Again, this can be accomplish by using the power iteration method. After we obtain the conditional state probabilities, we can use (16) to (20) to obtain the expected response time of the lower bound model.

## V. HOMOGENEOUS SERVERS

In this section we consider a system with  $K$  homogeneous servers having exponential service times with rate  $\mu$ . In this case, we can improve on the lower bound for the heterogeneous system as well as on the upper bound at high utilization. Here, the minimum expected delay policy becomes the classical *join the shortest queue (SQ)* policy.

We first describe the new upper bound model under very high system utilization. For the upper bound model  $M^u$  in Sec-

tion III we do not have a very tight upper bound under very *high* system utilization, since we put a constraint on the departure events based on the state of the system. Due to this constraint, the upper bound model saturates at a lower traffic intensity; if we can find an upper bound model that saturates at the same point as the original model, we can use the minimum of this model and  $M^u$  model as an upper bound. One simple upper bound for the homogeneous case which has the same saturation point as the original model is formed by assigning customers to servers in a cyclic fashion, [10]. In this case, each server in the system behaves as an  $E_k/M/1$ , and the mean response time of this system is well known [16]. Taking the minimum response time of this model and  $M^u$  provides a good upper bound over the entire range of traffic intensity.<sup>4</sup>

We now define the new lower bound model under the identical servers assumption. Let  $N(t) = (N_1(t), N_2(t), \dots, N_K(t))$  denote the joint queue lengths at time  $t > 0$  under SQ, and let  $N(t) = \sum_{k=1}^K N_k(t)$ . Let  $\hat{N}_k(t)$  denote the  $k$ th largest queue length,  $k = 1, 2, \dots, K$  at time  $t \geq 0$ . The new lower bound system operates as follows:

- Whenever  $N(t) < C_f = \sum_{i=1}^K C_i$ ,  $\hat{N}_1 - \hat{N}_K = d$ , and a departure would normally occur from the smallest queue, it is forced to occur instead from the next largest queue (i.e., if a departure would cause the system to exceed the maximum degree of imbalance  $d$ , then the departure is made to occur from the second shortest queue).
- Whenever  $N(t) \geq C_f$  and a departure occurs, it is taken from the largest queue.

Here  $C$  and  $d$  are parameters that can be tuned to provide a tight bound.

In order to describe in what sense this system is a lower bound, we introduce the concept of *majorization* [20]. Let  $X, Y \in IN^K$ .

DEFINITION 5.  $Y$  is said to majorize  $X$  (written  $X < Y$ ) iff

$$\begin{aligned} \sum_{l=1}^k \hat{X}_l &\leq \sum_{l=1}^k \hat{Y}_l, \quad k = 1, \dots, K-1, \\ \sum_{l=1}^K \hat{X}_l &= \sum_{l=1}^K \hat{Y}_l \end{aligned} \quad (21)$$

where  $\hat{X}_l(\hat{Y}_l)$  is the  $l$ -largest component of  $X(Y)$ . If we replace the equality in (21) by

$$\sum_{l=1}^K \hat{X}_l \leq \sum_{l=1}^K \hat{Y}_l,$$

we obtain a weaker ordering. In this case we say that  $Y$  weakly majorizes  $X$  (written  $X <_w Y$ ).

The following lemma states some properties regarding operations that can be performed on  $X$  and  $Y$  such that weak majorization is preserved.

<sup>4</sup> Note that this approach cannot be applied to the heterogeneous case since a cyclic assignment policy may not provide an upper bound response time.

LEMMA 5. Let  $X, Y \in IN^K$  such that  $X <_w Y$ , then

- 1)  $(\hat{X}_1, \dots, \hat{X}_k, \dots, \hat{X}_l + 1, \dots, \hat{X}_K) <_w (\hat{Y}_1, \dots, \hat{Y}_k + 1, \dots, \hat{Y}_l, \dots, \hat{Y}_K)$ ,  
for  $1 \leq k \leq l \leq K$
- 2)  $\hat{X}_1, \dots, (\hat{X}_k - 1)^+, \dots, \hat{X}_l, \dots, \hat{X}_K <_w (\hat{Y}_1, \dots, \hat{Y}_k, \dots, (\hat{Y}_l - 1)^+, \dots, \hat{Y}_K)$ ,  
for  $1 \leq k \leq l \leq K$

PROOF. The proof follows in a straightforward manner from the definition of " $<_w$ ." The reader is referred to [20] for a detailed proof.  $\square$

Before we define a stochastic comparison based on majorization, we introduce the notion of a *Schur-convex function*.

DEFINITION 6. A function  $\phi : IN \rightarrow IR$  is said to be *Schur-convex* iff

$$\phi(X) \leq \phi(Y), \quad \forall X, Y \in IN^K \quad \text{such that } X < Y.$$

DEFINITION 7. If  $X, Y \in IN^K$  are random variables, then we say  $X$  is smaller than  $Y$  in the sense of *Schur-convex order* (written  $X \leq_{scx} Y$ ) iff

$$\phi(X) \leq_{st} \phi(Y), \quad \forall \text{Schur-convex } \phi.$$

If the class of functions is restricted to be increasing *Schur-convex*, then we say that  $X$  is smaller than  $Y$  in the sense of *increasing Schur-convex order* ( $X \leq_{iscx} Y$ ).

Last, if  $\phi$  is a function such that  $\phi(X) \leq \phi(Y), \forall X, Y \in IN^K$  such that  $X <_w Y$ , then  $\phi$  can be shown to be an increasing *Schur-convex* function.

Let  $N^l(t)$  denote the joint queue length vector for the new lower bound system. We have the following result.

THEOREM 3. If  $N^l(0) = N(0)$ , then  $N^l(t) \leq_{iscx} N(t) \forall t > 0$ .

PROOF. Couple the initial queue lengths so that  $N^l(0) = N(0)$ .

Condition on the arrival times of the two systems. For the  $k$ th largest queue, we have an associated *service event process* which is a Poisson process with parameter  $\mu$ . Whenever a service event occurs associated with the  $k$ th largest queue, a departure occurs if there is one or more customer in the queue at the time of the event. Observe that the coupling of the the service event times at the different servers is only possible if the service times at the servers are all mutually independent sequences of i.i.d exponential random variables with the same parameter.

Let  $\{t_i\}_{i=0}^{\infty}$  be the sequence of times at which arrivals or service events occur ( $t_0 \equiv 0$ ). We will establish the relation  $N^l(t) <_w N(t)$  by induction on the event times. Clearly, if  $N^l(t_i) <_w N(t_i)$  then  $N^l(t) <_w N(t), t_i \leq t < t_{i+1}, i \geq 0$ .

*Basis step.* This follows from the coupling of the initial queue lengths.

*Inductive step.* Assume that  $N^l(t) <_w N(t)$  for  $t < t_i$ . We will establish it for  $t = t_i$ . There are two cases depending on whether the event is an arrival or a service event. For arrival event,  $N^l(t_i) <_w N(t_i)$  follows because arrivals are to the smallest queue, so Property 1 of Lemma 5 can be applied. For *service event*, there are two cases depending on whether  $N^l(t_i^-) < C_f$ . In either case, result follows from an application of Property 2 of Lemma 5.

This completes the inductive step and thus we have  $N^l(t) <_w N(t), t \geq 0$ . By the definition of weak majorization ( $<_w$ ), this implies that  $f(N^l(t)) \leq f(N(t))$  for any increasing *Schur-convex* function  $f(t)$ . Removing the conditioning on the arrival times and service times, we have  $N^l(t) \leq_{iscx} N(t) \forall t > 0$ .  $\square$

COROLLARY 1. If  $N^l(0) \leq_{iscx} N(0)$ , then  $N^l(t) \leq_{st} N(t)$ , for  $t \geq 0$ .

PROOF. This follows from the preceding theorem and the fact that

$$\phi(X) = \sum_{k=1}^K X_k \text{ is an increasing Schur-convex function. } \square$$

For the purpose of computing performance measures, let us define the following:

$S_0$  = set of states with  $0 \leq N_j \leq C$  and  $|N_i - N_j| \leq d, \forall i, j$ .

$s_0(C_0)$  = this is the only state in  $S_0$  that has a positive transition rate into it from states outside  $S_0$ .

$\tilde{\pi}(s_0(C_0))$  = conditional probability of state  $s_0(C_0)$ , given that the system is in  $S_0$ .

$s_0$  = aggregate state which represents all states of  $S_0$ .

$s_i$  = state which represents the system having  $C_f + i$  customers for  $i \geq 1$ .

$\pi^*(s_i)$  = steady state probability of state  $s_i$ .

$\tilde{N}(S_0)$  = mean number of customers given that the system is in  $S_0$ .

The mean number of customers and mean response time for this lower bound are:

$$\begin{aligned} N_l &= \tilde{N}(S_0)\pi^*(s_0) + \sum_{i=1}^{\infty} [C_f + i]\pi^*(s_i) \\ &= \tilde{N}(S_0)\pi^*(s_0) + C_f(1 - \pi^*(s_0)) + \lambda_0 \frac{K\mu}{(K\mu - \lambda)^2} \pi^*(s_0) \end{aligned} \quad (22)$$

and

$$R_l = \frac{1}{\lambda} \left[ \tilde{N}(S_0)\pi^*(s_0) + C_f(1 - \pi^*(s_0)) + \lambda_0 \frac{K\mu}{(K\mu - \lambda)^2} \pi^*(s_0) \right] \quad (23)$$

where:

$$\lambda_0 = \tilde{\pi}(s_0(C_0))\lambda \quad (24)$$

$$\pi^*(s_0) = \left[ 1 + \frac{\lambda_0}{K\mu - \lambda} \right]^{-1} \quad (25)$$

$$\pi^*(s_i) = \left[ 1 + \frac{\lambda_0}{K\mu - \lambda} \right]^{-1} \left( \frac{\lambda_0}{K\mu} \right) \left( \frac{\lambda}{K\mu} \right)^{i-1} \quad \text{for } i = 1, 2, \dots \quad (26)$$

Before we illustrate this application of this methodology, it is interesting to note that it is possible to reduce the state space further by lumpability as illustrated in [19].

## VI. APPLICATION AND NUMERICAL EXAMPLE

In this section, we present an example to illustrate the application of this methodology and the accuracy of the bounding algorithm. The example concerns the transmission policy in a computer network. For a packet switching system, there are basically two modes of data transmission. In one case, the vir-

tual circuit transmission mode, a path is first set up from the source node to the destination node. User packets then traverse the network following the path chosen during the initial connection setup. In this mode, user packets can arrive in sequence in which they were transmitted but the user has to pay for the overhead for the initial connection setup (ex: 3-way handshake in TCP/IP). The other transmission mode is the connectionless mode, where each individual packet or datagram independently traverses the network from source node to destination node. No initial connection is set up in this case and datagrams are forwarded through the network on an individual basis. Routing of each datagram is based on the destination address and the availability of output ports<sup>5</sup> at each intermediate nodes. Usually, connectionless transmission can yield better performance since there is no overhead for connection setup, but downside is that the user has to take care of packet resequencing and retransmission in case datagrams are dropped.

In this first example, we assume datagram transmission mode and we want to evaluate the transmission time of datagrams in each communication node (see Fig. 6).

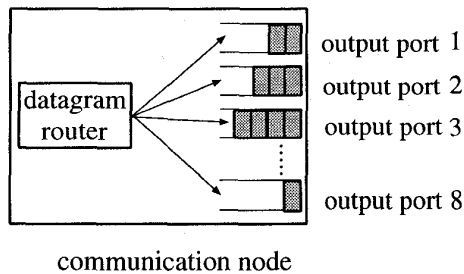


Fig. 6. Communication node for sending out datagram.

Assume that the communication node has eight output ports. We model the datagram arrival process of to be a Poisson process with rate  $\lambda$  and each output port as having an exponential service rate of  $\mu$ . We normalize the service rate to 1.0 and we increase the datagram arrival rate so that the system utilization will vary from 0.1 to 0.9. Table I illustrates the upper and lower bound on the mean response time of the datagram (the time to send the datagram across the output link of the communication node) as a function of system utilization. As can be seen by the percentage error,<sup>6</sup> the bounding methodology we developed provides a very good accuracy from the mean response time and at a moderate cost.

To illustrate the tradeoff between computational cost and accuracy of the bounds, we can vary the  $d$  and  $C$  parameters. By fixing the system utilization at 0.9 and increasing the number of states generated, we see the improvement obtained for the bounds on the mean response time. The results are illustrated in Table II.

5. Assuming these output ports can reach the destination node.

6. If the spread in the bounds is less than  $< 10^{-6}$ , we leave the percentage error entry blank.

TABLE I  
DATAGRAM RESPONSE TIME

System Utilization	Response Time Lower Bound	Response Time Upper Bound	Spread of Bounds	Percentage Error
0.1	0.1000252	0.1000252		
0.2	0.2000863	0.2000863		
0.3	0.3008306	0.3008306		
0.4	0.4052623	0.4052623		
0.5	0.5208155	0.5208162	0.0000007	$6.27 \times 10^{-5} \%$
0.6	0.6610700	0.6610820	0.0000120	$9.07 \times 10^{-4} \%$
0.7	0.8521012	0.8522784	0.0001772	0.0103 %
0.8	1.1640786	1.1652135	0.0011349	0.0487 %
0.9	1.9107856	1.9273843	0.0165987	0.4324 %

TABLE II  
COMPUTATIONAL COST VS. ACCURACY

$d$	$C$	States Generated	Response Time Upper Bound	Response Time Lower Bound	Spread of Bounds	Percentage Errors
4	8	1815	1.8521678	2.1078925	0.2557247	6.4576 %
4	10	2475	1.8973256	2.0013574	0.1040318	2.6684 %
5	12	6831	1.9107856	1.9273843	0.0165987	0.4324 %
6	13	15015	1.9123782	1.9261783	0.0138001	0.3591 %

## VII. CONCLUSION

The minimum expected delay routing policy is appealing to study not only due to its simplicity in implementation, but also due to the fact that it is theoretically difficult to analyze because the routing of arrivals is state dependent and no closed form solutions exist in general. Also, due to the fact that each server has an infinite capacity queue, the state space cardinality of the Markov model is infinite, and it becomes impossible to generate the entire state space to solve the Markov model numerically. We have presented an approach to bound the mean response time and the mean number of customers in the minimum expected delay routing policy, which is a generalization of the join the shortest queue routing policy. The algorithmic approach provides the flexibility to tradeoff computational resources and tighter bounds. There is ongoing work on the subject to determine  $d$  and  $C_i$  to obtain specified error bounds. We are also investigating the possibility of bounding the mean response time under more relaxed conditions, e.g., by allowing general service distributions.

## ACKNOWLEDGMENTS

The work of John C.S. Lui was supported by the Croucher Foundation and the Direct Grant. The work of Richard R. Muntz was supported by the National Science Foundation under grant CCR-9215064. The work of Don Towsley was supported by the National Science Foundation under grant NCR-9116183.

## REFERENCES

- [1] I.J.B.F. Adan, J. Wessels, and W.H.M. Zijm, "Analysis of the symmetric shortest queue problem," *Stochastic Models*, vol.6, pp. 691-713, 1990.
- [2] I.J.B.F. Adan, J. Wessels, and W.H.M. Zijm, "Analysis of the asymmetric shortest queue problem," *Queueing Systems*, vol. 8, pp. 1-58, 1991.

- [3] I. Adan, G.-J. van Houtum, and J. van der Wal, "Upper and lower bounds for the waiting time in the symmetric shortest queue system," Technical Report COSOR 92-09, Eindhoven Univ. of Technology.
- [4] N.R. Baker, "Optimal user search sequences and implications for information system operation," *J. Applied Probability*, vol.24, pp. 540-546, 1987.
- [5] J.P.C. Blanc, "A note on waiting times in systems with queues in parallel," *J. Applied Probability*, vol.24, pp. 540-546, 1987.
- [6] B.W. Conolly, "The autostrada queueing problem," *J. Applied Probability*, vol. 21, pp. 394-403, 1984.
- [7] J.W. Cohen and O.J. Boxma, *Boundary Value Problems in Queueing System Analysis*. North Holland, 1983.
- [8] P.J. Courtois, *Decomposability—Queueing and Computer System Applications*. New York: Academic Press, 1977.
- [9] P.J. Courtois and P. Semal, "Computable bounds for conditional steady-state probabilities in large Markov chains and queueing models," *IEEE J. Selected Areas in Communications*, vol 4, no. 6, Sept. 1986.
- [10] A. Ephremides, P. Varaiya, and J. Walrand, "A simple dynamic routing problem," *IEEE Trans. Automatic Control*, vol. 25, 1980.
- [11] L. Flatto and H.P. McKean, "Two queues in parallel," *Comm. Pure and Applied Mathematics*, vol. 30, pp. 255-263, 1977.
- [12] S. Halfin, "The shortest queue problem," *J. Applied Probability*, vol. 22, pp. 865-878, 1985.
- [13] W.K. Grassmann, "Transient and steady state results for two parallel queues," *Omega*, vol. 8, pp. 105-112, 1980.
- [14] J.F.C Kingman, "Two similar queues in parallel," *Annals of Mathematical Statistics*, vol 32, pp. 1,314-1,323, 1961.
- [15] C. Knessl, B.J. Matkowsky, Z. Schuss, and C. Tier, "Two parallel queues with dynamic routing," *IEEE Trans. Communications*, vol. 34, pp. 1,170-1,175, 1986.
- [16] L. Kleinrock, *Queueing Systems: Volume I: Theory*. New York: Wiley-Interscience Publication, 1975.
- [17] J.D.C Little, "A proof of the queueing formula  $L = \lambda W$ ," *Operations Research*, vol 9, pp. 383-387, 1967.
- [18] H.C. Lin and C.S. Raghavendra, "An analysis of the join the shortest queue policy," Electrical Eng. technical report, Univ. of Southern California, 1991.
- [19] J.C.S. Lui and R.R. Muntz, "Algorithmic approach to bounding the response time of a minimum expected delay routing system," *Proc. 1992 ACM SIGMETRICS/Performance '92 Conf.*, pp. 140-152.
- [20] A.W. Marshall and I. Olkin, *Inequalities: Theory of Majorization and Applications*. New York: Academic Press, 1979.
- [21] R.D. Nelson and T.K. Philips, "An approximation to the response time for shortest queue routing," *ACM SIGMETRICS*, vol 17, no 1, pp 181-189, 1989.
- [22] R.D. Nelson and T.K. Philips, "An approximation for the mean response time for shortest queue routing with general interarrival and service times," Technical Report RC15429, IBM T.J. Watson Research Lab, 1990.
- [23] M.F. Neuts, *Matrix-Geometric Solutions in Stochastic Models—An Algorithmic Approach*. Baltimore: Johns Hopkins Univ. Press, 1981.
- [24] B.M. Rao and M.J.M. Posner, "Algorithmic and approximate analysis of the shorter queue model," *Naval Research Logistics*, vol.34, pp. 381-398, 1987.
- [25] W.J. Stewart, "MARCA: Markov chain analyzer, a software package for Markov modeling," *Numerical Solution of Markov Chains*. Dekker Press, 1991.
- [26] D. Towsley and S. Chen, "Design and modeling policies for two server fork/join queueing systems," COINS Technical Report 91-39, Univ. of Massachusetts.
- [27] D. Towsley, P. Sparaggis, and C. Cassandras, "Stochastic ordering properties and optimal routing control for a class of finite capacity queueing systems," *IEEE Trans. Automatic Control*, vol. 37, no. 9, pp. 1,446-1,451, Septe. 1992.
- [28] Y.T. Wang and R.J.T. Morris, "Load sharing in distributed systems," *IEEE Trans. Computers*, vol. 34, no. 3, pp. 204-217, Mar. 1985.
- [29] W. Winston, "Optimality of the shortest line discipline," *J. Applied Probability*, vol 15, pp. 181-189, 1977.
- [30] Y. Zhao and W.K. Grassmann, "The shortest queue model with jockeying," *Naval Research Logistics*, vol. 37, pp. 773-787, 1990.
- [31] Y. Zhao and W.K. Grassmann, "A numerically stable algorithm for two server queue models," *Queueing Systems*, vol 8, pp. 59-80, 1991.



**John Chi-Shing Lui** received his PhD in computer science from UCLA in 1991. He then joined IBM San Jose and participated in the research and development of a parallel architecture project. He also participated in the parallel database project in IBM Yorktown Research Laboratory. He is currently an assistant professor of computer science at the Chinese University of Hong Kong. His research interests are in parallel and distributed system design, distributed multimedia systems, parallel I/O architecture, communication networks, mobile computing, and performance evaluation theory.



**Richard R. Muntz** received the BEE from Pratt Institute in 1963, the MEE from New York University in 1966, and the PhD in electrical engineering from Princeton University in 1969.

Dr. Muntz is a professor in the Computer Science Department, School of Engineering and Applied Science, University of California, Los Angeles. He is a member of the Board of Directors for SIGMETRICS. He was an associate editor for the *Journal of the ACM* from 1975 to 1980, and editor-in-chief for *ACM Computing Surveys* from 1992 to 1995. He is a member of Sigma Xi, Tau Beta Pi, IFIP WG7.3, the Association for Computing Machinery, and a fellow of the IEEE.

Dr. Muntz's current research interests are mass storage and storage hierarchies, parallel and distributed database systems, temporal/spatial data models and query processing, and computer performance evaluation.

**Don Towsley** received the BA degree in physics and the PhD degree in computer science from the University of Texas in 1971 and 1975, respectively. From 1976 to 1985, he was a member of the faculty of the Department of Electrical and Computer Engineering at the University of Massachusetts, Amherst. He is currently a professor of computer science at the University of Massachusetts. During 1982-1983, he was a visiting scientist at the IBM T.J. Watson Research Center, Yorktown Heights, New York, and during 1989-1990, he was a visiting professor at the Laboratoire MASI, Paris, France. His current interests include high speed networks and multimedia systems.

Dr. Towsley is currently and editor of *IEEE/ACM Transactions on Networking*, and is on the editorial boards of *Networks*, *Journal of Dynamic Discrete Event Systems*, and *Performance Evaluation*. He was a program co-chair of the joint ACM SIGMETRICS and PERFORMANCE '92 conference. He is a member of the ACM, IEEE, and ORSA, and is active in the IFIP Working Group 7.3 on Performance Modeling.