

Multi-server queueing systems with multiple priority classes

Mor Harchol-Balter* Takayuki Osogami† Alan Scheller-Wolf‡ Adam Wierman§

Abstract

We present the first near-exact analysis of an M/PH/ k queue with $m > 2$ preemptive-resume priority classes. Our analysis introduces a new technique, which we refer to as Recursive Dimensionality Reduction (RDR). The key idea in RDR is that the m -dimensionally infinite Markov chain, representing the m class state space, is recursively reduced to a 1-dimensionally infinite Markov chain, that is easily and quickly solved. RDR involves no truncation and results in only small inaccuracy when compared with simulation, for a wide range of loads and variability in the job size distribution.

Our analytic results are then used to derive insights on how multi-server systems with prioritization compare with their single server counterparts with respect to response time. Multi-server systems are also compared with single server systems with respect to the effect of different prioritization schemes – “smart” prioritization (giving priority to the smaller jobs) versus “stupid” prioritization (giving priority to the larger jobs). We also study the effect of approximating m class performance by collapsing the m classes into just two classes.

Keywords: M/GI/ k , M/PH/ k , multi-server queue, priority queue, matrix analytic methods, busy periods, multi-class queue, preemptive priority.

1 Introduction

Much of queueing theory is devoted to analyzing priority queues, where jobs (customers) are labeled and served in accordance with a priority scheme: high-priority (H) jobs preempt medium-priority (M) jobs, which in turn preempt low-priority (L) jobs in the queue. Priority queueing comes up in many applications. Sometimes the priority of a job is determined by the job’s owner via a Service Level Agreement (SLA), whereby certain customers have chosen to pay more so as to get high-priority access to some high-demand

*Supported by NSF Career Grant CCR-0133077, NSF Theory CCR-0311383, NSF ITR CCR-0313148, and IBM Corporation via Pittsburgh Digital Greenhouse Grant 2003.

Email: harchol@cs.cmu.edu; Web: <http://www.cs.cmu.edu/~harchol>; Address: Department of Computer Science, Carnegie Mellon University, 5000 Forbes Ave., Pittsburgh, PA, 15213; Phone: 412-268-7893; Fax: 412-268-5576.

†Email: osogami@cs.cmu.edu; Address: Department of Computer Science, Carnegie Mellon University, 5000 Forbes Ave., Pittsburgh, PA 15213; Phone: 412-268-3621.

‡Email: awolf@andrew.cmu.edu; Address: Tepper School of Business, Carnegie Mellon University, 5000 Forbes Ave., Pittsburgh, PA, 15213; Phone: 412-268-5066.

§Email: acw@cs.cmu.edu; Address: Department of Computer Science, Carnegie Mellon University, 5000 Forbes Ave., Pittsburgh, PA, 15213; Phone: 412-877-9455.

resource. Other times, the priority of a job is artificially created, so as to maximize a company's profit or increase system utilization. For example, an online store may choose to give high-priority to the requests of big spenders, so that those customers are less likely to go elsewhere, see [17].

Analyzing the mean response time (and higher moments of response time) for different classes of jobs is clearly an important problem.¹ While this problem has been well understood in the case of a single-server M/GI/1 queue since the 1950's [5], the problem becomes much more difficult when considered in the context of a multi-server M/GI/ k system, and even for an M/M/ k system when jobs have different completion rates. This is unfortunate since such multi-server systems are prevalent in many applications where prioritization is used, e.g., web server farms and super-computing centers.

The reason that priority queueing is so difficult to analyze in a multi-server setting is that jobs of different priorities may be in service (at different servers) at the same time, thus the Markov chain representation of the multi-class, multi-server queue appears to require tracking the number of jobs of each class. Hence one needs a Markov chain which is infinite in m dimensions, where m is the number of priority classes. While the analysis of a 1-dimensionally infinite Markov chain is easy, the analysis of an m -dimensionally infinite Markov chain ($m > 1$) is largely intractable.

Prior work

The number of papers analyzing multi-server priority queues is vast, however almost all are restricted to only *two priority classes*. Of those restricted to two priority classes, all assume *exponential service times*. The only papers *not* restricted to two priority classes are coarse approximations based on assuming that the multi-server behavior is related to that of a single server system [2] or approximations based on aggregating the many priority classes into two classes [19, 22].

Two priority classes

We start by describing the papers restricted to *two priority classes* and *exponentially-distributed service demands*. Techniques for analyzing the M/M/ k dual priority system can be organized into four types on which we elaborate below: (i) approximations via aggregation or truncation; (ii) matrix analytic methods; (iii) generating function methods; (iv) special cases where the priority classes have the same mean. Unless otherwise mentioned, preemptive-resume priorities should be assumed.

Nearly all analysis of dual priority M/M/ k systems involves the use of Markov chains, which with two priority classes grows infinitely in two dimensions (one dimension for each priority class). In order to overcome this, researchers have simplified the chain in various ways. Kao and Narayanan truncate the chain by either limiting the number of high priority jobs [12], or the number of low priority jobs [10]. Nishida

¹We will use the term *response time* throughout the paper to denote the time from when a job arrives until it is completed. We will also occasionally talk about the *delay* (wasted time), which we define as the job's response time minus its service requirement.

aggregates states, yielding an often rough approximation [22]. Kapadia, Kazmi and Mitchell explicitly model a finite queue system [13]. Unfortunately, aggregation or truncation is unable, in general, to capture the system performance as the traffic intensity grows large.

Although, in theory, the matrix analytic method can be used to directly analyze a 2D-infinite Markov chain (see for example [3]), the matrix analytic method is much simpler and more computationally efficient when it is applied to a 1D-infinite Markov chain. Therefore, most papers that use the matrix analytic method to analyze systems involving 2D-infinite Markov chains first reduce the 2D-infinite chain to a 1D-infinite chain by, for example, truncating the state space by placing an upper bound on the number of jobs [12, 10, 16, 21]. Miller [18] and Ngo and Lee [21] partition the state space into blocks and then “super-blocks,” according to the number of high priority customers in queue. This partitioning is quite complex and is unlikely to be generalizable to non-exponential job sizes. In addition, [18] experiences numerical instability issues when $\rho > 0.8$.

A third stream of research capitalizes on the exponential job sizes by explicitly writing out the balance equations and then finding roots via generating functions. In general these yield complicated mathematical expressions susceptible to numerical instabilities at higher loads. See King and Mitrani [19]; Gail, Hantler, and Taylor [8, 9]; Feng, Kowada, and Adachi [7]; and Kao and Wilson [11].

Finally there are papers that consider the special case where the multiple priority classes all having the same mean. These include Davis [6], Kella and Yechiali [14] (for non-preemptive priorities), and Buzen and Bondi [4].

The only work dealing with non-exponential service times is contained in a pair of papers, not yet published, by Sleptchenko et. al. [25, 26]. Both papers consider a two-priority, multi-server system where within each priority there may be a number of different classes, each with its own different exponential job size distribution. This is equivalent to assuming a *hyper-exponential job size distribution* for each of the *two priority classes*. The problem is solved via a combination of generating functions and the matrix analytic method. In theory, their technique may be generalizable to PH distributions, though they evaluate only hyper-exponential distributions due to the increased complexity necessary when using more general PH distributions.

More than two priority classes

For the case of *more than two priority classes*, there are only coarse approximations. The Bondi-Buzen (BB) approximation [2] is beautiful in its simplicity and usability. It is based on an intuitive observation that the “improvement” of priority scheduling over FCFS scheduling under k servers is similar to that for the case of one server with equal total capacity:

$$\frac{E[D^{M/GI/k/prio}]}{E[D^{M/GI/k/FCFS}]} \approx \frac{E[D^{M/GI/1/prio}]}{E[D^{M/GI/1/FCFS}]} = \text{scaling factor.} \quad (1)$$

Here $E[D^{M/GI/k/prio}]$ is the overall mean delay under priority scheduling with k servers of speed $1/k$, and $E[D^{M/GI/k/FCFS}]$ is defined similarly for FCFS. This relation is exact when job sizes are exponential with the same rate for all classes; however what happens when this is not the case has never been established.

The other approximation (which we denote by MK-N) which allows for *more than two priority classes* and *exponential* job sizes is due to Mitrani and King [19], and also used by Nishida [22] to extend the latter author’s analysis of two priority classes to $m > 2$ priority classes. The MK-N approximation analyzes the mean delay of the lowest priority class in an M/M/ k queue with $m \geq 2$ priority classes by *aggregating all the higher priority classes*. Thus, instead of aggregating all jobs into one class, as BB does, MK-N aggregates into two classes. The job size distribution of the aggregated class is then approximated with an exponential distribution by matching the first moment of the distribution.

Contributions of this paper

In Section 2, we introduce a new analytical approach that provides the first near-exact analysis of the M/PH/ k queue with $m \geq 2$ preemptive-resume priority classes. Our approach, which we refer to as Recursive Dimensionality Reduction (RDR), is very different from the prior approaches described above. RDR allows us to recursively reduce the m -dimensionally infinite state space, created by tracking the m priority classes, to a 1-dimensionally infinite state space, which is analytically tractable. The dimensionality reduction is done without any truncation; rather, we reduce dimensionality by introducing “busy period transitions” within our Markov chain, for various types of busy periods created by different job classes. The only approximation in the RDR method stems from the fact that we need to approximate these busy periods using Markovian (phase-type) PH distributions. We find that matching three moments of these busy periods is usually possible using a 2-phase Coxian distribution, and provides sufficient accuracy, within one or two percent of simulation, for all our experiments (our experiments span load ranging from $\rho = 0.05$ to $\rho = 0.95$ and job size variability ranging from $C^2 = 0$ to $C^2 = 128$). The accuracy of the RDR method can be increased arbitrarily by better approximating the busy periods.

In theory RDR can handle systems with any number of servers, any number of priority classes, and PH service times. In addition, RDR is quite efficient; for all the scenarios explored in this paper, the computation time under RDR is less than a few seconds. However, the complexity of the RDR method does increase with both the number of servers k and the number of classes m . Because RDR becomes less practical under high m and k , we develop a much simpler, but only slightly less accurate, approximation RDR-A (see Section 2.5). RDR-A simplifies calculations by approximating an m priority system with a two priority system, which is then solved using RDR.

In Section 3 we present results from both RDR and RDR-A for per-class mean response time for an M/PH/ k queue with multiple priority classes, and also discuss the computation of higher moments of response time. In Section 4, we use these results to obtain many interesting *insights* about priority queueing. First, in Section 4.1 we compare the performance of priority queueing in a multi-server system with k servers

each of speed $1/k$ versus a single server of speed 1. We find that the effect of priorities in a single server system can be very different than in a multi-server system of equal capacity. (A non-surprising consequence of this result is that the BB approximation, which relies on relating a multi-server system to a single server system, can exhibit large errors.) Next, in Section 4.2, we study the effect of priority policies that favor short jobs (“smart prioritization”) versus priority policies that favor long jobs (“stupid prioritization”) under systems with different numbers of servers. Understanding the effect of “smart” prioritization is important because many common scheduling policies are designed to give priority to short jobs. Lastly, in Section 4.3, we ask how effective class aggregation (aggregating $m > 2$ priority classes into just 2 priority classes) is as an approximation for dealing with systems having than two priority classes. We evaluate several types of class aggregation including that proposed by the MK-N approximation and that used in RDR-A to show when class aggregation serves as a reasonable approximation.

2 RDR analysis of M/PH/k with m priority classes

In this section we describe the RDR technique. We divide our explanation into three parts. As an introduction, in Section 2.1, we deal only with the simplest case of $m = 2$ priority classes and exponential job sizes, which we solve using the techniques in [24]. We then move to the difficult case of $m > 2$ priority classes, but still exponential service times, in Section 2.2. Here the techniques from [24] do not apply, so we introduce Recursive Dimensionality Reduction (RDR). The RDR approach uses the analysis of the $m - 1$ priority classes to analyze the m -th priority class. This is a non-trivial procedure for $m > 2$ since it involves evaluating many complex passage times (busy periods) in the chain representing the $m - 1$ priority classes, as these passage times now form transitions within the chain representing m priority classes. Finally in Section 2.3, we show how RDR can be applied to the most general case of $m > 2$ priority classes, with PH service times.

All the analysis up to through Section 2.3 deals with how to derive mean per-class response times. In Section 2.4 we illustrate how the RDR method can be extended to obtain *variance of response time* for each class. Finally, in Section 2.5, we introduce RDR-A, which is an approximation of RDR, allowing very fast (< 1 second) evaluation of high numbers of priority classes and servers, with small ($< 5\%$) error.

2.1 Simplest case: Two priority classes, exponential job sizes

Consider the simplest case of two servers and two priority classes, high (H) and low (L), with exponentially distributed sizes with rates μ_H and μ_L respectively. Figure 1(a) illustrates a Markov chain of this system, whose states track the number of high priority and low priority jobs; hence this chain grows infinitely in two dimensions. Observe that high priority jobs simply see an M/M/2 queue, and thus their mean response time is well-known. Low priority jobs, however, have access to either an M/M/2, M/M/1, or no server at all,

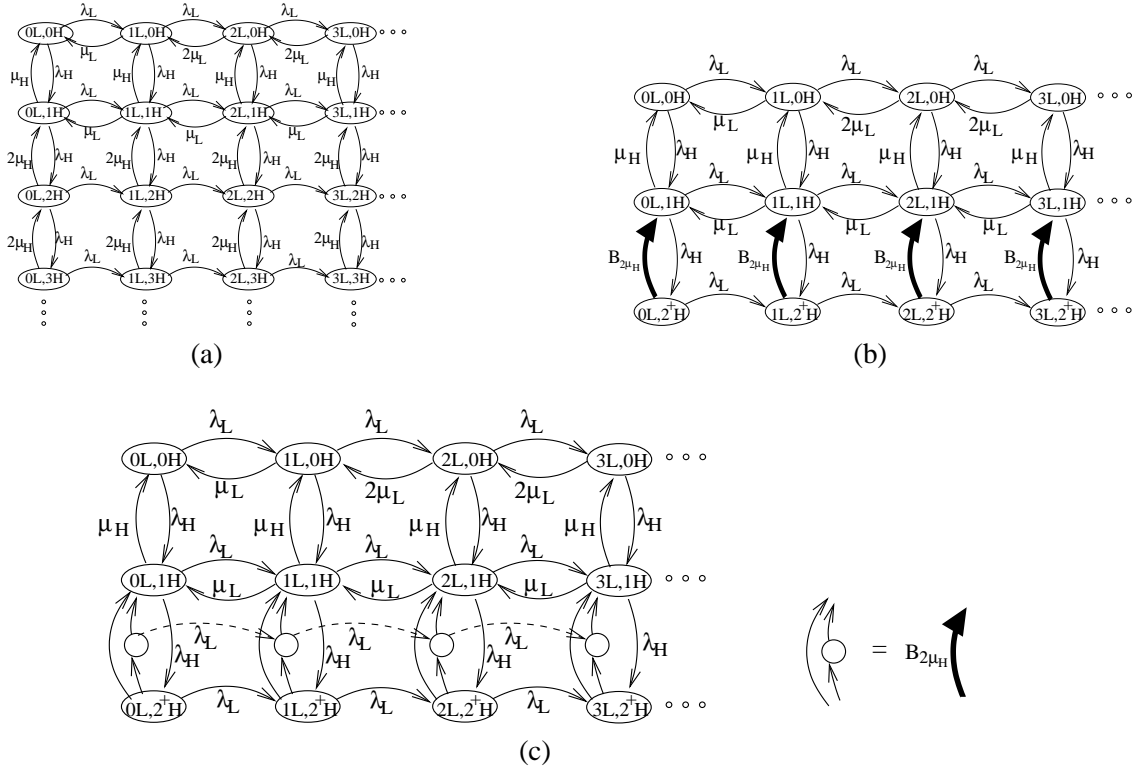


Figure 1: (a) Markov chain for a system with two servers and two priority classes where all jobs have exponential sizes. This Markov chain is infinite in two dimensions. Via the Dimensionality Reduction technique, we arrive at the chain in (b), which uses busy period transitions, and is only infinite in one dimension. In (b), the busy period is represented by a single transition. In (c), the busy period is represented by a two phase PH distribution (with Coxian representation), yielding a one-dimensionally infinite Markov chain.

depending on the number of high priority jobs. Thus their mean response time is more complicated, and this is where we focus our efforts.

Figure 1(b) illustrates the reduction of the 2D-infinite Markov chain to a 1D-infinite chain. The 1D-infinite chain tracks the number of low priority jobs exactly. For the high priority jobs, the 1D-infinite chain only differentiates between zero, one, and two-or-more high priority jobs. As soon as there are two-or-more high priority jobs, a *high priority busy period* is started. During the high priority busy period, the system only services high priority jobs, until the number of high priority jobs drops to one.² The length of time spent in this high priority busy period is exactly an M/M/1 busy period where the service rate is $2\mu_H$. We denote the duration of this busy period by the transition labeled $B_{2\mu_H}$.

The busy period $B_{2\mu_H}$ is *not* exponentially-distributed. Hence it is not clear how it should fit into a Markov model. We use a PH distribution (specifically a Coxian distribution) to match the first three

²Throughout the paper a “higher priority busy period” is defined as the time from when the system has k higher priority jobs until there are only $k - 1$ higher priority jobs.

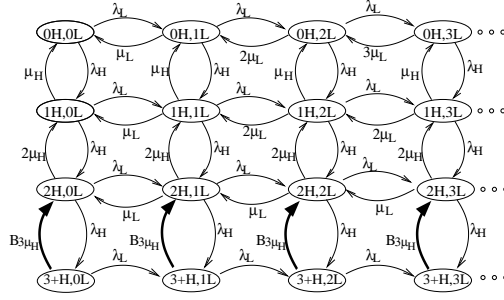


Figure 2: This chain illustrates the case of two priority classes and three servers. The busy period transitions are replaced by a Coxian phase-type distribution matching three moments of the busy period duration, as shown in Figure 1.

moments of the distribution of $B_{2\mu_H}$. Parameters of the PH distribution, whose first three moments match those of $B_{2\mu_H}$, are calculated via the closed form solutions provided in [23].

Figure 1(c) illustrates the same 1D-infinite chain as in Figure 1(b), except that the busy period transition is now replaced by a two phase Coxian distribution. The limiting probabilities in this 1D-infinite chain can be analyzed using the matrix analytic method, which yields the mean response time for low-priority jobs via Little's law. The only inaccuracy in the above approach is that only three moments of the high-priority busy period have been matched. We will see later that this suffices to obtain very high accuracy across a wide range of load and job size distributions.

Figure 2 shows the generalization to a three server system. We simply add one row to the chain shown in Figure 1, and now differentiate between 0, 1, 2, or 3-or-more high priority jobs. This can be easily extended to the case of $k > 3$ servers.

2.2 Harder case: m priority classes, exponential job sizes

We now turn to the more difficult case of $m > 2$ priority classes. We illustrate this for the case of two servers and three priority classes: high-priority (H), medium-priority (M), and low-priority (L). The mean response time for class H jobs and that for class M jobs are easy to compute. Class H jobs simply see an M/M/2 queue. Class M jobs see the same system that the low-priority jobs see in an M/M/2 queue having two priority classes. Replacing the L's by M's in the chain in Figure 1 yields the mean response time for the M class jobs.

The analysis of the class L jobs is the difficult part. The obvious approach would be to aggregate the H and M jobs into a single class, so that we have a 2-class system (H-M versus L jobs). Then we could apply the technique of the previous section, tracking exactly the number of low-priority jobs and maintaining limited state information on the H-M class. This is the approach that we follow in Section 2.5 in deriving our RDR-A approximation. However, this approach is imprecise because the duration of the busy periods in the H-M class depends on whether the busy period was started by 2H jobs, 1H and 1M job, or 2M jobs

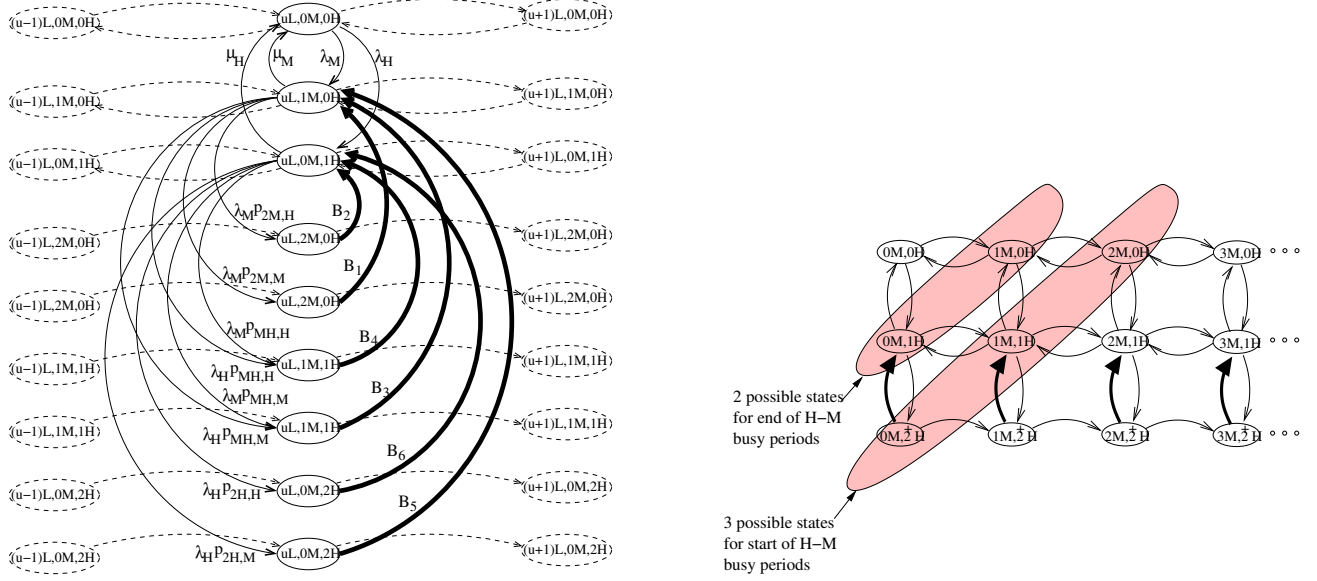


Figure 3: (Left) Portion of the 1D-infinite chain used to compute mean response time for low-priority jobs in the case of three priority classes and two servers, and all exponential service times. (Right) Chain used to compute moments of the durations of the six busy period transitions.

in service. By ignoring the priorities among H's and M's, we are ignoring the fact that some types of busy periods are more likely than others. Even given the information on who starts the busy period, this still does not suffice to determine its duration, because the duration is also affected by the prioritization within the aggregated H-M class.

Thus a precise response time analysis of class L requires maintaining more information. As before we want to exactly track the number of class L jobs. Given that there are two servers, we need to further differentiate between whether there are zero H and M jobs, one H or M job, or two or more H and M jobs. Whenever there are two or more H and M jobs, we are in an H-M busy period. For an M/M/2 with three priority classes, there are *six types of busy periods* possible, depending on the state at the start the busy period – $(0M, 2H)$, $(1M, 1H)$, or $(2M, 0H)$ – and the state in which the busy period ends – $(0M, 1H)$ or $(1M, 0H)$. We derive the busy period duration by conditioning on who starts and ends the busy period.

Figure 3 (left) shows the level of the 1D infinite chain in which the number of class L jobs is u . In state (uL, vM, wH) , v class M jobs and w class H jobs are in the system if $v + w < 2$; otherwise, the state (uL, vM, wH) denotes that we are in a H-M busy period that was started by v class M jobs and w class H jobs. Observe that there are six types of busy periods depicted, labeled B_1, B_2, \dots, B_6 ; the busy period is determined by the state in which it was started and the state in which it ends. Notice, for example, that both states in the fourth and fifth row are labeled $(uL, 2M, 0H)$, meaning that the busy period is started by two class M jobs; but these two states differ in the class of the job that is left at the end of the H-M busy period: In state $(uL, 2M, 0H)$ of the fourth row, the busy period ends leaving a class H job, whereas in state of the

fifth row, the busy period ends leaving a class M job. (Recall that the class of job left at the end of a busy period is probabilistically determined at the beginning of the busy period and the duration of the busy period is conditioned on the class of the job left at the end.) Here $p_{2M,H}$, for example, denotes the probability that the busy period started by two class M jobs ends leaving one class H job, whereas $p_{MH,M}$ denotes the probability that the busy period started by one class M and one class H job, ends leaving one class M job. The remaining probabilities are defined similarly.

It remains to derive the moments of the duration of busy periods, B_1, B_2, \dots, B_6 , and probabilities $p_{2M,M}, p_{2M,H}, p_{MH,M}, p_{MH,H}, p_{2H,M}$, and $p_{2H,H}$ in Figure 3(left). The trick to deducing these quantities is to observe that the six busy periods correspond to passage times between two “diagonal” (shaded) levels in the chain shown in Figure 3(right), which is the 1D-infinite chain that we used to analyze the class M performance. Note that the 3 states in the right shaded diagonal level correspond to the three possible “start” states for busy periods, and the two states in the left shaded diagonal level correspond to the two possible “end” states for the busy periods. Thus, for example, busy period B_1 in Figure 3(left) corresponds to the first passage time from state (2M,0H) to state (1M,0H) in the chain in Figure 3(right). Likewise, probability $p_{2M,M}$ corresponds to the probability that, in Figure 3(right), state (1M,0H) is the first state of the two possible “end” states that is reached, given that the “start” state is (2M,0H). Inter-level passage times and ending probabilities within the chain in Figure 3(right) can be calculated using techniques developed by Neuts in [20]. We provide a precise description of this in Appendix A. Observe that these computations are greatly facilitated by the fact that our chains are infinite in only one dimension.

The extension of RDR to $m > 3$ classes is straightforward. For example, for the case of $m = 4$ classes, we proceed as in Figure 3, where we first create a chain that tracks exactly the number of jobs in class 4, and creates busy periods for the aggregation of the three higher priority classes. Then to derive the busy periods for the three higher priority classes, we make use of the existing chain for three classes shown in Figure 3(left), and compute the appropriate passage times for that chain. For an M/M/ k with m priority classes, there are $\binom{m+k-2}{k} \binom{m+k-3}{k-1}$ possible busy periods. That is, the number of different types of busy periods is polynomial in k if m is constant ($\Theta(k^m)$), and it is polynomial in m if k is constant ($\Theta(m^k)$); however, it is exponential in k and m if neither k nor m is constant.³

³**Remark:** We note that in practice the number of busy periods can be reduced further, so that an M/M/ k with m priority classes has $\binom{m+k-3}{k-1}^2$ busy periods of class 1 to class $m - 1$ jobs. An advantage of this reduction is that the number of busy periods of class 1 to class $m - 1$ jobs becomes independent of the type of PH distributions that is used to approximate the busy period of class 1 to class $m - 2$ jobs.

The trick to reducing the number of busy periods is illustrated by considering the example of the M/M/2 with three classes, shown in Figure 3. Here, by taking the mixture of the six busy periods, B_1, B_2, \dots, B_6 , we can approximate the H-M busy period by *four* PH distributions. These four distributions of the H-M busy period are differentiated by the state from which we *enter* the H-M busy period (either (0M,1H) or (1M,0H)) and by the state we return to after the H-M busy period (either (0M,1H) or (1M,0H)).

We illustrate how using B_1 and B_3 , we can obtain (the moments of) the distribution of the conditional H-M busy period when we enter the H-M busy period from (1M,0H) and return to (1M,0H). When we are at state (1M,0H), an arrival of an H job or an M job starts an H-M busy period. When the arrival is an H job (respectively, an M job), the H-M busy period ends with an M job with probability $p_{MH,M}$ (respectively, $p_{2M,M}$), and the conditional duration of the H-M busy period is B_3 (respectively, B_1). Since the arrival processes are Poisson, this conditional H-M busy period, which ends with an M job, starts at state (1M,0H) with rate

Practically speaking, the RDR approach is fast for a small number of servers and a small number of priority classes. In examples we ran with an M/M/2 and 10 priority classes, the RDR algorithm yielded mean response times within tens of seconds.

2.3 General case: Analysis of M/PH/ k with m priority classes

In this section, we explicitly describe how RDR can be applied to analyze the case of PH job size distributions. We describe RDR for the case of two servers ($k = 2$) and two priority classes ($m = 2$), high (H) and low (L), where the class H jobs have a particular 2-phase PH job size distribution with Coxian representation, shown in Figure 4(a).⁴ Generalization to higher k 's and higher m 's is straightforward by applying the recursive algorithm introduced in Section 2.2.

Analyzing the performance of class H is trivial, since high-priority jobs simply see the mean response time in an M/PH/2 queue, which can be analyzed via standard matrix analytic methods. To analyze the class L jobs, as before, we create a 1D-infinite Markov chain tracking the class L jobs, and use busy period transitions to represent needed information regarding the class H jobs.

Observe that under the 2-phase Coxian job sizes distribution, we will need *four* different types of busy periods for high priority jobs, depending on the phases of the two jobs starting the busy period (1 & 1, or, 1 & 2) and the phase of the job left at the end of the busy period (1 or 2). To derive the durations of these busy periods, we observe that the busy periods correspond to passage times from shaded level 3 to shaded level 2 in the Markov chain shown in Figure 4(b). Figure 4(b) describes the behavior of class H jobs, where states track the number of high priority jobs in the system and the phases of the jobs being processed. Namely, at state (0H) there are no high priority jobs in the system; at state (1H, i), there is one high priority job in phase i ; at state (n H, i, j) there are n high priority jobs in the system and the two jobs are being processed are in phase i and j , respectively (jobs in the queue are all in phase 1). The first passage times in Figure 4 are computed again using techniques in [20].

Figure 4(c) shows a level of the chain that tracks the number of low priority jobs, where the number of low priority jobs is u . The low priority job sizes are assumed to be exponentially distributed, but this can be easily generalized to PH distributions. In state (u L,0H), no high priority jobs are in system. An arrival of a high priority job in state (u L,0H) triggers a transition to state (u L,1H,1). In state (i L,1H, j), one high priority job in phase j is in the system for $j = 1, 2$. An arrival of a high priority job in state (u L,1H, j) triggers a transition to state (i L,2⁺H,1, j) for $j = 1, 2$. In state (i L,2⁺H,1, j), at least two high priority jobs are in the system, and the two jobs that started the busy period were in phase one and j , respectively, for $j = 1, 2$. The four types of busy periods are labeled as B_1 , B_2 , B_3 , and B_4 , and the duration of these busy periods is

$\lambda_M \cdot p_{2M,M} + \lambda_H \cdot p_{MH,M}$. Thus, the duration of this conditional H-M busy period is B_3 with probability $\frac{\lambda_M \cdot p_{2M,M}}{\lambda_M \cdot p_{2M,M} + \lambda_H \cdot p_{MH,M}}$ and B_1 otherwise. The other three H-M busy periods can be analyzed analogously.

⁴Under the Coxian job size distribution, a job starts in phase one where it is processed for a time exponentially distributed with rate $\mu_H^{(1)}$, and then either completes (with probability $q_H = 1 - p_H$) or moves to phase two (with probability p_H).

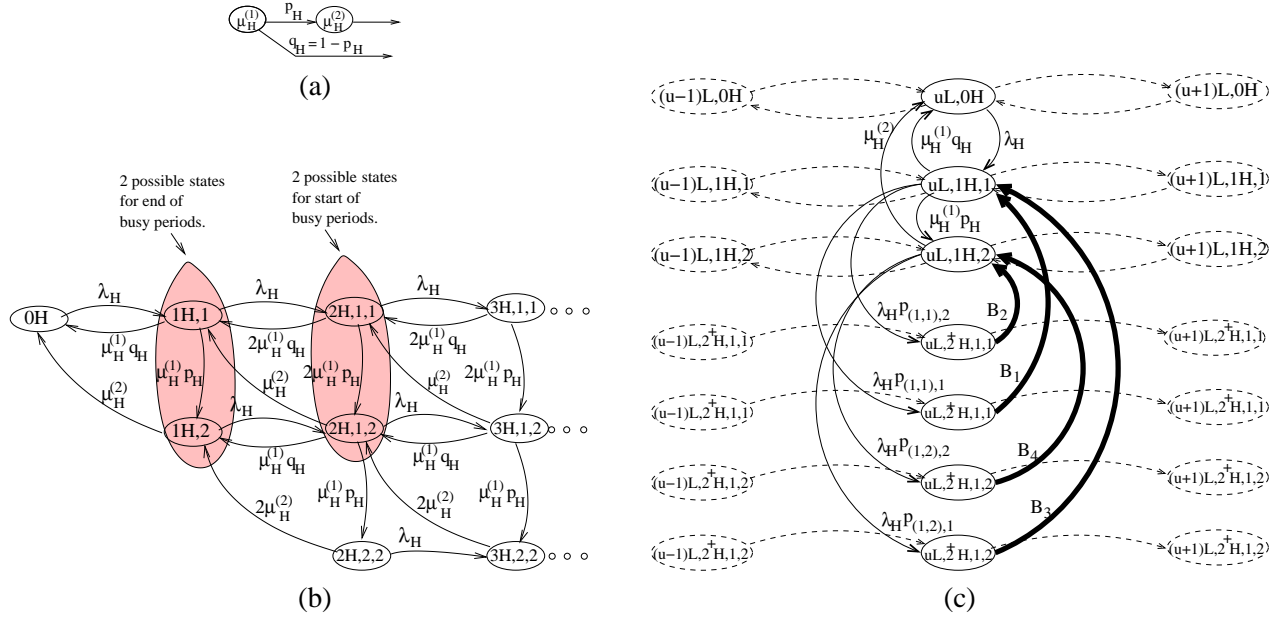


Figure 4: (a) A 2-phase PH distribution with Coxian representation. (b) Markov chain which will be used to compute the high-priority job busy periods, in the case where high-priority job size have a PH distribution with Coxian representation shown in (a). (c) Chain for a system with two servers and two priority classes where high priority jobs have Coxian service times.

approximated by PH distributions by matching the first three moments of the busy period distribution (note that the busy period cannot start with two jobs in phase two). Finally, $p_{(1,j),k}$ denotes the probability that a busy period started by two jobs in phases one and j , ends with a single job in phase k , for $j = 1, 2$, and $k = 1, 2$.

2.4 Computing variance of response time and higher moments

Throughout our discussion of RDR thus far, we have been concerned with computing the mean per-class response time. It turns out that computing higher moments of per-class response time is not much more difficult. Before we present our approach, we make two remarks. First, observe that it is trivial to derive all moments of the steady-state per-class *number of jobs* in the system, directly from the steady-state probabilities for the Markov chain, which we have already computed. Unfortunately, however, we cannot apply the beautiful generalization of Little's Law to higher moments (see [27, 1]) to immediately get the per-class higher moments of response time for free, since jobs do not necessarily leave our system in the order in which they arrive.

Below we sketch our approach for computing per-class variance in response time for the case of two servers, two priority classes (H and L), and exponential service times. We will refer to Figure 1(c) during

our discussion. For class H jobs, it is easy to compute the variance of their response time, via standard matrix analytic methods, since they are oblivious to class L jobs. Thus we will concentrate on class L jobs.

Consider the 1D-infinite Markov Chain shown in Figure 1(c) that tracks the number of class L jobs. Our approach thus far has been to compute the limiting probabilities, use those to derive the mean number of class L jobs in the system, and then apply Little's Law to yield mean response time for class L jobs. Now, we instead use the limiting probabilities to condition on what a class L arrival sees. Specifically, by PASTA (Poisson Arrivals See Time Averages) a class L arrival with probability $\pi_{(iL, jH)}$ will see state (iL, jH) when it arrives, and will cause the system state to change to $((i + 1)L, jH)$ at that moment.

To calculate the variance in response time seen by this class L arrival, we remove all the λ_L arcs from the Markov chain in Figure 1(c), so that there are no more class L arrivals. This enables us to view the response time for the class L arrival as the first passage time of this modified chain from state $((i + 1)L, jH)$ to the state where our class L arrival departs. The only complexity is in figuring out exactly in which state our class L arrival departs, where our class L arrival is the last class L job to enter the system.

The final class L arrival may depart the modified Markov chain the first time it hits $(1L, 0H)$ or $(1L, 1H)$, depending on the sample path the chain follows. We will compute the passage time to go from state $((i + 1)L, jH)$ to one of these states $\{ (1L, 0H) \text{ or } (1L, 1H) \}$. It is important to observe that the first time we hit a state with 1L, the state we hit cannot be $(1L, 2^+H)$, by virtue of the fact that the Markov chain doesn't have decreasing arcs in its bottom rows.

If $(1L, 1H)$ is the first state that we hit with 1L, then we know that we must have gotten there from $(2L, 1H)$, which means that the single L job remaining is in fact the last arrival. (We're assuming preemption is done "oldest first to be preempted"). Thus we need to now add on the passage time to go from $(1L, 1H)$ to $(0L, *)$ to get the full response time for the arrival.

If $(1L, 0H)$ is the first state that we hit with 1L, then we know that we got there from state $(2L, 0H)$. In this case, the remaining 1L is equally likely to be the last arrival or not. With probability half, the last arrival is already gone, in which case we add nothing to the response time. With probability half, this last arrival remains, in which case we now add on the passage time to go from $(1L, 0H)$ to $(0L, *)$ to get the full response time for the arrival.

Observe that computing the above passage times is straightforward, since all the λ_L arcs have been removed.

2.5 Introducing RDR-A

We have seen that the RDR method can become computationally intensive as the number of priority classes grows. This motivates us to introduce an approximation based on RDR called RDR-A. RDR-A applies to $m > 2$ priority classes and PH job size distributions.

The key idea behind RDR-A is that the RDR computation is far simpler when there are only two priority classes: H and L. In RDR-A, under m priority classes, we simply aggregate these classes into two priority

classes, where the $m - 1$ higher priority classes become the new aggregate H class and the m^{th} priority class becomes the L class. We define the H class to have a PH job size distribution that matches the first three moments of the aggregation of the $m - 1$ higher priority classes.

Observe that the RDR-A method is similar to the MK-N approximation. The difference is that in MK-N, both the H and L classes are exponentially-distributed. Thus under MK-N, the H class only matches the *first* moment of the aggregate $m - 1$ classes, whereas under RDR-A *three* moments are matched. The reason that we are able to match the first three moments, rather than just the first moment is that we have the RDR technique, which allows the analysis of multi-server priority queues with *PH* job size distributions, as described in Section 2.3.

3 Results and Validation

In this section we present results for per-class mean response times in M/M/ k and M/PH/ k queues with $m = 4$ priority classes, derived using RDR and RDR-A, respectively. To the best of our knowledge, these are the first such analytical results in the literature. We will validate our results against simulation, and show that their relative error is small. Furthermore the time required to generate our results is short, typically less than a second for each data point.

Figure 5 (top row) shows our results for per-class mean response times in an M/M/2 queue (left plot) and an M/PH/2 queue (right plot), both as a function of load ρ . The PH distribution used is a 2-phase PH distribution with squared coefficient of variation, $C^2 = 8$. All job classes have the same distribution, and the load is distributed evenly between the four classes. The left plot is derived using RDR and the right plot using RDR-A. Observe that the M/PH/2 queue (right plot) results in higher mean response time than the M/M/2 queue (left plot), as expected. In both cases the mean response time of the lower-priority classes dwarfs that of the higher-priority classes.

Figure 5 (bottom row) shows the relative per-class error for our results, when compared with simulation. Throughout the paper we always show error in *delay* (queueing time) rather than response time (sojourn time), since the error in delay is proportionally greater. We define relative error as

$$\text{error} = 100 \times \frac{(\text{mean delay by RDR or RDR-A}) - (\text{mean delay by simulation})}{(\text{mean delay by simulation})} \quad (\%).$$

We only show the error for classes 2, 3, and 4, since our analysis is virtually exact for class 1 (solved via matrix-analytic methods). We see that the relative error in the mean delay of RDR and RDR-A compared to simulation is within 2% for all classes and typically within 1%, for all ρ 's. We will see later that this error increases only slightly when we move to the case of priority classes with different means.

Figure 6 (left) again uses RDR-A to calculate per-class mean response time in the M/PH/2 queue with four classes, but this time as a function of C^2 , the squared coefficient of variation of the job size distribution.

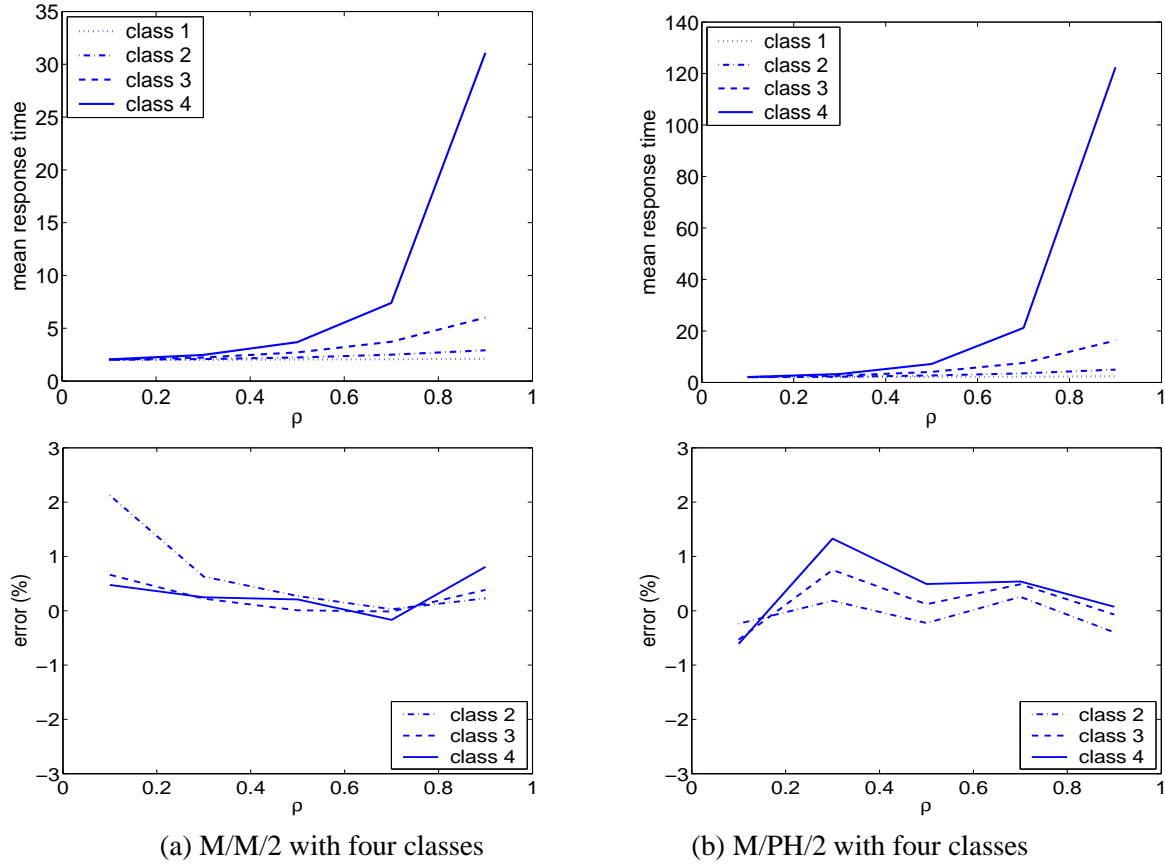


Figure 5: Top row shows per-class mean response time for M/M/2 (left) and M/PH/2 (right) with four priority classes. Left graph is derived using RDR and right graph is derived using RDR-A. Bottom row shows the error in our analytically-derived mean delay relative to simulation results, for the corresponding graphs in the top row.

(Again, all classes have the same job size distribution). As we see from the figure, the per-class mean response time increases nearly linearly with C^2 . Figure 6 (right) shows the relative error in mean delay when the results of the RDR-A analysis in the left plot are compared with simulation. Again the error is under 2%. We will see later that this error increases only slightly when we move to the case of priority classes with different means.

Finally, we note that in the above computations RDR is much more computationally efficient than simulation. Simulation requires tens of minutes to generate each figure, since the simulation is run 30 times, and in each run 1,000,000 events are generated. By comparison our analysis takes only a few seconds for each figure. Further, if we try to reduce the number of events in the simulation to 100,000 events, to speed it up, we see five times as much variation in the simulation around our analytical values. Thus, it is possible that as we increase the number of events in simulation, the difference in our analysis and the simulation may decrease even further.

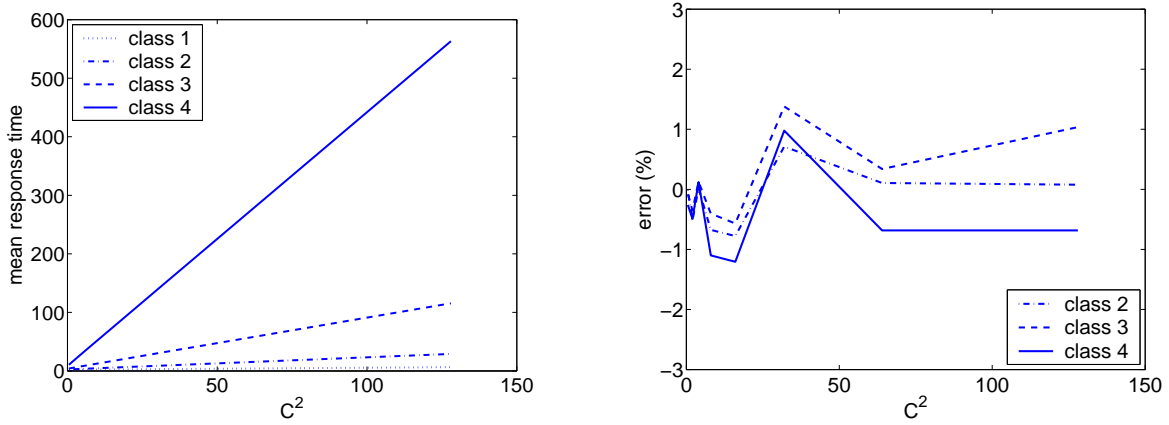


Figure 6: (Left) Per-class mean response times for $M/PH/2$ with four priority classes, derived via RDR-A analysis. (Right) Relative error in analysis of mean delay compared with simulation.

4 Comparisons and Insights

In this section we apply RDR and RDR-A to answer fundamental questions on prioritization in multi-server systems. In Section 4.1 we compare the behavior of multi-server versus single server systems under prioritization. In this context, we also evaluate the BB approximation, which approximates the effect of prioritization in a multi-server system by that in a single server system. In Section 4.2 we evaluate the effect of prioritization schemes which favor short jobs in multi-server systems. Finally, in Section 4.3 we study the effect of aggregating multiple priority classes into just two classes, so as to significantly speed up the analysis. In this context we also evaluate the MK-N approximation discussed earlier.

4.1 Comparing multi-server versus single server performance under prioritization

In this section we compare systems with different numbers of servers. It is important to note that throughout these comparisons, we *hold the total system capacity fixed*. That is, we compare a single server of unit speed with a 2-server system, where each server has speed half, with a 4-server system, where each server has speed one-fourth, etc.

Figure 7 considers an $M/PH/k$ system with two priority classes where k is one (left) then two (middle) then four (right), but the total system capacity is held fixed, and load is fixed at $\rho = 0.8$. The low-priority jobs are exponentially-distributed. The high-priority jobs follow a Coxian distribution where the squared coefficient of variation for high priority jobs, C_H^2 , is varied. The means of the two classes are the same and the load is split evenly between the two classes. The plots show per-class mean response time as a function of C_H^2 . All results are computed using RDR.

The first thing to observe is that the response times in the case of one server appear very different from the response times in the case of two servers, or four servers. The effect of prioritization in a single server system offers little (quantitative) insight into the effect of prioritization in a multi-server system, aside from

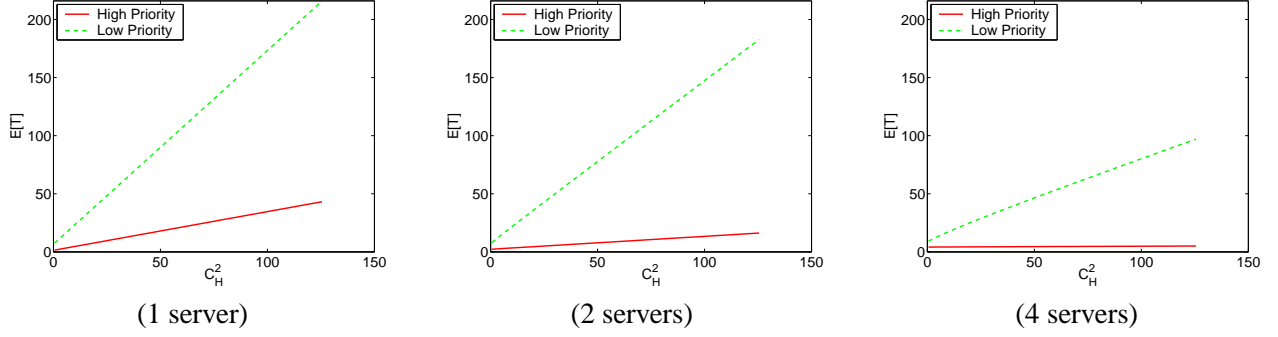


Figure 7: *Contrasting per-class mean response time under one server (left), two servers (middle) and four servers (right) for an $M/PH/k$ with two priority classes. Total system capacity is fixed throughout, and $\rho = 0.8$. Results are obtained using RDR.*

the fact that in all cases the response times appear to be a nearly linear function of C_H^2 .

Figure 7 also illustrates some other interesting points. We see that as we increase the number of servers, under *high* C_H^2 , the performance of both high-priority and low-priority jobs improves. By contrast, under *low* C_H^2 , the performance can get worse as we increase the number of servers. To understand this phenomenon, observe that when C_H^2 is high, short jobs can get stuck behind long jobs, and increasing the number of servers can allow the short jobs to get a chance to serve. By contrast when C_H^2 is low, all jobs are similar in size, so we don't get the benefit of allowing short jobs to jump ahead of long jobs when there are more servers. However we do get the negative effect of increasing the number of servers, namely the underutilization of system resources when there are few jobs in the system, since each of the k servers only has speed $1/k$. The behavior under low C_H^2 , where more servers lead to worse performance, is more prominent under lower load ρ .

Figure 7 already implies that the effect of prioritization on mean response time in a multi-server system may be quite different from that in a single server system. In Figure 8 we investigate this phenomena more closely, by evaluating when the BB approximation [2], which is based on this assumption of similar behavior in single and multi-server priority queues, is accurate. Looking at Figure 8, we see that the error in the BB approximation appears to increase for higher C^2 (right graph) and for more classes. With four classes and two servers, the error is already 10% when $C^2 = 8$ and higher for higher C^2 . By contrast, for the same 4-class case as shown in Figure 8, the error in RDR is always $< 2\%$ independent of C^2 and the number of servers (we have omitted this graph). In the above graphs all classes were statistically identical. In the case where the classes have different means, the error in BB can be much higher, whereas RDR-A is insensitive to this.

4.2 The effect of biasing toward short jobs in multi-server versus single server systems

Until now, we have assumed that all job classes are statistically equivalent. In this section and the next section, we remove this assumption. In this section we consider the effect of priority schemes which favor

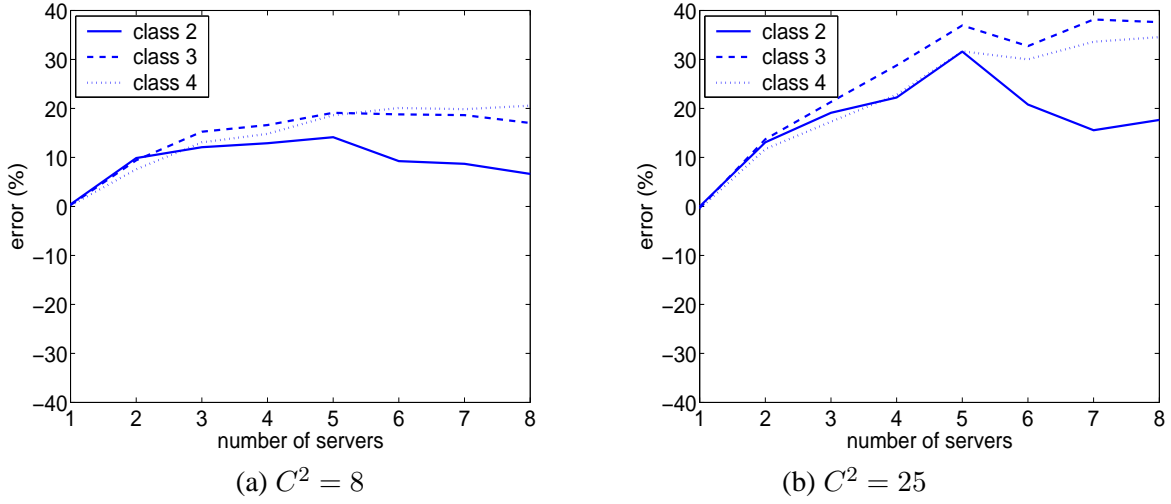


Figure 8: Error in predicting mean delay using the BB approximation (compared with simulation) for an M/PH/2 with four classes where $C^2 = 8$ (left) or $C^2 = 25$ (right) and $\rho = 0.8$.

short jobs in multi-server systems. Biasing towards short jobs is a common method for improving mean response time in any system. We use RDR to understand how the benefit of favoring short jobs in a single server system compares to that for a multi-server system.

Figure 9 considers a job size distribution comprised of an exponential of mean 1, representing jobs which are “short” in expectation, and an exponential of mean 10, representing jobs which are “long” in expectation (where job sizes are measured in a single-server system). The probability of each type of job is chosen to split load evenly between the “short” and “long” jobs. The SMART scheduling policy assigns high priority to the “short” jobs, and the STUPID scheduling policy assigns high priority to the “long” jobs (possibly due to economic reasons). Figure 9 shows the results for a (a) one server, (b) two server, and (c) four server system.

Looking at Figure 9, the SMART and STUPID policies are the same when load ρ is low. At low load, the response time for both policies converges to simply the mean job size, which in these figures is $\frac{20}{11}$ for the single server system, $\frac{40}{11}$ for the 2-server system, and $\frac{80}{11}$ for the 4-server system (recall that in a system with k servers, each server runs at $1/k$ th the speed).

The most interesting observation is that more servers lead to less differentiation between SMART and STUPID schemes. For example, at load $\rho = 0.6$, there is a factor of five differentiation between SMART and STUPID with one server and only a 25% difference between SMART and STUPID with four servers. The effect appears more prominent under lighter load. This can be explained by recalling our earlier observation that multi-server systems allow short jobs a chance to jump ahead of long jobs, hence the negative effects of the STUPID scheme are mitigated.

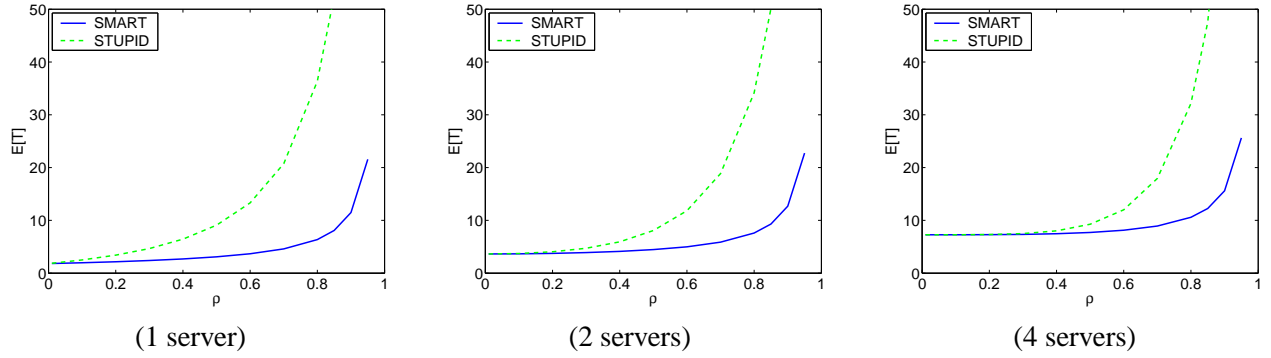


Figure 9: Illustration of mean response time under SMART versus STUPID prioritization in a 2-class system, where the classes are exponentially-distributed with means one and ten respectively, for the case of one server, two servers, and four servers.

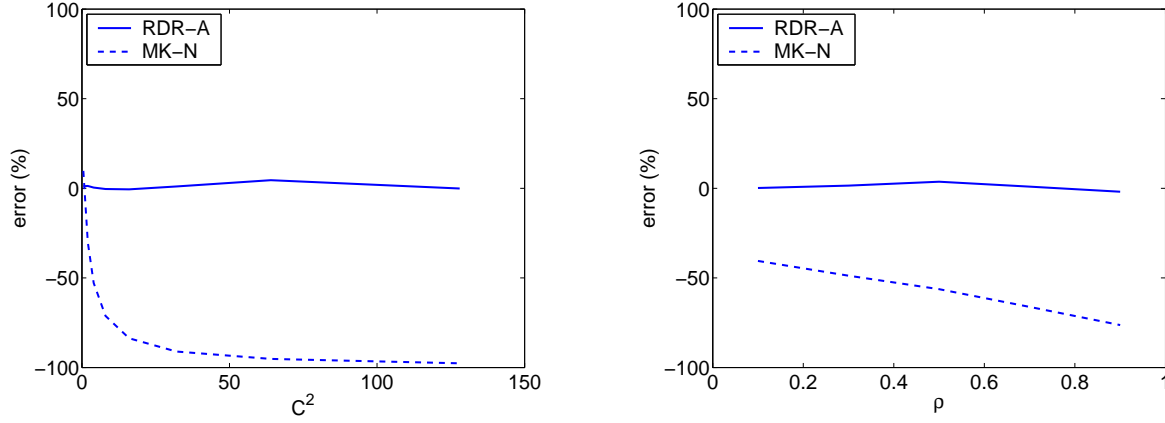


Figure 10: Effect of aggregation. Graphs show error in mean delay of the 4th (lowest priority) class in the MK-N and RDR-A approximations for an M/PH/2 with SMART prioritization. On the left as a function of C^2 where $\rho = 0.8$, and on the right as a function of ρ where $C^2 = 8$. The classes all have a 2-phase Coxian distribution with squared coefficient of variation C^2 and means: 1, 2, 4, and 8.

4.3 How effective is class aggregation: RDR-A

In the early 80's Mitrani and King (later followed by Nishida in the early 90's) proposed analyzing prioritization in a multi-server system via aggregation as follows: To obtain the mean response time of the m^{th} class, simply aggregate classes 1 through $m - 1$ into a single high-priority class, and let class m represent the low-priority class – then analyze the remaining two class system. The above MK-N approximation required further approximating the single aggregate class by an *exponential* job size distribution, since it was not known how to analyze even a two class multi-server system with non-exponential job size distributions.

Since RDR enables the analysis of multi-server priority queues with general PH job size distributions, we can reapply the MK-N aggregation idea, but where now we are able to capture the higher moments of the aggregated class. We call this approximation RDR-A, since it combines the use of RDR together with

aggregation.

To understand the effect of aggregation, we consider a two server system with four priority classes. All the classes have a two phase PH distribution, with varying squared coefficient of variation (C^2). The classes differ however in their mean, having means 1, 2, 4, and 8, respectively, and are prioritized according to the SMART scheme; classes with lower means have higher priority. (STUPID prioritization yields similar insights.) Figure 10 examines the error in the mean delay of the 4th class under RDR-A and under MK-N as a function of C^2 (left) and as a function of ρ (right).

We see that the error in RDR-A is never more than 5% regardless of C^2 or ρ . By contrast, the error in MK-N is almost never less than 50%, and gets worse under higher load and C^2 .

What this tells us is that “aggregation into two classes” is a good method for approximating prioritization in multi-server systems where the number of classes is $m > 2$. However, the aggregation needs to be done carefully – the distribution of the aggregate class must be modeled more closely than can be captured by an exponential distribution. Thus another benefit of RDR is revealed; by allowing for PH job size distributions it enables more accurate approximations of multi-class systems via aggregation.

5 Conclusion

This paper introduces the RDR technique, providing the first near-exact analysis of an M/PH/ k queue with $m \geq 2$ priority classes. The RDR algorithm is efficient (requiring only a second or two for each data point in the paper) and accurate (resulting in $< 2\%$ error for all cases that we studied). Furthermore, RDR appears to maintain its accuracy across a wide range of loads and job size variability (in this paper we studied load ρ , ranging from 0.05 to 0.95 and studied squared coefficient of variation, C^2 , ranging from 0 to 128).

Although the RDR algorithm is efficient when the number of priority classes is small, it becomes less practical when the number of priority classes grows (e.g., for an M/M/2 with 10 priority classes, the running time can get as high as tens of seconds). Hence we also introduce the RDR-A approximation, which works by aggregating the $m > 2$ priority classes into only two priority classes. The distribution of each aggregate class is then captured by a PH distribution, and the resulting 2 class system (with PH job sizes) is solved using RDR. The RDR-A algorithm is extremely efficient (< 1 second for a data point, regardless of the number of classes), since its running time is that of the RDR algorithm for only two classes. Furthermore, the RDR-A algorithm has high accuracy ($< 5\%$ error) across all loads and C^2 .

We use our analysis to obtain insights about priority queueing in multi-server systems. We start by comparing multi-server systems with single server systems of equal capacity. We find that the effect of prioritization in multi-server systems cannot be predicted by considering a comparable single server system. The reason is that adding servers creates complex effects not present in a single server. For example, multiple servers provide a strong benefit in dealing with highly variable job sizes, however they also hinder performance under lighter load. We also compare multi-server with single server systems, by evaluating the

error in the Bondi-Buzen (BB) approximation which is based on relating multi-server performance under prioritization to single-server performance. We find that the error in BB can get quite high, when C^2 grows or the number of classes grows.

We next consider the effect of “smart” prioritization, where classes of jobs with smaller means are given priority over those with larger means. We find that “smart” prioritization has a much stronger effect in a single-server system, than in a multi-server system of equal capacity. This can be explained in part by the observation that multiple servers inherently help out short jobs by allowing them to jump ahead of long jobs.

Lastly, we evaluate the effect of class aggregation as an approximation method for analyzing a high number of classes. We find that aggregation when done carefully – by capturing several moments of the aggregated class – works surprisingly well, resulting in very low error. However, when the aggregate class is approximated only with respect to its first moment (by just an exponential distribution), aggregation can be very poor, resulting in error of well over 50%. The fact that RDR allows the first analysis of classes with PH job size distributions enables this good aggregation approximation.

In this paper we have focused on the problem of multi-server queues with $m > 2$ priorities. What makes this problem difficult is the fact that its Markov chain representation grows infinitely in m dimensions, and there are dependencies between those dimensions (the class M jobs depend on the class H jobs, and the class L jobs depend on both the class M and class H jobs). The RDR algorithm greatly simplifies this problem by reducing the dimensionality of the Markov chain to just one. There are many other problems that also exhibit high dimensionality in their Markov chain representation, and it is possible that the RDR method introduced here may be applicable to those problems as well.

References

- [1] D. Bertsimas and D. Nakazato. The distributional Little’s Law and its applications. *Operations Research*, 43(2):298–310, 1995.
- [2] A. Bondi and J. Buzen. The response times of priority classes under preemptive resume in M/G/m queues. In *ACM Sigmetrics*, pages 195–201, August 1984.
- [3] L. Bright and P. Taylor. Calculating the equilibrium distribution in level dependent quasi-birth-and-death processes. *Stochastic Models*, 11:497–514, 1995.
- [4] J. Buzen and A. Bondi. The response times of priority classes under preemptive resume in M/M/m queues. *Operations Research*, 31:456–465, 1983.
- [5] A. Cobham. Priority assignment in waiting line problems. *Operations Research*, 2:70–76, 1954.
- [6] R. Davis. Waiting-time distribution of a multi-server, priority queueing system. *Operations Research*, 14:133–136, 1966.
- [7] W. Feng, M. Kawada, and K. Adachi. Analysis of a multiserver queue with two priority classes and (M,N)-threshold service schedule ii: preemptive priority. *Asia-Pacific Journal of Operations Research*, 18:101–124, 2001.
- [8] H. Gail, S. Hantler, and B. Taylor. Analysis of a non-preemptive priority multiserver queue. *Advances in Applied Probability*, 20:852–879, 1988.

- [9] H. Gail, S. Hantler, and B. Taylor. On a preemptive Markovian queues with multiple servers and two priority classes. *Mathematics of Operations Research*, 17:365–391, 1992.
- [10] E. Kao and K. Narayanan. Modeling a multiprocessor system with preemptive priorities. *Management Science*, 2:185–97, 1991.
- [11] E. Kao and S. Wilson. Analysis of nonpreemptive priority queues with multiple servers and two priority classes. *European Journal of Operational Research*, 118:181–193, 1999.
- [12] E. P. C. Kao and K. S. Narayanan. Computing steady-state probabilities of a nonpreemptive priority multiserver queue. *Journal on Computing*, 2(3):211 – 218, 1990.
- [13] A. Kapadia, M. Kazumi, and A. Mitchell. Analysis of a finite capacity nonpreemptive priority queue. *Computers and Operations Research*, 11:337–343, 1984.
- [14] O. Kella and U. Yechiali. Waiting times in the non-preemptive priority M/M/c queue. *Stochastic Models*, 1:257 – 262, 1985.
- [15] G. Latouche and V. Ramaswami. *Introduction to Matrix Analytic Methods in Stochastic Modeling*. ASA-SIAM, 1999.
- [16] H. Leemans. *The Two-Class Two-Server Queue with Nonpreemptive Heterogeneous Priority Structures*. PhD thesis, K.U.Leuven, 1998.
- [17] D. McWherter, B. Schroeder, N. Ailamaki, and M. Harchol-Balter. Priority mechanisms for OLTP and transactional web applications. In *Proceedings of the 20th International Conference on Data Engineering (ICDE 2004)*, pages 535–546, April 2004.
- [18] D. Miller. Steady-state algorithmic analysis of M/M/c two-priority queues with heterogeneous servers. In R. L. Disney and T. J. Ott, editors, *Applied probability - Computer science, The Interface, volume II*, pages 207–222. Birkhauser, 1992.
- [19] I. Mitrani and P. King. Multiprocessor systems with preemptive priorities. *Performance Evaluation*, 1:118–125, 1981.
- [20] M. Neuts. Moment formulas for the Markov renewal branching process. *Advances in Applied Probabilities*, 8:690–711, 1978.
- [21] B. Ngo and H. Lee. Analysis of a pre-emptive priority M/M/c model with two types of customers and restriction. *Electronics Letters*, 26:1190–1192, 1990.
- [22] T. Nishida. Approximate analysis for heterogeneous multiprocessor systems with priority jobs. *Performance Evaluation*, 15:77–88, 1992.
- [23] T. Osogami and M. Harchol-Balter. A closed-form solution for mapping general distributions to minimal PH distributions. In *Performance TOOLS*, pages 200–217, 2003.
- [24] T. Osogami, M. Harchol-Balter, and A. Scheller-Wolf. Analysis of cycle stealing with switching cost. In *ACM Sigmetrics 2003*, pages 184–195, 2003.
- [25] A. Sleptchenko. Multi-class, multi-server queues with non-preemptive priorities. Technical Report 2003-016, EURANDOM, Eindhoven University of Technology, 2003.
- [26] A. Sleptchenko, A. van Harten, and M. van der Heijden. An exact solution for the state probabilities of the multi-class, multi-server queue with preemptive priorities, 2003 – Manuscript.
- [27] W. Whitt. A review of $L = \lambda W$ and extensions. *Queueing Systems*, 9:235–268, 1991.

A Moments of busy periods in nonhomogeneous QBD processes

(Can omit in final version if space constraints)

Neuts' algorithm [20] is an efficient algorithm that calculates the moments of various types of busy periods in very general processes, i.e. M/G/1 type semi-Markov processes. Because of its generality, however, the description of the algorithm in [20] is sophisticated, and thus non-trivial to understand or implement. Since Neuts' algorithm can be applied to the performance analysis of many computer and communication systems, it is a shame that it has not been used more frequently in the literature.

The purpose of this section is therefore to make Neuts' algorithm more accessible by re-describing his algorithm restricted to the first three moments of particular types of busy periods in QBD processes. In [20], the analysis is limited to homogeneous (level independent) processes. However, it is trivially extended to nonhomogeneous (level dependent) QBD process which we need in this paper. Therefore, we describe Neuts' algorithm for nonhomogeneous QBD processes that repeats after level k . We omit all proofs, which are provided in detail in [20], instead we focus on intuition and interpretation. We include everything needed to apply Neuts' algorithm within our solution framework, so that readers who wish to apply our methodology can do so.

In our paper, we consider a QBD process with state space $E = \{(i, j) | i \geq 0, 1 \leq j \leq P_i\}$, which has generator matrix \mathbf{Q} :

$$\mathbf{Q} = \begin{pmatrix} \mathbf{A}_1^{(0)} & \mathbf{A}_0^{(0)} & & & \\ \mathbf{A}_2^{(1)} & \mathbf{A}_1^{(1)} & \mathbf{A}_0^{(1)} & & \\ & \mathbf{A}_2^{(2)} & \mathbf{A}_1^{(2)} & \ddots & \\ & & \ddots & \ddots & \ddots \end{pmatrix}$$

where $\mathbf{A}_i^{(l)}$ is a $P_l \times P_l$ matrices. We assume that our QBD process has a repeating structure: there exists $k < \infty$ such that $\mathbf{A}_i^{(l)} = \mathbf{A}_i^{(k)}$ for $i = 0, 1, 2$ and for all $l \geq k$. Figure 11 shows a particular QBD process with $P_l = 2$ for all l .

We define level i as denoting the set of states of the form (i, j) for $j = 1, \dots, P_i$. Our goal can be roughly stated as deriving the passage time required to get from state (l, j) to level $l - 1$ conditioned on the particular state first reached in level $l - 1$. More precisely, we seek the first three moments of the distribution of the time required to get from state (l, j) to state $(l - 1, k)$, given that state $(l - 1, k)$ is the first state reached in level $l - 1$ for any $l \geq 1$, $1 \leq j \leq P_l$, and $1 \leq k \leq P_{l-1}$. In Section A.1, we introduce more notation. In Section A.2, we show how Neuts algorithm can be applied to derive the first three moments of the conditional passage time in the repeating part ($l > k$). In Section A.3, we extend the analysis to nonrepeating part ($l \leq k$). In Section A.4, we summarize other generalization that Neuts' algorithm allows.

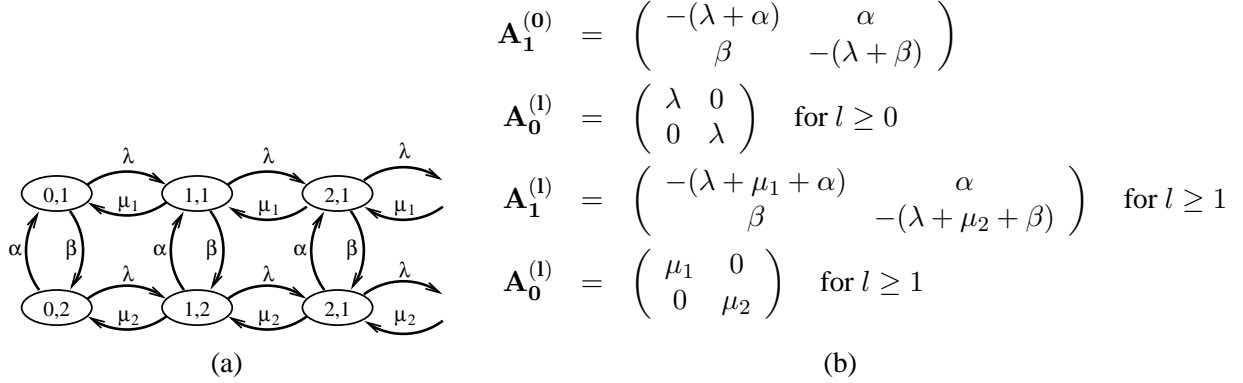


Figure 11: An example of a QBD process: (a) states and transition rates, and (b) submatrices of the generator matrix.

A.1 Notation

We define the transition probability matrix, $\mathbf{P}(x)$, as

$$\mathbf{P}(x) = \begin{pmatrix} \alpha_{0,1}(x) & \alpha_{0,0}(x) & & & \\ \alpha_{1,2}(x) & \alpha_{1,1}(x) & \alpha_{1,0}(x) & & \\ & \alpha_{2,2}(x) & \alpha_{2,1}(x) & \ddots & \\ & & \ddots & \ddots & \ddots \end{pmatrix}$$

where the (s, t) element, $\mathbf{P}_{st}(x)$, is the probability that the sojourn time at state s is $\leq x$ and the first transition out of state s is to state t . Observe that $\alpha_{l,1}(x)$ is a $P_l \times P_l$ submatrices for $l \geq 0$.

Next, we define the r -th moment of submatrices $\alpha_{l,i}(x)$ as

$$\alpha_{l,i}^{(r)} = \int_0^\infty x^r d\alpha_{l,i}(x)$$

for $i = 0, 1, 2$, $r = 1, 2, 3$, and $l \geq 0$, where an integral of a matrix M is a matrix of the integrals of the elements in M .

Example

Consider the QBD process shown in Figure 11. Let $\gamma_1 = \lambda + \mu_1 + \alpha$ and $\gamma_2 = \lambda + \mu_2 + \beta$. Then, $\alpha_{l,i}(x)$'s and their moments $\alpha_{l,i}^{(r)}$ for $r = 1, 2, 3$ look as follows:

$$\begin{aligned}\alpha_{l,0}(x) &= \begin{pmatrix} 1 - e^{-\gamma_1 x} & 0 \\ 0 & 1 - e^{-\gamma_2 x} \end{pmatrix} \begin{pmatrix} \frac{\lambda}{\gamma_1} & 0 \\ 0 & \frac{\lambda}{\gamma_2} \end{pmatrix} \\ \alpha_{l,1}(x) &= \begin{pmatrix} 1 - e^{-\gamma_1 x} & 0 \\ 0 & 1 - e^{-\gamma_2 x} \end{pmatrix} \begin{pmatrix} 0 & \frac{\alpha}{\gamma_1} \\ \frac{\beta}{\gamma_2} & 0 \end{pmatrix} \\ \alpha_{l,2}(x) &= \begin{pmatrix} 1 - e^{-\gamma_1 x} & 0 \\ 0 & 1 - e^{-\gamma_2 x} \end{pmatrix} \begin{pmatrix} \frac{\mu_1}{\gamma_1} & 0 \\ 0 & \frac{\mu_2}{\gamma_2} \end{pmatrix}\end{aligned}$$

and

$$\begin{aligned}\alpha_{l,0}^{(r)} &= \begin{pmatrix} \frac{r!}{\gamma_1^r} & 0 \\ 0 & \frac{r!}{\gamma_2^r} \end{pmatrix} \begin{pmatrix} \frac{\lambda}{\gamma_1} & 0 \\ 0 & \frac{\lambda}{\gamma_2} \end{pmatrix} \\ \alpha_{l,1}^{(r)} &= \begin{pmatrix} \frac{r!}{\gamma_1^r} & 0 \\ 0 & \frac{r!}{\gamma_2^r} \end{pmatrix} \begin{pmatrix} 0 & \frac{\alpha}{\gamma_1} \\ \frac{\beta}{\gamma_2} & 0 \end{pmatrix} \\ \alpha_{l,2}^{(r)} &= \begin{pmatrix} \frac{r!}{\gamma_1^r} & 0 \\ 0 & \frac{r!}{\gamma_2^r} \end{pmatrix} \begin{pmatrix} \frac{\mu_1}{\gamma_1} & 0 \\ 0 & \frac{\mu_2}{\gamma_2} \end{pmatrix}\end{aligned}$$

for $l \geq 2$.

Finally, let $\mathbf{G}_l(x)$ be an $P_l \times P_l$ matrix whose (j, k) element, $(\mathbf{G}_l)_{jk}(x)$, is the probability that the time to visit level $l - 1$ is at most x and the first state visited in level $l - 1$ is $(l - 1, k)$ given that we started at (l, j) . Also, let $(\mathbf{G}_l^{(r)})_{jk}$ be the r -th moment of $(\mathbf{G}_l)_{jk}(x)$; namely, $\mathbf{G}_l^{(r)} = \int_0^\infty x^r d\mathbf{G}_l(x)$ for $r = 1, 2, 3$.

Matrix $\mathbf{G}_l = \lim_{x \rightarrow \infty} \mathbf{G}_l(x)$ is a fundamental matrix used in the matrix analytic method, and various algorithms to calculate \mathbf{G}_l have been proposed [15]. The most straightforward (but slow) algorithm for the repeating part ($l \geq k$) is to iterate

$$\mathbf{G}_l = -(\mathbf{A}_1^{(k)})^{-1} \mathbf{A}_0^{(k)} (\mathbf{G}_l)^2 - (\mathbf{A}_1^{(k)})^{-1} \mathbf{A}_2^{(k)} \quad (2)$$

until it converges. Once \mathbf{G}_k is calculated, \mathbf{G}_l for $l < k$ can be calculated recursively:

$$\mathbf{G}_l = -\left(\mathbf{A}_1^{(l)} + \mathbf{A}_0^{(l)} \mathbf{G}_{l+1}\right)^{-1} \mathbf{A}_2^{(l)}. \quad (3)$$

Notice that $(\mathbf{G}_l)_{jk}(x)$ is not a proper distribution function and $(\mathbf{G}_l)_{jk} = \lim_{x \rightarrow \infty} (\mathbf{G}_l)_{jk}(x)$, which is the probability that the first state in level $l - 1$ is state $(l - 1, k)$ given that we start at state (l, j) , can be less

than 1. Therefore, $(\mathbf{G}_1^{(r)})_{jk}$ is not a proper moment rather a conditional moment: “the r -th moment of the distribution of the first passage time to level $l - 1$ given that the first state in level $l - 1$ is state $(l - 1, k)$ and given that we start at state (l, j) ” multiplied by “the probability that the first state reached in level $l - 1$ is state $(l - 1, k)$ given that we start at state (l, j) .”

A.2 Moments of passage time in the repeating part

The quantities that we need in our methodology are (a) the probability that the first state reached in level $l - 1$ is state $(l - 1, t)$ given that we start at state (l, s) and (b) the r -th moment of the distribution of the first passage time to level $l - 1$ given that the first state in level $l - 1$ is state $(l - 1, t)$ and given that we start at state (l, s) . Quantity (a) is given by $(\mathbf{G}_1)_{s,t}$, and quantity (b) is given by $\frac{(\mathbf{G}_1^{(r)})_{s,t}}{(\mathbf{G}_1)_{s,t}}$. Matrix \mathbf{G}_1 is obtained by an iterative substitution of (2) and (3) for any $l \geq 1$. Therefore, our goal is to derive matrices $\mathbf{G}_1^{(r)}$ for $r = 1, 2, 3$. In this section, we derive the repeating part, $\mathbf{G}_1^{(r)} = \mathbf{G}_1^{(k)}$ for $l \geq k$. In Section A.3 we derive the nonrepeating part, $\mathbf{G}_1^{(r)}$ for $l < k$.

Once \mathbf{G}_k is obtained, matrix $\mathbf{G}_k^{(1)}$ is obtained by iterating

$$\mathbf{G}_k^{(1)} = \alpha_2^{(1)} + \alpha_1^{(1)} \mathbf{G}_k + \alpha_1 \mathbf{G}_k^{(1)} + \alpha_0^{(1)} \mathbf{G}_k \mathbf{G}_k + \alpha_0 \mathbf{G}_k^{(1)} \mathbf{G}_k + \alpha_0 \mathbf{G}_k \mathbf{G}_k^{(1)} \quad (4)$$

Similarly, matrix $\mathbf{G}_k^{(2)}$ is obtained by iterating

$$\begin{aligned} \mathbf{G}_k^{(2)} = & \alpha_2^{(2)} + \alpha_1^{(2)} \mathbf{G}_k + 2\alpha_1^{(1)} \mathbf{G}_k^{(1)} + \alpha_1 \mathbf{G}_k^{(2)} \\ & + \alpha_0^{(2)} \mathbf{G}_k \mathbf{G}_k + 2\alpha_0^{(1)} (\mathbf{G}_k^{(1)} \mathbf{G}_k + \mathbf{G}_k \mathbf{G}_k^{(1)}) + \alpha_0 (\mathbf{G}_k^{(2)} \mathbf{G}_k + 2\mathbf{G}_k^{(1)} \mathbf{G}_k^{(1)} + \mathbf{G}_k \mathbf{G}_k^{(2)}) \end{aligned} \quad (5)$$

and matrix $\mathbf{G}_k^{(3)}$ is obtained by iterating

$$\begin{aligned} \mathbf{G}_k^{(3)} = & \alpha_2^{(3)} + \alpha_1^{(3)} \mathbf{G}_k + 3\alpha_1^{(2)} \mathbf{G}_k^{(1)} + 3\alpha_1^{(1)} \mathbf{G}_k^{(2)} + \alpha_1 \mathbf{G}_k^{(3)} \\ & + \alpha_0^{(3)} \mathbf{G}_k \mathbf{G}_k + 3\alpha_0^{(2)} (\mathbf{G}_k^{(1)} \mathbf{G}_k + \mathbf{G}_k \mathbf{G}_k^{(1)}) + 3\alpha_0^{(1)} (\mathbf{G}_k^{(2)} \mathbf{G}_k + 2\mathbf{G}_k^{(1)} \mathbf{G}_k^{(1)} + \mathbf{G}_k \mathbf{G}_k^{(2)}) \\ & + \alpha_0 (\mathbf{G}_k^{(3)} \mathbf{G}_k + 3\mathbf{G}_k^{(2)} \mathbf{G}_k^{(1)} + 3\mathbf{G}_k^{(1)} \mathbf{G}_k^{(2)} + \mathbf{G}_k \mathbf{G}_k^{(3)}). \end{aligned} \quad (6)$$

We now give intuition behind expressions (4)-(6). The right hand side of (4) can be divided into three parts: [0] $\alpha_2^{(1)}$, [1] $\alpha_1^{(1)} \mathbf{G}_k + \alpha_1 \mathbf{G}_k^{(1)}$, and [2] $\alpha_0^{(1)} \mathbf{G}_k \mathbf{G}_k + \alpha_0 \mathbf{G}_k^{(1)} \mathbf{G}_k + \alpha_0 \mathbf{G}_k \mathbf{G}_k^{(1)}$. For $h = 0, 1, 2$, the (s, t) element of part [h] gives “the first moment of the distribution of the time to get from state (k, s) to state $(k - 1, t)$ given that the first transition out of state (k, s) is to level $k + h - 1$ and the first state reached in level $k - 1$ is $(k - 1, t)$ ” multiplied by “the probability that the first transition out of state (k, s) is to level $k + h - 1$ and the first state reached in level $k - 1$ is $(k - 1, t)$.” Part [1] consists of two terms. The first term, $\alpha_1^{(1)} \mathbf{G}_k$, is the contribution of the time to the first transition, and the second term, $\alpha_1 \mathbf{G}_k^{(1)}$, is the contribution of the time it takes to reach $(k - 1, t)$ after the first transition. Similarly, part [2] consists

of three terms. The first term, $\alpha_0^{(1)} G_k G_k$, is the contribution of the time to the first transition, the second term, $\alpha_0 G_k^{(1)} G_k$, is the contribution of the time it takes to come back from level $k + 1$ to level k after the first transition, and the third term, $\alpha_0 G_k G_k^{(1)}$, is the contribution of the time it takes to go from level k to level $k - 1$.

The right hand sides of (5) and (6) can similarly be divided into three parts: part [0] consists of terms containing α_2 or $\alpha_2^{(r)}$; part [1] consists of terms containing α_1 or $\alpha_1^{(r)}$; part [2] consists of terms containing α_0 or $\alpha_0^{(r)}$. The three parts of (5) and (6) can be interpreted exactly the same way as the three parts of (4) except that “the first moment” in (4) must be replaced by “the second moment” and “the third moment” in (5) and (6), respectively. The three terms in part [1] of (5) can be interpreted as follows. Let T_α be the time to the first transition and let T_G be the time it takes from level k to level $k - 1$. Then, the second moment of the distribution of these two times is

$$E[(T_\alpha + T_G)^2] = E[(T_\alpha)^2] + 2E[T_\alpha]E[T_G] + E[(T_G)^2],$$

since T_α and T_G are independent. Roughly speaking, $\alpha_1^{(2)} G_k$ corresponds to $E[(T_\alpha)^2]$, $2\alpha_1^{(1)} G_k^{(1)}$ corresponds to $2E[T_\alpha]E[T_G]$, and $\alpha_1 G_k^{(2)}$ corresponds to $E[(T_G)^2]$. The other terms can be interpreted in the same way.

A.3 Extension to nonrepeating part

For $l < k$, $G_l^{(r)}$ is calculated recursively.

$$\begin{aligned} G_1^{(1)} &= \alpha_{l,2}^{(1)} + \alpha_{l,1}^{(1)} G_l + \alpha_{l,1} G_l^{(1)} + \alpha_{l,0}^{(1)} G_{l+1} G_l + \alpha_{l,0} G_{l+1}^{(1)} G_l + \alpha_{l,0} G_{l+1} G_l^{(1)} \\ &= (\mathbf{I} - \alpha_{l,1} - \alpha_{l,0} G_{l+1})^{-1} \left(\alpha_{l,2}^{(1)} + \alpha_{l,1}^{(1)} G_l + \alpha_{l,0}^{(1)} G_{l+1} G_l + \alpha_{l,0} G_{l+1}^{(1)} G_l \right) \end{aligned}$$

$$\begin{aligned} G_1^{(2)} &= \alpha_{l,2}^{(2)} + \alpha_{l,1}^{(2)} G_l + 2\alpha_{l,1}^{(1)} G_l^{(1)} + \alpha_{l,1} G_l^{(2)} \\ &\quad + \alpha_{l,0}^{(2)} G_{l+1} G_l + 2\alpha_{l,0}^{(1)} (G_{l+1}^{(1)} G_l + G_{l+1} G_l^{(1)}) + \alpha_{l,0} (G_{l+1}^{(2)} G_l + 2G_{l+1}^{(1)} G_l^{(1)} + G_{l+1} G_l^{(2)}) \\ &= (\mathbf{I} - \alpha_{l,1} - \alpha_{l,0} G_{l+1})^{-1} \\ &\quad \left(\alpha_{l,2}^{(2)} + \alpha_{l,1}^{(2)} G_l + 2\alpha_{l,1}^{(1)} G_l^{(1)} \right. \\ &\quad \left. + \alpha_{l,0}^{(2)} G_{l+1} G_l + 2\alpha_{l,0}^{(1)} (G_{l+1}^{(1)} G_l + G_{l+1} G_l^{(1)}) + \alpha_{l,0} (G_{l+1}^{(2)} G_l + 2G_{l+1}^{(1)} G_l^{(1)}) \right) \end{aligned}$$

$$\begin{aligned}
\mathbf{G}_1^{(3)} &= \alpha_{l,2}^{(3)} + \alpha_{l,1}^{(3)} \mathbf{G}_l + 3\alpha_{l,1}^{(2)} \mathbf{G}_l^{(1)} + 3\alpha_{l,1}^{(1)} \mathbf{G}_l^{(2)} + \alpha_{l,1} \mathbf{G}_l^{(3)} \\
&\quad + \alpha_{l,0}^{(3)} \mathbf{G}_{l+1} \mathbf{G}_l + 3\alpha_{l,0}^{(2)} (\mathbf{G}_{l+1}^{(1)} \mathbf{G}_1 + \mathbf{G}_{l+1} \mathbf{G}_1^{(1)}) + 3\alpha_{l,0}^{(1)} (\mathbf{G}_{l+1}^{(2)} \mathbf{G}_1 + 2\mathbf{G}_{l+1}^{(1)} \mathbf{G}_1^{(1)} + \mathbf{G}_{l+1} \mathbf{G}_1^{(2)}) \\
&\quad + \alpha_{l,0} (\mathbf{G}_{l+1}^{(3)} \mathbf{G}_1 + 3\mathbf{G}_{l+1}^{(2)} \mathbf{G}_1^{(1)} + 3\mathbf{G}_{l+1}^{(1)} \mathbf{G}_1^{(2)} + \mathbf{G}_{l+1} \mathbf{G}_1^{(3)}). \\
&= (\mathbf{I} - \alpha_{l,1} - \alpha_{l,0} \mathbf{G}_{l+1})^{-1} \\
&\quad \left(\alpha_{l,2}^{(3)} + \alpha_{l,1}^{(3)} \mathbf{G}_l + 3\alpha_{l,1}^{(2)} \mathbf{G}_l^{(1)} + 3\alpha_{l,1}^{(1)} \mathbf{G}_l^{(2)} \right. \\
&\quad + \alpha_{l,0}^{(3)} \mathbf{G}_{l+1} \mathbf{G}_l + 3\alpha_{l,0}^{(2)} (\mathbf{G}_{l+1}^{(1)} \mathbf{G}_1 + \mathbf{G}_{l+1} \mathbf{G}_1^{(1)}) + 3\alpha_{l,0}^{(1)} (\mathbf{G}_{l+1}^{(2)} \mathbf{G}_1 + 2\mathbf{G}_{l+1}^{(1)} \mathbf{G}_1^{(1)} + \mathbf{G}_{l+1} \mathbf{G}_1^{(2)}) \\
&\quad \left. + \alpha_{l,0} (\mathbf{G}_{l+1}^{(3)} \mathbf{G}_1 + 3\mathbf{G}_{l+1}^{(2)} \mathbf{G}_1^{(1)} + 3\mathbf{G}_{l+1}^{(1)} \mathbf{G}_1^{(2)}) \right)
\end{aligned}$$

A.4 Generalization allowed

Finally, we mention some generalizations that Neuts' algorithm allows. (1) We restricted ourselves to the first three moments, but this can be generalized to any higher moments. (2) We restricted ourselves to the first passage time from level l to level $l - 1$, but this can be generalized to level l from level $l - i$. (3) We restricted to QBD processes, but this can be generalized to M/G/1 type semi-Markov processes. (4) We restricted ourselves to the moments of the distribution of the duration of busy periods, but this can be generalized to the moments of the joint distribution of the duration of a busy period and the number of transitions during the busy period.