

**Computer Science Department Technical Report
University of California
Los Angeles, CA 90024-1596**

**BOUNDING THE MEAN RESPONSE TIME OF A MINIMUM
EXPECTED DELAY ROUTING SYSTEM: AN ALGORITHMIC
APPROACH**

**J. C.-S. Lui
R. R. Muntz
D. Towsley**

**July 1993
CSD-930022**

Bounding the Mean Response Time of a Minimum Expected Delay Routing System: An Algorithmic Approach

John C.S. Lui
Computer Science Department
The Chinese University of Hong Kong

Richard R. Muntz¹
UCLA Computer Science Department

Don Towsley²
University of Massachusetts
Computer Science Department

¹The work of this author was supported in part by the National Science Foundation under grant CCR-9215064.

²The work of this author was supported in part by the National Science Foundation under grant NCR-9116183.

Abstract

We study a heterogeneous multi-server queueing system in which the minimum expected delay routing policy is used, i.e., an arriving customer is assigned to the server 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. 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. Examples are presented which show the excellent relative accuracy achievable with modest computational cost.

Index terms: Load balancing, shortest queue routing, bounds, Markov models.

1 Introduction

Routing policies can have tremendous effect on the performance of a multi-server system where each server has its own queue. The minimum expected delay routing policy, although not optimal, 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 approach to 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 multi-server queueing systems in which the assignment of customers to servers is chosen at arrival instants using the minimum expected delay routing (MED) policy (a natural generalization of the join-the-shortest-queue (SQ) policy for homogeneous servers). 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 μ_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 general arrival process with a mean rate of λ . 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] and has been shown to minimize the queue length vector in the sense of Schur-convex ordering, [27]. This latter fact carries the implication that the total number of customers in the system is stochastically minimized by the SQ policy as is the mean stationary response time (when it exists).

Of more interest to us is the literature dealing with the performance evaluation of the SQ and MED policies. In the case of the SQ policy for two identical servers, numerous authors have provided exact, though not necessarily efficiently computable, solutions Kingman [14], Flatto and McKean [11], Zhao and Grassmann [31], and Cohen and Boxma [6], Adan [1], et al. Several authors have provided similar solutions to the heterogeneous server problem; e.g., Knessl, et al. [15] and Adan, et al. [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.

Numerous authors have proposed approximations for the SQ and MED policies. These include Conolly [5], Rao and Posner [25], 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 [24]. The last paper, [26], is also noteworthy in that it provides upper and 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 [4] 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 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 [19], Nelson and Philips [22, 23], 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. Avritzer [3] studied a dynamic load balancing algorithm which used a threshold policy in an asymmetric distributed system. The result was only applicable to two distinct types of servers and a small class of threshold sizes, no formal proof was given on how to obtain performance bounds. None of the work cited above treated more than two heterogeneous servers and simultaneously provided error bounds. Lui and Muntz [20] were the first to propose a methodology to bound the mean response time of a minimum expected delay routing system. This paper differs from [20] in several ways. First, we derive improved bounds for the homogeneous servers case, and secondly, we use sample path analysis to prove the bounds, which yields 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 2, we formally define the queueing model we are analyzing. In Sections 3 and 4, we present the modified models and prove that these models do provide bounds. In Section 5, we provide a methodology for obtaining tighter bounds in the special case of homogeneous servers. Section 6 shows

that using lumpability, we can reduce the state space further and thereby obtain better bounds at the same cost. In Section 7, we present example applications, and conclusions are given in Section 8.

2 Model

We consider a queueing system, as depicted in Figure 1, with K heterogeneous servers with associated queues being fed by a general arrival process with mean rate λ . The service times at servers form mutually independent sequences of exponential random variables with rates $\mu_1 \geq \mu_2 \geq \dots \geq \mu_K$. and are independent of arrival times. 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. We call this the *minimum expected delay routing policy*. When all service rates are equal, this policy reduces to the classic join-the-shortest queue (*SQ*) routing algorithm.

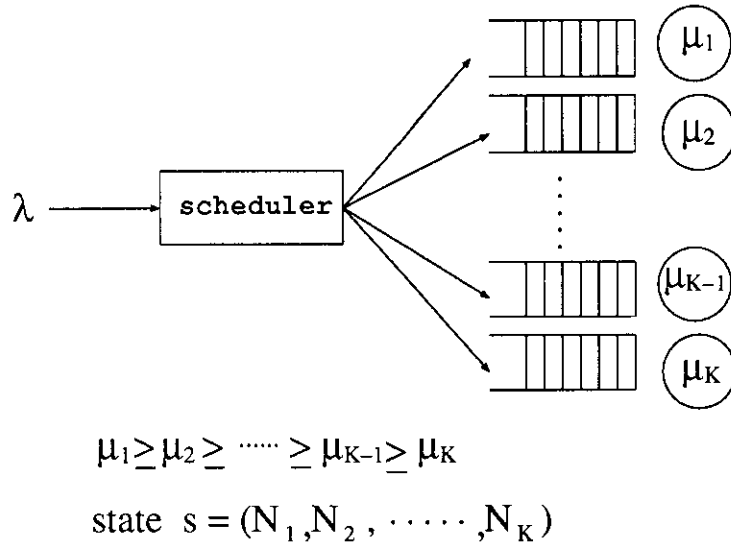


Figure 1: Minimum Expected Delay Routing Policy Queueing Model

We can construct a Markov model, M , for this queueing system with state space:

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

Assume the system is stable; that is $\lambda < \sum_{i=1}^K \mu_i$. The unique steady state probability vector for this continuous-time Markov model satisfies the following system of linear equations:

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

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

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

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 (a) the modified system has an efficient solution and (b) 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.

3 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 general with mean rate λ and K servers with service rates $\mu_i, i = 1, 2, \dots, K$, where $\mu_1 \geq \mu_2 \geq \dots \geq \mu_K$ ¹. The upper bound model M^u has two additional parameters. The first parameter we term the *artificial capacity vector* $\overline{C} = (C_1, \dots, C_K)$. The second parameter is a threshold setting d , which is the maximum allowable difference between the longest queue and the shortest queue in 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 construct the upper bound model M^u .

Definition 1 Let $\overline{C}^* = (C_1^*, \dots, C_K^*)$ 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 \overline{C} is an integer multiple of \overline{C}^* , i.e., $\overline{C} = j\overline{C}^*$ for some $j \geq 1$.

For example, if $\mu_1 = \frac{3}{2}$ and $\mu_2 = 1$, then $C = 3$ and $\overline{C}^* = (3, 2)$. So the artificial capacity vector can be $\overline{C} = j\overline{C}^*$ 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 in the i^{th} server. Define $q(s)$ 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

¹We 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.

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 (i.e., become active again). 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 *repetitive 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 $\cup_{i=0}^{\infty} \mathcal{S}_i \dots$ where:

$$\begin{aligned} \mathcal{S}_0 &= \{(N_1, \dots, N_K) | 0 \leq N_j \leq C_j \text{ for } j = 1, \dots, K\} \\ \mathcal{S}_i &= \{(N_1, \dots, N_K) | iC_j \leq N_j \leq (i+1)C_j \text{ for } j = 1, \dots, K\} - \{(iC_1, \dots, iC_K)\} \quad i \geq 1 \end{aligned}$$

Due to the routing of arrivals and the constraint on departures, we can show that all transitions from \mathcal{S}_i to \mathcal{S}_{i+1} are actually from one state in \mathcal{S}_i and the transitions from \mathcal{S}_{i+1} to \mathcal{S}_i can only go to one state in \mathcal{S}_i . As will be shown later, this property allows us to efficiently solve the model via *exact* decomposition based on the partition $\{\mathcal{S}_0 \cup \mathcal{S}_1 \cup \dots\}$. 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 let $C_i = 10$, for $i = 1, 2, 3, 4$. It is easy to see that \mathcal{S}_0 consists of all states for which each queue has between 0 and 10 customers; \mathcal{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 \mathcal{S}_0 to \mathcal{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 \mathcal{S}_1 to \mathcal{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 \mathcal{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.

3.1 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. Since M^u and M have the same mean arrival rate it follows that if M^u has a larger mean number in system than M , the mean response time of M^u is an upper bound on the mean response time for M . We therefore concentrate on the mean number in system in the remainder of this section.

Let p denote a policy that routes an arriving customer to a server on the basis of the server queue lengths. Let $l_p^*(\mathbf{N})$ denote the identity of the queue to which the customer is routed under policy p when the queue length vector is \mathbf{N} . When we are interested in the joint queue length at time t under a specific policy, we will denote it as $\mathbf{N}^p(t)$. We assume that p is stationary. We start by defining an auxiliary concept that will be useful in the proof.

Definition 4 *A policy p is a proper policy if $\mathbf{N} \leq \mathbf{N}'$ (here " \leq " is taken to mean componentwise) implies that $\mathbf{N} + \mathbf{e}_{l_p^*(\mathbf{N})} \leq \mathbf{N}' + \mathbf{e}_{l_{p'}^*(\mathbf{N}')}$ where \mathbf{e}_k is the vector of all 0's 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 time instants 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 μ_k . (Note that a service event is also 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 customer service times are still i.i.d. exponential with the same mean.

Consider the i -th event. Let $\mathbf{N}_i = (N_{i,1}, \dots, N_{i,K})$ be the joint queue lengths immediately after the i -th event. Let \mathbf{N}_0 denote the initial queue lengths. We have the following evolution equations. If the $(i+1)$ -st event corresponds to an arrival,

$$N_{i+1,k} = N_{i,k} + 1\{l_p^*(\mathbf{N}_i) = k\}, \quad 1 \leq k \leq K \quad (2)$$

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 $\mathbf{N}^u(t)$ 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,

$$N_{i+1,k}^u = N_{i,k}^u + 1\{l_p^*(\mathbf{N}_i^u) = k\}, \quad 1 \leq k \leq K \quad (4)$$

If the $(i+1)$ -st event 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 $\mathbf{N}(0) \leq_{st} \mathbf{N}^u(0)$ and p is a proper routing policy, then*

$$\mathbf{N}(t) \leq_{st} \mathbf{N}^u(t), \quad t \geq 0.$$

Proof. Couple the initial queue lengths so that $\mathbf{N}(0) \leq \mathbf{N}^u(0)$. 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:

$$\mathbf{N}_i \leq \mathbf{N}_i^u, \quad i \geq 0.$$

For $i = 0$, $\mathbf{N}(0) \leq \mathbf{N}^u(0)$.

Assume $\mathbf{N}_i \leq \mathbf{N}_i^u$ 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:

$$N(t) \leq_{st} N^u(t), \quad \text{for } t \geq 0$$

■

Let $N_i = \lim_{t \rightarrow \infty} N_i(t)$ when it exists, $1 \leq i \leq K$ and let $N = \sum_{i=1}^K N_i$. Based on the above 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,

$$E[R] \leq E[R_u]$$

3.2 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 ² We partition the state space of M^u , $\mathcal{S}^u = \cup_{i=0}^{\infty} \mathcal{S}_i$ and $\mathcal{S}_i \cap \mathcal{S}_j = \emptyset$, $\forall i \neq j$, where:

$$\begin{aligned} \mathcal{S}_0 &= \{(N_1, \dots, N_K) | 0 \leq N_j \leq C_j \text{ for } j = 1, \dots, K\} \\ \mathcal{S}_i &= \{(N_1, \dots, N_K) | iC_j \leq N_j \leq (i+1)C_j \text{ for } j = 1, \dots, K\} - \{(iC_1, \dots, iC_K)\} \\ Q_{\mathcal{S}_i, \mathcal{S}_j} &= \text{transition rate matrix from states in } \mathcal{S}_i \text{ to states in } \mathcal{S}_j. \end{aligned}$$

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

This is a block tridiagonal transition rate matrix and therefore represents a quasi-birth-death process. By aggregating each partition \mathcal{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 \mathcal{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 .

²Although M^u yields an upper bound on the mean response time for arbitrary arrival process, the computation algorithm described in this section is for Poisson arrival processes only.

$$Q^u = \begin{bmatrix} Q_{\mathcal{S}_0, \mathcal{S}_0} & Q_{\mathcal{S}_0, \mathcal{S}_1} & 0 & 0 & 0 & \cdots \\ Q_{\mathcal{S}_1, \mathcal{S}_0} & Q_{\mathcal{S}_1, \mathcal{S}_1} & Q_{\mathcal{S}_1, \mathcal{S}_2} & 0 & 0 & \cdots \\ 0 & Q_{\mathcal{S}_2, \mathcal{S}_1} & Q_{\mathcal{S}_2, \mathcal{S}_2} & Q_{\mathcal{S}_2, \mathcal{S}_3} & 0 & \cdots \\ 0 & 0 & Q_{\mathcal{S}_3, \mathcal{S}_2} & Q_{\mathcal{S}_3, \mathcal{S}_3} & Q_{\mathcal{S}_3, \mathcal{S}_4} & \cdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix}$$

Figure 2: Transition Rate Matrix for the Upper Bound Model.

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

$$s_i(C_0) = (N_1, N_2, \dots, N_K) \in \mathcal{S}_i \text{ where } N_j = (i+1)C_j \quad \forall j = 1, 2, \dots, K$$

This follows from the rule used to assign an arriving customer to a server. Also, there are K states from \mathcal{S}_i that have non-zero transition rates to a state in \mathcal{S}_{i-1} where $i \geq 1$. Each corresponds to a state in which a particular server is the last to complete its “active” (non-suspended) customer. Let us call these states $s_i(l)$, $1 \leq l \leq K$, $i \geq 1$. These states are:

$$\begin{aligned} s_i(l) &= (N_1, N_2, \dots, N_K) \in \mathcal{S}_i \quad l = 1, \dots, K \\ \text{where } N_l &= iC_l + 1 \text{ and} \\ N_j &= iC_j \text{ for } j \neq l \text{ and } j = 1, 2, \dots, K \end{aligned}$$

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 $s_{i+1}(l)$, $l = 1, 2, \dots, K$:

$$\begin{aligned} s_i(C_0) &\rightarrow s_{i+1}(l) && 1\{l_p^*(s_i(C_0)) = l\}\lambda \\ s_{i+1}(l) &\rightarrow s_i(C_0) && \mu_l \text{ for } l = 1, 2, \dots, K \end{aligned}$$

The second important observation is that the submatrices $Q_{\mathcal{S}_i, \mathcal{S}_i}$, for $i \geq 1$, are all identical. The conditional state probabilities $P\{s \in \mathcal{S}_i | \mathcal{S}_i\}$ can now be computed exactly using the following lemma from [8]:

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 sub-matrix 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:

$$\begin{aligned}\vec{\pi}_{|A} [Q_{A,A} + Q_{A,B} \underline{e} \underline{e}_i^T] &= \vec{0} \\ \vec{\pi}_{|A} \underline{e} &= 1\end{aligned}$$

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

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

Theorem 1 Let $\tilde{Q}_{\mathcal{S}_i, \mathcal{S}_i}$ be the transition rate matrix which is equal to $Q_{\mathcal{S}_i, \mathcal{S}_i}$, except for the following modifications:

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

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

The solution to the following system of linear equations:

$$\vec{\pi} \tilde{Q}_{\mathcal{S}_i, \mathcal{S}_i} = \vec{0} \quad \text{and} \quad \vec{\pi} \underline{e} = 1$$

is the conditional steady state probability vector for states in \mathcal{S}_i , that is:

$$\tilde{\pi}(s) = \frac{\pi(s)}{\sum_{s \in \mathcal{S}_i} \pi(s)} \quad \forall s \in \mathcal{S}_i$$

Proof: Let us partition the state space $\mathcal{S}^u = \{\mathcal{S}'_i \cup \mathcal{S}''_i\}$ where $\mathcal{S}'_i = \cup_{j=0}^i \mathcal{S}_j$ and $\mathcal{S}''_i = \{\mathcal{S}^u - \mathcal{S}'_i\}$. There is only a single return state in \mathcal{S}'_i , which is $s_i(C_0)$, from the states in \mathcal{S}''_i . Based on Lemma 2, the modification of Equation (6) provides the conditional steady state probability, given the system is in \mathcal{S}'_i . Now partition the state space $\mathcal{S}'_i = \{\mathcal{S}^1_i \cup \mathcal{S}_i\}$ where $\mathcal{S}^1_i = \cup_{j=0}^{i-1} \mathcal{S}_j$. Based on Lemma 2 and the definition of the MED routing policy, using the modification given in Equation (7) we obtain the conditional state probability vector, given the system is in \mathcal{S}_i . ■

Since we can compute the conditional state probabilities for each partition \mathcal{S}_i exactly, we can exactly aggregate each \mathcal{S}_i into a single state $s_i, i \geq 0$. The aggregated process is depicted in Figure 3 where, λ_0, λ_{agg} and μ_{agg} are:

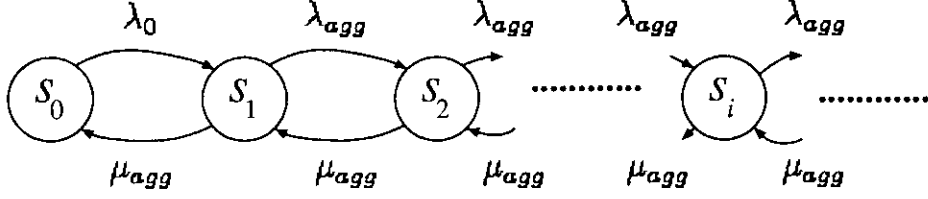


Figure 3: Aggregate Process for the Upper Bound Model.

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

Solving this chain, we have:

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

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

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

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

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

$$\vec{\pi} \dot{Q}_{\mathcal{S}_i, \mathcal{S}_i} = \vec{0} \quad \text{and} \quad \vec{\pi} \underline{e} = 1$$

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:

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

From Little's result [18], the upper bound mean system response time R_u is:

$$R_u = \frac{1}{\lambda} \left[\hat{N}(s_i)\pi^*(s_0) + \hat{N}(s_1)(1 - \pi^*(s_0)) + C_0\lambda_0 \frac{\mu_{agg}}{(\mu_{agg} - \lambda_{agg})^2} \pi^*(s_0) \right] \quad (12)$$

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.

4 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, we assume that the arrival process is general with rate λ 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.

4.1 Proof that M^l Provides a Lower Bound for M

The proof that M^l provides a lower bound on the mean response time of the system for all proper routing policies p 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 $\mathbf{N}(0) \leq_{st} \mathbf{N}'(0)$ and p is a proper routing policy, then*

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

Proof. Without loss of generality, we can couple the systems, so that $\mathbf{N}(0) \leq \mathbf{N}'(0)$. Now 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 the fact that p is a proper policy suffices to establish that $\mathbf{N}(t) \leq \mathbf{N}'(t)$. Removal of the conditioning yields the desired result. ■

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 $N^p(t) = \sum_{i=1}^K N_i^p$ 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 $\gamma_n^p : \mathbf{I} \rightarrow \mathbf{I}$ 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 variable's,

- $\{A_n\}$ is a sequence of random variable 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 random variable 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 random variable 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:

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

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

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.

Induction 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. Therefore if the k -th fastest server is idle in the full service model then there are less than k customers in the full service model. On the other hand, if the k -th fastest server is busy in the normal service model there must be at least k customers in this model. 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 systems, 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 $B_n \leq \mu_k/\mu$, then $B_n \leq \mu_j/\mu$. Hence it follows from the evolution equations, that 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. ■

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*

$$N^l(t) \leq_{st} N^p(t), \quad t \geq 0,$$

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. ■

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

4.2 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:

- \mathcal{S}_0 = set of states with $0 \leq N_j \leq C_j, j = 1, 2, \dots, K$. and such that the threshold d is respected
 \mathcal{G}_1 = $\{\mathcal{S}_0 - (0, 0, \dots, 0)\}$.
 a_i = a state, in the complement of \mathcal{S}_0 , in which the system contains i customers.
 $Q_{\mathcal{G}_1, a_i}$ = transition rate matrix between \mathcal{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 Figure 4. (Note that some of the $Q_{\mathcal{G}_1, a_i} = 0$ but this will not effect the development that follows.)

$$\begin{array}{c|cccccc}
 Q_{a_0, a_0} & Q_{a_0, \mathcal{G}_1} & 0 & 0 & 0 & 0 & \dots \\
 Q_{\mathcal{G}_1, a_0} & Q_{\mathcal{G}_1, \mathcal{G}_1} & Q_{\mathcal{G}_1, a_1} & Q_{\mathcal{G}_1, a_2} & Q_{\mathcal{G}_1, a_3} & Q_{\mathcal{G}_1, a_4} & \dots \\
 \hline
 g_{a_1, a_0}^* & 0 & g_{a_1, a_1}^* & \lambda & 0 & 0 & \dots \\
 0 & 0 & g_{a_2, a_1}^* & g_{a_2, a_2}^* & \lambda & 0 & \dots \\
 0 & 0 & 0 & g_{a_3, a_2}^* & g_{a_3, a_3}^* & \lambda & \dots \\
 \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots
 \end{array}$$

Figure 4: Transition rate matrix for lower bound model.

Since \mathcal{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 \mathcal{S}_0 (where $\mathcal{S}_0 = \{a_0 \cup \mathcal{G}_1\}$), then we can aggregate \mathcal{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 \mathcal{S}_0 from all states outside \mathcal{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 \mathcal{S}_0 can be obtained by solving the following system of linear equations:

$$\begin{aligned}
 \vec{\pi}(\mathcal{S}_0) \left[Q_{\mathcal{S}_0, \mathcal{S}_0} + \sum_{i=1}^{C_j+1} Q_{\mathcal{S}_0, a_i} \underline{e}_0^T \right] &= \vec{0} \\
 \vec{\pi}(\mathcal{S}_0) \underline{e} &= 1
 \end{aligned}$$

where $\vec{\pi}(\mathcal{S}_0)$ is the steady state probability vector, given that the system is in \mathcal{S}_0 . We can now apply exact aggregation; the aggregated process is depicted in Figure 5.

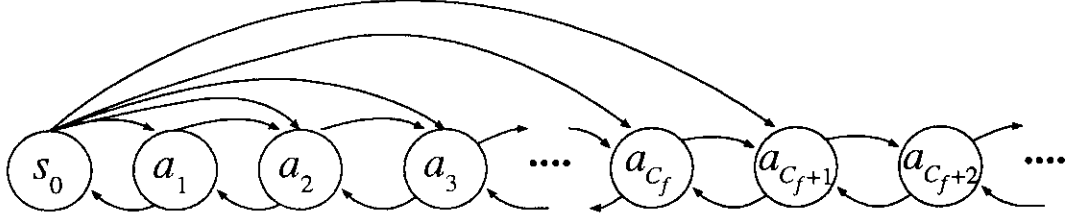


Figure 5: Aggregate Process for the Lower Bound Model.

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:

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

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

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

To obtain the mean number of customers in the system, N_l , and the mean response time, R_l , let

$$\tilde{N}(\mathcal{S}_0) = \sum_{s \in \mathcal{S}_0} r(s) \tilde{\pi}(s)$$

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

$$N_l = \tilde{N}(\mathcal{S}_0) \pi^*(s_0) + \sum_{i=1}^{C_f} i \pi^*(a_i) + \sum_{i=C_f+1}^{\infty} i \pi^*(a_i)$$

$$= \hat{N}(\mathcal{S}_0)\pi^*(s_0) + \sum_{i=1}^{C_f} i\pi^*(a_i) + \sum_{i=C_f+1}^{\infty} i\pi^*(a_{C_f+1}) \left(\frac{\lambda}{\mu^*}\right)^{i-C_f-1}$$

After simplifying, the mean number of customers N_l is:

$$N_l = \hat{N}(\mathcal{S}_0)\pi^*(s_0) + \sum_{i=1}^{C_f} i\pi^*(a_i) + \frac{C_f \pi^*(a_{C_f+1})\mu^*}{\mu^* - \lambda} + \frac{\pi^*(a_{C_f+1})\mu^*}{(\mu^* - \lambda)^2} \quad (19)$$

From Little's result [18], the lower bound mean response time is:

$$R_l = \frac{1}{\lambda} \left[\hat{N}(\mathcal{S}_0)\pi^*(s_0) + \sum_{i=1}^{C_f} i\pi^*(a_i) + \frac{C_f \pi^*(a_{C_f+1})\mu^*}{\mu^* - \lambda} + \frac{\pi^*(a_{C_f+1})\mu^*}{(\mu^* - \lambda)^2} \right] \quad (20)$$

5 Homogeneous Servers

In this section we consider a system with K homogeneous servers having exponential service times with rate μ . 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 the Section 3, 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 [17]. Taking the minimum response time of this model and M^u provides a good upper bound over the entire range of traffic intensity³.

We now define the new lower bound model under the identical servers assumption. Let $\mathbf{N}(t) = (N_1(t), N_2(t), \dots, N_K(t))$ 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:

³Note that this approach cannot be applied to the heterogeneous case since a cyclic assignment policy may not provide an upper bound response time.

- 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, then 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, which is a form of jockeying).
- 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* [21]. Let $\mathbf{X}, \mathbf{Y} \in \mathbb{N}^K$.

Definition 5 \mathbf{Y} is said to majorize \mathbf{X} (written $\mathbf{X} \prec \mathbf{Y}$) iff

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

where \hat{X}_l (\hat{Y}_l) is the l -largest component of \mathbf{X} (\mathbf{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 \mathbf{Y} *weakly majorizes* \mathbf{X} (written $\mathbf{X} \prec_w \mathbf{Y}$).

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

Lemma 5 Let $\mathbf{X}, \mathbf{Y} \in \mathbb{N}^K$ such that $\mathbf{X} \prec_w \mathbf{Y}$, then

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

Proof. The proof follows in a straightforward manner from the definition of “ \prec_w ”. The reader is referred to [21] for a detailed proof. ■

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

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

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

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

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

One property of these orderings that will be of use to us is expressed in the following Lemma:

Lemma 6 Let $\mathbf{X}, \mathbf{Y} \in IN^K$ be random variable's such that:

$$\mathbf{X} \leq_{scx} \mathbf{Y} \quad (\mathbf{X} \leq_{iscx} \mathbf{Y})$$

There exist two random variable's $\tilde{\mathbf{X}}$ and $\tilde{\mathbf{Y}}$ such that:

$$\tilde{\mathbf{X}} =_{st} \mathbf{X}, \quad \tilde{\mathbf{Y}} =_{st} \mathbf{Y} \quad \text{and} \quad \tilde{\mathbf{X}} \prec \tilde{\mathbf{Y}} \quad (\tilde{\mathbf{X}} \prec_w \tilde{\mathbf{Y}})$$

almost surely.

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

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

Proof. Couple the initial queue lengths so that $\mathbf{N}^l(0) \prec \mathbf{N}(0)$. 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 μ . 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 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 $\{t_i\}_{n=0}^{\infty}$ be the sequence of times at which arrivals or service events occur ($t_0 \equiv 0$). We will establish the relation $\mathbf{N}^l(t) \prec_w \mathbf{N}(t)$ by induction on the event times. Clearly, if $\mathbf{N}^l(t_i) \prec_w \mathbf{N}(t_i)$ then $\mathbf{N}^l(t) \prec_w \mathbf{N}(t)$, $t_i \leq t < t_{i+1}$, $i = 0, 1, \dots$

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

Inductive step. Assume that $\mathbf{N}^l(t) \prec_w \mathbf{N}(t)$ for $t < t_i$. We will establish it now for $t = t_i$. There are two cases depending on whether the event is an arrival or a service event.

Arrival. $\mathbf{N}^l(t_i) \prec_w \mathbf{N}(t_i)$ follows because arrivals are always to the smallest queue, so property 1 of Lemma 5 can be applied.

Service event. There are two cases depending on whether $N^l(t_i^-) < C_f$. In either case, the result follows from an application of property 2 of Lemma 5.

This completes the inductive step and thus we have $\mathbf{N}^l(t) \prec_w \mathbf{N}(t)$, $t \geq 0$. By the definition of weak majorization (\prec_w), this implies that $f(\mathbf{N}^l(t)) \leq f(\mathbf{N}(t))$ for any increasing Shur-convex function $f(t)$. Removing the conditioning on the arrival times and service times, we have:

$$N^l(t) \leq_{iscx} N(t) \forall t > 0$$

■

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

Proof. This follows from the preceding theorem and the fact that $\phi(\mathbf{X}) = \sum_{k=1}^K X_k$ is an increasing Schur-convex function. ■

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

$$\mathcal{S}_0 = \text{set of states with } 0 \leq N_j \leq C \text{ and } |N_i - N_j| \leq d, \forall i, j$$

- $s_o(C_o)$ = this is the only state in \mathcal{S}_0 that has a positive transition rate into it from states outside \mathcal{S}_0 .
- $\tilde{\pi}(s_o(C_o))$ = conditional probability of state $s_o(C_o)$, given that the system is in \mathcal{S}_0 .
- s_0 = aggregate state which represents all states of \mathcal{S}_0 .
- s_i = state which represents the system having $C_f + i$ customers, where $i = 1, 2, \dots$
- $\pi^*(s_i)$ = steady state probability of state s_i .
- $\tilde{N}(\mathcal{S}_0)$ = mean number of customers given that the system is in \mathcal{S}_0 .

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

$$\begin{aligned}
N_l &= \tilde{N}(\mathcal{S}_0)\pi^*(s_0) + \sum_{i=1}^{\infty} [C_f + i] \pi^*(s_i) \\
&= \tilde{N}(\mathcal{S}_0)\pi^*(s_0) + C_f(1 - \pi^*(s_0)) + \lambda_0 \frac{K\mu}{(K\mu - \lambda)^2} \pi^*(s_0) \quad \text{and} \quad (22)
\end{aligned}$$

$$R_l = \frac{1}{\lambda} \left[\tilde{N}(\mathcal{S}_0)\pi^*(s_0) + C_f(1 - \pi^*(s_0)) + \lambda_0 \frac{K\mu}{(K\mu - \lambda)^2} \pi^*(s_0) \right] \quad (23)$$

where:

$$\lambda_0 = \tilde{\pi}(s_o(C_o))\lambda \quad (24)$$

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

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

6 State Space Reduction by Lumpability

In the previous section, we discussed a methodology for constructing an upper bound model M^u and a lower bound model M^l . The computational costs in solving the models can be broken down into:

1. obtaining the conditional state probabilities in \mathcal{S}_0 and \mathcal{S}_1 ,
2. obtaining the steady state probabilities of the aggregated process and,
3. obtaining the performance measure, e.g., expected response time or expected number of customers.

The larger the state space cardinality of \mathcal{S}_i , the more accurate are the results obtained. In this section, we discuss how we can reduce the state space of \mathcal{S}_i by lumping *similar* states.

Kemeny and Snell [16] studied the conditions under which an aggregated process is still Markovian. The condition for a Markov process to be lumpable with respect to a partition $\{\mathcal{P}_0 \cup \mathcal{P}_1 \cup \dots\}$, where $\mathcal{P}_i \cap \mathcal{P}_j = \emptyset$, for $i \neq j$, is that for every pair of sets \mathcal{P}_i and \mathcal{P}_j , r_{k,\mathcal{P}_j} has the same value for every state $k \in \mathcal{P}_i$ where

$$r_{k,\mathcal{P}_j} = \sum_{l \in \mathcal{P}_j} q_{k,l} \quad \text{for } k \in \mathcal{P}_i \quad (27)$$

We can apply this notion to our minimum expected delay routing problem.

Let J be the number of distinct types of servers in the model where two servers are of the same type if and only if they have the same service rate. For any state s define the following mapping:

$$f : s \rightarrow \{l_i | i = 1, 2, \dots, J\} \quad (28)$$

where:

l_i = a set of tuples $(\alpha_{ij}, \beta_{ij})$

α_{ij} = a queue length for a server of type i that appears in state s

β_{ij} = the number of servers of type i that have queue length α_{ij} in state s

We define a partition of the state space \mathcal{S}_u (\mathcal{S}_l) by specifying that $s_1, s_2 \in \mathcal{S}_u$ (\mathcal{S}_l) are in the same partition if and only if $f(s_1) = f(s_2)$.

For example, assume that we have a four server system with $\mu_1 = \mu_2 = 4$, $\mu_3 = 3$ and $\mu_4 = 2$. There are three distinct types of servers, so $J = 3$. We can group states such as $s_1 = (3, 4, 2, 1)$ and $s_2 = (4, 3, 2, 1)$ into the same partition since the $l_i, i = 1, 2, 3$, for both states are:

$$l_1 = \{(4, 1), (3, 1)\}; l_2 = \{(2, 1)\}; l_3 = \{(1, 1)\}$$

It is not difficult to see that the condition for lumpability is satisfied and we can often significantly reduce the state space of the model.

7 Numerical Examples

In this section, we present two examples in order to illustrate the bounding algorithm.

The system we consider in our first example consists of eight homogeneous servers. To vary the system utilization ρ from 0.1 to 0.9, we fix the input arrival rate at 8.0 and

vary the service rates for all servers. For $\rho = 0.1$ to 0.5 , we set $d = 4, C_i = 8$, for $\rho = 0.6$ to 0.7 , we set $d = 4, C_i = 10$, and for $\rho = 0.8$ to 0.9 , we set $d = 5, C_i = 12$. Table 1 illustrates the upper and lower bound on mean response time as a function of system utilization. Percentage error⁴ is defined to be $\frac{R_u - R_l}{R_u + R_l} \times 100\%$. Note that the bounds are very tight.

The second system we consider has eight heterogeneous servers with $\mu_1 = \mu_2 = \mu_3 = \mu_4 = 3, \mu_5 = \mu_6 = 2$, and $\mu_7 = \mu_8 = 1$. To vary the system utilization from 0.1 to 0.9 , we fix the service rates for all servers and vary the input arrival rate. For $\rho = 0.1$ to $\rho = 0.5$, we set $d = 4, \vec{C} = (6, 6, 6, 6, 4, 4, 2, 2)$, for $\rho = 0.6$ to $\rho = 0.8$, we set $d = 5, \vec{C} = (9, 9, 9, 9, 6, 6, 3, 3)$, and for $\rho = 0.8$ to $\rho = 0.9$, we set $d = 6, \vec{C} = (12, 12, 12, 12, 8, 8, 4, 4)$. Table 2 illustrates the upper and lower bound on the mean response time and the tightness of the bounds.

To illustrate the tradeoff between computational cost and accuracy of the bounds consider the homogeneous queueing system in the first example. By fixing the system utilization at 0.9 and increasing the number of states generated, we see the improvement of the bounds on the mean response time. The results are illustrated in Table 3.

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 1: Homogeneous Servers System.

⁴If the spread in the bounds is less than $< 10^{-6}$, we leave the entries for the spread of the bounds and the percentage error blank.

System Utilization	Response Time Lower Bound	Response Time Upper Bound	Spread of Bounds	Percentage Error
0.1	0.3337116	0.3337280	0.0000164	0.00247 %
0.2	0.3380923	0.3381092	0.0000169	0.00250 %
0.3	0.3486201	0.3486987	0.0000786	0.01127 %
0.4	0.3672009	0.3672981	0.0000972	0.01320 %
0.5	0.3978023	0.3979084	0.0001061	0.01333 %
0.6	0.4471098	0.4472873	0.0001775	0.01985 %
0.7	0.5239872	0.5249875	0.0010003	0.09536 %
0.8	0.6506371	0.6609821	0.0103450	0.78872 %
0.9	0.9788093	1.0237158	0.0449065	2.242494%

Table 2: Heterogeneous Servers System.

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 %

Table 3: Computational Cost vs. Accuracy.

8 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 a priori determine d and C_i in order 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.

References

- [1] I.J.B.F. Adan, J. Wessels, W.H.M. Zijm. *Analysis of the Symmetric Shortest Queue Problem*. Stochastic Models, Vol.6, 691-713, 1990.
- [2] I.J.B.F. Adan, J. Wessels, W.H.M. Zijm. *Analysis of the Asymmetric Shortest Queue Problem*. Queueing Systems, Vol. 8, 1-58, 1991.
- [3] Alberto Avritzer. *Dynamic Load Sharing Algorithms in Asymmetric Distributed Systems*. UCLA Computer Science Technical Report, CSD-900023.
- [4] J.P.C. Blanc. *A Note on Waiting Times in Systems with Queues in Parallel*. Journal of Applied Probability, Vol.24, 540-546, 1987.
- [5] B.W. Conolly. *The Autostrada Queueing Problem*. Journal of Applied Probability. Vol. 21, 394-403, 1984.
- [6] J.W. Cohen, O.J. Boxma. *Boundary Value Problems in Queueing System Analysis*. North Holland, 1983.
- [7] P. J. Courtois. *Decomposability — queueing and computer system applications*. Academic Press, New York, 1977.
- [8] P. J. Courtois, P. Semal. *Computable Bounds for Conditional Steady-State Probabilities in Large Markov Chains and Queueing Models*. IEEE JSAC, Vol 4, number 6, September, 1986.
- [9] Nico M. van Dijk. *The Importance of Bias-terms for Error Bounds and Comparison Results*. First International Conference on Numerical Solution of the Markov Chains, January 1990.
- [10] A. Ephremides, P. Varaiya, J. Walrand. *A Simple Dynamic Routing Problem*. IEEE Trans. on Auto. Control, Vol. 25, 1980.
- [11] L. Flatto, H.P McKean. *Two Queues in Parallel*. Communication on Pure and Applied Mathematics, Vol. 30, 255-263, 1977.
- [12] S. Halfin. *The Shortest Queue Problem*. Journal of Applied Probability, Vol. 22, 865-878, 1985.
- [13] W.K. Grassmann. *Transient and Steady State Results for Two Parallel Queues*. Omega, 8, 105-112, 1980.
- [14] J.F.C Kingman. *Two Similar Queues in Parallel*. Annals of Mathematical Statistics, Vol 32, 1314-1323, 1961.

- [15] C. Knessl, B.J. Matkowsky, Z. Schuss, C. Tier. *Two Parallel Queues with Dynamic Routing*. IEEE Trans. on Communications, Vol. 34, 1170-1175, 1986.
- [16] J.G. Kemeny and J.L. Snell. *Finite Markov Chains*. Van Nostrand Company, 1960.
- [17] L. Kleinrock. *Queueing Systems: Volume I: Theory* Wiley-Interscience Publication . New York. 1975.
- [18] J.D.C Little. *A Proof of the Queueing Formula $L = \lambda W$* , Operations Research, Vol 9, 383-387, 1967.
- [19] Hwa-Chun and C.S. Raghavendra. *An Analysis of the Join the Shortest Queue Policy* Electrical Engineering Technical Report, University of Southern California, 1991.
- [20] J.C.S. Lui, R.R. Muntz. *Algorithmic Approach to Bounding the Response Time of a Minimum Expected Delay Routing System*, Proc. 1992 ACM SIGMETRICS/Performance'92 Conference, 140-152.
- [21] Marshall and Olkin. *Inequalities: Theory of Majorization and Applications*, Academic Press, New York, 1979.
- [22] R.D. Nelson, T.K. Philips. *An Approximation to the Response Time for Shortest Queue Routing*. ACM SIGMETRICS Vol 17, No 1. 1989, pp 181-189.
- [23] R.D. Nelson, T.K. Philips. *An Approximation for the Mean Response Time for Shortest Queue Routing with General Interarrival and Service Times*. IBM T.J. Watson Research Lab, Technical Report RC15429, 1990.
- [24] M.F. Neuts. *Matrix-Geometric Solutions in Stochastic Models - An Algorithmic Approach*, The Johns Hopkins University Press, Baltimore, 1981.
- [25] B.M. Rao, M.J.M. Posner. *Algorithmic and Approximate Analysis of the Shorter Queue Model*. Naval Research Logistics, Vol.34, 381-398, 1987.
- [26] D. Towsley, S. Chen. *Design and Modeling Policies for Two Server Fork/Join Queueing Systems*. University of Mass. COINS Technical Report:91-39.
- [27] D. Towsley, P. Sparaggis, C. Cassandras, "Stochastic ordering properties and optimal routing control for a class of finite capacity queueing systems", *IEEE Transactions on Automatic Control*, **37**, 9 1446-1451, September 1992.
- [28] Y.T. Wang and Robert J. T. Morris. *Load Sharing in Distributed Systems*. IEEE Transactions on Computers, Vol. C-34, No 3, 204-217. March 1985.
- [29] W. Winston. *Optimality of the Shortest Line Discipline*, Journal of Applied Probability. Vol 15, 181-189, 1977.

- [30] Y. Zhao and W.K. Grassmann. *The Shortest Queue Model with Jockeying*. Naval Research Logistics. Vol. 37. 773-787, 1990.
- [31] Y. Zhao and W.K. Grassmann. *A Numerically Stable Algorithm for Two Server Queue Models*. Queueing Systems, Vol 8, 59-80, 1991.