

REAL-TIME DELAY ESTIMATION BASED ON DELAY HISTORY

by

Rouba Ibrahim and Ward Whitt

IEOR Department
Columbia University
{rei2101, ww2040}@columbia.edu

Abstract

Motivated by interest in making delay announcements to arriving customers who must wait in call centers and related service systems, we study the performance of alternative real-time delay estimators based on recent customer delay experience. The main estimators considered are: (i) the delay of the last customer to enter service (LES), (ii) the delay experienced so far by the customer at the head of the line (HOL), and (iii) the delay experienced by the customer to have arrived most recently among those who have already completed service (RCS). We compare these delay-history estimators to the standard estimator based on the queue length (QL), commonly used in practice, which requires knowledge of the mean interval between successive service completions in addition to the queue length. We characterize performance by the mean squared error (MSE). We do analysis and conduct simulations for the standard $GI/M/s$ multi-server queueing model, emphasizing the case of large s . We obtain analytical results for the conditional distribution of the delay given the observed HOL delay. An approximation to its mean value serves as a refined estimator. For all three candidate delay estimators, the MSE relative to the square of the mean is asymptotically negligible in the many-server and classical heavy-traffic limiting regimes.

Keywords: delay estimation, real-time delay estimation, delay prediction, delay announcements, many-server queues, call centers, heavy traffic.

May 30, 2007; Revision: April 9, 2008

1. Introduction

In this paper, we study alternative ways to estimate the delay (before entering service) of an arriving customer in a service system. These delay estimates may be used to make delay announcements to arriving customers, especially when the delay will be relatively long. Such real-time delay announcements can be very helpful with invisible queues, as in call centers, where service requests are made by telephone; see Gans et al. (2003) for background on call centers.

Since the steady-state waiting-time distribution tends to be quite highly variable (e.g., often exponential or approximately so), good real-time delay estimation necessarily relies on state information; see Whitt (1999). From the perspective of statistical precision, for a single-number estimate we would ideally want to use the conditional expected delay given all information available at the arrival epoch, but complexity leads to considering more elementary alternatives.

The Standard Queue-Length (QL) Delay Estimator. The standard state-dependent delay estimator, commonly used in practice (assuming service from a queue in first-come first-served order, but without any other specific model assumptions), is the *queue-length (QL) delay estimator*, defined as

$$\theta_{QL}(t) \equiv \frac{Q(t) + 1}{r(t)}, \quad (1.1)$$

where the notation \equiv means “defined as,” t is the current time (time of the arrival for which the announcement is made), $Q(t)$ is the queue length (number of customers waiting) and $r(t)$ is the rate at which customers enter service (typically not known precisely). If the number of servers is $s(t)$, and can be assumed to remain at that level in the near future, with each server serving a single customer without interruption, and the current average service time is $m(t)$, then the rate customers enter service may be approximated by $r(t) = s(t)/m(t)$. Furthermore, when the mean service time is stable, we can replace $m(t)$ by a long-run average service time m . The QL delay estimator then becomes $\theta_{QL}(t) \equiv m(Q(t) + 1)/s(t)$, which requires knowledge of only $s(t)$, the number of servers, and $Q(t)$, the queue length, at each time t , which is information that usually is readily available.

Estimators Based on Delay History. In this paper, we examine alternative estimators based on the delays actually experienced by recent customers, in particular: (i) the delay of the last customer to enter service (LES), (ii) the delay experienced so far by the customer at

the head of the line (HOL), (iii) the delay experienced by the customer to have arrived most recently among those that have completed service (RCS).

These delay estimators based on recent delay history are appealing because they are easy to interpret, and because they are simple and robust, applying to a broad range of models, without requiring knowledge of the model or its parameters. If somehow the queue length, $Q(t)$, or the rate at which customers enter service, $r(t)$, is unknown or incorrect, then we would have difficulties with the standard QL estimator. With any prediction system, it is good to monitor its performance, but that is often not possible for the customer. A delay-history delay estimator has the advantage that the basis for the prediction is evident.

The HOL delay estimator was used as an announcement in an Israeli bank studied by Mandelbaum et al. (2000) and is mentioned as a candidate delay announcement by Nakibly (2002) in her study of delay predictions. Something similar to LES or RCS is used by the U. S. Citizenship and Immigration Service (USCIS); they publish the arrival time of recently completed applications in order to give an idea about upcoming delays. In this study, we are motivated in part by recent work by Armony et al. (2006), who studied delay announcements in many-server queues with customer abandonment, focusing on customer response to the announcements, leading to balking and new abandonment behavior. They developed ways to approximately describe the equilibrium system performance using LES delay announcements. Armony et al. (2006) discuss the motivation for the LES delay estimator and other delays estimators based on recent delay history.

Quantifying the Effectiveness. We quantify the effectiveness of the delay estimators through the mean squared error (MSE), which we approximate analytically and estimate via simulation. To illustrate, let $W_{LES}(w)$ denote the random delay of a new arrival, conditional on that customer having to wait and an observed LES delay of w (under specified conditions, e.g., in steady state). Let $\theta_{LES}(w)$ be a candidate estimator based on this information. We will primarily be concerned with the direct estimator $\theta_{LES}^d(w) \equiv w$, the refined estimator $\theta_{LES}^r(w) \equiv E[W_{LES}(w)]$ and approximations of the refined estimator, since the refined estimator is difficult to determine. The MSE of such an estimator is

$$MSE \equiv MSE(\theta_{LES}(w)) \equiv E \left[(W_{LES}(w) - \theta_{LES}(w))^2 \right] . \quad (1.2)$$

For the refined estimator $\theta_{LES}^r(w)$, the MSE coincides with the variance $Var(W_{LES}(w))$. It is well known that the mean minimizes the MSE (using that information).

To estimate these MSE's via simulation, we use the average squared error (ASE), defined by

$$ASE \equiv \frac{1}{n} \sum_{j=1}^n (a_j - e_j)^2, \quad (1.3)$$

where a_j is the actual delay and e_j is the estimated delay for appropriate customers. For example, if we want to estimate the performance of LES when the observed delay is $w = 0.40$, then we consider all arrivals who must wait ($a_j > 0$) for which the LES delay e_j falls in an interval such as $[0.39, 0.41]$. On the other hand, if we wish to consider the overall average performance of LES, then we consider all j such that $a_j > 0$.

Study in an Idealized Setting. In this paper, we study the performance of the delay-history delay estimators and compare them to the standard QL delay estimator in the relatively simple idealized setting of the $GI/M/s$ queueing model, which has a renewal arrival process, s homogeneous servers working in parallel, unlimited waiting space, a FCFS service discipline and i.i.d. exponential service times with mean m , which are independent of the arrival process. For this $GI/M/s$ model, the QL estimator $\theta_{QL}(t) \equiv m(Q(t) + 1)/s$ is an ideal delay estimator. Indeed, there are no serious competitors, as far as statistical precision is concerned (provided that we have no information about remaining service times). Given the queue length, the future evolution of the system is independent of the past. (This even remain true for more general arrival processes.) Consequently, $\theta_{QL}(t)$ is the conditional mean delay given all information available at time t , so that it minimizes the MSE.

We study the alternative delay-history delay estimators in this simple context in order to gain insight about the relative performance of alternative estimators in more complex scenarios (which are much more difficult to analyze directly). We know that the QL estimator will have superior performance for the $GI/M/s$ model, but we want to understand by how much. That knowledge will help us understand the advantage of the QL estimator over these alternative delay estimators when the QL estimator is appropriate, and will provide useful background when considering these alternative delay estimators for more complicated systems for which these alternative estimators may be preferred.

Motivation for Considering Alternative Delay Estimators. Whenever the actual service system is well modelled by a $GI/M/s$ queueing model and the system state is known accurately at each time, then there is little motivation for considering other delay estimators besides the standard QL estimator. However, real service systems rarely are as simple as the

$GI/M/s$ model. First, the service-time distribution might well be non-exponential, as shown for call centers by Brown et al. (2005). Second, the number of servers and mean service times often change over time, in part because the servers are humans who serve in different shifts and may well have different service-time distributions. Third, the queue length may not be directly observable. That is nicely illustrated by the ticket queues studied by Xu et al. (2007). Upon arriving at a ticket queue, each customer is issued a numbered ticket. The number currently being served is displayed. The queue length is not known to ticket-holding customers or even to system managers, because they do not observe customer abandonments.

Finally, the system is often much more complicated: For one example, there may be multiple customer classes and multiple service pools with some form of skill-based routing (SBR); see Gans et al. (2003). For a second example, with web chat, servers may serve several customers simultaneously, different servers may participate in a single service, and there may be interruptions in the service times, as the customers explore material on the web in between conversations with agents. For a third example, when delays are large – which is when we most want to make delay announcements – customers often abandon from queue. In these more complicated settings, the queue length is typically known, but the rate customers enter service is often not known and/or difficult to estimate reliably. That causes problems for the QL estimator.

When the $GI/M/s$ model is not appropriate for one of these reasons, the QL estimator may not perform well.

Example (non-exponential service times). To dramatically illustrate the possible difficulties with the QL delay estimator in the presence of a non-exponential service-time distribution (without trying to be realistic), we consider a limiting hyperexponential (H_2) distribution, in which each service time is either an exponential with mean 10, with probability $1/10$, or the deterministic value 0, with probability $9/10$. Thus the service time has mean 1, but busy servers will only be serving customers with the exponential distribution. Let $s = 100$ and suppose that an arrival finds the queue empty but all the servers busy. Then the QL delay estimate for this new arrival is $1/s = 1/100$, but the actual delay is exponentially distributed with mean $1/10$ (the minimum of 100 exponential random variables, each with mean 10). Hence, the actual mean delay is ten times greater than predicted by the QL delay estimator. Consistent with this extreme example, we have found that our alternative delay estimators actually outperform the QL delay estimator in the $D/H_2/100$ model with moderately variable

H_2 distributions. ■

Similarly, when there is a large amount of customer abandonment, the QL estimator will tend to overestimate the potential delay (the delay assuming that the customer has infinite patience), because many customers in queue may abandon before entering service, and the standard QL estimator fails to take that into account. As discussed in Whitt (1999), the QL estimator can be revised to provide an accurate estimate of delays with abandonments when the time-to-abandon distribution is exponential. However, as discussed in Whitt (2006), the performance measures in the overloaded $M/M/s + GI$ model, with non-exponential time-to-abandon distribution, depend strongly on the time-to-abandon distribution beyond its mean. Since the time-to-abandon distribution has been found to be non-exponential in practice, see Brown et al. (2005), there also are potential difficulties with the generalized QL estimator based on the $M/M/s + M$ model. We investigate alternative delay estimators in the presence of abandonments in a sequel to this paper, Ibrahim and Whitt (2008). There we give examples with non-exponential distributions in which both the standard QL estimator and the refinement for the $M/M/s + M$ model are outperformed by delay estimators based on recent delay history.

From the above discussion, we conclude that other estimators besides the standard QL estimator are worth considering; we do not conclude that the standard QL estimator or other estimators based on the queue length are necessarily bad. Indeed, we will show advantages of the QL estimator when it can be used.

This Study. Here we study the performance of the delay estimators based on delay history in the relatively simple idealized setting of the $GI/M/s$ model. Motivated by call centers, we are especially interested in the case of large s , but we consider all possible s .

For this more elementary $GI/M/s$ model, we obtain strong analytical results and make comparisons through computer simulations. Unlike Armony et al. (2006), here we do not consider customer response and we do not consider balking or customer abandonment, although we recognize that those phenomena are important. Moreover, here we are not concerned with what to announce, for which we should consider interpretation and response, but only with the effectiveness of the candidate delay estimators in predicting the actual delay encountered (assuming no customer response).

We find that the conditional distribution of the delay to be estimated, given the observed past delay, is often approximately normally distributed, implying that the conditional distribution is approximately characterized by its mean and variance. The observed delay is the

natural *direct estimator* of the delay to be encountered by the new arrival, while the mean of the conditional distribution of the delay of the new arrival, given that observed delay, is a natural *refined estimator* based on the same information. (In general, these are different!) The refined estimator depends on the model and its parameters. Since the conditional mean is complicated, we develop approximations for it.

For the $GI/M/s$ model, we will show that the QL estimator does indeed perform better than the alternative estimators based on recent delays, and we will quantify the difference. Roughly, the MSE differs by the constant factor $c_a^2 + 1$, where c_a^2 is the squared coefficient of variation (SCV, variance divided by the square of the mean) of an interarrival time. Thus, the MSE's of the delay-history estimators are about the same as the MSE of the QL estimator when the arrival-process variability is low, but considerably greater when the arrival-process variability is high.

Related Literature. There is a large body of related literature with somewhat different goals. We are doing statistical inference for queues, but as in Avramidis et al. (2004), Brown et al. (2005) and Glynn and Whitt (1989), most statistical inference for queues aims to estimate the *model* or the *steady-state performance*. There is an interesting stream of literature related to estimating *past* performance in a partially observed system from transactional data, stemming from Larson (1990). There has been much interesting recent inference work, including delay estimation, related to the *Internet*, as surveyed by Coates et al. (2002), but our setting and time scales tend to be very different. In addition to Whitt (1999), delay estimation for real-time delay prediction is investigated by Ward and Whitt (2000) and Nakibly (2002); these focus on processor-sharing and priority disciplines, respectively. Our real-time focus is in the spirit of real-time queueing, as in Doytchinov et al. (2000) and references therein.

Organization of the paper. We start in §2 by defining alternative delay estimators based on recent delay history and giving some expressions for them for the $GI/M/s$ model. We present results of initial simulation experiments in §3. We establish properties of two basic delay estimators – LES and the Head-of-the-Line (HOL) estimator – in §4. We present confirming simulations related to those analytical results in §5. We discuss insights from heavy-traffic limits in §6. Finally, we draw conclusions in §7. We present additional material in the e-companion, including more experimental results, more heavy-traffic limits and a cautionary example showing the possible pitfalls of the LES and HOL delay estimators for non-exponential

service-time distributions. We present even more experimental results in an online supplement available on the authors' web pages, Ibrahim and Whitt (2007).

2. Alternative Estimators

The $GI/M/s$ Model. We now specify the $GI/M/s$ model: The service times are independent and identically distributed (i.i.d.) exponential random variables V_n with mean 1. The interarrival times are i.i.d. positive random variables U_n with a non-lattice cumulative distribution function (cdf) F . (We will also consider the deterministic arrival process, which violates this condition; consequently, it will require slightly different analysis.) We omit the subscripts from U and V when the specific index is not important. Let F have finite third moment, characterized by $\nu_3^a \equiv E[U^3]/(E[U])^3$. Then F necessarily has finite first and second moments. Assume that $E[U] = 1/(s\rho)$, where s is the number of servers and $\rho \equiv E[V]/(sE[U])$ is the traffic intensity. Let F have SCV $c_a^2 \equiv Var(U)/(E[U])^2$. Let $A \equiv \{A(t) : t \geq 0\}$ be the renewal counting process (arrival process) associated with U_n , defined by

$$A(t) \equiv \max \{n \geq 0 : U_1 + \dots + U_n \leq t\}, \quad t \geq 0. \quad (2.1)$$

The $GI/M/s$ system is well known to be stable, and have a proper limiting steady-state behavior, if and only if $\rho < 1$. All our simulation results are for the $GI/M/s$ model in steady state, even though the estimation procedures apply more generally.

The No-Information (NI) Steady-State Estimator. The candidate delay estimators differ depending on the information used. If no information at all is used beyond the model, then it is natural to use the steady-state distribution. In particular, with W_∞ denoting the steady-state waiting time before beginning service, the no-information (NI) steady-state delay estimator for a customer that must wait before beginning service is $\theta_{NI} \equiv E[(W_\infty | W_\infty > 0)]$. It serves as a useful reference point. Any other estimator exploiting additional real-time information should do at least as well to be worth serious consideration.

For the $GI/M/s$ model, it is well known that $(W_\infty | W_\infty > 0)$ has an exponential distribution – see §XII.3 of Asmussen (2003) – so that the SCV is 1. Since the SCV is 1, the NI estimator is quite highly variable, and so necessarily has low predictive power. For the $M/M/s$ special case, the mean is $1/s(1 - \rho)$, so that $MSE = Var((W_\infty | W_\infty > 0)) = 1/s^2(1 - \rho)^2$.

The Full-Information Queue-Length (QL) Delay Estimator. The other extreme would be full-information at the arrival epoch, which we take to mean that we know: (i) the queueing

model, (ii) the number of customers in the system at that arrival epoch and (iii) the elapsed service times of all customers in service. If we knew the remaining service times as well, then we could compute the exact delay, but we assume that the remaining service times are unknown. Of course, for exponential service times, the elapsed service times give no useful information about the remaining service times because of the lack-of-memory property of the exponential distribution. Thus the (full-information) queue-length (QL) estimator for the $GI/M/s$ model only exploits the queue-length $Q(t)$ and knowledge of the model.

Let $W_Q(n)$ represent a random variable with the conditional distribution of the delay of a new arriving customer at some time t , given that the arriving customer must wait before starting service and given that the queue length at that time (not counting the new arrival) is $Q(t) = n$. (For $n \geq 1$, the customer must necessarily wait; for $n = 0$ our conditioning implies that all servers are busy but the queue length is 0.) For the $GI/M/s$ model, the random variable $W_Q(n)$ can be represented as

$$W_Q(n) \equiv \sum_{i=1}^{n+1} (V_i/s) , \quad (2.2)$$

when $Q(t) = n$. The natural QL delay estimator, based on the observed queue length $Q(t) = n$, is the mean $\theta_{QL}(n) \equiv E[W_Q(n)] = (n+1)/s$. The QL estimator requires knowledge of s and the mean service time $E[V]$ (here taken to be 1) as well as $Q(t)$.

We have the division by s in (2.2) because the times between successive service completions when all servers are busy are i.i.d. random variables distributed as the minimum of s exponential random variables, each with mean 1, which makes the minimum exponential with mean $1/s$. It is significant that this estimator is independent of the arrival process and thus also of the traffic intensity. It applies equally well to steady-state and transient settings.

As discussed in Whitt (1999), $W_Q(n)$ has the desirable property that the estimation gets relatively more accurate as the observed queue length n increases:

$$E[W_Q(n)] = \frac{n+1}{s}, \quad Var[W_Q(n)] = \frac{n+1}{s^2} \quad \text{and} \quad c_{W_Q(n)}^2 \equiv \frac{Var[W_Q(n)]}{(E[W_Q(n)])^2} = \frac{1}{n+1}, \quad (2.3)$$

so that $c_{W_Q(n)}^2 \rightarrow 0$ as $n \rightarrow \infty$.

Thus, whenever the queue length is large, the QL estimator $E[W_Q(n)]$ will be relatively accurate. If we consider heavy-traffic regimes, where the queue length approaches infinity, as we will do later, then this QL delay estimator will perform well. For example, the halfwidth of a 95% confidence interval is about $2/\sqrt{n}$, which is about 20% of a mean conditional waiting time when $n = 100$. Such a large value of n is not uncommon when s too is large.

For the $M/M/s$ model, there is a simple expression for the average MSE in steady state, which helps judge the performance of other estimators; the MSE's for the other delay estimators should all fall between the QL estimator (best possible) and the NI estimator (worst possible, knowing the model). Let Q_∞^w be a random variable with the conditional distribution of the steady-state queue length upon arrival given that the customer must wait before beginning service. In the $M/M/s$ model, $Q_\infty^w + 1$ has a geometric distribution with mean $1/(1 - \rho)$. That is easily deduced from the time reversibility of the $M/M/s$ model, which implies that Q_∞^w has the steady state distribution of the number in system in an $M/M/1$ queue with traffic intensity ρ ; e.g., see Proposition 5.6.3 of Ross (1996). Hence,

$$E[MSE(W_Q(Q_\infty^w))] \equiv \sum_{n=0}^{\infty} MSE(W_Q(n))P(Q_\infty^w = n) = E[Var(W_Q(Q_\infty^w))] = \frac{1}{s^2(1 - \rho)}, \quad (2.4)$$

so that the ratio between the worst possible NI MSE and the best possible QL MSE is

$$\frac{MSE(\theta_{NI})}{MSE(\theta_{QL}(Q_\infty^w))} = \frac{Var(W_\infty | W_\infty > 0)}{E[Var(W_Q(Q_\infty^w))]} = \frac{1/s^2(1 - \rho)^2}{1/s^2(1 - \rho)} = \frac{1}{1 - \rho}. \quad (2.5)$$

For example, a case of principle interest for call centers has $s = 100$ and $\rho = 0.95$. Then the average MSE for NI is 20 times greater than the average MSE for QL. We will show that the delay-history estimators produce a corresponding ratio of approximately $c_a^2 + 1 = 2$.

The Last Customer to Enter Service (LES). The first candidate *direct* delay estimator is the delay (before starting service) of the last customer to *enter* service (LES). The direct LES estimator is appealing because it is relatively easy to obtain and interpret, but there also are a variety of *refined* LES estimators we can consider; all are based on the LES observation.

To a large extent, the alternative refined LES delay estimators (and others as well) are obtained by replacing the known queue length n in (2.2) by random variables that estimate the queue length, based on the available delay history. Let $W_{LES}(w, d)$ be the delay of a new arrival, given that the new arrival must wait before starting service and given that the last customer to enter service experienced delay w before entering service and there was elapsed time d since that customer entered service. Let t_a be the arrival epoch of the new customer and t_e be the time the last customer entered service prior to t_a . (Throughout this paper we use the fact that, almost surely, no two events – arrivals or service completions – will occur simultaneously.) Necessarily, $d = t_a - t_e$ and $t_e - w$ is the arrival epoch of the customer entering service at t_e . A key observation is that the queue length at time t_e must be distributed as $A(w)$, because customers enter service from the queue in order of arrival. However, $W_{LES}(w, d)$ has

a relatively complicated exact distribution, because we do not know precisely what happens in the interval $[t_e, t_a]$.

If we impose an extra condition, then this random variable $W_{LES}(w, d)$ has a relatively simple distribution. The **extra condition** is that the epoch t_e is also simultaneously the last service completion prior to t_a . That extra condition will necessarily hold if at least one customer remains in the queue at time t_e . In turn, that sufficient condition is very likely to be satisfied if w is relatively large (the case of primary interest). Under the extra condition that t_e is also the last service completion before t_a , we have the simple representation

$$W_{LES}(w, d) \equiv \sum_{i=1}^{A(w+d)+1} (V_i/s), \quad (2.6)$$

where the summands are i.i.d. and independent of $A(w+d)$, because the queue length seen by the new arrival at time t_a will be $A(w+d)$, the number of arrivals in the interval of length $w+d$ preceding the arrival epoch t_a . Formula (2.6) allows us to characterize the distribution of $W_{LES}(w, d)$, under the assumed extra condition. Just like (2.2), (2.6) requires knowledge of s and the mean service time as well as w . Here we also require knowledge of the renewal arrival process or, equivalently, the interarrival-time distribution.

An important reference point for the refined LES estimator in (2.6) is the $D/M/s$ model, with a deterministic arrival process, having constant interarrival times, because under the extra condition leading to (2.6), we then have $W_{LES}(w, d) = W_Q(Q(t_a))$, since $A(w+d) = Q(t_a)$, making (2.2) coincide with (2.6). Thus we see that the loss of efficiency in going from QL to LES (direct or refined) is primarily due to the variability in the arrival process.

We assume that the experienced LES waiting time w is always available, but we might not know d , so that we might want to consider as an alternative refined estimator the mean of the random variable $W_{LES}(w)$, which assumes d is unavailable, but dropping d makes the distribution even more complicated. If we can assume that $w \gg d$, then there should be negligible difference. In general, we have the natural approximations based on (2.6):

$$W_{LES}(w) \approx \sum_{i=1}^{A(w+(V_0/s))+1} (V_i/s) \approx \sum_{i=1}^{A(w+(1/s))+1} (V_i/s), \quad (2.7)$$

where V_0 is an exponential random variable with mean 1 independent of V_i for $i \geq 1$, because the time between successive service completions when all servers are busy is distributed as V_0/s . (Assuming that the queue is nonempty at time t_e , that time is a service completion epoch. Then d is the age of the Poisson all-servers-busy departure process with rate s under

Poisson inspection by the arrival process.) The second approximation is obtained by inserting the expected value. It is also based on the extra condition, which will hold approximately for large w .

The Head-Of-The-Line (HOL) Estimator. A second candidate direct delay estimator, which is closely related to the direct LES estimator, is the elapsed waiting time of the customer at the *head of the line* (HOL) (queue), assuming that there is at least one customer waiting at the new arrival epoch. The direct HOL delay estimator was used as an announcement in an Israeli bank studied by Mandelbaum et al. (2000) and mentioned as a candidate delay announcement by Nakibly (2002). It is appealing compared to LES because the conditional distribution of the delay to be estimated is more tractable given the HOL information.

The customer at the head of the line will enter service after the next service completion. That remaining time is exponential with mean $1/s$. Let $W_{HOL}(w)$ be a random variable with the conditional distribution of the waiting time (before starting service) of a new arrival given that the new arrival must join the queue, given that there already is at least one customer in queue, and given that the customer at the head of the line has already spent time w in queue. The random variable $W_{HOL}(w)$ is closely related to the random variable $W_{LES}(w, d)$, but has the advantage that we do not need to use d . Moreover, we do not need to impose the extra condition that we made for $W_{LES}(w, d)$, but instead we need to impose a new one: The **extra condition** now is the assumption that there is at least one customer in queue at the arrival epoch t_a ; otherwise there would be no customer at the head of the line. We propose the random variable $W_{HOL}(w)$ as an approximation for the random variable $W_{LES}(w)$ where we omit the lag d , as well as for its own sake. Closely paralleling the previous formulas, we have

$$W_{HOL}(w) \equiv \sum_{i=1}^{A(w)+2} (V_i/s) . \quad (2.8)$$

The Delay of the Last Customer to Complete Service (LCS). A third candidate direct delay estimator is the delay of the last customer to *complete* service (LCS). We naturally would want to consider this alternative estimator if we only learn customer delay experience after they complete service. That might be the case for customers and outside observers.

Let $W_{LCS}(w, v, d)$ be the delay of a new arrival, given that the new arrival must wait before starting service and given that the last customer to complete service experienced delay w before entering service, had individual service time v , and there was elapsed time d since that customer completed service. As before, let t_a be the arrival epoch of the new customer;

let t_c be the time the last customer completed service prior to t_a . The mean of the random variable $W_{LCS}(w, v, d)$ is a natural refined estimator, but this random variable has a relatively complicated distribution. Some data may be unavailable, so that we may want to consider as alternative refined estimators the means of the random variables $W_{LCS}(w, d)$, which assumes v is unavailable, and $W_{LCS}(w)$, which assumes that neither v nor d is available. Dropping v or the pair (v, d) makes the representation even more complicated.

The Delay of the Most Recent Arrival to Complete Service (RCS). Under some circumstances, the LCS and LES direct estimators will be similar, but they actually can be very different when s is large, because the last customer to complete service may have experienced his waiting time much before the last customer to enter service. We emphasize that customers need not depart in order of arrival. Indeed, with exponential service times, when all s servers are busy, each of the s servers is equally likely to generate the next service completion. Thus, for large s the LCS estimator is not really a viable alternative, as we will show. Consequently, we propose other candidate delay estimators based on the delay experience of customers that have already completed service. The first is the delay experienced by the customer that arrived most recently (and thus entered service most recently) among those customers who have already completed service (RCS). We find that RCS is far superior to LCS when s is large.

Among the Last $c\sqrt{s}$ Customers to Complete Service (RCS- $c\sqrt{s}$). A disadvantage of the RCS estimator is that we must analyze a lot of data, going arbitrarily far back in the past. From heavy-traffic analysis in §6 and the e-companion, we deduce that the most recent arrival time of a customer that has completed service is very likely to occur among the last $c\sqrt{s}$ customers when s is large (and the system is normally loaded). So we introduce a new estimator, which requires less information processing: Let RCS- $c\sqrt{s}$ be the delay of the customer to have arrived most recently among the last $c\sqrt{s}$ customers who have already completed service. Clearly, these last three estimators LCS, RCS and RCS- $c\sqrt{s}$ are complicated, so that we primarily rely on simulation to evaluate their relative performance. Through extensive simulation experiments, we found that the average squared error of RCS- $c\sqrt{s}$ is essentially identical to that of RCS when $c = 4$, differs by at most 1% when $c = 2$ and differs by at most 10% when $c = 1$.

Averages. Our main estimators are individual delays experienced by a recent customer, rather than an average over many past delays. Only the no-information steady-state estimator

($W_\infty | W_\infty > 0$) can be said to use averages. We can extend the LES, LCS, RCS and RCS- $c\sqrt{s}$ estimators to get LES- k , LCS- k , RCS- k and RCS- $c\sqrt{s} - k$ by averaging over the last k experienced delays. With the exception of LCS with large s (which does not have desirable properties), we have found that averages do not help, when the delays are relatively large (the case of primary interest to us). There is a simple explanation: When delays are large, the delays change relatively slowly compared to the size of the delays. Theoretically, this can be explained by the heavy-traffic snapshot principle; see Section 6. In this setting it is better to use recent information than to eliminate noise by averaging.

3. Initial Simulation Experiments: Comparing the Estimators

In this section we present initial simulation experiments, aiming to compare the alternative estimators defined in §2. We focus on the *average squared error* (ASE) of the estimator, defined in (1.3). For large samples, the ASE should agree with the MSE in steady state.

Table 1 shows the ASE's for seven different delay estimators in the $GI/M/s$ model with $s = 100$. We consider three categories of estimators: (i) the two reference estimators QL and NI , (ii) the direct delay estimators LES and HOL, and (iii) the three estimators based on delays of customers who have already completed service - LCS, RCS and RCS- \sqrt{s} . We consider three interarrival-time distributions - M , D and H_2 - and four values of the traffic intensity ρ - 0.98, 0.95, 0.93 and 0.90. The H_2 distribution has SCV $c_a^2 = 4$ and balanced means (the two component exponential distributions contribute equally to the mean). We performed 10 independent replications of long runs in each case. The half width of the 95% confidence interval is shown below each estimate. Corresponding results for other values of s - 1, 10, 400 and 900 - are contained in the online supplement, Ibrahim and Whitt (2007). The cases $s = 10$ and $s = 1$ are shown in the e-companion.

These estimators appear in Table 1 with the better performance toward the left; i.e., in terms of efficiency (low ASE), the estimators are ordered by

$$QL > LES \approx HOL > RCS \approx RCS - \sqrt{s} > LCS > NI . \quad (3.1)$$

As expected, the full-information QL estimator performs best, while the no-information NI estimator performs worst. The performance of LES and HOL are very close, while the performance of RCS and RCS- \sqrt{s} are very close. The QL estimator is significantly better than LES; LES is slightly better than RCS; RCS is significantly better than LCS; and LCS is significantly better than NI. Very roughly, $ASE(LES)/ASE(QL) \approx (c_a^2 + 1)/\rho$, so LES performs nearly as

*Estimated ASE in units of 10^{-3}
 $M/M/s$ model with $s = 100$*

ρ	QL	LES	HOL	RCS	RCS- \sqrt{s}	LCS	NI
0.98	5.03 ± 0.02	10.2 ± 0.05	10.2 ± 0.05	12.5 ± 0.05	12.9 ± 0.05	26.7 ± 0.06	255 ± 36
0.95	2.04 ± 0.02	4.3 ± 0.05	4.3 ± 0.05	6.4 ± 0.05	6.7 ± 0.05	16.5 ± 0.06	41.8 ± 2.7
0.93	1.44 ± 0.002	3.07 ± 0.003	3.08 ± 0.003	5.06 ± 0.003	5.32 ± 0.003	13.1 ± 0.13	20.8 ± 1.2
0.90	0.99 ± 0.003	2.2 ± 0.006	2.2 ± 0.006	3.9 ± 0.008	4.2 ± 0.009	9.4 ± 0.27	9.7 ± 0.7

$D/M/s$ model with $s = 100$

ρ	QL	LES	HOL	RCS	RCS- \sqrt{s}	LCS	NI
0.98	2.48 ± 0.05	2.62 ± 0.05	2.62 ± 0.05	3.77 ± 0.05	3.94 ± 0.05	10.3 ± 0.11	61.5 ± 3.9
0.95	1.01 ± 0.02	1.15 ± 0.02	1.15 ± 0.02	2.20 ± 0.03	2.34 ± 0.03	6.38 ± 0.12	10.1 ± 0.40
0.93	0.73 ± 0.02	0.87 ± 0.02	0.87 ± 0.02	1.85 ± 0.03	1.96 ± 0.03	4.90 ± 0.13	5.20 ± 0.32
0.90	0.52 ± 0.015	0.67 ± 0.016	0.66 ± 0.017	1.54 ± 0.035	1.63 ± 0.037	3.44 ± 0.15	2.68 ± 0.23

$H_2/M/s$ model with $s = 100$

ρ	QL	LES	HOL	RCS	RCS- \sqrt{s}	LCS	NI
0.98	12.4 ± 0.70	60.4 ± 3.2	60.4 ± 3.2	66.1 ± 3.2	67.0 ± 3.2	103.4 ± 34.0	1505 ± 226
0.95	4.82 ± 0.095	22.5 ± 0.46	22.5 ± 0.47	27.7 ± 0.45	28.4 ± 0.45	56.3 ± 0.58	243.3 ± 22.7
0.93	3.44 ± 0.094	15.5 ± 0.44	15.5 ± 0.44	20.4 ± 0.49	21.1 ± 0.50	44.5 ± 1.02	121.4 ± 10.2
0.90	2.35 ± 0.040	10.2 ± 0.21	10.2 ± 0.21	14.6 ± 0.24	15.2 ± 0.24	33.1 ± 0.53	55.4 ± 2.9

Table 1: A comparison of the efficiency of different real-time delay estimators for the $GI/M/100$ queue as a function of the traffic intensity ρ and the interarrival-time distribution (M , D and H_2). Only the direct estimators are considered. Estimates of the average squared error ASE are shown together with the half width of the 95% confidence interval. The units are 10^{-3} throughout.

well as QL for low-variability arrival processes such as the D arrival process, but much worse for high-variability arrival processes such as the H_2 arrival process.

It is instructive to look at the relative average squared error (RASE), which is obtained by dividing the ASE by $E[W_\infty|W_\infty > 0]^2$, because the associated steady-state relative mean squared error (RMSE), defined as $MSE/E[W_\infty|W_\infty > 0]^2$, is *linear* as a function of ρ for the QL estimator: $RMSE(QL) = (1 - \rho)$. (The RMSE is identically 1 for the NI estimator.) We show the RASE plots for the $D/M/100$ model in Figure 1. With the D arrival process, LES and HOL are virtually identical (with the plots lying on top of each other), so we only show LES. Both LES and HOL are nearly as good as QL and much better than RCS; LCS is so bad that it is not even shown. Corresponding plots for other interarrival-time distributions and other s appear in the online supplement. The plots for the $M/M/100$ and $H_2/M/100$ models are in the e-companion.

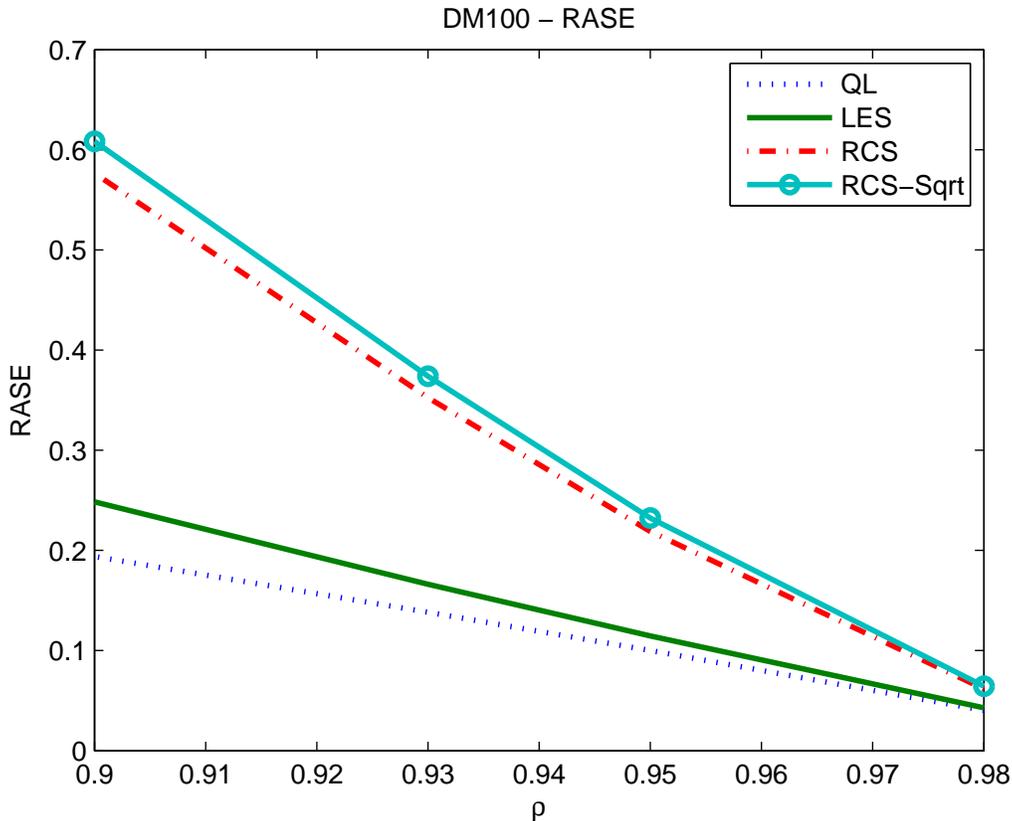


Figure 1: The relative average squared error ($RASE$) for the $D/M/100$ model.

Experience shows that the NI estimator performs especially poorly in very heavy traffic, while LCS performs especially poorly with large s in light traffic. For large s and small ρ , LCS

Conditional ASE for the $M/M/100$ model in units of 10^{-3}
Observed delays in between $4E[W|W > 0]$ and $6E[W|W > 0]$

ρ	QL	LES	HOL	RCS	RCS- \sqrt{s}	LCS	NI
0.99	49.4 ± 7.0	86.6 ± 6.9	86.3 ± 6.9	89.4 ± 7.2	90.1 ± 7.2	108.8 ± 10.2	11,586 ± 1250
0.98	24.8 ± 1.8	47.5 ± 3.1	47.3 ± 3.0	50.1 ± 3.1	50.6 ± 3.1	69.6 ± 3.7	3,542 ± 431
0.95	10.5 ± 0.23	20.4 ± 0.63	20.1 ± 0.62	23.5 ± 0.82	24.0 ± 0.80	50.4 ± 3.3	564 ± 27
0.93	7.54 ± 0.20	15.2 ± 0.31	14.9 ± 0.29	18.7 ± 0.43	19.3 ± 0.45	52.0 ± 3.2	286 ± 8.0
0.90	5.62 ± 0.21	11.1 ± 0.38	10.7 ± 0.38	15.3 ± 0.61	16.1 ± 0.66	50.9 ± 25.2	137.4 ± 6.7

Table 2: A comparison of the efficiency of different real-time delay estimators conditional on the level of delay experienced for the $M/M/100$ model as a function of the traffic intensity ρ . Actual delays are considered that fall in the interval $(4E[W|W > 0], 6E[W|W > 0])$. Estimates of the conditional average squared error ASE are shown together with the half width of the 95% confidence interval. The units are 10^{-3} throughout.

even performs worse than the NI estimator. There is only one case in Table 2; more cases can be seen when $s = 400$ and $s = 900$ in the supplement.

Since delay estimates are more relevant when the observed delays in the system are longer, it is natural to consider the behavior of the estimators for larger delays. We have complemented the experiments described above by considering how the delay estimators perform when we only consider actual delays that fall in one of the intervals: $(E[W|W > 0], 2E[W|W > 0])$, $(2E[W|W > 0], 4E[W|W > 0])$, $(4E[W|W > 0], 6E[W|W > 0])$ and $(6E[W|W > 0], \infty)$. Table 2 illustrates the results for the $M/M/100$ model when the observed delays fall in the interval $(4E[W|W > 0], 6E[W|W > 0])$. Other cases appear in the online supplement. The performance of the estimators for these larger delays is approximately as in Table 1. As should be expected, the NI estimator fares even worse in this comparison.

4. Analysis of the HOL and LES Estimators

The representation (2.8) allows us to characterize the probability distribution of the random variable $W_{HOL}(w)$, which we do both for its own sake and as an approximation for the random variables $W_{LES}(w)$ and $W_{RCS}(w)$. When we use the HOL estimator, we assume that there is at least one customer in queue at the new arrival epoch t_a . Very similar formulas hold for the LES estimator based on formula (2.6), under the extra assumption given there. Since the formulas are virtually identical, we do not display separate results for LES.

We emphasize that the random variable $W_{HOL}(w)$ applies to both transient and steady-state scenarios. We can have arbitrary traffic intensity ρ , including $\rho > 1$, under which there is no proper steady state. We assume that the renewal arrival process $\{A(t) : t \geq 0\}$ and the traffic intensity ρ are specified and unchanging in the interval $[t_a - w, t_a]$, which is the relevant system history for our estimation at time t_a .

We start by showing that the distribution of $W_{HOL}(w)$ depends on s in a relatively simple way. For that purpose, we introduce an extra subscript s to indicate the dependence upon s , getting $W_{HOL,s}(w)$. Let $\stackrel{d}{=}$ denote equality in distribution.

Theorem 4.1. (*dependence upon s*) *For the GI/M/s model,*

$$W_{HOL,s}(w) \stackrel{d}{=} \frac{W_{HOL,1}(sw)}{s} \quad (4.1)$$

for all ρ , w and s .

Proof. We show the equality in distribution by establishing equality w.p.1 for a special construction. We construct a convenient family of systems indexed by s . For each s , let the service times be exponential random variables V_n with mean 1 as before. Start by defining interarrival times U_n with mean $1/\rho$ to use for the case of $s = 1$. Then in the system with $s > 1$, let the n^{th} interarrival time be U_n/s . Let $\{A_s(t) : t \geq 0\}$ be the renewal counting process in system s , having interarrival times U_n/s . Then $A_s(w/s) = A_1(w)$ for all s and w ; since we have re-scaled the interarrival times, we just re-scale time in the associated renewal counting process. This construction yields equality for the random variables in (4.1) and all $w \geq 0$. Since the distribution is independent of the construction, that implies the claimed relation (4.1). ■

We now show that we get relatively simple asymptotic expressions characterizing the distribution of $W_{HOL,s}(w)$ when $sw \rightarrow \infty$. That applies when $w \rightarrow \infty$ for fixed s , but it also can apply when $s \uparrow \infty$ and $w \downarrow 0$, as occurs in the QED many-server heavy-traffic limiting regime, to be discussed in §6; then $w = O(1/\sqrt{s})$ so that $sw \rightarrow \infty$ while $w \rightarrow 0$.

Let $N(m, \sigma^2)$ denote a normally distributed random variable with mean m and variance σ^2 . Let \Rightarrow denote convergence in distribution.

Theorem 4.2. (*distribution of $W_{HOL,s}(w)$*) *Consider the GI/M/s queue with traffic intensity ρ operating in the time interval $[t_a - w, t_a]$. (a) For any $\rho > 0$, $s \geq 1$ and $w > 0$,*

$$E[W_{HOL,s}(w)] = \frac{E[A(w)] + 2}{s} \quad (4.2)$$

and

$$\text{Var}[W_{HOL,s}(w)] = E[A(w) + 2]\text{Var}(V/s) + \text{Var}(A(w) + 2)(E[V/s])^2. \quad (4.3)$$

(b) If the arrival process is Poisson, then

$$E[W_{HOL,s}(w)] = \rho w + \frac{2}{s} \quad (4.4)$$

and

$$\text{Var}[W_{HOL,s}(w)] = \frac{2\rho w}{s} + \frac{2}{s^2}, \quad (4.5)$$

so that

$$c_{W_{HOL,s}(w)}^2 = \frac{2}{\rho s w} - \frac{6}{(\rho s w)^2} + O\left(\frac{1}{(\rho s w)^3}\right) \quad \text{as } s w \rightarrow \infty. \quad (4.6)$$

(c) For a general renewal arrival processes with a non-lattice interrenewal-time distribution, if $s w \rightarrow \infty$, then

$$sE[W_{HOL,s}(w)] - \rho s w \rightarrow \frac{(c_a^2 + 3)}{2}, \quad (4.7)$$

$$\frac{W_{HOL,s}(w)}{w} \rightarrow \rho \quad \text{w. p. 1} \quad \text{and} \quad \frac{E[W_{HOL,s}(w)]}{w} \rightarrow \rho, \quad (4.8)$$

$$s^2 \text{Var}(W_{HOL,s}(w)) - \rho s w (c_a^2 + 1) \rightarrow \left(\frac{5(c_a^2 + 1)^2}{4} - \frac{2\nu_a^3}{3} + 1 \right), \quad (4.9)$$

$$s^2 E[(W_{HOL,s}(w) - \rho w)^2] - \rho s w (c_a^2 + 1) \rightarrow K, \quad (4.10)$$

$$s^2 E[(W_{HOL,s}(w) - w)^2] - (s w)^2 (1 - \rho)^2 - s w [(2\rho - 1)c_a^2 + 4\rho - 3] \rightarrow K, \quad (4.11)$$

where

$$K \equiv K(c_a^2, \nu_a^3) \equiv \left(\frac{3c_a^4}{2} + 4c_a^2 + \frac{9}{2} - \frac{2\nu_a^3}{3} \right), \quad (4.12)$$

$$s w c_{W_{HOL,s}(w)}^2 \rightarrow \frac{c_a^2 + 1}{\rho} \quad \text{and} \quad \frac{W_{HOL,s}(w) - \rho w}{\sqrt{\rho w (c_a^2 + 1)/s}} \Rightarrow N(0, 1). \quad (4.13)$$

Proof. Since $W_{HOL}(w)$ in (2.8) is a random sum of i.i.d. random variables, where $A(w)$ is independent of the summands V_i/s , we have (4.2). Formula (4.3) follows from the conditional variance formula, e.g., p. 51 of Ross (1996). For (4.6), we use elementary operations on series, as in 3.6.22 in Abramowitz and Stegun (1972). When we let $s w$ increase, we first apply Theorem 4.1 to reduce the analysis to the case $s = 1$. Henceforth assume that $s = 1$. When we restrict attention to $s = 1$, it suffices to simply let $w \rightarrow \infty$. When we let w increase,

$$E[A(w) + 2] - \rho w \rightarrow \frac{(c_a^2 + 1)}{2} + 1 \quad \text{as } w \rightarrow \infty, \quad (4.14)$$

see Corollary 3.4.7 of Ross (1996) or (2.7) and (2.8) of Whitt (1982), which review a classic result. Combining (4.14) and (4.2) gives (4.7), which immediately implies the second limit in

(4.8). For the w.p.1 limit in (4.8), we apply the strong law of large numbers for the partial sums of V_n and the renewal arrival process $A(w)$: With probability one,

$$\frac{\sum_{i=1}^n V_i}{n} \rightarrow E[V] = 1 \quad \text{and} \quad \frac{A(w) + 2}{w} \rightarrow \frac{1}{E[U]} = \rho, \quad (4.15)$$

so that

$$\frac{\sum_{i=1}^{A(w)+2} V_i}{w} = \frac{A(w) + 2}{w} \times \frac{\sum_{i=1}^{A(w)+2} V_i}{A(w) + 2} \rightarrow \rho \quad \text{w. p. 1.} \quad (4.16)$$

The asymptotic variance formula (4.9) follows from (4.3) and the asymptotic form of the variance for a renewal process, e.g., as in (2.7) and (2.8) of Whitt (1982):

$$\text{Var}(A(w)+2) = \text{Var}(A(w)) = \rho w c_a^2 + \frac{5(c_a^2 + 1)^2}{4} - \frac{2\nu_a^3}{3} - \frac{(c_a^2 + 1)}{2} + o(1) \quad \text{as } w \rightarrow \infty. \quad (4.17)$$

The associated limits (4.10) and (4.11) follow from (4.9). For (4.10), we use

$$\begin{aligned} E[(W_{HOL,s}(w) - \rho w)^2] &= \text{var}(W_{HOL,s}(w) - \rho w) + (E[W_{HOL,s}(w) - \rho w])^2 \\ &= \text{var}(W_{HOL,s}(w)) + (E[W_{HOL,s}(w) - \rho w])^2. \end{aligned} \quad (4.18)$$

The calculation for (4.11) is similar. The first limit in (4.13) follows immediately from (4.7) and (4.9). The central limit theorem in (4.13) follows from the central limit theorem for renewal-reward processes, e.g., Theorem 7.4.1 of Whitt (2002). We use the convergence-together theorem, Theorem 11.4.7 of Whitt (2002), to justify neglecting the asymptotically negligible terms. ■

Remark 4.1. (exact values by numerical inversion) It is possible to exploit (4.2) and (4.3) in order to compute the exact means and variances. To do so, we can exploit numerical transform inversion of Laplace transforms, as discussed in §13 of Abate and Whitt (1992). The Laplace transform of $E[A(t)]$ is $\hat{m}_1(s) \equiv \hat{f}(s)/[s(1 - \hat{f}(s))]$, where $\hat{f}(s)$ is the Laplace transform of the density function of the interarrival-time cdf F (here assumed to exist). The associated Laplace transform of $E[A(t)^2]$ is $2\hat{m}_1(s)^2 - \hat{m}_1(s)$, as can be seen from exercise XI.13 on p. 386 of Feller (1971). Since we are interested in estimation for relatively large delays, we will rely on the asymptotic approximations. ■

Remark 4.2. (nonhomogeneous Poisson arrival process) We can also analyze the random variable $W_{HOL,s}(w)$ in the case of a nonhomogeneous Poisson arrival process with intensity function $\{\lambda(t) : t \geq 0\}$. The exact relations (4.4) and (4.5) have natural extensions to that

case. We again have representation (2.8), but now with $A(w)$ being a Poisson random variable having mean

$$m_a(w) \equiv \int_{t_a-w}^{t_a} \lambda(t) dt, \quad (4.19)$$

which depends on the arrival time t_a and the intensity function as well as the experienced waiting time w . Unless we specify how the intensity function behaves, we have no simple asymptotic story as w increases, though. ■

Theorem 4.2 shows that the first-order asymptotic behavior of the random variable $W_{HOL,s}(w)$ as sw increases depends on the general interarrival-time distribution F only through its first two moments or, equivalently, through the mean $E[U] = 1/\rho s$ and the SCV c_a^2 . Equations (4.9) and (4.13) show that both the variance $Var(W_{HOL,s}(w))$ and the SCV $c_{W_{HOL,s}(w)}^2$ are approximately proportional to $c_a^2 + 1$ for large sw .

Theorem 4.2 shows that it may be useful to consider various refined estimators instead of the direct estimator $\theta_{HOL}^d \equiv w$. We would want to use the refined estimator $\theta_{HOL}^r \equiv E[W_{HOL,s}(w)]$, because the mean necessarily minimizes the MSE, but we do not have a convenient formula for the mean. Theorem 4.2 leads us to consider two other refined estimators: the *simple refined estimator* $\theta_{HOL}^{sr} \equiv \rho w$ and the *asymptotic refined estimator* $\theta_{HOL}^{ar} \equiv \rho w + (c_a^2 + 3)/(2s)$, based on the the limit (4.7) as $sw \rightarrow \infty$. Note that the formulas for the mean and variance for Poisson arrivals in (4.4) and (4.5) are exact, whereas the formulas for non-Poisson formulas are only approximations.

For fixed $\rho < 1$, the three refined estimators $\theta_{HOL}^r(w)$, $\theta_{HOL}^{sr}(w)$ and $\theta_{HOL}^{ar}(w)$ are all relatively consistent and asymptotically relatively efficient as $sw \rightarrow \infty$, whereas the direct HOL estimator w has neither of these properties. By *relatively consistent*, we mean that the ratio of the estimator to the quantity being estimated (here $W_{HOL,s}(w)$) converges to 1; by *asymptotically relatively efficient*, we mean that the relative mean squared error (RMSE $\equiv MSE/Mean^2$) converges to 0.

At first glance, the simple refined estimator looks very appealing, because it combines simplicity with good asymptotic properties. However, we found that the direct estimator consistently outperforms the simple refined estimator in experiments evaluating the steady-state performance for typical parameter values. Evidently, the extra constant term in θ_{HOL}^{ar} helps. The following (somewhat loosely stated) theorem supports that empirical observation. Let $MSE(\theta_{HOL}(W_\infty))$ denote the steady-state MSE of the estimator $\theta_{HOL}(w)$ when w is averaged with respect to the conditional delay $(W_\infty | W_\infty > 0)$, where W_∞ is the steady-state

delay.

Theorem 4.3. (*comparison of alternative HOL estimators*) Consider the $GI/M/s$ queue with traffic intensity $\rho < 1$ in steady state. If the arrival process is Poisson or if we take the limit in (4.7) as the exact mean, then the steady-state MSE's are ordered by

$$MSE(\theta_{HOL}^{ar}(W_\infty)) < MSE(\theta_{HOL}^d(W_\infty)) < MSE(\theta_{HOL}^{sr}(W_\infty)). \quad (4.20)$$

Moreover,

$$\begin{aligned} MSE(\theta_{HOL}^d(W_\infty)) - MSE(\theta_{HOL}^{ar}(W_\infty)) &= E \left[\left((1 - \rho)(W_\infty | W_\infty > 0) - \frac{(c_a^2 + 3)}{2s} \right)^2 \right] \\ &< \frac{(c_a^2 + 3)^2}{4s^2} = MSE(\theta_{HOL}^{sr}(W_\infty)) - MSE(\theta_{HOL}^{ar}(W_\infty)). \end{aligned} \quad (4.21)$$

Proof. The MSE formulas in (4.21) are obtained by directly adding and subtracting the mean inside the MSE formula, with the mean here regarded as being given by (4.7). The key inequality in (4.21) follows from a bound on the mean steady-state waiting time in the $GI/M/1$ queue. The conditional delay $(W_\infty | W_\infty > 0)$ in the $GI/M/s$ model has the same exponential distribution as in the $GI/M/1$ model; e.g., see p. 398 of Wolff (1989). Its mean is $(1 - \omega)^{-1}$, where ω is the root of the transform equation $\hat{f}(1 - \omega) = \omega$, where $\hat{f}(s)$ is the Laplace-Stieltjes transform of the interarrival-time cdf. However, it is known that $1 - \omega > 2(1 - \rho)/(c_a^2 + 1)$; e.g., apply Theorem 2 of Whitt (1984), noting that in the $D/M/1$ queue $1 - \omega > 2(1 - \rho)$, which follows from elementary inequalities for the exponential function: $e^{-2(1-\rho)} \geq 1 - 2(1 - \rho)$. From (4.21), we see that $MSE(\theta_{HOL}^d(W_\infty)) < MSE(\theta_{HOL}^{sr}(W_\infty))$ if and only if

$$E \left[\left((1 - \rho)(W_\infty | W_\infty > 0) - \frac{(c_a^2 + 3)}{2s} \right)^2 \right] < \frac{(c_a^2 + 3)^2}{4s^2}, \quad (4.22)$$

which, upon expanding the quadratic and using the fact that the second moment is twice the square of the first moment, holds if and only if

$$E[W_\infty | W_\infty > 0] < \frac{c_a^2 + 3}{s(1 - \rho)}, \quad (4.23)$$

which is implied by the delay bound. ■

To illustrate, we show numerical results in Table 3 for the candidate delay estimators θ_{HOL}^d , θ_{HOL}^{sr} and θ_{HOL}^{ar} in the $H_2/M/s$ model with $s = 100$ and $s = 1$. We display the values of their approximate MSE's in steady state predicted by formulas (4.11), (4.10) and (4.9), and we show the contributing terms, displayed in the order given in Theorem 4.2. In each case, one term

grows without bound as ρ increases while the other terms remains constant or nearly constant. We take the expected value of each MSE formula, where w is distributed randomly as the steady-state conditional delay ($W_\infty | W_\infty > 0$). We use the simulation estimates of the first two moments of the conditional delay. Table 3 is consistent with Theorem 4.3. As a consequence of Theorem 4.3, we suggest using the asymptotic refined estimator θ_{HOL}^{ar} .

We remark that the limit in (4.13) implies that $W_{HOL,s}(w)$ should be approximately normally distributed when sw is not too small. Our simulation experiments show that all the random variables $W_{HOL,s}(w)$, $W_{LES,s}(w)$ and $W_{RCS,s}(w)$ tend to be normally distributed when sw is not too small.

We can combine (4.13) and (2.3) to compare the efficiency of the QL and refined HOL estimators under high congestion. Let $W(t)$ be the virtual waiting time at time t , the time an arrival at time t would have to wait before beginning service. Since

$$W(t) = \sum_{i=1}^{Q(t)+1} (V_i/s), \quad (4.24)$$

the law of large numbers implies that $W(t)/Q(t) \rightarrow 1/s$ as $Q(t) \rightarrow \infty$. Thus, when $Q(t)$ is large, we have $W(t) \approx Q(t)/s$ (even if $W(t)$ itself is not large). Assuming that n is large with $w \approx n/s$ in (4.13) and (2.3), we have both sw and n large and

$$\frac{c_{W_{HOL,s}(w)}^2}{c_{W_{Q,s}(n)}^2} \approx \frac{(c_a^2 + 1)/\rho sw}{1/(n + 1)} \approx \frac{c_a^2 + 1}{\rho}. \quad (4.25)$$

Since we have introduced HOL partly as an approximation for LES, it is interesting to consider the difference between the HOL and LES observed delays and the difference between the random variables $W_{HOL,s}(w)$ and $W_{LES,s}(w, d/s)$. (We let $t_a - t_e = d/s$ because it should be proportional to $1/s$ with s servers.) First note that if at least one customer remains in queue after the last customer to enter service at time t_e , then the HOL customer at time t_e (after the customer entered service) will remain the HOL customer at time t_a . As a consequence, the HOL customer arrived immediately after the LES customer. Thus the HOL customer waits more than the LES customer by the time $t_a - t_e$ but less by the single interarrival time between them. Clearly these differences should become asymptotically negligible in the appropriate scaling.

We now compare the random variables $W_{HOL,s}(w)$ and $W_{LES,s}(w, d)$. We establish a stochastic bound between these random variables. Let \leq_{st} denote ordinary stochastic order; see §9.1 of Ross (1996). The following bound shows that the difference between $W_{HOL,s}(w)$ and $W_{LES,s}(w, d)$ is stochastically bounded and thus asymptotically negligible compared to w

Evaluating the alternative HOL estimators
 Approximations in the $H_2/M/100$ model

ρ	0.88	0.92	0.96	0.98
$E[W W > 0]$	0.1902	0.2964	0.6114	1.307
conf. int.	± 0.0030	± 0.0067	± 0.029	± 0.17
$E[W^2 W > 0]$	0.07205	0.1761	0.7446	3.436
conf. int.	± 0.0022	± 0.0095	± 0.060	± 0.67
$MSE(\theta^d)$	0.00826	0.0135	0.0293	0.0640
term 1	0.00103	0.00113	0.00119	0.00137
term 2	0.00677	0.0120	0.0276	0.0622
term 3	0.00045	0.00045	0.00045	0.00045
$MSE(\theta^{sr})$	0.00882	0.00141	0.00298	0.0645
term 1	0.00837	0.0136	0.0293	0.0640
term 2	0.00045	0.00045	0.00045	0.00045
$MSE(\theta^{ar})$	0.00759	0.0129	0.0286	0.0632
term 1	0.00837	0.0136	0.0293	0.0640
term 2	-0.000775	-0.000775	-0.000775	-0.000775

Approximations in the $H_2/M/1$ model

ρ	0.85	0.90	0.95	0.98
$E[W W > 0]$	15.01	23.50	48.64	115.7
conf. int.	± 0.18	± 0.42	± 1.6	± 8.80
$E[W^2 W > 0]$	446.2	1105.7	4707.1	25650.5
conf. int.	± 8.03	± 39.2	± 263.2	± 3280
$MSE(\theta^d)$	62.59	104.9	230.3	565.7
term 1	10.04	11.06	11.76	10.26
term 2	48.04	89.3	214.0	550.9
term 3	4.5	4.5	4.5	4.5
$MSE(\theta^{sr})$	68.31	110.3	235.5	571.6
term 1	63.81	105.8	231.0	567.1
term 2	4.5	4.5	4.5	4.5
$MSE(\theta^{ar})$	56.06	98.02	223.3	559.3
term 1	63.81	105.8	231.0	567.1
term 2	-7.75	-7.75	-7.75	-7.75

Table 3: Evaluation of the MSE approximations for the estimators θ_{HOL}^d , θ_{HOL}^{sr} and θ_{HOL}^{ar} in steady-state using (4.11), (4.9) and (4.10) together with simulation estimates of the first two moments of the conditional delay $E[W_\infty|W_\infty > 0]$. The $H_2/M/s$ model is considered as a function of the traffic intensity ρ for $s = 100$ and $s = 1$.

and these individual random variables as $sw \rightarrow \infty$. We say that a family of random variables $\{X(w) : w > 0\}$ is *stochastically bounded* if for any $\epsilon > 0$ there exists a positive constant $K(\epsilon)$ such that $P(|X(w)| > K(\epsilon)) < \epsilon$. By Markov's inequality, for nonnegative random variables it suffices to have the means $E[X(w)]$ uniformly bounded: $P(|X(w)| > K(\epsilon)) \leq E[X(w)]/K(\epsilon)$.

Theorem 4.4. (*bound on the difference between $W_{HOL,s}(w)$ and $W_{LES,s}(w, d/s)$)* Consider the GI/M/s model. Assume that there is at least one customer in queue at the new arrival epoch, so that (2.8) is valid for HOL and (2.6) is valid for LES. Then

$$W_{LES,s}(w, d/s) - X(s, w, d) \leq_{st} W_{HOL,s}(w) \leq_{st} W_{LES,s}(w, d/s) + X(s, w, d) , \quad (4.26)$$

where $X(s, w, d)$ is distributed as

$$X(s, w, d) \equiv \sum_{i=1}^{A(w+(d/s))-A(w)+1} (V_i/s) . \quad (4.27)$$

As $w \rightarrow \infty$ for fixed s , $E[X(s, w, d)] \rightarrow (\rho d + 1)/s$; as $sw \rightarrow \infty$, $E[X(s, w, d)]/w \rightarrow 0$. so that

$$\frac{|W_{HOL}(w) - W_{LES}(w, d)|}{w} \rightarrow 0 \quad \text{as} \quad sw \rightarrow \infty . \quad (4.28)$$

For the M/M/s model,

$$X(s, w, d) = \sum_{i=1}^{A(d/s)+1} (V_i/s) , \quad (4.29)$$

so that

$$E[X(s, w, d)] = (\rho d + 1)/s \quad \text{and} \quad \text{Var}(X(s, w, d)) = (2\rho d + 1)/s^2 . \quad (4.30)$$

Proof. Without altering the individual distributions of $W_{HOL,s}(w)$ and $W_{LES,s}(w, d/s)$, we can make a special construction in which we use exactly the same exponential random variables V_i/s for the two estimators. The random numbers of summands differ by $A(w + (d/s)) - A(w) - 1$, which is bounded above by $A(w + (d/s)) - A(w) + 1$, which we use in (4.27). Since the renewal process A has rate ρs , we can then apply Blackwell's renewal theorem, p. 155 of Asmussen (2003), to get $E[A(w + d/s) - A(w)] \rightarrow \rho d$ as $sw \rightarrow \infty$. Recall that we have assumed that the interarrival time cdf F is non-lattice. Hence we get $E[X(s, w, d)]/w \rightarrow 0$ as $sw \rightarrow \infty$, which implies (4.28). ■

5. Simulations Related to Theorem 4.2

Based on (4.11) in Theorem 4.2, we approximate the MSE of the direct HOL, LES and RCS estimators by

$$\text{MSE}(\theta_{HOL}^d(w)) \approx (1 - \rho)^2 w^2 + \frac{((2\rho - 1)c_a^2 + 4\rho - 3)w}{s} + \frac{K}{s^2} , \quad (5.1)$$

for K in (4.12). As above, let $MSE(\theta_{HOL}^d(W_\infty))$ denote the MSE in steady state, i.e., when we replace w in (5.1) by $(W_\infty|W_\infty > 0)$. We obtain

$$MSE(\theta_{HOL}^d(W_\infty)) \approx (1 - \rho)^2 E[W_\infty^2 | W_\infty > 0] + \frac{((2\rho - 1)c_a^2 + 4\rho - 3)E[W_\infty | W_\infty > 0]}{s} + \frac{K}{s^2}, \quad (5.2)$$

where W_∞ is the steady-state delay.

We have compared the ASE for HOL, LES and RCS to $MSE(\theta_{HOL}^d(W_\infty))$ and found close agreement, with the agreement being slightly better for HOL and LES than for RCS. In making this comparison, we substitute the simulation estimates of the two moments $E[W_\infty | W_\infty > 0]$ and $E[W_\infty^2 | W_\infty > 0]$ into (5.2). We must calculate or approximate these conditional moments in order to have a full approximation, but we do not consider that step here. We obtain good results comparing approximation (5.2) to the ASE for the cases of exponential (M), hyperexponential (H_2 with $c_a^2 = 4$) and Erlang (E_2) interrenewal-time distributions. We did experiments for $s = 1, 10, 100, 400, 900$, each for four values of ρ , increasing with s in order to represent typical cases. The errors were consistently less than 5% for HOL and LES in these experiments, as illustrated by the results for LES with M and H_2 interarrival-time distributions in Table 4.

Testing the MSE(HOL_∞) Approximations
in the $GI/M/100$ model

ρ	M	% diff.	D	% diff.	H_2	% diff.
0.98	10.20	-0.3%	2.67	-1.9%	62.8	-3.9%
0.95	4.20	1.4%	1.20	-4.1%	22.9	-1.9%
0.93	3.06	0.4%	0.92	-5.8%	15.9	-2.1%
0.90	2.20	-1.5%	0.72	-7.5%	10.5	-3.2%

Table 4: Evaluation of the approximations for the steady-state MSE of HOL in (5.2) and (5.4) by comparing to simulation estimates of the ASE for LES in the $GI/M/100$ model as a function of the interarrival-time distribution and the traffic intensity ρ . The simulation estimates appear in Table 1. The approximations in units of 10^{-3} and the relative percent differences are shown here.

We found that the approximation in (5.2) does not perform nearly as well for the case of a deterministic (D) arrival process, which should not be surprising, because the deterministic interrenewal-time distribution is a lattice distribution not covered by Theorem 4.2. Instead of (5.1), we propose the following approximation for the direct estimator with a D arrival process:

$$MSE(\theta_{HOL,D}^d(w)) \approx (1 - \rho)^2 w^2 + \frac{\rho w + (2/s)}{s}, \quad (5.3)$$

which is obtained by making the simple approximation $A(w) \approx \rho s w$. We then obtain the

following analog of the steady-state approximation (5.2):

$$MSE(\theta_{HOL,D}^d(W_\infty)) \approx (1 - \rho)^2 E[W_\infty^2 | W_\infty > 0] + \frac{\rho E[W_\infty | W_\infty > 0] + (2/s)}{s}. \quad (5.4)$$

Approximation (5.4) performs much better than approximation (5.2) with $c_a^2 = 0$, yielding errors of about 5% (ranging up to 11%), instead of about 5 – 25%, as shown in Table 4. For the refined estimator, we would also change the mean estimator to (4.4) instead of (4.7).

In order to evaluate the approximations for a specified observed delay w , we consider data from the simulation where the observed *HOL* delay falls in a small interval about $w \equiv 2E[W_\infty | W_\infty > 0]$. (We choose interval widths to make roughly reasonable, comparable sample sizes.) Table 5 shows the results of such an experiment for the *GI/M/100* model with $\rho = 0.95$. (The width of the sampling interval in each case was chosen to have roughly comparable sample sizes.) Table 5 shows that the approximations for the *HOL* conditional mean and variance are remarkably accurate approximations for all three estimators: *HOL*, *LES* and *RCS*, with the variance being slightly higher for *RCS*. We found that the estimated distribution of the actual delay is approximately normally distributed in each case, as predicted by the limit in (4.13).

Testing the Approximations (4.7) and (4.9)
with observed w in a small interval about $2E[W_\infty | W_\infty > 0]$

interarrival-time dist.	M	D	H_2
$2E[W_\infty W_\infty > 0]$	0.40	0.20	0.96
selected <i>HOL</i> w interval	[0.39, 0.41]	[0.19, 0.21]	[0.94, 0.98]
sample size	128,287	99,747	151,556
sample mean observed	0.3998	0.2000	0.9597
$E[W_{HOL}(w)]$ est.	0.4003	0.1996	0.9625
$Var(W_{HOL}(w))$ est.	0.0080	0.0020	0.0448
$E[W_{LES}(w)]$ est.	0.3996	0.1995	0.9617
$Var(W_{LES}(w))$ est.	0.0081	0.0021	0.0450
$E[W_{RCS}(w)]$ est.	0.3938	0.1929	0.9586
$Var(W_{RCS}(w))$ est.	0.0103	0.0029	0.0507
Predicted mean by (4.7)	0.400	0.205	0.947
Pred. variance by (4.9)	0.0076	0.0021	0.0455

Table 5: Comparing the approximations for $E[W_{HOL}(w)]$ and $Var(W_{HOL}(w))$ for fixed w following from (4.7) and (4.9) with simulation estimates of the mean and variance of the *HOL*, *LES* and *RCS* estimators in the *GI/M/100* model with $\rho = 0.95$ as a function of the interarrival-time distribution. Data are collected for observed waiting times contained in a small interval about $2E[W_\infty | W_\infty > 0]$. The resulting sample sizes are shown.

6. Insights from Heavy-Traffic Limits

We can gain additional insight about the performance of the different estimators by considering heavy-traffic limits for the $GI/M/s$ model. To do so, we consider a family of models indexed by the parameter ρ , so we introduce a second subscript ρ in addition to s . We let the service times remain unchanged. We assume that we start with interarrival times U_n having mean $1/s$. In system (s, ρ) , we use interarrival times U_n/ρ , so that they have mean $1/s\rho$. That makes the traffic intensity in model ρ be ρ .

We consider both the classical heavy-traffic (HT) regime in which $\rho \uparrow 1$ for fixed s and the Quality-and-Efficiency-Driven (QED) many-server heavy-traffic (HT) regime in which both $\rho \uparrow 1$ and $s \rightarrow \infty$ with $((1 - \rho)\sqrt{s} \rightarrow \beta$ for $0 < \beta < \infty$; see Chapters 5, 9 and 10 of Whitt (2002) for background. The queue length tends to be of order $1/(1 - \rho)$ in both limiting regimes, but the delays behave differently. The delay are of order $1/(1 - \rho)$ in the classical HT regime, but are of order $1 - \rho$ or $1/\sqrt{s}$ in the QED HT regime.

The Heavy-Traffic Snapshot Principle Just as in the application of heavy-traffic limits to plan queueing simulations reviewed in §5.8 of Whitt (2002), the time scaling in the heavy-traffic stochastic-process limits provides important insight. In particular, we can apply the celebrated *heavy-traffic snapshot principle*, see Reiman (1982) and p. 187 of Whitt (2002), which in our context tells us that the waiting times (of other customers) tend to change negligibly during the time a customer spends waiting when the system is in heavy traffic. In other words, the snapshot principle immediately implies that the LES and HOL estimators are asymptotically exact in heavy-traffic limits (specifically, the ratio converges to one). It also shows that, asymptotically in the heavy-traffic limit, there is no advantage in averaging over delays of past customers.

Since we are primarily concerned with waiting times, it is appropriate to focus on the virtual waiting time stochastic process, which describes the waiting time of a potential arrival who would come at time t . We first consider the classical HT regime. Let $W_{s,\rho}(t)$ be the virtual waiting time at time t in model (s, ρ) . The waiting time of the k^{th} arrival at time $A_{k,s,\rho}$ is just $W_{s,\rho}(A_{k,s,\rho}-)$, where $g(t-)$ is the left limit of the function g at time t .

The classical heavy-traffic stochastic-process limit for the virtual waiting time process states that

$$(1 - \rho)W_{s,\rho}((1 - \rho)^{-2}t) \Rightarrow RBM(t) \quad \text{as } \rho \uparrow 1, \quad (6.1)$$

where the limit stochastic process $RBM(t)$ is a reflected Brownian motion, which has continuous sample paths, and the convergence in distribution is for the entire stochastic process with sample paths in the function space D ; see Whitt (2002). The space scaling in (6.1) implies that the waiting times will be of order $O(1/(1-\rho))$, while the time scaling in (6.1) implies that the waiting times will only change significantly over time intervals of length of order $O(1/(1-\rho)^2)$. As a consequence, we conclude that the HOL and LES estimators are relatively consistent in the classical HT regime.

A similar story holds in the QED HT regime. The stochastic-process limit for the virtual waiting time process in the QED regime is obtained by Puhalskii and Reiman (2000). Let $W_{s,\rho}(t)$ be the virtual waiting time at time t in model (s, ρ) . Paralleling (6.1), in the QED regime we have the stochastic-process limit

$$\sqrt{s}W_{s,\rho}(t) \Rightarrow Y(t) \quad \text{as } \rho \uparrow 1, \quad (6.2)$$

where the limit process $Y(t)$ is no longer RBM but again is a diffusion process with continuous sample paths and again the convergence in distribution is for the entire stochastic process with sample paths in the function space D .

The time and space scaling in (6.2) is drastically different from (6.1), but we nevertheless obtain the same conclusions about our estimators. Now the waiting times are getting small instead of large, being of order $O(1/\sqrt{s})$, but there is no time scaling at all, so that the waiting times will only change significantly over time intervals of length of order $O(1)$. As a consequence, we conclude that the HOL and LES estimators are also relatively consistent in the QED HT regime. Again, we conclude that there will be no advantage to averaging the delays experienced over past customers.

Steady-State Heavy-Traffic Limits In the e-companion we also establish heavy-traffic limits in both regimes for steady-state random variables. We focus on the HOL estimator; by Theorem 4.4, the LES estimator behaves the same. We see what happens “on average” to the random variable $W_{HOL,s,\rho}(w)$ (where the observed delay w has the steady-state distribution). From the steady-state HT limits, we deduce that both the direct QL and HOL estimators are (weakly) relatively consistent: the ratio of the estimator to the random quantity being estimated converges to 1. We also establish limits establishing the asymptotic efficiency of the different estimators (comparing MSE’s). In these HT limits the direct and refined estimators have asymptotically the same efficiency, while the QL estimator is asymptotically more efficient

than these delay-history estimators by the constant factor $c_a^2 + 1$, consistent with Theorem 4.2. Since associated heavy-traffic stochastic-process limits have been established for other models, the estimators should have similar nice properties for other models.

7. Conclusions

Insights that can be Generalized. Even though we are primarily interested in service systems that are more complex than the $GI/M/s$ queueing model, in this paper we studied the performance of alternative delay estimators in this relatively simple idealized $GI/M/s$ setting. Our goal has been to gain insight into how the estimators will perform in more complex settings. Our results for the $GI/M/s$ model indicate what to expect more generally. Although it remains to be verified in each specific context, we anticipate that many of the performance conclusions for the $GI/M/s$ model (reviewed below) will extend to other settings. At a minimum, the results here serve as a basis for comparison in further examination of delay estimation.

Performance of the Estimators. An important reference point for the delay estimators based on delay history is the standard QL estimator based on the observed queue length, defined in (1.1). For QL, the only source of uncertainty is the remaining service times of the customers ahead of the arrival. That uncertainty can be reduced if the remaining service times can be reliably estimated, as emphasized by Whitt (1999).

As can be seen from formulas (2.6)-(2.8), to a large extent, the LES and HOL estimators can be regarded as the QL estimator modified by replacing the known queue length by an estimate of that queue length. Since the queue length is equal (or approximately equal) to the number of arrivals during the observed waiting time, the queue length is estimated by the expected number of arrivals during the observed waiting time. Thus the increase in MSE in going from QL to the LES, HOL and RCS estimators is primarily due to variability in the arrival process. The MSE tends to be larger for LES and HOL than QL by the constant factor $(c_a^2 + 1)$, where c_a^2 is the SCV of an interarrival time, a common measure of variability for a renewal arrival process; see Whitt (1982).

As a consequence, the delay estimators based on delay history will perform about the same as the QL estimator when the arrival process has very low variability, but the relative performance will degrade as that arrival-process variability increases. From the perspective of statistical precision, the QL estimator should be preferred to the delay-history estimators if it

is available, unless there is negligible arrival-process variability. The delay-history estimators offer the advantage of transparency, but that is obtained at the expense of statistical precision. This insight should apply very broadly.

Overall, we conclude that the greatest source of estimation uncertainty is the remaining service times. After that, it is the arrival-process variability, as partially characterized by the SCV c_a^2 . We conclude that the estimators $\theta_{QL}(n)$, $\theta_{LES}^d(w)$, $\theta_{HOL}^d(w)$ and $\theta_{RCS}^d(w)$ can be very useful, but they are not extraordinarily accurate. The refined estimators for HOL, LES and RCS can remove all or nearly all of the bias, but non-negligible variance remains. The greatest hope for more reliable estimation seems to lie in being able to better predict the remaining service times, which is certainly possible if the service times are actively *controlled*, and is possible to some extent if either the service-time distribution is non-exponential or if it is possible to classify the customers, as discussed in Whitt (1999). An important direction for further research is to develop more sophisticated estimators that exploit much more of the information. Nevertheless, there may always be a role for the transparent delay estimators based on recent delay history considered here.

We considered several different delay estimators based on recent delay history, notably LES, HOL and RCS. Through analysis and extensive simulation experiments, we conclude that the LES and HOL delay estimators are very similar, with both being more accurate than the others based on delay history, but less accurate than the full-information queue-length (QL) estimator. For large s , RCS is far superior to the delay of the last customer to complete service (LCS), because customers need not complete service in the same order they arrive. For low traffic intensities with large s , LCS was even outperformed by the no-information estimator (NI). The reason is that the LCS customer may have arrived too long ago. We conclude that RCS should only be preferred to HOL and LES if delay information is not available until after customers complete service, but the MSE is not much greater for RCS than for LES and HOL.

For the $GI/M/s$ model, the random delay $W_{HOL}(w)$ given the HOL observation w is remarkably tractable, as can be seen from the representation (2.8). Theorem 4.2 gives the exact mean and variance of $W_{HOL}(w)$ for Poisson arrivals. It is significant that the mean $E[W_{HOL}(w)]$ is not simply w , but instead is a linear function of it: $\rho w + (2/s)$ with Poisson arrivals. That mean serves as a refined estimate, which has lower MSE than the direct estimator, but it requires extra information. Bias in the direct estimators can be expected more generally.

For general renewal arrivals, Theorem 4.2 establishes asymptotic results that generate sim-

ple approximations, which may well describe the behavior of these estimators in other settings. As sw increases, the random variable $W_{HOL}(w)$ is asymptotically normally distributed with explicit mean and variance (§4), which has been substantiated by simulation, as discussed in §5. From (4.13), we see that the squared coefficient of variation $c_{W_{HOL},s(w)}^2$ is asymptotically proportional to $(c_a^2 + 1)/\rho sw$ as $sw \rightarrow \infty$. That implies very accurate prediction when sw is large. These properties of $W_{HOL}(w)$ (and $W_{LES}(w)$) can be expected to hold more generally.

In §6 and the e-companion we showed that heavy-traffic limits provide important insight. The heavy-traffic snapshot principle provides strong support for all these delay-history estimation procedures, and shows that there should be little benefit from averaging over past customer delays, under heavy loads. The relative errors of the LES and HOL estimators are asymptotically negligible in both the classical and many-server heavy-traffic regimes. The MSE relative to the mean is asymptotically negligible for all the candidate delay estimators based on delay history. The QL estimator is asymptotically more efficient than HOL and LES by the constant factor $c_a^2 + 1$ in both heavy-traffic regimes. Since similar heavy-traffic limits have already been established for much more general models, these heavy-traffic properties can be expected to hold more generally.

Possible New Applications. For call centers as well as other service systems (e.g., delays in receiving new products or getting an application processed by the INS), there may be new applications of these alternative delay estimators based on recent delay history. They can also be used by customers and third parties that do not have access to all the state information available to the service provider. This might work as follows: Large groups of customers might voluntarily route their delay experience electronically to a centralized consumer-group monitor that makes this information available to its customer base in real time. The customers in turn could have their communication equipment set up to simultaneously query the monitor whenever the customer contacts the service provider. In this way, the flow of critical information could take place in milliseconds, which is far shorter than a short telephone call. This is not beyond current technology.

In the same spirit, the LES delay estimator could be used by outside parties to verify that the service provider is providing accurate delay estimates. The service provider could agree to publish its delay estimates, providing extra coded information giving the customer identification for each observed LES delay. Customers or authorized third parties could then verify that the delays, appropriately recorded, coincided with that same delay when it was

quoted as an LES delay. The information available to each customer would not go beyond its own delay experience, and yet, collectively, customers could verify the accuracy of the delay predictions. Such verification might well be regarded as a legitimate customer concern. And service providers might want to offer the verification as a way to provide better service.

Acknowledgments. The reported research was supported by NSF grant DMI-0457095.

References

- Abate, J. and W. Whitt. 1992. The Fourier-series method for inverting transforms of probability distributions. *Queueing Systems* 10, 5–88.
- Abramowitz, M. and I. A. Stegun. 1972. *Handbook of Mathematical Functions*, National Bureau of Standards, U. S. Dept. of Commerce, Washington, D.C.
- Armony, M., N. Shimkin and W. Whitt. 2006. The impact of delay announcements in many-server queues with abandonments. *Operations Research*, forthcoming.
Available at <http://columbia.edu/~ww2040>.
- Asmussen, S. 2003. *Applied Probability and Queues*, second edition, Springer, New York.
- Avramidis, A. N., A. Deslauriers and P. L’Ecuyer. 2004. Modeling daily arrivals to a telephone call center. *Management Sci.* 50, 896–908.
- Brown, L., N. Gans, A. Mandelbaum, A. Sakov, H. Shen, S. Zeltyn and L. Zhao. 2005. Statistical analysis of a telephone call center: a queueing-science perspective. *J. Amer. Statist. Assoc.* 100, 36–50.
- Coates, M., A. O. Hero, III, R. Nowak and B. Yu. 2002. Internet tomography. *IEEE Signal Processing Magazine* 19, 47–65.
- Doytchinov, B., J. Lehoczy and S. Shreve. 2001. Real-time queues in heavy traffic with earliest-deadline-first queue discipline. *Ann. Appl. Probab.* 11, 332–378.
- Feller, W. 1971. *An Introduction to Probability Theory and its Applications*, vol. II, second ed., Wiley, New York.
- Gans, N., G. Koole and A. Mandelbaum. 2003. Telephone call centers: tutorial, review and research prospects. *Manufacturing and Service Opns. Mgmt.* 5, 79–141.
- Glynn, P. W. and W. Whitt. 1989. Indirect estimation via $L = \lambda W$. *Operations Research* 37, 82–103.
- Ibrahim, R. E. and W. Whitt. 2007. Supplement to “Real time delay estimation based on delay history.” IEOR Department, Columbia University, New York, NY. Available at <http://columbia.edu/~ww2040>.

- Ibrahim, R. E. and W. Whitt. 2008. Real-time delay announcements in service systems with customer abandonment, in preparation. Available at <http://columbia.edu/~ww2040>.
- Larson, R. C. 1990. The queue inference engine: deducing queue statistics from transactional data. *Management Sci.* 36, 586–601.
- Mandelbaum A., A. Sakov and S. Zeltyn. 2000. Empirical analysis of a call center. Technical Report, Faculty of Industrial Engineering and Management, The Technion, Israel.
- Nakibly, E. 2002. *Predicting Waiting Times in Telephone Service Systems*, MS thesis, the Technion, Haifa, Israel.
- Puhalskii, A. A. and M. I. Reiman. 2000. The multiclass $GI/PH/N$ queue in the Halfin-Whitt regime. *Adv. Appl. Prob.* 32, 564–595.
- Reiman, M. I. 1982. The heavy traffic diffusion approximation for sojourn times in Jackson networks. In *Applied Probability – Computer Science, the Interface, II*, R. L. Disney and T. J. Ott (eds.), Birhauser, Boston, 409–422.
- Ross, S. M. 1996. *Stochastic Processes*, second edition, Wiley, New York.
- Ward, A. W. and W. Whitt. 2000. Predicting response times in processor-sharing queues. In *Analysis of Communication Networks: Call Centres, Traffic and Performance*, D. R. McDonald and S. R. E. Turner (eds.), Fields Institute Communications 28, American Math. Society, Providence, RI, 1-29.
- Whitt, W. 1982. Approximating a point process by a renewal process: two basic methods. *Operations Research* 30, 125–147.
- Whitt, W. 1999. Predicting queueing delays. *Management Sci.* 45, 870–888.
- Whitt, W. 2002. *Stochastic-Process Limits*, Springer, New York.
- Whitt, W. 2004a. Efficiency-driven heavy-traffic approximations for many-server queues with abandonments. *Management Science* 50, 1449–1461.
- Whitt, W. 2006. Fluid Models for Multiserver Queues with Abandonments. *Operations Research* 54, 37–54.
- Wolff, R. W. 1989. *Stochastic Modeling and the Theory of Queues*, Prentice Hall, Englewood Cliffs, NJ.

Xu, S. H., L. Gao and J. Ou. 2007. Service performance analysis and improvement for a ticket queue with balking customers. *Management Sci.* 53, 971–990.

e-Companion

to

REAL-TIME DELAY ESTIMATION BASED ON DELAY HISTORY

by

Rouba Ibrahim and Ward Whitt

IEOR Department
Columbia University
{rei2101, ww2040}@columbia.edu

A. Introduction

We present additional material in this e-companion. First, in §B we present additional experimental results; we present many more in an online supplement available on the authors' web pages. Next, in §C we establish steady-state heavy-traffic limits for these estimators. At the end of the section, we show that the bad performance of the LCS estimator for large s can be explained in part by its behavior in the QED many-server heavy-traffic limiting regime. Unlike the LES, HOL and RCS delay estimators, the LCS delay estimator is *not* asymptotically consistent in this limiting regime. Finally, in §D we present a cautionary example showing the possible pitfalls of the LES and HOL delay estimators.

B. Additional Tables and Figures

Paralleling Table 1 in §3, which displays the ASE's for seven different estimators in the $GI/M/100$ model for the M , D and H_2 arrival processes, we display the corresponding estimated ASE's for the same estimators for the $GI/M/s$ models with $s = 10$ and $s = 1$ in Tables 6 and 7 below. The estimator LCS fares better as s decreases. The ASE's of LCS and RCS do not differ greatly for $s = 10$ and are identical for $s = 1$.

Paralleling Figure 1 in §3, where we display plots of the relative average squared errors (RASE's) for several of the estimators in the $D/M/100$ model, we display the RASE's for the $M/M/100$ and $H_2/M/100$ models in Figures 2 and 3. Again we see linear or near-linear performance as a function of ρ . The advantage of QL over LES increases as c_a^2 increases. Again the HOL and LES values fall on top of each other, so we only show LES.

Paralleling Table 3 in §4, where we compare the approximations for the MSE's of the three estimators θ_{HOL}^d , θ_{HOL}^{ar} and θ_{HOL}^{sr} in the $H_2/M/s$ model with $s = 100$ and $s = 1$, we show

*Estimated ASE in units of 10^{-1}
M/M/s model with $s = 10$*

ρ	QL	LES	HOL	RCS	RCS- \sqrt{s}	LCS	NI
0.98	4.95 ± 0.23	10.1 ± 0.42	10.1 ± 0.41	10.8 ± 0.41	10.9 ± 0.42	11.9 ± 0.41	257.2 ± 48.1
0.95	1.98 ± 0.025	4.16 ± 0.040	4.17 ± 0.042	4.83 ± 0.039	4.94 ± 0.041	5.87 ± 0.041	39.61 ± 2.3
0.93	1.42 ± 0.013	3.03 ± 0.032	3.05 ± 0.037	3.67 ± 0.036	3.77 ± 0.033	4.62 ± 0.036	20.01 ± 0.66
0.9	1.00 ± 0.017	2.19 ± 0.033	2.20 ± 0.042	2.79 ± 0.036	2.88 ± 0.035	3.63 ± 0.036	10.10 ± 0.49
0.85	0.661 ± 0.0032	1.50 ± 0.0076	1.53 ± 0.012	2.04 ± 0.0092	2.11 ± 0.0085	2.69 ± 0.0097	4.41 ± 0.083

D/M/s model with $s = 10$

ρ	QL	LES	HOL	RCS	RCS- \sqrt{s}	LCS	NI
0.98	2.49 ± 0.084	2.63 ± 0.083	2.63 ± 0.086	2.99 ± 0.085	3.05 ± 0.086	3.57 ± 0.086	59.3 ± 10.2
0.95	1.01 ± 0.018	1.16 ± 0.018	1.16 ± 0.020	1.50 ± 0.019	1.55 ± 0.019	2.00 ± 0.019	10.1 ± 0.83
0.93	0.730 ± 0.010	0.876 ± 0.011	0.877 ± 0.013	1.21 ± 0.012	1.26 ± 0.011	1.66 ± 0.012	5.24 ± 0.29
0.9	0.518 ± 0.0058	0.663 ± 0.0057	0.663 ± 0.0091	0.977 ± 0.0077	1.02 ± 0.0066	1.37 ± 0.0078	2.66 ± 0.12
0.85	0.352 ± 0.0025	0.494 ± 0.0026	0.494 ± 0.0057	0.779 ± 0.0047	0.814 ± 0.0028	1.06 ± 0.0047	1.24 ± 0.0053

H_2 /M/s model with $s = 10$

ρ	QL	LES	HOL	RCS	RCS- \sqrt{s}	LCS	NI
0.98	12.8 ± 0.69	62.6 ± 4.0	62.6 ± 4.1	64.4 ± 4.1	65.1 ± 4.1	67.3 ± 5.6	1594 ± 258
0.95	4.81 ± 0.081	22.3 ± 0.47	22.3 ± 0.48	23.9 ± 0.47	24.6 ± 0.47	26.5 ± 0.81	229 ± 9.1
0.93	3.42 ± 0.069	15.4 ± 0.35	15.4 ± 0.37	17.0 ± 0.35	17.5 ± 0.35	19.4 ± 0.35	115 ± 6.8
0.9	2.34 ± 0.036	10.1 ± 0.18	10.1 ± 0.20	11.6 ± 0.19	11.8 ± 0.18	13.7 ± 0.18	54.4 ± 2.9
0.85	1.50 ± 0.022	6.00 ± 0.12	6.02 ± 0.13	7.25 ± 0.12	7.50 ± 0.13	8.97 ± 0.076	22.8 ± 1.37

Table 6: A comparison of the efficiency of different real-time delay estimators for the $GI/M/10$ queue as a function of the traffic intensity ρ and the interarrival-time distribution (M , D and H_2). Only the direct estimators are considered. Estimates of the average squared error ASE are shown together with the half width of the 95% confidence interval. The units are 10^{-1} throughout.

Estimated ASE
M/M/s model with s = 1

ρ	QL	LES	HOL	RCS	RCS- \sqrt{s}	LCS	NI
0.95	20.1 ± 0.42	42.2 ± 0.77	42.4 ± 0.79	44.1 ± 0.78	44.1 ± 0.78	44.1 ± 0.78	405.0 ± 23.4
0.93	14.4 ± 0.19	30.6 ± 0.37	30.7 ± 0.39	32.4 ± 0.37	32.4 ± 0.37	32.4 ± 0.37	207.5 ± 10.4
0.9	9.99 ± 0.084	21.8 ± 0.19	22.0 ± 0.21	23.5 ± 0.19	23.5 ± 0.19	23.5 ± 0.19	100.6 ± 3.4
0.85	6.68 ± 0.043	15.1 ± 0.093	15.4 ± 0.095	16.6 ± 0.010	16.6 ± 0.010	16.6 ± 0.010	44.9 ± 0.88

D/M/s model with s = 1

ρ	QL	LES	HOL	RCS	RCS- \sqrt{s}	LCS	NI
0.95	10.1 ± 0.15	11.6 ± 0.15	11.6 ± 0.16	12.6 ± 0.15	12.6 ± 0.15	12.6 ± 0.15	101.1 ± 7.2
0.93	7.32 ± 0.081	8.79 ± 0.078	8.79 ± 0.086	9.73 ± 0.080	9.73 ± 0.080	9.73 ± 0.080	52.7 ± 2.4
0.9	5.19 ± 0.038	6.64 ± 0.037	6.65 ± 0.041	7.56 ± 0.040	7.56 ± 0.040	7.56 ± 0.040	26.8 ± 0.94
0.85	3.53 ± 0.018	4.96 ± 0.018	4.95 ± 0.020	5.82 ± 0.020	5.82 ± 0.021	5.82 ± 0.020	12.4 ± 0.36

H₂/M/s model with s = 1

ρ	QL	LES	HOL	RCS	RCS- \sqrt{s}	LCS	NI
0.95	48.7 ± 1.13	226.4 ± 5.14	226.5 ± 5.23	231.1 ± 5.15	231.1 ± 5.15	231.1 ± 5.15	2339 ± 425
0.93	34.3 ± 0.63	154.4 ± 2.9	154.4 ± 2.9	158.9 ± 3.0	158.9 ± 3.0	158.9 ± 3.0	1151 ± 181
0.9	23.48 ± 0.37	101.3 ± 2.3	101.4 ± 2.4	105.5 ± 2.4	105.5 ± 2.4	105.5 ± 2.4	552.9 ± 103
0.85	14.95 ± 0.104	60.0 ± 0.52	60.2 ± 0.53	63.9 ± 0.51	63.9 ± 0.51	63.9 ± 0.51	224.4 ± 6.2

Table 7: A comparison of the efficiency of different real-time delay estimators for the $GI/M/1$ queue as a function of the traffic intensity ρ and the interarrival-time distribution (M , D and H_2). Only the direct estimators are considered. Estimates of the average squared error ASE are shown together with the half width of the 95% confidence interval.

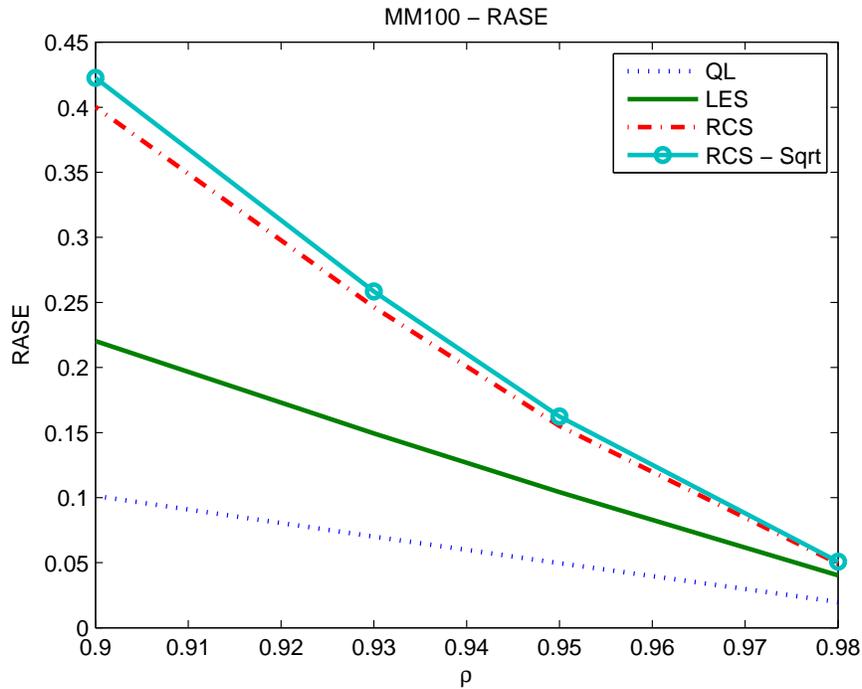


Figure 2: The relative average squared error ($RASE$) for the $M/M/100$ model.

corresponding results for the $M/M/s$ model with $s = 100$ and $s = 1$ in Table 8. We have used simulation to estimate all quantities here, even though we could compute them analytically. This case thus provides a crosscheck on both our analytic formulas and the simulations.

Finally, we present one table illustrating our study of the number of past customers we need to consider for RCS, as discussed at the end of §2. Table 9 presents simulation results for the $H_2/M/100$ model as a function of ρ . These results support the conclusion that $RCS - c\sqrt{s}$ is virtually identical to RCS itself when $c = 4$, and that small errors are observed when $c = 2$ and $s = 1$. These conclusions held uniformly over all interarrival-time distributions and all $s \geq 100$.

Evaluating the alternative HOL estimators
 Approximations in the $M/M/s$ model for $s = 100$ and $s = 1$

ρ	0.85	0.90	0.93	0.95	0.98	0.99
$E[W W > 0]$	0.0666	0.0993	0.1435	0.2012	0.500	0.901
conf. int.	± 0.0018	± 0.0027	± 0.0018	± 0.0019	± 0.037	± 0.059
$E[W^2 W > 0]$	0.0089	0.0196	0.0414	0.0811	0.500	1.53
conf. int.	± 0.0006	± 0.0012	± 0.0016	± 0.0026	± 0.097	± 0.24
$MSE(\theta^d)$	0.00153	0.00219	0.00307	0.00422	0.01020	0.01823
term 1	0.00020	0.00020	0.00020	0.00020	0.00020	0.00015
term 2	0.00073	0.00139	0.00227	0.00342	0.00940	0.01748
term 3	0.00060	0.00060	0.00060	0.00060	0.00060	0.00060
$MSE(\theta^{sr})$	0.00173	0.00239	0.00327	0.00442	0.01040	0.01844
term 1	0.00113	0.00179	0.00267	0.00382	0.00980	0.01784
term 2	0.00060	0.00060	0.00060	0.00060	0.00060	0.00060
$MSE(\theta^{ar})$	0.00133	0.00199	0.00287	0.00402	0.01000	0.01804
term 1	0.00113	0.00179	0.00267	0.00382	0.00980	0.01784
term 2	0.00020	0.00020	0.00020	0.00020	0.00020	0.00020

Approximations in the $M/M/1$ model

ρ	0.80	0.85	0.90	0.95	0.96	0.98
$E[W W > 0]$	5.01	6.68	9.98	20.04	24.80	50.70
conf. int.	± 0.03	± 0.04	± 0.08	± 0.36	± 0.33	± 2.4
$E[W^2 W > 0]$	50.3	89.6	200.3	806.6	1211	5290
conf. int.	± 0.69	± 1.36	± 5.1	± 37.4	± 45	640
$MSE(\theta^d)$	12.02	15.36	21.98	42.08	51.58	103.4
term 1	2.01	2.01	2.00	2.02	1.94	2.11
term 2	4.01	7.35	13.98	34.07	43.64	95.25
term 3	6.00	6.00	6.00	6.00	6.00	6.00
$MSE(\theta^{sr})$	14.02	17.35	23.97	44.07	53.61	105.31
term 1	8.02	11.35	17.97	38.07	47.61	99.31
term 2	6.00	6.00	6.00	6.00	6.00	6.00
$MSE(\theta^{ar})$	10.02	13.35	19.97	40.07	49.61	101.31
term 1	8.02	11.35	19.97	38.07	47.61	99.31
term 2	2.00	2.00	2.00	2.00	2.00	2.00

Table 8: Evaluation of the MSE approximations for the estimators θ_{HOL}^d , θ_{HOL}^{sr} , and θ_{HOL}^{ar} in steady-state using (4.11), (4.9) and (4.10) together with simulation estimates of the first two moments of the conditional delay $E[W_\infty|W_\infty > 0]$. The $M/M/s$ model is considered as a function of the traffic intensity ρ for $s = 100$ and $s = 1$.

<i>ASE in the $H_2/M/s$ model with $s = 100$</i>						
ρ	$ASE(RCS)$	$ASE(RCS - s)$	$ASE(RCS - 4\sqrt{(s)})$	$ASE(RCS - 2\sqrt{(s)})$	$ASE(RCS - \sqrt{(s)})$	$ASE(RCS - \log(s))$
0.98	2.439×10^{-2} $\pm 4.84 \times 10^{-4}$	2.439×10^{-2} $\pm 4.84 \times 10^{-4}$	2.439×10^{-2} $\pm 4.84 \times 10^{-4}$	2.442×10^{-2} (0.123) $\pm 4.88 \times 10^{-4}$	2.511×10^{-2} (2.95) $\pm 4.92 \times 10^{-4}$	3.724×10^{-2} (52.7) $\pm 6.81 \times 10^{-4}$
0.97	2.229×10^{-2} $\pm 4.70 \times 10^{-4}$	2.229×10^{-2} $\pm 4.70 \times 10^{-4}$	2.229×10^{-2} $\pm 4.70 \times 10^{-4}$	2.229×10^{-2} (0.141) $\pm 4.73 \times 10^{-4}$	2.367×10^{-2} (3.28) $\pm 4.73 \times 10^{-4}$	3.566×10^{-2} (55.6) $\pm 5.80 \times 10^{-4}$
0.95	1.989×10^{-2} $\pm 3.67 \times 10^{-4}$	1.989×10^{-2} $\pm 3.67 \times 10^{-4}$	1.989×10^{-2} $\pm 3.67 \times 10^{-4}$	1.992×10^{-2} (0.136) $\pm 3.67 \times 10^{-4}$	2.058×10^{-2} (3.48) $\pm 3.64 \times 10^{-4}$	3.175×10^{-2} (59.6) $\pm 5.45 \times 10^{-4}$
0.93	1.715×10^{-2} $\pm 3.56 \times 10^{-4}$	1.715×10^{-2} $\pm 3.56 \times 10^{-4}$	1.715×10^{-2} $\pm 3.56 \times 10^{-4}$	1.718×10^{-2} (0.150) $\pm 3.54 \times 10^{-4}$	1.780×10^{-2} (3.78) $\pm 3.60 \times 10^{-4}$	2.800×10^{-2} (63.2) $\pm 5.89 \times 10^{-4}$
0.90	1.344×10^{-2} $\pm 4.90 \times 10^{-4}$	1.344×10^{-2} $\pm 4.90 \times 10^{-4}$	1.344×10^{-2} $\pm 4.90 \times 10^{-4}$	1.347×10^{-2} (0.182) $\pm 4.89 \times 10^{-4}$	1.399×10^{-2} (4.06) $\pm 4.99 \times 10^{-4}$	2.233×10^{-2} (66.3) $\pm 8.61 \times 10^{-4}$

Table 9: A comparison of the efficiency of the candidate RCS- $f(s)$ delay estimators for the $H_2/M/s$ queue with $s = 100$ as a function of the traffic intensity ρ . Below each point estimates for ASE is shown with the half width of the 95-percent confidence interval. Also included in parentheses are the values of the relative percent difference.

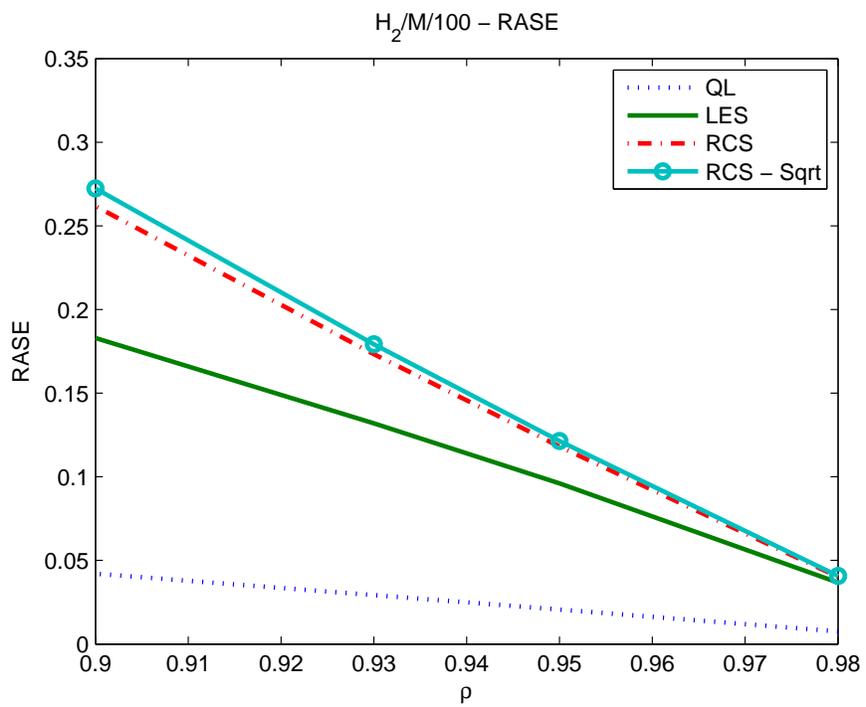


Figure 3: The relative average squared error ($RASE$) for the $H_2/M/100$ model.

C. Heavy-Traffic Limits

In this section we present additional heavy-traffic limits, extending the discussion in §6. We start by establishing heavy-traffic limits for the steady-state random variables. We see what happens “on average” to the random variable $W_{HOL,s,\rho}(w)$. We consider both the classical heavy-traffic regime in which $\rho \uparrow 1$ for fixed s and the QED (many-server heavy-traffic limiting) regime in which both $\rho \uparrow 1$ and $s \rightarrow \infty$ with $((1 - \rho)\sqrt{s} \rightarrow \beta$ for $0 < \beta < \infty$; see Chapters 5, 9 and 10 of Whitt (2002) for background. For more on the QED regime for $GI/G/s$ queues, see Halfin and Whitt (1981), Puhalskii and Reiman (2000), Jelenkovic et al. (2004) and Whitt (2004b, 2005).

The Classical Heavy-Traffic Regime. We start with the classic heavy-traffic (HT) regime in which $\rho \uparrow 1$ with fixed s . We look at the distribution of $W_{HOL,s}(w)$, assuming that the observed waiting time w experienced by the customer at the head of the line is a random variable $W_{\infty,s,\rho}^h$, assumed to be the steady-state delay in model (s, ρ) experienced by a customer at the head of the line at an arrival epoch, conditional on there being at least one customer in the queue. Thus let $W_{HOL,s,\rho}(W_{\infty,s,\rho}^h)$ denote a random variable with the distribution

$$P(W_{HOL,s,\rho}(W_{\infty,s,\rho}^h) \leq x) \equiv \int_0^\infty P(W_{HOL,s,\rho}(w) \leq x) dP(W_{\infty,s,\rho}^h \leq w), \quad (3.1)$$

in model (s, ρ) , where in this subsection s is held fixed. This means that $E[W_{HOL,s,\rho}(W_{\infty,s,\rho}^h)] \equiv E[E[W_{HOL,s,\rho}(W_{\infty,s,\rho}^h)|W_{\infty,s,\rho}^h]]$. The random variable $W_{\infty,s,\rho}^h$ is not quite distributed as the steady-state waiting time at the arrival epoch, $W_{\infty,s,\rho}$, or the conditional steady-state waiting time, $(W_{\infty,s,\rho}|W_{\infty,s,\rho} > 0)$, but it is asymptotically equivalent to both of these in the heavy-traffic limit.

In order to relate the HOL and QL estimators, it is important to exploit the joint convergence of the steady-state queue length and waiting time. Such joint convergence is discussed extensively for the single-server queue in Chapter 9 of Whitt (2002); it was also used in Iglehart and Whitt (1970), which treated more general models. Let $(Q_{\infty,s,\rho}, W_{\infty,s,\rho})$ be a random vector with the limiting steady-state distribution of $(Q_{k,s,\rho}, W_{k,s,\rho})$, where $Q_{k,s,\rho}$ is the queue length and $W_{k,s,\rho}$ is the delay just before $A_{k,s,\rho}$, where $A_{k,s,\rho}$ is the k^{th} arrival epoch, all in model (s, ρ) .

Here we will use the following established steady-state heavy-traffic limit:

$$(1 - \rho)(Q_{\infty,s,\rho}, W_{\infty,s,\rho}) \Rightarrow (L, L/s) \quad \text{as } \rho \uparrow 1, \quad (3.2)$$

where $L \stackrel{d}{=} \text{Exp}(c_a^2 + 1)/2$ with $\text{Exp}(m)$ denoting a random variable having an exponential distribution with mean m . We give a detailed proof in a subsection below starting from the known steady-state distribution for $Q_{\infty,s,\rho}$. The joint convergence follows from the limit for $Q_{\infty,s,\rho}$ and the law of large numbers, using the representation

$$(Q_{\infty,s,\rho}, W_{\infty,s,\rho}) = \left(Q_{\infty,s,\rho}, (Q_{\infty,s,\rho} + 1) \left(\frac{\left[\sum_{i=1}^{Q_{\infty,s,\rho}+1} (V_i/s) \right]}{(Q_{\infty,s,\rho} + 1)} \right) \right). \quad (3.3)$$

We can apply (3.2) and previous results to get the following limits for our estimators. Let $\text{RMSE} \equiv \text{MSE}/\text{Mean}^2$ be the relative mean squared error. Let $c_{W_{Q,s,\rho}}^2(Q_{\infty,s,\rho})$ be the random variable assuming the value $c_{W_{Q,s,\rho}}^2(n)$ with probability $P(Q_{\infty,s,\rho} = n)$ for $n \geq 0$. Let other random variables involving c^2 and RMSE be defined analogously. We prove the following theorem in a subsection below.

Theorem C.1. (*classical heavy-traffic limit*) *If $\rho \uparrow 1$ in the family of GI/M/s models indexed by (s, ρ) with fixed s , then*

$$\frac{W_{Q,s,\rho}(Q_{\infty,s,\rho})}{E[W_{Q,s,\rho}(Q_{\infty,s,\rho})|Q_{\infty,s,\rho}]} = \frac{W_{Q,s,\rho}(Q_{\infty,s,\rho})}{(Q_{\infty,s,\rho} + 1)/s} \Rightarrow 1, \quad (3.4)$$

$$\frac{W_{\infty,s,\rho}}{W_{\infty,s,\rho}^h} \Rightarrow 1 \quad \text{and} \quad \frac{W_{HOL,s,\rho}(W_{\infty,s,\rho}^h)}{W_{\infty,s,\rho}^h} \Rightarrow 1, \quad (3.5)$$

from which we can deduce that

$$(1 - \rho)(Q_{\infty,s,\rho}, W_{\infty,s,\rho}, W_{\infty,s,\rho}^h, W_{Q,s,\rho}(Q_{\infty,s,\rho}), W_{HOL,s,\rho}(W_{\infty,s,\rho}^h)) \Rightarrow (L, L/s, L/s, L/s, L/s) \quad (3.6)$$

and

$$(1 - \rho)^{-1}(c_{W_{Q,s,\rho}(Q_{\infty,s,\rho})}^2, c_{W_{HOL,s,\rho}(W_{\infty,s,\rho}^h)}^2, \text{RMSE}(W_{\infty,s,\rho}^h)) \Rightarrow (1/L, (c_a^2 + 1)/L, (c_a^2 + 1)/L) \quad (3.7)$$

where $L \stackrel{d}{=} \text{Exp}((c_a^2 + 1)/2)$ as above, so that

$$\frac{W_{HOL,s,\rho}(W_{\infty,s,\rho}^h)}{W_{Q,s,\rho}(Q_{\infty,s,\rho})} \Rightarrow 1, \quad \frac{c_{W_{HOL,s,\rho}(W_{\infty,s,\rho}^h)}^2}{c_{W_{Q,s,\rho}(Q_{\infty,s,\rho})}^2} \Rightarrow c_a^2 + 1, \quad (3.8)$$

$$\frac{\text{RMSE}(W_{\infty,s,\rho}^h)}{c_{W_{HOL,s,\rho}(W_{\infty,s,\rho}^h)}^2} \Rightarrow 1 \quad \text{and} \quad \frac{\text{RMSE}(W_{\infty,s,\rho}^h)}{c_{W_{Q,s,\rho}(Q_{\infty,s,\rho})}^2} \Rightarrow c_a^2 + 1. \quad (3.9)$$

The limits in (3.4) and (3.5) show that the direct QL and HOL estimators are (weakly) relatively consistent in the classical heavy-traffic limit, while the limits in (3.7)–(3.9) compare

the asymptotic efficiency of the different estimators. In this heavy traffic limit, the direct and refined HOL estimators have asymptotically the same efficiency, while the QL estimator is asymptotically more efficient by the constant factor $c_a^2 + 1$.

We conjecture (but have not yet proved) that there is appropriate uniform integrability, so that the moments of these random variables converge as well as distributions, see p. 31 of Billingsley (1999). Then from (3.7) and (3.8) we obtain associated convergence of the moments:

$$E \left[\frac{c_{W_{HOL,s,\rho}^h}^2(W_{\infty,s,\rho}^h)}{c_{W_{Q,s,\rho}}^2(Q_{\infty,s,\rho})} \right] \rightarrow c_a^2 + 1 \quad \text{and} \quad \frac{E[c_{W_{HOL,s,\rho}^h}^2(W_{\infty,s,\rho}^h)]}{E[c_{W_{Q,s,\rho}}^2(Q_{\infty,s,\rho})]} \rightarrow c_a^2 + 1, \quad (3.10)$$

and similarly for the direct estimator. These limits supplement the previous limits, implying that the QL delay estimator is asymptotically more efficient than the HOL and LES delay estimators by the constant factor $c_a^2 + 1$ in the classical heavy-traffic limit.

The QED Many-Server Heavy-Traffic Regime. We now consider the QED HT regime, in which both $\rho \uparrow 1$ and $s \uparrow \infty$ with $(1 - \rho)\sqrt{s} \rightarrow \beta$ for some positive constant β .

This alternative QED regime is appealing because, unlike the classical HT regime, the probability that a customer is delayed approaches a nondegenerate limit, strictly between 0 and 1:

$$P(W_{\infty,s,\rho} > 0) \rightarrow \alpha \quad \text{and} \quad P(Q_{\infty,s,\rho} > 0) \rightarrow \alpha, \quad 0 < \alpha < 1, \quad (3.11)$$

where $\alpha \equiv \alpha(\beta/\sqrt{c_a^2 + 1})$ for $\alpha(x) \equiv [1 + x\Phi(x)/\phi(x)]^{-1}$, where ϕ is the cdf and ϕ is the probability density function (pdf) of the standard normal $N(0, 1)$; see (1.1) of Whitt (2004b).

With minor modifications, the story is the same as for the classical HT regime, so we will be brief. A major difference is that the queue length is of order $O(\sqrt{s}) = O(1/(1 - \rho))$, while the waiting time is of order $O(1/\sqrt{s}) = O((1 - \rho))$. As before, the ratio $W_{\infty,s,\rho}/Q_{\infty,s,\rho}$ is of order $O(1/s)$, but now $s \rightarrow \infty$.

Paralleling (3.2), we have the joint limit

$$(Q_{\infty,s,\rho}/\sqrt{s}, (1 - \rho)Q_{\infty,s,\rho}, \sqrt{s}W_{\infty,s,\rho}, W_{\infty,s,\rho}/(1 - \rho)) \Rightarrow (Z, \beta Z, Z, Z/\beta), \quad (3.12)$$

where $P(Z > 0) = \alpha$ for the same $\alpha \equiv \alpha(\beta/\sqrt{c_a^2 + 1})$ defined above and $(Z|Z > 0) \stackrel{d}{=} L \stackrel{d}{=} \text{Exp}((c_a^2 + 1)/2)$. The limit for $Q_{\infty,s,\rho}$ was established by Halfin and Whitt (1981), but Whitt (2004b) corrects an error in the expression for α when the arrival process is non-Poisson. The joint limit with $W_{\infty,s,\rho}$ can be established as in (3.3). Paralleling (3.39), here we have

$$\begin{aligned} & ((1 - \rho)(Q_{\infty,s,\rho}|Q_{\infty,s,\rho} > 0), (W_{\infty,s,\rho}|W_{\infty,s,\rho} > 0)/(1 - \rho), W_{\infty,s,\rho}^h/(1 - \rho), (1 - \rho)A(W_{\infty,s,\rho}^h)) \\ \Rightarrow & (\beta L, L/\beta, L/\beta, \beta L), \end{aligned} \quad (3.13)$$

where again $L \stackrel{d}{=} (Z|Z > 0) \stackrel{d}{=} \text{Exp}(c_a^2 + 1)/2$; as before, the important point is that the same random variable L appears in all four components on the right.

We now state the theorem, omitting the proof.

Theorem C.2. (*QED heavy-traffic limit*) *If $\rho \uparrow 1$ and $s \uparrow \infty$ so that $(1 - \rho)\sqrt{s} \rightarrow \beta$ for $0 < \beta < \infty$ in the family of $GI/M/s$ models indexed by ρ and s , then*

$$\frac{W_{Q,s,\rho}(Q_{\infty,s,\rho})}{(Q_{\infty,s,\rho} + 1)/s} \Rightarrow 1 \quad \text{and} \quad \frac{W_{HOL,s,\rho}(W_{\infty,s,\rho}^h)}{W_{\infty,s,\rho}^h} \Rightarrow 1 . \quad (3.14)$$

$$(1 - \rho)^{-1}(W_{Q,s,\rho}(Q_{\infty,s,\rho}), W_{HOL,s,\rho}(W_{\infty,s,\rho}^h)) \Rightarrow (L/\beta, L/\beta) \quad (3.15)$$

and

$$(1 - \rho)^{-1}(c_{W_{Q,s,\rho}(Q_{\infty,s,\rho})}^2, c_{W_{HOL,s,\rho}(W_{\infty,s,\rho}^h)}^2, RMSE(W_{\infty,s,\rho}^h)) \Rightarrow (1/\beta L, (c_a^2 + 1)/\beta L, (c_a^2 + 1)/\beta L) \quad (3.16)$$

where $L \stackrel{d}{=} \text{Exp}((c_a^2 + 1)/2)$ as above, so that

$$\frac{W_{HOL,s,\rho}(W_{\infty,s,\rho}^h)}{W_{Q,s,\rho}(Q_{\infty,s,\rho})} \Rightarrow 1, \quad \frac{c_{W_{HOL,s,\rho}(W_{\infty,s,\rho}^h)}^2}{c_{W_{Q,s,\rho}(Q_{\infty,s,\rho})}^2} \Rightarrow c_a^2 + 1 . \quad (3.17)$$

$$\frac{RMSE(W_{\infty,s,\rho}^h)}{c_{W_{HOL,s,\rho}(W_{\infty,s,\rho}^h)}^2} \Rightarrow 1 \quad \text{and} \quad \frac{RMSE(W_{\infty,s,\rho}^h)}{c_{W_{Q,s,\rho}(Q_{\infty,s,\rho})}^2} \Rightarrow c_a^2 + 1 . \quad (3.18)$$

Just as in the classical HT regime, we conjecture that there is appropriate uniform integrability, so that the moments converge as well as distributions. Then we will obtain associated convergence of the moments, just as in (3.10).

Heavy-Traffic Detail: Proof of (3.2). In this section we prove the classical heavy-traffic limit for the steady-state joint distribution of the queue length and waiting time at arrival epochs stated in (3.2):

$$(1 - \rho)(Q_{\infty,\rho}, W_{\infty,\rho}) \Rightarrow (L, L/s) \quad \text{as} \quad \rho \uparrow 1 , \quad (3.19)$$

where $L \stackrel{d}{=} \text{Exp}(c_a^2 + 1)/2$ with $\text{Exp}(m)$ denoting a random variable that is exponentially distributed with mean m . We consider this a known result, but we cannot point to a place where a proof is given.

We draw on well-known properties of the steady-state distribution of the $GI/M/s$ queue. The key initial result is the fact that the conditional distribution of the queue length at an arrival epoch, given that the arrival must wait, is a geometric distribution, i.e.,

$$P(Q_{\infty,\rho} = j | W_{\infty,\rho} > 0) = (1 - \omega)\omega^j, \quad j \geq 0 , \quad (3.20)$$

where the single parameter ω in (3.20) is the unique root of the equation

$$\omega = \int_0^\infty e^{-(1-\omega)sx} dF(x) \equiv \hat{f}((1-\omega)s), \quad (3.21)$$

where \hat{f} is the Laplace-Stieltjes transform of the cdf F , i.e.,

$$\hat{f}(z) \equiv \int_0^\infty e^{-zx} dF(x); \quad (3.22)$$

see (14.10), (14.11), (14.12) and (14.19) of Cooper (1982). This property was used in the proof of Theorem 4.3.

The key then is the way that the root $\omega \equiv \omega(\rho)$ depends on the traffic intensity ρ as $\rho \uparrow 1$. Anticipating that we should have $\omega(\rho) \uparrow 1$ as $\rho \uparrow 1$, we see that the argument of the Laplace-Stieltjes transform should approach 0 in the limit. It should thus come as no surprise that we can rigorously establish the desired result by expanding the Laplace transform $\hat{f}(z)$ in a Taylor series about $z = 0$; see p. 435 of Feller (1971) for supporting theory. As was first observed by Smith (1953, p. 461), it follows that

$$\frac{1 - \omega(\rho)}{1 - \rho} \rightarrow \frac{2}{c_a^2 + 1} \quad \text{as } \rho \uparrow 1. \quad (3.23)$$

The expansion appears in a more general context in formula (17) of Abate and Whitt (1994). In the special case of the $GI/M/s$ queue, equation (7) there reduces to equation (3.21) here. An alternative approach involving upper and lower bounds is given in Whitt (1984); that focuses on the more elementary $GI/M/1$ model, but the key root has the same structure. The equation differs only by the constant factor s appearing in the equation (3.21). Additional theoretical results about characterizing roots for queues appears in Neuts (1986), Choudhury and Whitt (1994) and Glynn and Whitt (1994).

It is well known – see pages 1-2 of Feller (1971) – that if X_m is a random variable with a geometric distribution having mean m , then

$$\frac{X_m}{cm} \Rightarrow \text{Exp}(1/c) \quad \text{as } m \rightarrow \infty. \quad (3.24)$$

By (3.20), $(Q_{\infty,\rho} | W_{\infty,\rho} > 0)$ has a geometric distribution with mean $1/(1 - \omega(\rho))$. Thus we can combine (3.20), (3.23) and (3.24) to obtain

$$(1 - \rho)(Q_{\infty,\rho} | W_{\infty,\rho} > 0) \Rightarrow \text{Exp}((c_a^2 + 1)/2) \quad \text{as } \rho \uparrow 1. \quad (3.25)$$

It is also known that

$$P(W_{\infty,\rho} > 0) = \frac{A}{1 - \omega} \quad \text{where } A = \left[\frac{1}{1 - \omega} + X \right]^{-1}, \quad (3.26)$$

with $X \equiv X(\rho) \rightarrow X(1)$, $0 < X(1) < \infty$, as $\rho \uparrow 1$; see (14.14)–(14.17) of Cooper (1982). Hence

$$P(W_{\infty,\rho} > 0) = [1 + (1 - \omega(\rho))X(\rho)]^{-1} \rightarrow 1 \quad \text{as } \rho \uparrow 1. \quad (3.27)$$

Combining (3.25) and (3.27), we obtain the first part of (3.19):

$$(1 - \rho)Q_{\infty,\rho} \Rightarrow L \stackrel{d}{=} \text{Exp}((c_a^2 + 1)/2) \quad \text{as } \rho \uparrow 1. \quad (3.28)$$

Given that

$$W_{\infty,\rho} \stackrel{d}{=} \sum_{i=1}^{Q_{\infty,\rho}+1} (V_i/s), \quad (3.29)$$

we have

$$\frac{W_{\infty,\rho}}{Q_{\infty,\rho} + 1} \Rightarrow \frac{1}{s} \quad \text{as } \rho \uparrow 1 \quad (3.30)$$

by the weak law of large numbers, since $Q_{\infty,\rho} \Rightarrow \infty$ as a consequence of (3.28). We then apply Theorem 11.4.5 of Whitt (2002) to write the joint limit

$$((1 - \rho)Q_{\infty,\rho}, W_{\infty,\rho}/(Q_{\infty,\rho} + 1)) \Rightarrow (L, (1/s)). \quad (3.31)$$

We then can apply the continuous mapping theorem with the function $h : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ defined by $h(x, y) = (x, xy)$ to get

$$h(((1 - \rho)Q_{\infty,\rho}, W_{\infty,\rho}/(Q_{\infty,\rho} + 1))) \Rightarrow h(L, (1/s)) = (L, L/s), \quad (3.32)$$

but

$$h(((1 - \rho)Q_{\infty,\rho}, W_{\infty,\rho}/(Q_{\infty,\rho} + 1))) = \left((1 - \rho)Q_{\infty,\rho}, (1 - \rho)W_{\infty,\rho} \frac{Q_{\infty,\rho}}{Q_{\infty,\rho} + 1} \right). \quad (3.33)$$

Since $Q_{\infty,\rho} \Rightarrow \infty$,

$$\frac{Q_{\infty,\rho}}{Q_{\infty,\rho} + 1} \Rightarrow 1 \quad \text{as } \rho \uparrow 1. \quad (3.34)$$

Hence,

$$|h(((1 - \rho)Q_{\infty,\rho}, W_{\infty,\rho}/(Q_{\infty,\rho} + 1))) - (1 - \rho)(Q_{\infty,\rho}, W_{\infty,\rho})| \Rightarrow 0 \quad \text{as } \rho \uparrow 1. \quad (3.35)$$

Thus we can combine (3.32), (3.35) and the convergence-together theorem, Theorem 11.4.7 of Whitt (2002), to complete the proof of (3.19).

Proof of Theorem C.1. First we show that $W_{\infty,s,\rho}^h \Rightarrow \infty$ as $\rho \uparrow 1$. As a consequence of the limit in (3.2), we must have $W_{\infty,s,\rho} \Rightarrow \infty$ as $\rho \uparrow 1$. Suppose that we do *not* have $W_{\infty,s,\rho}^h \Rightarrow \infty$. Then there must exist a subsequence $\{\rho_k\}$ with $\rho_k \uparrow 1$ as $k \rightarrow \infty$, a constant K and a positive constant $\epsilon > 0$ such that $P(W_{\infty,s,\rho_k}^h > K) > \epsilon$ for all k . Since

$$W_{\infty,s,\rho} \stackrel{d}{=} \sum_{i=1}^{A(W_{\infty,s,\rho}^h)+2} (V_i/s) , \quad (3.36)$$

conditional on $W_{\infty,s,\rho} > 0$, which holds with probability 1 in the limit, there must exist a new constant K' such that $P(W_{\infty,s,\rho_k} > K') > \epsilon/2$ for all k as well, but that contradicts the established limit $W_{\infty,s,\rho} \Rightarrow \infty$ as $\rho \uparrow 1$. Hence we must have $W_{\infty,s,\rho}^h \Rightarrow \infty$ as $\rho \uparrow 1$, as claimed above.

Given that $\rho \uparrow 1$ and $W_{\infty,s,\rho}^h \Rightarrow \infty$, we get $A(W_{\infty,s,\rho}^h)/W_{\infty,s,\rho}^h \Rightarrow s$ and

$$\frac{W_{HOL,s,\rho}(W_{\infty,s,\rho}^h)}{W_{\infty,s,\rho}^h} = \left(\frac{\sum_{i=1}^{A(W_{\infty,s,\rho}^h)+2} (V_i/s)}{A(W_{\infty,s,\rho}^h) + 2} \right) \left(\frac{A(W_{\infty,s,\rho}^h) + 2}{W_{\infty,s,\rho}^h} \right) \Rightarrow (1/s) \times s = 1 , \quad (3.37)$$

by the law of large numbers for partial sums and renewal processes. Similarly, by (3.2), we also have $Q_{\infty,s,\rho} \Rightarrow \infty$, so that

$$\frac{W_{Q,s,\rho}(Q_{\infty,s,\rho})}{Q_{\infty,s,\rho} + 1} = \frac{\sum_{i=1}^{Q_{\infty,s,\rho}+1} (V_i/s)}{Q_{\infty,s,\rho} + 1} \Rightarrow 1/s . \quad (3.38)$$

The limits (3.37) and (3.38) imply (3.4) and (3.5).

Since the limits in (3.37) and (3.38) are deterministic, we can apply Theorem 11.4.5 of Whitt (2002) to obtain joint convergence of all these with the limits in (3.2):

$$\begin{aligned} & \left((1-\rho)Q_{\infty,s,\rho}, (1-\rho)W_{\infty,s,\rho}, (1-\rho)W_{\infty,s,\rho}^h, \frac{W_{Q,s,\rho}(Q_{\infty,s,\rho})}{Q_{\infty,s,\rho} + 1}, \frac{W_{HOL,s,\rho}(W_{\infty,s,\rho}^h)}{W_{\infty,s,\rho}^h} \right) \\ & \Rightarrow \left(L, \frac{L}{s}, \frac{L}{s}, \frac{1}{s}, 1 \right) . \end{aligned} \quad (3.39)$$

We next apply the continuous mapping theorem, see Section 3.4 of Whitt (2002), with the function $h : \mathbb{R}^5 \rightarrow \mathbb{R}^5$ defined by $h(v, w, x, y, z) = (v, w, x, vy, xz)$ to get (3.6) from (3.39).

To continue, we next consider the random variable $c_{W_{HOL,s,\rho}(W_{\infty,s,\rho}^h)}^2$. Starting from the limit in (3.6), we can apply the Skorohod representation theorem, Theorem 3.2.2 on p. 78 of Whitt (2002), to get random variables $\tilde{W}_{\infty,s,\rho}^h$ with the same probability law as $W_{\infty,s,\rho}^h$ but for which we have the convergence $(1-\rho)\tilde{W}_{\infty,s,\rho}^h \rightarrow \tilde{L}/s$ as $\rho \uparrow 1$ w.p.1, where $\tilde{L} \stackrel{d}{=} L \stackrel{d}{=} \text{Exp}((c_a^2 + 1)/2)$. Next note that $c_{W_{HOL,s,\rho}(w)}^2/c_{W_{HOL,s,1}(w)}^2 \rightarrow 1$ w.p.1 as $\rho \uparrow 1$ and $w \rightarrow \infty$ in any order. Then,

by (4.13),

$$\frac{c_{W_{HOL,s,\rho}^h(\tilde{W}_{\infty,s,\rho}^h)}^2}{1-\rho} = \left(\frac{c_{W_{HOL,s,\rho}^h(\tilde{W}_{\infty,s,\rho}^h)}^2}{c_{W_{HOL,s,1}^h(\tilde{W}_{\infty,s,\rho}^h)}^2} \right) \left(\frac{\tilde{W}_{\infty,s,\rho}^h c_{W_{HOL,s,1}^h(\tilde{W}_{\infty,s,\rho}^h)}^2}{(1-\rho)\tilde{W}_{\infty,s,\rho}^h} \right) \rightarrow \frac{(c_a^2+1)/s}{\tilde{L}/s} \quad (3.40)$$

as $\rho \uparrow 1$ w.p.1. Essentially the same reasoning applies to the random variable RMSE $(W_{\infty,s,\rho}^h)$, giving the same limit. The equality in distribution then implies the associated convergence in distribution for the last two components of the original random vector in (3.7). We now treat the first component. Since $(Q_{\infty,s,\rho}+1)c_{W_{Q,s,\rho}(Q_{\infty,s,\rho})}^2 = 1$, a deterministic quantity, by (2.2), we can apply (4.13) to get

$$\begin{aligned} \frac{c_{W_{HOL,s,\rho}^h(W_{\infty,s,\rho}^h)}^2}{c_{W_{Q,s,\rho}(Q_{\infty,s,\rho})}^2} &= \left(\frac{Q_{\infty,s,\rho}+1}{W_{\infty,s,\rho}^h} \right) \left(\frac{W_{\infty,s,\rho}^h c_{W_{HOL,s,\rho}^h(W_{\infty,s,\rho}^h)}^2}{(Q_{\infty,s,\rho}+1)c_{W_{Q,s,\rho}(Q_{\infty,s,\rho})}^2} \right) \\ &= \left(\frac{Q_{\infty,s,\rho}+1}{W_{\infty,s,\rho}^h} \right) W_{\infty,s,\rho}^h c_{W_{HOL,s,\rho}^h(W_{\infty,s,\rho}^h)}^2 \Rightarrow s \times \frac{c_a^2+1}{s} = c_a^2+1 \end{aligned} \quad (3.41)$$

We then reason as before in establishing (3.39), first to express this limit jointly with the last two components of (3.7) and then to apply the continuous mapping theorem to complete the proof of (3.7) itself. Finally, (3.8) and (3.9) follow from the previous results. ■

Customers Who Have Completed Service. In this final subsection, supplementing the application of the snapshot principle in §6, we consider the estimators based on the delays experienced by previous customers to *complete* service. Unlike for the LES and HOL estimators, we find that the LCS estimator behaves very differently in the classical and QED HT regimes. The way to see this is to observe that the LCS customer completed service a full service time in the past. That LCS customer arrived a waiting time plus a service time in the past.

In both heavy-traffic regimes, the service time is an exponential random variable with mean 1. In the classical HT regime, the waiting times are exploding in heavy traffic, so that a service time is negligible compared to the waiting time. Thus we see that LCS will be asymptotically equivalent to LES and HOL in the classical HT regime, for any fixed number of servers. The LCS estimator will be consistent as well in the classical heavy-traffic regime.

However, the story is very different in the QED HT regime. The service times remain unchanged, but now the waiting times become smaller, being of order $O(1/\sqrt{s})$. Now the service time is the same order as the time scaling. The stochastic-process limit in (6.2) describes the waiting time experience of each customer, but for the last customer to complete service at time t , we have a different limit. Let $A_{s,\rho}^L(t)$ denote the arrival time of the last customer to

complete service at time t in model (s, ρ) . The relevant limit now will be

$$\sqrt{s}W_{s,\rho}(A_{s,\rho}^L(t)) \Rightarrow Y(t - V) \quad \text{as } \rho \uparrow 1, \quad (3.42)$$

where $Y(t)$ is the limit process in (6.2) and V is a service time, an exponential random variable with mean 1. In other words, the waiting time at time t is approximately $Y(t)/\sqrt{s}$, while the waiting time of the last customer to complete service immediately prior to time t is approximately $Y(t - V)/\sqrt{s}$. Thus, in the QED HT limit the LCS estimator is *not* consistent. The effectiveness of the LCS estimator depends on the difference between $Y(t - V)$ and $Y(t)$. However, we do not attempt to do further analysis; here we are content to observe that the LCS estimator has inferior asymptotic performance in the QED HT regime. That is consistent with our simulation results, which show that the LCS estimator performs poorly for large s .

Fortunately, there is better information that we can obtain from customers who have already completed service in the QED HT regime. Other customers who have completed service are very likely to have arrived much more recently than the last customer to complete service. The minimum service time among the last m customers to complete service is $1/m$. Since the waiting times are of order $1/\sqrt{s}$, it is natural to consider $m = O(\sqrt{s})$; then the minimum service time among these customers also will be of order $O(1/\sqrt{s})$.

As a bound, first consider the customer among the last $c\sqrt{s}$ customers to complete service with the minimum service time. That customer's service time is exponentially distributed with mean $1/c\sqrt{s} = O(1/\sqrt{s})$. By (6.2), the customer's waiting time is also of order $O(1/\sqrt{s})$. Since the times between successive service completions are i.i.d. exponential random variables with mean $1/s$, the last $c\sqrt{s}$ service completions occur over a time interval having mean $c/\sqrt{s} = O(1/\sqrt{s})$. Hence this customer arrived $O(1/\sqrt{s})$ in the past. Hence we deduce that if we consider the customer among the last $c\sqrt{s}$ customers to complete service with the minimum service time, then that delay estimator is consistent in the QED HT regime.

Even better will be the RCS and RCS- $c\sqrt{s}$ estimators, because those customers necessarily arrive at least as recently. We summarize these conclusions in the following theorem. To state the theorem, let $W_{\infty,s,\rho}^{RCS}$ and $W_{\infty,s,\rho}^{RCS-c\sqrt{s}}$ be the steady-state RCS and RCS- $c\sqrt{s}$ delays in model (s, ρ) ; and let $W_{RCS,s,\rho}(w)$ and $W_{RCS-c\sqrt{s},s,\rho}(w)$ be the associated random variables having the conditional distribution of the delay to be estimated given the observed RCS and RCS- $c\sqrt{s}$ delays.

Theorem C.3. (*performance of LCS, RCS and RCS- $c\sqrt{s}$ in the QED HT regime*) If $\rho \uparrow 1$ and $s \uparrow \infty$ so that $(1 - \rho)\sqrt{s} \rightarrow \beta$ for $0 < \beta < \infty$ in the family of GI/M/s models indexed by

s and ρ , then the RCS and RCS- $c\sqrt{s}$ estimators are relatively consistent, i.e.,

$$\frac{W_{RCS,s,\rho}(W_{\infty,s,\rho}^{RCS})}{W_{\infty,s,\rho}^{RCS}} \Rightarrow 1 \quad \text{and} \quad \frac{W_{RCS-c\sqrt{s},s,\rho}(W_{\infty,s,\rho}^{RCS-c\sqrt{s}})}{W_{\infty,s,\rho}^{RCS-c\sqrt{s}}} \Rightarrow 1, \quad (3.43)$$

but the LCS estimator is not relatively consistent.

In this relatively crude sense, the estimators LES, HOL, RCS and RCS- $c\sqrt{s}$ are all asymptotically equivalent in the QED regime, but LCS is not. However, it remains to describe the asymptotic efficiency of RCS and RCS- $c\sqrt{s}$, paralleling the results for the HOL (and LES) estimator SCV's in (3.16) and (3.17).

D. A Pathological Example for LES

We have drawn very positive conclusions about the LES delay estimator $W_{LES}(w)$ in the $GI/M/s$ queue. To provide some balancing perspective, in this section we demonstrate potential weaknesses of the estimator $W_{LES}(w)$ for other service-time distributions. To illustrate the possible deficiencies of the LES estimator, we consider a specific stable $D/G/1$ queueing model with non-exponential service-time distribution in light traffic. Let the arrival process be deterministic with interarrival times 1.

We deliberately choose a difficult service-time distribution: let the service-time distribution be a two-point probability distribution, which usually assumes a very small value ϵ , but occasionally takes a very large value M ; specifically, let

$$P(V = M \gg 1) = \delta = 1 - P(V = \epsilon \ll 1), \quad (4.1)$$

where the traffic intensity

$$\rho \equiv E[V]/E[U] = E[V] = \delta M + (1 - \delta)\epsilon. \quad (4.2)$$

We suppose that δ is very small, so that ρ itself is very small and the service time is only equal to the large value M very rarely. If δ is sufficiently small, relatively few customers will have to wait in queue before starting service, but occasionally a customer will have one of the very long service times.

To see the deficiencies of the LES estimator, we will consider an epoch at which a customer with service time M arrives at an empty system. If δ is small enough, then with high probability the customer with the large service time M will not have to wait before starting service, but he will remain in service for a long time, precisely M . Thus the following M customers will

all have to wait before starting service. For each of them, however, the last served customer to have entered service – the customer with service time M – will have not had to wait at all.

To quantify the effect, let us call the customer with service time M customer 0. Then, assuming that these following M customers themselves all have ϵ service times (which has high probability), customer k will have to wait precisely $M - k + (k - 1)\epsilon$ before starting service. Customer number M will have to wait only $(M - 1)\epsilon$. But, for all M customers with positive waiting times, the last served customer will have waited 0 before starting service.

To go further, suppose that ϵ is very small, so that $(M - 1)\epsilon$ is itself less than 1. Then customer M will have to wait less than 1 before starting service, so that $M + 1$ will not have to wait at all before starting service. We thus have the strange estimation phenomenon: *The delay of the last served customer is 0 for all customers that themselves experience positive delays.* Thus, whenever an estimation needs to be made (because the customer must wait in queue), the estimated delay will be 0. Moreover, the actual delays of these customers who have to wait may be quite large: as large as $M - 1$ and averaging about $M/2$ for all customers forced to wait. This example allows arbitrarily large M , but after choosing M , we must choose ϵ and δ suitably small.

We have only described one possible scenario. The story we have described breaks down when two or more customers with large service time M interact, but by choosing δ sufficiently small, this deviation from the story can be made to occur relatively rarely. Thus the phenomenon we have described will hold for the vast majority of the customers that are delayed.

We can make the situation described above apply w.p.1 if we abandon the condition of i.i.d. service times. If we instead assume that customers $2kM$ have service times M , while all other customers have service time ϵ with $\epsilon < 1/M$ (e.g., $\epsilon = 0$), then we obtain the scenario above w.p.1. In addition, the average delay is approximately $M/4$, so the average delay can be made arbitrarily large by choosing M large. Thus this scenario does not only apply in very light traffic. Nevertheless, we regard this example as pathological. We are thinking of situations in which the delay of a new arrival should not be too different from the delay of the last customer to enter service.

For this example, the HOL estimator would fare somewhat better, but it would not do so great either. Given the scenario described above, when the customer at the head of the line has waited $w = k$, the random variable $W_{HOL}(w)$ depicting the delay of this new arrival is very likely to take the value $M - k + (k - 1)\epsilon$ instead of w .

References

- Abate, J. and W. Whitt. 1994. A heavy-traffic expansion for asymptotic decay rates of tail probabilities in multichannel queues. *Operations Res. Letters* 15, 223–230.
- Choudhury, G. L. and W. Whitt. 1994. Heavy-traffic asymptotic expansions for the asymptotic decay rates in the BMAP/G/1 Queue. *Stochastic Models* 10, 453–498.
- Cooper, R. B. 1981. *Introduction to Queueing Theory*, second edition, North-Holland, New York.
- Glynn, P. W. and W. Whitt. 1994. Logarithmic Asymptotics for steady-state tail probabilities in a single-server queue. *J. Appl. Prob.* 31A, 131–156 (also called *Studies in Applied Probability, Papers in Honour of Lajos Takacs*, J. Galambos and J. Gani (eds.), Applied Probability Trust, Sheffield, England).
- Halfin, S. and W. Whitt. 1981. Heavy-traffic limits for queues with many exponential servers. *Operations Research* 29, 567–588.
- Iglehart, D. L. and W. Whitt. 1970. Multiple channel queues in heavy traffic II: sequences, networks, and batches. *Advances in Applied Probability* 2, 355–369.
- Jelenkovic P., A. Mandelbaum A. and P. Momcilovic. 2004. Heavy traffic limits for queues with many deterministic servers. *Queueing Systems* 47, 53–69.
- Neuts, M. F. 1986. The caudal characteristic curve of queues. *Adv. Appl. Probab.* 18, 221–254.
- Smith, W. L. 1953. On the distribution of queueing times. *Proc. Camb. Phil. Soc.* 49, 449–461.
- Whitt, W. 1984. On approximations for queues, I: extremal distributions. *AT&T Bell Lab. Tech. J.* 63, 115–138.
- Whitt, W. 2004b. A diffusion approximation for the G/GI/n/m queue. *Operations Research* 52, 922–941.
- Whitt, W. 2005. Heavy-traffic limits for the G/H2*/n/m queue. *Math. Oper. Res.* 30, 1–27.