

Real-Time Delay Estimation Based on Delay History in Many-Server Service Systems with Time-Varying Arrivals

Rouba Ibrahim, Ward Whitt

Industrial Engineering & Operations Research, Columbia University, New York, NY 10027, USA
 rei2101@columbia.edu, ww2040@columbia.edu

Motivated by interest in making delay announcements in service systems, we study real-time delay estimators in many-server service systems, both with and without customer abandonment. Our main contribution here is to consider the realistic feature of time-varying arrival rates. We focus especially on delay estimators exploiting recent customer delay history. We show that time-varying arrival rates can introduce significant estimation bias in delay-history-based delay estimators when the system experiences alternating periods of overload and underload. We then introduce refined delay-history estimators that effectively cope with time-varying arrival rates together with non-exponential service-time and abandonment-time distributions, which are often observed in practice. We use computer simulation to verify that our proposed estimators outperform several natural alternatives.

Key words: delay estimation; delay announcements; time-varying arrival rates; simulation
History: Received: April 2009; Accepted: June 2010 by Michel Pindo, after 2 revisions.

1. Introduction

We investigate alternative ways to estimate, in real time, the delay (before entering service) of an arriving customer in a service system with time-varying arrival rates. We consider time-varying arrival rates because arrival processes to service systems, in real life, typically vary significantly over time.

Our delay estimators may be used to make delay announcements. Delay announcements may be especially helpful when delays are sometimes long, as in a hospital emergency department (ED). In many cases waiting customers are unable to accurately estimate their own delay, and would therefore gain from delay announcements. That is typically true with invisible queues, as occur in call centers; see Aksin et al. (2007) for background on call centers.

1.1. Delay-History-Based Estimators

In this paper, we examine alternative estimators based on recent customer delay history in the system. As in Armony et al. (2009), a candidate delay estimator based on recent customer delay history is the delay of the last customer to have entered service, before our customer’s arrival at time t , denoted by the last customer to enter service (LES). That is, letting w be the delay of the last customer to have entered service, the corresponding LES delay estimate is $\theta_{LES}(t, w) \equiv w$. Armony et al. (2009) 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.

Closely related to LES is the elapsed waiting time of the customer at the head of the line (HOL), assuming that there is at least one customer waiting at the new arrival epoch. The HOL delay estimator was mentioned as a candidate delay announcement by Nakibly (2002). For a detailed discussion of the HOL and LES estimators, see Ibrahim and Whitt (2009a,b). Experience indicates that the LES and HOL estimators have very similar performance. In complex systems, the LES delay is more likely to be observable than the HOL delay, because arrival and service completion times are more likely to be known than the experience of customers who have not yet completed their service; e.g., customers may have abandoned and that might not be known. Nevertheless, here we focus on HOL, because it is easier to analyze. However, we do so with the understanding that similar results will hold for LES.

1.2. Motivation For Delay-History-Based Estimators

We now briefly explain why it is important to study the performance of delay-history-based estimators; for more discussion, see section 1 of Ibrahim and Whitt (2009a). First, delay-history-based estimators are

currently used in service systems. For one example, the US Citizenship and Immigration Service (USCIS) publishes the arrival time of the most recently completed application to give an idea about upcoming delays. For another example, the HOL estimator was used as an announcement in an Israeli bank studied by Mandelbaum et al. (2000).

Second, delay-history-based estimators are appealing for complicated service systems. For one example, there may be multiple customer classes with multiple service pools. For another example, with web chat, servers typically 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 yet another example, consider ticket queues studied by Xu et al. (2007). Upon arrival at a ticket queue, each customer is issued a numbered ticket. The number currently being served is displayed. The queue length (QL) is not known to ticket-holding customers or even to system managers, because they do not observe customer abandonments. Even in systems with no customer abandonment, we may not know the QL in the system at a new arrival epoch. In a ticket queue (as at a supermarket), a ticketed customer may elect to go and do other shopping and plan to come back later to get in line. (Customers may also abandon, but that does not have to be the case.) Customers with tickets could return to the queue at some point in time and “preempt” customers who are already in line (e.g., if they have a lower numbered ticket). Now, suppose that there is a new arrival at the station. It is unclear whether ticketed customers (currently doing some other shopping) will return quickly enough to be inserted before that new arrival. Therefore, the QL cannot be determined at the new arrival epoch. Nevertheless, it is possible to determine who the LES (or HOL) customer is, and to know his/her delay.

Delay-history-based estimators are appealing, from a practical perspective, whenever the QL is not known, but also because they do not depend on the model and use very little information about the system. They are robust because they respond automatically to changes in system parameters (e.g., number of servers, mean service time, and arrival rate).

To fully understand a complex service system, we need to study it in detail. However, to help develop a service science, we are systematically studying various delay estimators in controlled environments, i.e., in structured models, starting with $GI/M/s$ and extending to $GI/GI/s$ (non-exponential service times), $GI/GI/s+GI$ (abandonment with non-exponential patience distributions) in Ibrahim and Whitt (2009a,b) and now $M_t/GI/s$ and $M_t/GI/s+GI$ (time-varying arrival rates).

1.3. The Case of a Stationary Arrival Process

In Ibrahim and Whitt (2009a,b), we studied the performance of the LES and HOL delay estimators in many-server systems, both with and without customer abandonment, by studying conventional stationary queueing models. In Ibrahim and Whitt (2009a), we studied the performance of HOL in the $GI/M/s$ queueing model, which has a renewal arrival process, s homogeneous servers, an unlimited waiting room, and the first-come-first-served service discipline. The service times are independent of the arrival process, and independent and identically distributed (i.i.d.) exponential random variables.

We showed that HOL is an effective estimator in the $GI/M/s$ model. As a frame of reference, we considered the classical delay estimator based on the QL which multiplies the QL plus one times the mean interval between successive service completions, ignoring customer abandonment. For this special idealized model with i.i.d. exponential service times and no customer abandonment, the QL estimator is provably the most effective estimator, under the mean squared error (MSE) criterion; see section 4. The HOL estimator performs worse than QL, because it does not exploit QL information. Nevertheless, we showed that the difference in performance need not be too great, particularly when the arrival process has low variability. Because the model is highly structured, we were able to obtain analytical results.

In Ibrahim and Whitt (2009b), we considered the $GI/GI/s+GI$ model, which includes independent sequences of i.i.d. service times and abandonment times with general distributions. As one would expect, QL can overestimate customer delay when there is significant customer abandonment in the system. We showed that QL performs poorly in a heavily loaded $GI/GI/s+GI$ model, while HOL remains an effective estimator.

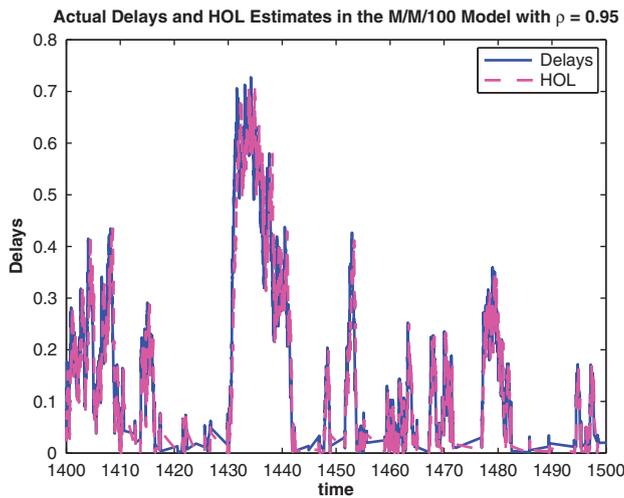
When customer abandonment is a serious issue, it is possible to refine the QL-based delay estimator by using the exact expected conditional delay, given the QL, in the $G/M/s+M$ model; we denote this by QL_m . However, for non-exponential service-time and abandonment distributions, the delay-history-based estimators can also outperform this refined QL-based estimator QL_m even when the QL and the model are known; e.g., see figures 1–4 of Ibrahim and Whitt (2009b).

However, we do not mean to suggest that the QL does not provide useful information when it is known. Indeed, our best estimator for the $GI/GI/s+GI$ model is an approximation-based estimator, referred to as QL_{ap} , which exploits the QL as well as model parameters; we also will make use of QL_{ap} here for the $M_t/GI/s+GI$ model in section 8.

1.4. Time-Varying Arrival Rates

In this paper, we study the performance of the HOL estimator with time-varying arrival rates. We do so

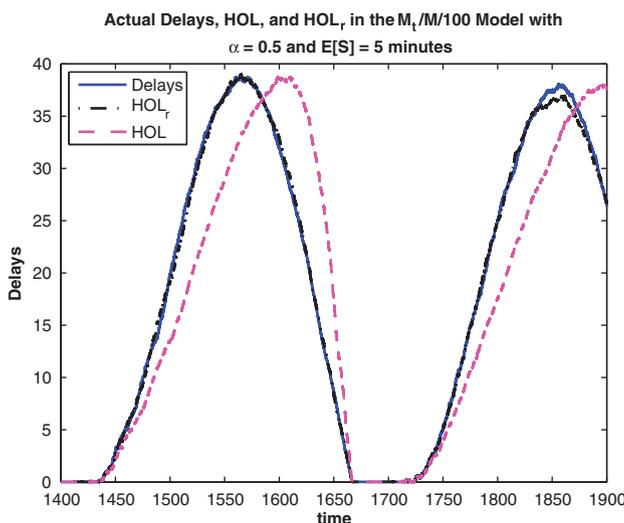
Figure 1 Sample Paths of Actual Delays and Head of the Line (HOL) Delay Estimates with Constant Arrival Rate



primarily because arrival rates typically vary significantly over time in real-life service systems.

The HOL estimator can perform poorly when the delays vary systematically over time, as can occur when there are alternating periods of significant overload and underload. Then the delay of a new arrival may not be like the HOL delay. To demonstrate potential problems with the HOL estimator, we plot simulation sample paths of HOL delay estimates given, and actual delays observed, as a function of time, in simulation runs from two different heavily loaded many-server systems. In Figure 1, we consider the stationary $M/M/100$ model with traffic intensity $\rho = 0.95$ and mean service time 5 minutes; in Figure 2, we consider the $M_t/M/100$ model with sinusoidal arrival rates, again with traffic intensity $\rho = 0.95$, but now defined as the long-run average, and mean service time 5 minutes. We consider a daily

Figure 2 Sample Paths of Actual Delays and Delay Estimates Using Head of the Line (HOL) and HOL_r with Sinusoidal Arrival Rate



cycle, so that there is one peak during the day. We let the relative amplitude be $\alpha = 0.5$. (The ratio of the peak arrival rate to the average arrival rate is $1 + \alpha$.) We measure time and, thus, the delays in units of mean service times. The overall plotted time interval of length 500 mean service times is slightly less than 2 days, so we see two peaks.

For Figure 2, we deliberately chose an extreme case in which the system alternates between extreme overload and underload, while the number of servers remains fixed. In that setting, the maximum delays themselves are about 40 mean service times or 200 minutes, about 60 times greater than in the stationary environment. Delay estimation tends to be especially important with such large delays. Figure 2 shows that, with time-varying arrival rates, the HOL curve is clearly shifted to the right of the actual-delay curve; i.e., there is a time lag between the HOL estimates and the actual delays observed, leading to big errors.

Figure 2 also shows a third plot, the plot of a refined HOL estimator, denoted by HOL_r , which we develop in section 4. Clearly, it eliminates the time lag; visually the HOL_r plot falls on top of the actual delays. The ratio of the average squared errors $ASE(HOL)/ASE(HOL_r)$, defined in section 3, is about 95 in Figure 2. (If we would reduce the relative amplitude from 0.5 to 0.1, then the ratio would be only 1.3; it then requires careful analysis to see the improvement provided by HOL_r over HOL; see Ibrahim and Whitt [2009c] for the plot.)

In this paper, we not only show that HOL may not be an effective estimator with time-varying arrivals, particularly when the system alternates between phases of underload and overload, but we also develop refinements of the HOL estimator that remain effective for time-varying arrival rates. Through analysis and simulation, we show that these new estimators perform remarkably well with time-varying arrival rates, far better than HOL.

However, the improved performance of the refined HOL estimators comes at the expense of exploiting more information about the system, such as the arrival rate, the number of servers, and the mean service time. That requirement greatly reduces the advantage over QL-based delay estimators. Indeed, our strategy for obtaining the refined HOL estimators involves two steps: (i) representing or approximating the expected conditional delay given the QL and (ii) estimating the QL, given the observed HOL delay and the model parameters. Hence, the refined HOL estimators are valuable only when the QL is not known. However, such cases are not uncommon, as in web chat and ticket queues, when we directly observe arrivals and service completions, but not the queue, because we do not observe customer abandonments.

Because our refined estimators exploit more information about the system, we also investigate (i) how our refined estimators perform if the extra information

is known imperfectly, because it too must be estimated, and (ii) how this additional information can be estimated in real time. We propose estimation procedures for alternative system parameters, and quantify the estimation error resulting from those procedures. These additional experiments show that the refined estimators can be useful in practice.

1.5. Literature Review, Contributions, and Organization

The literature on delay announcements is large and growing. In broad terms, there are two main areas of research. The first area studies the effect of delay announcements on system dynamics; e.g., see Whitt (1999b), Armony and Maglaras (2004), Guo and Zipkin (2007), Armony et al. (2009), Allon et al. (2009), and references therein. The second area studies alternative ways of estimating customer delay in service systems; e.g., see Nakibly (2002), Whitt (1999a), Jouini et al. (2007), and Ibrahim and Whitt (2009a, b). For a more detailed review, see section 2 of Jouini et al. (2007).

This paper falls in the second main area of research. Our main contributions are the following: (i) to show that time-varying arrival rates can cause estimation bias for delay-history-based delay estimators, (ii) to propose new and easily implementable delay estimators, based on the history of delays in the system, that effectively cope with time-varying arrivals and general service-time and abandon-time distributions, (iii) to provide analytical results quantifying the performance of some delay estimators, and (iv) to describe results of a wide range of simulation experiments evaluating alternative delay estimators, with time-varying arrivals.

The rest of this paper is organized as follows: In section 2, we describe the modeling framework. In section 3, we describe measures quantifying the performance of our candidate delay estimators. In section 4, we introduce a new delay estimator for the $M_t/GI/s$ model. In section 5, we provide analytical results for the performance of this estimator in the $M_t/M/s$ model. In section 6, we present simulation results showing that it is effective in the $M_t/GI/s$ model. In section 7, we propose ways of obtaining the additional system information required for implementing the new delay estimator of section 4. In section 8, we develop a new delay estimator for the $M_t/GI/s+GI$ model. In section 9, we present simulation results showing that it is effective. We make concluding remarks in section 10. Additional material appears in Ibrahim and Whitt (2009c), available on the authors' web pages.

2. The Framework

We consider many-server queueing models with time-varying arrival rates, both with and without customer abandonment. We model the arrival process as a

non-homogeneous Poisson process, which is the accepted model for capturing time-varying arrivals. It is completely characterized by its deterministic arrival-rate function $\lambda \equiv \{\lambda(u); -\infty < u < \infty\}$. There is statistical evidence suggesting that a non-homogeneous Poisson process is a good fit for the arrival process to a call center; see Brown et al. (2005). We adopt this model for arrivals, although we recognize its shortcomings. For example, this model does not reproduce an essential feature of call center arrivals, which is the overdispersion of the number of arrivals relative to the Poisson distribution (i.e., the variance is larger than the mean); see Avramidis et al. (2004). Moreover, the arrival rate in a real-life system is often not known with certainty. Therefore, it could be assumed to be a random variable; see Jongbloed and Koole (2001). It is natural, however, to begin an investigation in a relatively tractable setting, for which we are able to obtain analytical results. Our results provide useful background for similar studies in even more complicated settings.

In sections 4–6, we consider the $M_t/GI/s$ model, which has a non-homogeneous Poisson arrival process, i.i.d. service times distributed as a random variable S with a general distribution, having mean $E[S] = \mu^{-1}$ and no customer abandonment. Motivated by large service systems, we are primarily interested in the case of large s , which we take to be fixed. It is possible to choose appropriate time-varying staffing (making s a function of time) so that delays are stabilized at low levels; e.g., see Green et al. (2007). However, in practice there often is not adequate flexibility in setting staffing levels. Our fixed staffing assumption captures the spirit of such situations. We leave to future research the important extension to time-varying staffing levels.

Our delay estimators apply to arbitrary arrival-rate functions, but to analyze the performance of these estimators we restrict attention to periodic arrival-rate functions, under which the queueing system has a dynamic steady state, provided that the average arrival rate, denoted by $\bar{\lambda}$, is strictly less than the maximum possible service rate, $s\mu$; e.g., see Heyman and Whitt (1984). For our analysis, both analytically and by simulation, we further restrict attention to the special case of sinusoidal arrival rates. That is commonly done in studies of queues with time-varying arrivals; e.g., see Green et al. (2007) and references therein. Sinusoidal arrival rates capture the spirit of daily cycles.

In sections 8 and 9 we consider the $M_t/GI/s+GI$ model, which adds customer abandonment. The abandonment times are i.i.d. with mean v^{-1} and a general cumulative distribution function (cdf) F . As in Ibrahim and Whitt (2009b), we see that the abandonment distribution has a significant impact.

3. Performance Measures for the Delay Estimators

In this section, we indicate how we evaluate the performance of our candidate delay estimators. We use computer simulation to do the actual estimation. In our simulation experiments, we quantify the performance of a delay estimator by computing the *average squared error* (ASE), defined by

$$ASE \equiv \frac{1}{k} \sum_{i=1}^k (p_i - e_i)^2, \quad (1)$$

where $p_i > 0$ is the potential waiting time of delayed customer i , e_i is the delay estimate given to customer i , and k is the number of customers in our sample. In our simulation experiments, we measure p_i for both served and abandoning customers. For abandoning customers, we compute the delay experienced, had the customer not abandoned, by keeping him “virtually” in queue until he would have begun service. Such a customer does not affect the waiting time of any other customer in queue. Because we measure time in units of mean service times, the ASE is given in units of mean service time squared per customer.

As discussed in Ibrahim and Whitt (2009a,b), the ASE approximates the expected MSE for a system in steady state with a constant arrival rate, but the situation is more complicated with time-varying arrivals. We regard ASE as directly meaningful, but now we indicate how it relates to the MSE. Let $W_{HOL}(t, w)$ represent a random variable with the conditional distribution of the potential delay of an arriving customer, given that this customer must wait before starting service, and given that the elapsed delay of the customer at the HOL at the time of his arrival, t , is equal to w . Let $\theta_{HOL}(t, w)$ be some given single-number delay estimate which is based on the HOL delay, w , and the time of arrival, t . Then, the MSE of the corresponding delay estimator is given by

$$MSE(\theta_{HOL}(t, w)) \equiv E[(W_{HOL}(t, w) - \theta_{HOL}(t, w))^2], \quad (2)$$

which is a function of w and t . In order to obtain the overall MSE of HOL at time t , we average with respect to the *unconditional* distribution of the HOL waiting time at time t , $W_{HOL}(t)$, i.e.,

$$MSE(t) \equiv E[MSE(\theta_{HOL}(t, W_{HOL}(t)))]. \quad (3)$$

Finally, in order to relate the ASE in (1) to the MSE, we need to average $MSE(t)$ defined in (3) appropriately over time, but because the ASE represents a customer average instead of a time average, we need to use a weighted time average of the time-dependent MSE in (2) in order to relate it to the ASE.

In particular, if T is the cycle length, then

$$ASE \approx \frac{\int_0^T \lambda(u) MSE(u) du}{\int_0^T \lambda(u) du}, \quad (4)$$

where $MSE(t)$ is defined in (3); for supporting theory see the appendix of Massey and Whitt (1994).

In addition to the ASE, we quantify the performance of a delay estimator by computing the *root relative average squared error* (RRASE), defined by

$$RRASE \equiv \frac{\sqrt{ASE}}{(1/k) \sum_{i=1}^k p_i}, \quad (5)$$

using the same notation as in (1). The denominator in (5) is the average potential waiting time of customers who must wait. The RRASE is useful because it measures the effectiveness of an estimator relative to the average potential waiting time, given that the customer must wait.

4. Delay Estimators for the $M_t/GI/s$ Model

In this section, we propose a new refined HOL-based delay estimator, HOL_r for the $M_t/GI/s$ model. Our idea is to use the refined estimator $\theta_{HOL}^r(t, w) \equiv E[W_{HOL}(t, w)]$ instead of the HOL estimator $\theta_{HOL}(t, w) \equiv w$, because the mean necessarily minimizes the MSE based on this information. However, this mean is difficult to compute, so we propose an approximation. We approximate the mean in the given $M_t/GI/s$ model by its exact value in the corresponding $M_t/GI/s$ model, with exponential service time having the given mean $E[S]$.

For the $M_t/M/s$ model, we have the representation

$$W_{HOL}(t, w) \equiv \sum_{i=1}^{A(t)-A(t-w)+2} S_i/s, \quad (6)$$

where $\{A(t): t \geq 0\}$ denotes the arrival (counting) process. We have division by s in (6) 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 rate μ , which makes the minimum exponential with rate $s\mu$. The random variable $A(t) - A(t - w)$ has a Poisson distribution with mean $\int_{t-w}^t \lambda(u) du$. Since $W_{HOL}(t, w)$ in (6) is a random sum of i.i.d. random variables, where $A(t) - A(t - w)$ is independent of the summands S_i/s , we can easily compute this mean. Hence our refined HOL estimator for the $M_t/GI/s$ model is this mean

$$\begin{aligned} \theta_{HOL_r}(t, w) &\equiv E[W_{HOL, M_t/M/s}(t, w)] \\ &= \frac{1}{s\mu} \left(2 + \int_{t-w}^t \lambda(u) du \right). \end{aligned} \quad (7)$$

In general, with a non-exponential service-time distribution, $\theta_{HOL_r}(t, w)$ in (7) need not equal $E[W_{HOL}(t, w)]$, because many remaining service times at time t are residual service times for service times begun before time t . Consequently, these service times have a different distribution than the original service time. However, we can make stochastic comparisons. A cdf G of a non-negative random variable is said to be new better (worse) than used—NBU (NWU)—if $G_t^c(x) \equiv G^c(t+x)/G^c(t) \leq (\geq) G^c(x)$ for all $t \geq 0$ and $x \geq 0$, where $G^c(x) \equiv 1 - G(x)$; see Barlow and Proschan (1975, p. 159). In the parlance of survival analysis, a cdf is NBU (NWU) if the probability of surviving for an additional x time units, given survival up to time t , decreases (increases) with t .

PROPOSITION 1. *If the service-time cdf is NBU (NWU), then $\theta_{HOL_r}(t, w) \geq (\leq) E[W_{HOL}(t, w)]$.*

PROOF. The NBU and NWU condition means that the residual service times are stochastically ordered compared with the original service times. Intuitively, approximating an NBU (NWU) distribution by an exponential leads to overestimating (underestimating) the residual service times, and thus the overall delay. Given the elapsed times, the remaining service times are mutually independent. The minimum (the time until the next departure) is thus stochastically ordered compared with the minimum of mutually independent original service-time distributions. The random variable $W_{HOL}(t, w)$ is the sum of several of those intervals between successive departures. Even though those intervals may be dependent, the mean of the sum is the sum of the means. Hence the means are ordered, as claimed. \square

More importantly, simulation shows that HOL_r provides a good approximation even when the service-time distribution is not nearly exponential; see section 6.

We conclude this section by reviewing the QL estimator, previously considered in Ibrahim and Whitt (2009a, b). Let $W_Q(t, n)$ represent a random variable with the conditional distribution of the delay of an arriving customer, given that this customer must wait before starting service, and given that the QL seen upon arrival, at time t , is equal to n . Again, the QL estimator is obtained by using the exact expected value $E[W_Q(t, n)]$ for the corresponding $M_t/M/s$ model with the same mean service time.

In the $M_t/M/s$ model, $W_Q(t, n)$ is the sum of $n+1$ i.i.d. exponential random variables, each with rate $s\mu$. The QL estimate given to a customer who finds n other customers in queue upon arrival is $\theta_{QL}(t, n) \equiv E[W_Q(t, n)] = (n+1)/s\mu$, which depends on t only through n , which is directly observable. The optimal delay estimator, conditional on the number of customers, n , seen in line at time t , using the MSE criterion, is the one announcing the mean, $E[W_Q(t, n)]$. That is why the QL

estimator is the optimal delay estimator, under the MSE criterion, in the $M_t/M/s$ model.

By essentially the same reasoning as for Proposition 1, we can obtain bounds for the mean delay compared with $\theta_{QL}(t, n)$ when the service-time cdf is NBU or NWU.

PROPOSITION 2. *If the service-time cdf is NBU (NWU), then $\theta_{QL}(t, n) \geq (\leq) E[W_Q(t, n)]$.*

Fortunately, again simulation shows that QL remains effective in the $M_t/GI/s$ model, even when the service-time distribution is not nearly exponential; see section 6. For the $M_t/M/s$ model, we obtain analytical results quantifying the difference in performance between QL and HOL_r in the next section.

5. Analytical Expressions for the $M_t/M/s$ Model

The QL estimator has the desirable property that the estimation obtains relatively more accurate as the observed QL n increases. For the conditional waiting time at time t based on an observed QL of n , we have the representation

$$W_Q(t, n) \equiv \sum_{i=1}^{n+1} S_i/s. \quad (8)$$

The expectation, variance, and squared coefficient of variation (SCV, equal to the variance divided by the square of the mean) of $W_Q(t, n)$ are given by

$$\begin{aligned} E[W_Q(t, n)] &= \frac{n+1}{s\mu}, & \text{Var}[W_Q(t, n)] &= \frac{n+1}{s^2\mu^2}, \\ c_{W_Q(t, n)}^2 &\equiv \frac{\text{Var}[W_Q(t, n)]}{(E[W_Q(t, n)])^2} = \frac{1}{n+1}, \end{aligned} \quad (9)$$

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

To treat HOL_r , we use the representation in (6), which allows us to characterize the probability distribution of the random variable $W_{HOL}(t, w)$, in the $M_t/M/s$ model.

PROPOSITION 3. *For the $M_t/M/s$ model,*

$$\text{Var}[W_{HOL}(t, w)] = \frac{2}{s^2\mu^2} \left(1 + \int_{t-w}^t \lambda(u) du \right), \quad (10)$$

which, combined with (7), yields

$$c_{W_{HOL}(t, w)}^2 = \frac{\text{Var}[W_{HOL}(t, w)]}{(E[W_{HOL}(t, w)])^2} = 2 \frac{1 + \int_{t-w}^t \lambda(u) du}{(2 + \int_{t-w}^t \lambda(u) du)^2}. \quad (11)$$

PROOF. Formula (10) follows from the conditional variance formula, e.g., Ross (1996, p. 51). Formula (11) immediately follows from (7) and (10). \square

Since $\theta_{HOL_r}(t, w) \equiv E[W_{HOL}(t, w)]$ and $\theta_{QL}(t, n) \equiv E[W_Q(t, n)]$, we can compare the performance of HOL_r and QL by comparing the respective SCV's in (9) and (11). (When the delay estimate equals the conditional mean, the MSE coincides with the variance.)

To obtain further results, we consider a sinusoidal arrival-rate function

$$\lambda(u) = \bar{\lambda} + \beta \sin(\gamma u) \equiv \bar{\lambda} + \bar{\lambda}\alpha \sin(2\pi u/\Gamma) \quad (12)$$

for $-\infty < u < \infty$,

where $\bar{\lambda}$ is the average arrival rate, α is the relative amplitude, and Γ is the cycle length. (We define $\beta \equiv \bar{\lambda}\alpha$ and $\gamma \equiv 2\pi/\Gamma$.) Given the cycle length, Γ , we can deduce the place where any time u falls within the cycle, in dynamic steady state. Henceforth, we focus solely on the interval $0 \leq u \leq \Gamma$, which describes a full cycle.

With sinusoidal arrival rates, we obtain analytical results comparing the performance of the QL and HOL_r estimators. We determine the limit of the ratio of the SCV's as $n \rightarrow \infty$. Formula (13) coincides with formula (4.25) of Ibrahim and Whitt (2009a) for the stationary $GI/M/s$ model. As before, the condition $n \rightarrow \infty$ arises naturally in heavy traffic, either with fixed s or as $s \rightarrow \infty$; e.g., see Garnett et al. (2002). (When $s \rightarrow \infty$ along with the arrival rate, the QL is of order s and \sqrt{s} in the ED and QED regimes.) Recall that $\rho \equiv \bar{\lambda}/s\mu$.

PROPOSITION 4. For the $M_t/M/s$ model with sinusoidal arrival rates,

$$\frac{c_{W_{HOL}(t,w)}^2}{c_{W_Q(n)}^2} \rightarrow \frac{2}{\rho} \text{ as } n \rightarrow \infty, \quad (13)$$

for all t , provided that $w/n \rightarrow 1/s\mu$.

PROOF. Using Equations (7), (10)–(12), we obtain the following expressions for the mean, variance, and SCV of $W_{HOL}(t, w)$, in the $M_t/M/s$ model with sinusoidal arrivals:

$$E[W_{HOL}(t, w)] = \frac{2 + \bar{\lambda}w + (\beta/\gamma)(\cos(\gamma t - \gamma w) - \cos(\gamma t))}{s\mu} \quad (14)$$

and

$$\text{Var}[W_{HOL}(t, w)] = 2 \frac{1 + \bar{\lambda}w + (\beta/\gamma)(\cos(\gamma t - \gamma w) - \cos(\gamma t))}{s^2\mu^2}, \quad (15)$$

which yields

$$\begin{aligned} c_{W_{HOL}(t,w)}^2 &= \frac{\text{Var}[W_{HOL}(t, w)]}{(E[W_{HOL}(t, w)])^2} \\ &= 2 \frac{1 + \bar{\lambda}w + (\beta/\gamma)(\cos(\gamma t - \gamma w) - \cos(\gamma t))}{[2 + \bar{\lambda}w + (\beta/\gamma)(\cos(\gamma t - \gamma w) - \cos(\gamma t))]^2} \end{aligned} \quad (16)$$

for $0 \leq t \leq \Gamma$. Using (16), and recalling that $-1 \leq \cos(u) \leq 1$ for all u , we obtain the following bounds for the SCV of $W_{HOL}(t, w)$:

$$\begin{aligned} \frac{2 + 2\bar{\lambda}w - 4\beta/\gamma}{(2 + \bar{\lambda}w + 2\beta/\gamma)^2} &\leq c_{W_{HOL}(t,w)}^2 \\ &\leq \frac{2 + 2\bar{\lambda}w + 4\beta/\gamma}{(2 + \bar{\lambda}w - 2\beta/\gamma)^2}. \end{aligned} \quad (17)$$

Let $W(t)$ be the potential waiting time at time t , the time that an arrival at t would have to wait before beginning service. Since

$$W(t) = \sum_{i=1}^{Q(t)+1} S_i/s, \quad (18)$$

where $Q(t)$ is the number of customers waiting in queue upon arrival at t , the law of large numbers implies that $W(t)Q(t) \rightarrow 1/s\mu$ as $Q(t) \rightarrow \infty$. Thus, when $Q(t)$ is large, we have $W(t) \approx Q(t)/s\mu$. Assuming that n in (9) is large with $w = n/s\mu + o(n)$ as $n \rightarrow \infty$, where $o(n)$ denotes a quantity that is asymptotically negligible when divided by n , and combining that with (17), for large n we obtain

$$\begin{aligned} \frac{(2 + 2\rho(n + o(n)) - 4\beta/\gamma)(n + 1)}{(2 + \rho(n + o(n)) + 2\beta/\gamma)^2} &\leq \frac{c_{W_{HOL}(t,w)}^2}{c_{W_Q(n)}^2} \\ &\leq \frac{(2 + 2\rho(n + o(n)) + 4\beta/\gamma)(n + 1)}{(2 + \rho(n + o(n)) - 2\beta/\gamma)^2} \end{aligned} \quad (19)$$

for all t . By a sandwiching argument, (19) yields (13) as $n \rightarrow \infty$. \square

6. Simulations Experiments for the $M_t/GI/s$ Model

In this section, we present simulation results for the $M_t/GI/s$ model, quantifying the performance of QL , HOL , and HOL_r with sinusoidal arrival rates. For the service-time distribution, we consider M (exponential), D (deterministic), and $LN(1, 4)$ (lognormal with mean equal to 1 and variance equal to 4). The $LN(1, 4)$ (D) distribution exhibits high (low) variability, relative to M . We consider a lognormal distribution because there is statistical evidence suggesting a good fit of the service-time distribution to the lognormal distribution in call centers; see Brown et al. (2005).

6.1. Description of the Experiments

We fix the number of servers, $s = 100$, because we are interested in large service systems. We consider non-homogeneous Poisson arrival processes with the sinusoidal arrival-rate functions in (12). We vary $\bar{\lambda}$ to obtain alternative values of ρ , for fixed s . We consider values of ρ ranging from 0.90 to 0.98. These values of ρ are chosen to let our systems alternate between periods of overload and underload. We consider two

Table 1 The Relative Frequency, γ , as a Function of the Mean Service Time $E[S]$ for a Daily Cycle

Relative frequency γ	Mean service time $E[S]$
0.0220	5 minutes
0.0436	10 minutes
0.131	30 minutes
0.262	1 hour
1.571	6 hours
3.14	12 hours
6.28	24 hours
12.6	48 hours

The relative frequency is the frequency computed with measuring units so that $E[S] = 1$.

values of the relative amplitude: $\alpha = 0.1$ and $\alpha = 0.5$. Simulation point and 95% confidence interval estimates are based on 10 independent replications of five million events each, where an event is either an arrival or a service completion. That is, each simulation run terminates when the sum of the number of arrivals and the number of service completions is equal to five million. Here, we show a sample of our simulation results; see Ibrahim and Whitt (2009c) for more.

The parameters of the arrival-rate intensity function, $\lambda(u)$ in (12), should be interpreted relative to the mean service time, $E[S]$. As in section 1.4, we measure time in units of mean service times; hence $\mu = 1$. Then, we refer to γ in (12) as the relative frequency. Table 1 displays values of the relative frequency as a function of $E[S]$, assuming a daily cycle. For interpretation, we also will specify the associated mean service time in minutes, given a daily cycle.

Here, we consider two different values of γ . First, we consider $\gamma = 0.131$, which corresponds to $E[S] = 30$ minutes, assuming a daily cycle. This choice of $E[S]$ could be used to describe the experience of waiting customers in a call center, for example. Second, we consider $\gamma = 1.57$, which corresponds to $E[S] = 6$ hours. This choice of $E[S]$ could be used to describe the experience of waiting patients in a crowded hospital ED. With $E[S] = 30$ minutes and $\alpha = 0.1$ ($E[S] = 6$ hours and $\alpha = 0.5$), and daily cycles, the arrival rate varies relatively slowly (rapidly) with respect to the service times.

In Table 2, we present simulation (point and 95% confidence interval estimates) quantifying the performance of QL, HOL_r , and HOL in the $M_t/GI/s$ model

Table 2 A Comparison of the Efficiency of QL, HOL_r , and HOL in the $M_t/GI/100$ Model, as a Function of the Traffic Intensity, ρ

ρ	$M_t/M/100, \alpha = 0.1, E[S] = 30$ minutes			$M_t/M/100, \alpha = 0.5, E[S] = 6$ hours		
	QL	HOL_r	HOL	QL	HOL_r	HOL
0.9	2.26 ± 0.051	4.29 ± 0.088	4.61 ± 0.098	2.24 ± 0.023	4.27 ± 0.033	9.01 ± 0.015
0.93	3.77 ± 0.10	7.29 ± 0.21	8.04 ± 0.26	2.83 ± 0.029	5.45 ± 0.063	14.1 ± 0.25
0.95	5.08 ± 0.072	10.1 ± 0.15	11.7 ± 0.20	3.49 ± 0.033	6.82 ± 0.073	21.4 ± 0.28
0.97	7.16 ± 0.098	14.1 ± 0.20	17.5 ± 0.24	4.82 ± 0.12	9.46 ± 0.22	39.0 ± 1.5
0.98	9.14 ± 0.30	18.0 ± 0.59	23.9 ± 1.0	6.77 ± 0.32	13.3 ± 0.62	63.3 ± 3.9
ρ	$M_t/LN(1, 4)/100, \alpha = 0.1, E[S] = 30$ minutes			$M_t/LN(1, 4)/100, \alpha = 0.5, E[S] = 6$ hours		
	QL	HOL_r	HOL	QL	HOL_r	HOL
0.9	4.36 ± 0.25	7.30 ± 0.34	7.78 ± 0.36	2.08 ± 0.13	3.60 ± 0.19	7.79 ± 0.33
0.93	6.89 ± 0.15	11.3 ± 0.34	12.8 ± 0.34	3.48 ± 0.18	5.90 ± 0.27	14.0 ± 0.49
0.95	9.82 ± 0.28	15.9 ± 0.42	19.0 ± 0.56	5.70 ± 0.14	9.52 ± 0.22	22.5 ± 0.38
0.97	17.2 ± 0.81	27.0 ± 1.3	35.1 ± 2.1	9.92 ± 0.60	15.9 ± 0.89	34.2 ± 1.1
0.98	23.2 ± 0.94	35.8 ± 1.4	48.9 ± 2.4	20.1 ± 2.2	31.0 ± 3.3	52.1 ± 3.2
ρ	$M_t/D/100, \alpha = 0.1, E[S] = 30$ minutes			$M_t/D/100, \alpha = 0.5, E[S] = 6$ hours		
	QL	HOL_r	HOL	QL	HOL_r	HOL
0.9	0.972 ± 0.025	2.31 ± 0.034	2.47 ± 0.036	3.02 ± 0.023	4.14 ± 0.039	7.35 ± 0.054
0.93	1.23 ± 0.024	3.84 ± 0.063	4.18 ± 0.078	3.71 ± 0.027	5.01 ± 0.026	8.91 ± 0.045
0.95	1.31 ± 0.027	5.19 ± 0.041	6.01 ± 0.041	4.33 ± 0.038	5.84 ± 0.051	10.5 ± 0.068
0.97	1.35 ± 0.026	7.26 ± 0.065	9.29 ± 0.038	5.41 ± 0.086	7.54 ± 0.075	15.5 ± 0.14
0.98	1.34 ± 0.042	8.29 ± 0.057	11.3 ± 0.069	6.01 ± 0.075	8.84 ± 0.076	21.1 ± 0.49

Point and 95% confidence interval estimates of the average squared error (ASE) are shown (in units of mean service time squared per customer). Estimated ASEs are in units of 10^{-3} .

HOL, head of the line; QL, queue length.

with M , $LN(1,4)$, and D service-time distributions. We discuss these results next.

6.2. Comparing HOL_r and HOL

Table 2 shows that, for $\alpha = 0.1$ and $E[S] = 30$ minutes, HOL_r performs better than HOL , particularly for high values of ρ . We obtain consistent results with M , $LN(1,4)$, and D service times: $ASE(HOL)/ASE(HOL_r)$ is roughly equal to 1 for $\rho = 0.9$, and roughly equal to 1.4 for $\rho = 0.98$. The case with high ρ corresponds to extreme fluctuations between phases of underload and overload, in which case HOL performs relatively poorly.

With $\alpha = 0.5$, and $E[S] = 6$ hours, the difference in performance between HOL and HOL_r is significant, for all ρ considered. For example, with D service times, $ASE(HOL)/ASE(HOL_r)$ ranges from about 1.8 for $\rho = 0.9$ to about 2.4 for $\rho = 0.98$. With M service times, $ASE(HOL)/ASE(HOL_r)$ ranges from about 2.1 for $\rho = 0.9$ to about 4.8 for $\rho = 0.98$. The HOL_r estimator is also relatively more accurate than HOL . For example, with $LN(1,4)$ service times, $RRASE(HOL_r)$ ranges from about 27% for $\rho = 0.9$ to about 15% for $\rho = 0.98$. In this case, $RRASE(HOL)$ ranges from about 38% for $\rho = 0.9$ to about 20% for $\rho = 0.98$.

6.3. Comparing HOL_r and QL

In the $M_t/M/s$ model, QL is provably the optimal estimator given the observed QL upon arrival, under the MSE criterion; see section 4. With $\alpha = 0.1$, $E[S] = 30$ minutes, and M service times, Table 2 shows that $RRASE(QL)$ ranges from about 21% for $\rho = 0.9$ to about 10% for $\rho = 0.98$. With non-exponential service times, QL remains the most effective estimator, under the MSE criterion. It is relatively accurate, in all models considered. For example, with $\alpha = 0.5$, $E[S] = 6$ hours, and $LN(1,4)$ service times, $RRASE(QL)$ ranges from about 20% for $\rho = 0.9$ to about 12% for $\rho = 0.98$.

Consistent with section 5, the approximation for the ratio of the SCV's in (13) provides a remarkably accurate approximation for the ratio of the ASE's with M service times, particularly for high values of ρ , as we would expect. (The distortion caused by the customer average in (4) is evidently minor.) For example, with $E[S] = 30$ minutes and $\alpha = 0.1$, Table 2 shows that the relative error between simulation point estimates for $ASE(HOL_r)/ASE(QL)$ and numerical values given by (13) is less than 3% for $\rho = 0.98$.

With $LN(1,4)$ service times, $E[S] = 30$ minutes, and $\alpha = 0.1$, Table 2 shows that $ASE(HOL_r)/ASE(QL)$ ranges from about 1.7 for $\rho = 0.9$ to about 1.5 for $\rho = 0.98$, which is less than predicted by (13). Similarly, with D service times, $E[S] = 6$ hours, and $\alpha = 0.5$, Table 2 shows that $ASE(HOL_r)/ASE(QL)$ is approximately equal to 1.5 for all ρ .

7. Estimating the Required Additional Information for HOL_r

We have shown, both analytically and using simulation, that the HOL estimator can perform poorly when the arrival rate varies considerably over time while the staffing is fixed. We showed that the new refined HOL estimator, HOL_r , performs remarkably better than HOL in the $M_t/GI/s$ queueing model, with time-varying arrival rates; see section 6.

However, the statistical accuracy of HOL_r is obtained at the expense of ease of implementation. In addition to the HOL delay, w , HOL_r depends on the arrival-rate function, $\lambda(t)$, and the mean time between successive service completions (which equals $1/s\mu$ with s simultaneously busy servers and i.i.d. exponential service times with rate μ); see (7). In practice, the implementation of HOL_r requires knowledge of those system parameters, which may require estimation from data. Any estimation procedure inevitably produces some estimation error, which would affect the performance of HOL_r .

In this section, we propose estimation procedures for the arrival rate and the mean time between successive service completions in real-life service systems. Further, we quantify the estimation error resulting from those procedures, and its impact on the performance of HOL_r ; see Table 3. We show that the HOL_r estimator remains effective even with imperfect information about system parameters.

To estimate the arrival-rate function, $\lambda(t)$, we propose relying on forecasts relying on data from previous days, and observations over the current day, up to date. For $\theta_{HOL_r}(t, w)$ in (7), we need estimates of the arrival-rate function over the interval $[t - w, t]$. Here, we assume that the arrival process is a non-homogeneous Poisson process with an integrable arrival-rate function. As we observe customer arrival times, but not the arrival rates, we need to forecast future rates based on historical call volumes. For ways of forecasting future arrival rates, we refer the reader to recent work on forecasting arrival rates to service systems such as call centers. For one example, Shen and Huang (2008) propose an approach to forecast the time series of an inhomogeneous Poisson process by first building a factor model for the arrival rates, and then forecasting the time series of factor scores. As another example, Aldor-Noiman (2006) propose an arrival count model, which is based on a mixed Poisson process approach incorporating day-of-week, periodic, and exogenous effects. For other related work, see Avramidis et al. (2004), Brown et al. (2005), and references therein.

We might also rely on historical data from previous days to estimate the mean time between successive service completions, combined with real-time data

Table 3 Performance of $HOL_r(x)$ Delay Estimators, as a Function of the Traffic Intensity, ρ , and Alternative x , in the $M_t/M/100$ Queueing Model with $\alpha = 0.5$ and $E[S] = 5$ minutes

$M_t/M/100, \alpha = 0.5, E[S] = 5$ minutes									
ρ	$x = 0.1$		$x = 0.05$		$x = 0.02$		$HOL_r(x)$		HOL
	Sample size	Estimation interval							
0.9	385	(20 minutes)	1537	(77 minutes)	9604	(480 minutes)	9604	385	
0.93	385	(20 minutes)	1537	(77 minutes)	9604	(480 minutes)	9604	1537	
0.95	385	(20 minutes)	1537	(77 minutes)	9604	(480 minutes)	9604	1537	
0.97	385	(20 minutes)	1537	(77 minutes)	9604	(480 minutes)	9604	1537	
0.98	385	(20 minutes)	1537	(77 minutes)	9604	(480 minutes)	9604	1537	
	385	(20 minutes)	1537	(77 minutes)	9604	(480 minutes)	9604	385	

Sample sizes needed and length of estimation intervals required are also included. Estimates of the average squared errors are given in units of mean service time squared per customer. HOL_r , head of the line; QL, queue length.

over the recent past. However, we consider a procedure based on real-time estimation alone, and investigate its feasibility. As a real-time estimator, we propose computing the sample average, \hat{m} , of (recent) time intervals between successive service completions in the system. In doing so, as an approximation, we assume (i) that all servers are simultaneously busy and (ii) that the times between successive service completions are i.i.d. (As we are interested in systems that are heavily loaded, the assumption of busy servers is not too restrictive. The second assumption is exact for exponential service times, but not more generally.) Given that assumption, we can apply elementary statistics to compute the sample size, $n(x)$, needed to obtain a desired margin of relative error, x , at a given confidence level. (Specifically, the half width of a confidence interval is a function of the number of observations used. Therefore, we can obtain a desired margin of relative error by changing the number of observations, thus leading to a different half width.) The error, x , measures the relative error between the actual mean and the sample mean.

To illustrate, consider the $M_t/M/100$ model with exponential service times. Then, $n(0.05) \approx 1540$ at the 95% confidence level. That is, the sample size required to obtain a relative error margin of $x = 0.05$ is roughly equal to 1540, at the 95% confidence level. It is important to get a sense of how long it would take to get a total of 1540 service completions in the system. For example, suppose that the mean service time is equal to 5 minutes. The length of the estimation interval is roughly equal to 77 minutes. Indeed, each service request requires, on average, 5 minutes to process, and there are 100 servers working in parallel. This numerical example illustrates that the computational burden of obtaining estimates of system parameters that are within a relative error margin of $x = 0.05$ of their actual values is not unreasonable.

There remains to study the effect of the estimation error, x , on the performance of the HOL_r estimator. To that end, we consider modified HOL_r delay estimators, denoted by $HOL_r(x)$, depending on the relative error, x , in estimating $1/s\mu$. That is, the $HOL_r(x)$ estimators use the following delay estimate:

$$\theta_{HOL_r}(t, x, w) = \frac{1+x}{s\mu} \left(2 + \int_{t-w}^t \lambda(u)du \right),$$

where $-1 < x < 1$, and $(1+x)/s\mu$ is our estimate of the mean time between successive service completions, including a relative error x . We study the performance of $HOL_r(x)$ for alternative small values of x . Clearly, the performance of $HOL_r(x)$ should degrade as $|x|$ increases, but we would like to know by how much.

In Table 3, we study the performance of $HOL_r(x)$ as a function of the traffic intensity, ρ , in the $M_t/M/100$ queueing model, with $\alpha = 0.5$ and $E[S] = 5$ minutes. We

also include the sample sizes needed to obtain system parameter estimates within that error margin and, in parentheses, the corresponding required length of the estimation interval (under our model assumptions). We consider values of x between -0.1 and 0.1 . For these values, we find that HOL_r still performs considerably better than HOL . For example, for $x = 0.05$, the ratio $ASE(HOL)/ASE(HOL_r(x))$ ranges from about 14 to about 23 for values of ρ between 0.9 and 0.98. For $x = -0.05$, $ASE(HOL)/ASE(HOL_r(x))$ ranges from about 16 to about 27 for ρ between 0.9 and 0.98. That is, simulation shows that HOL_r remains remarkably more effective than HOL , even with imperfect information about system parameters, as would commonly occur in practice.

Additional simulation results are presented in the online supplement to the main paper. There, we consider lognormal and deterministic service times, and alternative arrival-rate parameters. We find that $HOL_r(x)$ usually performs better than HOL when the relative error, x , is at most 5%. For example, in the $M_t/H_2/100$ model with $\alpha = 0.5$, $E[S] = 6$ hours, and $x = -0.05$, the ratio $ASE(HOL)/ASE(HOL_r(x))$ ranges from 2.4 to 2.8.

8. Delay Estimators for the $M_t/GI/s+GI$ Model

In this section, we propose a new delay estimator for the $M_t/GI/s+GI$ model, based on the HOL delay observed upon arrival to the system. In section 9 we show that this new estimator, QL_{hr} , performs remarkably well. In particular, QL_{hr} effectively copes with both time-varying arrivals and non-exponential abandonment-time distributions. As a frame of reference, we also consider a delay estimator based on the QL seen upon arrival to the system. This estimator, QL_{mr} , was previously considered in Whitt (1999a) and Ibrahim and Whitt (2009b).

8.1. Actual and Potential Waiting Times

As in Garnett et al. (2002), we need to distinguish between the *actual* and *potential* waiting times of a given delayed customer in a queueing model with customer abandonment. A customer's actual waiting time is the amount of time that this customer spends in queue, until he either abandons or joins service, whichever comes first. A customer's potential waiting time is the delay he would experience, if he had infinite patience (his patience is quantified by his abandon time). For example, the potential waiting time of a delayed customer who finds n other customers waiting ahead in queue upon arrival is the amount of time needed to have $n+1$ consecutive departures from the system. (Departures from the system are either service completions or abandonments from the queue.) Our delay estimators, described next, estimate the potential waiting times of delayed customers.

8.2. The Approximation-Based QL -Based Delay Estimator (QL_{ap})

In Ibrahim and Whitt (2009b), we introduced an approximation-based QL -based delay estimator, QL_{ap} , which exploits established approximations for performance measures in the $M/GI/s+GI$ model, developed by Whitt (2005). We showed that QL_{ap} consistently outperforms all other estimators considered in the $GI+GI+s+GI$ model, with a stationary arrival process. Here, we propose an analog of QL_{ap} that uses the observed HOL delay, and effectively copes with time-varying arrival rates. We begin by briefly reviewing the QL_{ap} estimator for the $GI/GI/s+GI$ model; a more complete description can be found in section 3.5 of Ibrahim and Whitt (2009b) and Whitt (2005).

The QL_{ap} estimator approximates the $GI+GI+s+GI$ model by the corresponding $GI/M/s+M(n)$ model, with state-dependent Markovian abandonment rates. In particular, we assume that a customer who is j th from the *end* of the queue has an exponential abandonment time with rate ψ_j , where ψ_j is given by

$$\psi_j \equiv h(j/\lambda), \quad 1 \leq j \leq k, \quad (20)$$

where k is the current QL , λ is the arrival rate (assumed constant), and h is the abandonment-time hazard-rate function, defined as $h(t) \equiv f(t)/(1-F(t))$, $t \geq 0$, where f is the corresponding density function (assumed to exist). Here is how (20) is derived: If we knew that a given customer had been waiting for time t , then the rate of abandonment for that customer, at that time, would be $h(t)$. We therefore need to estimate the elapsed waiting time of that customer, given the available state information. Assuming that abandonments are relatively rare compared with service completions, it is reasonable to act as if there have been j arrival events because our customer arrived. As a simple rough estimate for the time between successive arrival events is the reciprocal of the arrival rate, $1/\lambda$, the elapsed waiting time is approximated by j/λ and the corresponding abandonment rate by (20).

For the $GI/M/s+M(n)$ model, we need to make further approximations in order to describe the potential waiting time of a customer who finds n other customers waiting in line, upon arrival. Let $W_Q(n)$ represent a random variable with the conditional distribution of the potential delay of an arriving customer, given that this customer must wait before starting service, and given that the QL seen upon arrival, is equal to n . We have the approximate representation:

$$W_Q(n) \approx \sum_{i=0}^n X_i, \quad (21)$$

where X_{n-i} is the time between the i th and $(i+1)$ th departure events. As the distribution of the X_i 's is complicated, we assume that successive departure

events are either service completions, or abandonments from the HOL. We also assume that an estimate of the time between successive departures is $1/\lambda$. Under our first assumption, after each departure, all customers remain in line except the customer at the HOL. The elapsed waiting time of customers remaining in line increases, under our second assumption, by $1/\lambda$. Let X_{n-l} , which is the time between the l th and $(l+1)$ th departure events, have an exponential distribution with rate $s\mu + \delta_n - \delta_l$, where $\delta_k = \sum_{j=1}^k \psi_j = \sum_{j=1}^k h(j/\lambda)$, $k \geq 1$, and $\delta_0 \equiv 0$. That is the case because X_{n-l} is the minimum of s exponential random variables with rate μ (corresponding to the remaining service times of customers in service), and $n-l$ exponential random variables with rates ψ_i , $l+1 \leq i \leq n$ (corresponding to the abandonment times of the customers waiting in line).

The QL_{ap} delay estimate given to a customer who finds n customers in queue upon arrival is

$$\theta_{QL_{ap}}(n) = \sum_{i=0}^n \frac{1}{s\mu + \delta_n - \delta_{n-i}}, \quad (22)$$

that is, $\theta_{QL_{ap}}(n)$ approximates the mean of the potential waiting time, $E[W_Q(n)]$.

8.3. The QL_h Estimator

We are now ready to propose a new delay estimator for the $M_t/GI/s+GI$ model, which we refer to as QL_h . This estimator requires knowledge of the abandonment-time hazard-rate function, h . That is convenient from a practical point of view, because it is relatively easy to estimate hazard rates from system data; see Brown et al. (2005).

We proceed in two steps: (i) we use the observed HOL delay, w , to estimate the QL seen upon arrival and (ii) we use this QL estimate to implement a new delay estimator, paralleling (22). Unlike QL_{ap} , QL_h exploits the HOL delay, and does not assume knowledge of the QL seen upon arrival.

For step (i), let $N_w(t)$ be the number of arrivals in the interval $[t-w, t]$ who do not abandon. That is, $N_w(t)+1$ is the number of customers seen in the queue upon arrival at time t , given that the observed HOL delay at t is equal to w . It is significant that N_w has the structure of the number in system in an $M_t/GI/\infty$ infinite-server system, starting out empty in the infinite past, with arrival rate $\lambda(u)$ identical to the original arrival rate in $[t-w, t]$ (and equal to 0 otherwise). The individual service-time distribution is identical to the abandonment-time distribution in our original system. Thus, $N_w(t)$ has a Poisson distribution with mean

$$m(t, w) \equiv E[N_w(t)] = \int_{t-w}^t \lambda(s)(1 - F(t-s))ds, \quad (23)$$

where F is the abandonment-time cdf.

For step (ii), we use $m(t, w)+1$ as an estimate of the QL seen upon arrival, at time t . In (20), we replace λ by $\hat{\lambda}$, where $\hat{\lambda}$ is defined as the average arrival rate over the interval $[t-w, t]$, i.e., $\hat{\lambda} \equiv (1/w) \int_{t-w}^t \lambda(s)ds$. We do so because we now have a non-stationary arrival process instead of a stationary arrival process. Paralleling (22), the QL_h delay estimate given to a customer such that the observed HOL delay, at his time of arrival, t , is equal to w , is given by

$$\theta_{QL_h}(t, w) \equiv \sum_{i=0}^{m(t, w)+1} \frac{1}{s\mu + \hat{\delta}_n - \hat{\delta}_{n-i}} \quad (24)$$

for $m(t, w)$ in (23), $\hat{\delta}_k = \sum_{j=1}^k h(j/\hat{\lambda})$, and $\hat{\delta}_0 = 0$. If we actually know the QL, then we can replace $m(t, w)$ by $Q(t)$, i.e., we can use QL_{ap} . There remains to investigate ways of estimating the abandonment-time distribution needed to implement QL_h . We envision that such estimates will be based on long-term estimates of customer time-to-abandon distribution, instead of real-time information about customer abandonment times. Providing additional details relating to this estimation is outside the scope of this paper, and is left for future research.

9. Simulation Results for the $M_t/M/s+GI$ Model

In this section, we present simulation results for the $M_t/M/s+GI$ model with sinusoidal arrival rates. For the abandonment-time distribution, we considered M (exponential), E_{10} (Erlang, sum of 10 exponentials) and H_2 (hyperexponential with SCV equal to four and balanced means), but here we only discuss the first two cases; see Ibrahim and Whitt (2009c) for a discussion of the H_2 case. We consider the QL_{mr} , QL_h , and HOL delay estimators. In this section, we show plots of the simulation results. Corresponding tables with estimates of 95% confidence intervals, in addition to more simulation results, appear in Ibrahim and Whitt (2009c).

9.1. Description of the Experiments

We vary the number of servers, s , but consider only relatively large values ($s \geq 100$), because we are interested in large service systems. We let the service rate, μ , be equal to 1. For the arrival-rate function, $\lambda(u)$ in (12), we fix the relative frequency, $\gamma = 1.571$. This value of γ corresponds to a mean service time $E[S] = 6$ hours, for daily arrival-rate cycles; see Table 1.

We consider a relative amplitude $\alpha = 0.5$, and an average arrival rate $\bar{\lambda} = 140$. The instantaneous offered load in the system, at time t , is given by $\lambda(t)/s\mu$. With $\alpha = 0.5$, the offered load varies between 0.7 and 2.1. Because of customer abandonment, the congestion is not extraordinarily high when the system is significantly overloaded. We let the abandonment rate, $\nu = 1$, because that seems to be a

representative value. Simulation results for all models are based on 10 independent replications of length 1 month each, assuming a daily cycle.

9.2. Results for the $M_t/M/s+M$ Model

Consistent with theory in section 8, Figure 3 shows that QL_m is the best possible estimator, under the MSE criterion. The RRASE of QL_m ranges from about 14% for $s = 100$ to about 4% when $s = 1000$. Figure 3 shows that $s \times ASE(QL_m)$, the ASE of QL_m multiplied by the number of servers s , is nearly constant for all values of s considered. This shows that QL_m is asymptotically correct as s increases, i.e., $ASE(QL_m)$ approaches 0 as s increases.

The QL_h estimator is the second best estimator for this model. The RRASE of QL_h ranges from about 20% for $s = 100$ to about 6% for $s = 1000$. That is, QL_h is relatively accurate for this model. The difference in performance between QL_h and QL_m is not too great: $ASE(QL_h)/ASE(QL_m)$ is close to 1.6 for all s . Moreover, Figure 3 shows that QL_h is asymptotically correct: $s \times ASE(QL_h)$ is also roughly equal to a constant for all s .

The HOL estimator performs much worse than QL_m and QL_h . For example, the ratio $ASE(HOL)/ASE(QL_h)$ ranges from about 3 for $s = 100$ to about 20 for $s = 1000$. The RRASE of HOL ranges from about 33% for $s = 100$ to about 27% for $s = 1000$. That is, we do not see a considerable improvement in the performance of HOL, as s increases. That is confirmed by Figure 3, where we see that $s \times ASE(HOL)$ increases linearly, as s increases.

9.3. Results for the $M_t/M/s+E_{10}$ Model

The QL_h estimator is the most effective estimator, under the MSE criterion, for this model. The RRASE of QL_h ranges from about 11% for $s = 100$ to about 4% for

Figure 3 $E[S] = 6$ hours, $\alpha = 0.5$

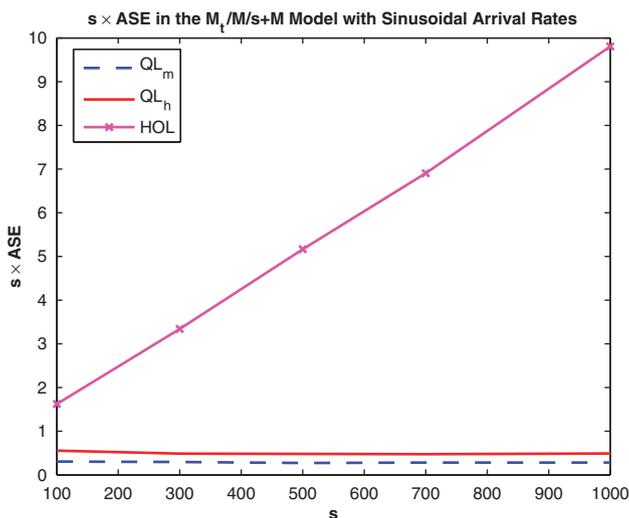
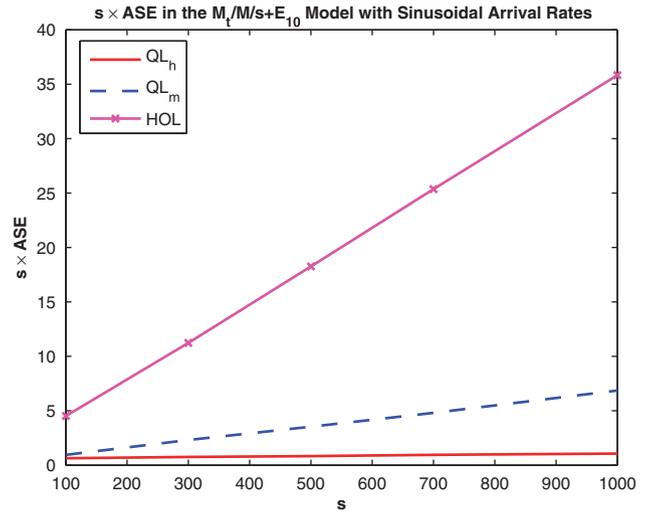


Figure 4 $E[S] = 6$ hours, $\alpha = 0.5$



$s = 1000$. That is, QL_h is relatively accurate for this model. Figure 4 shows that QL_h is asymptotically correct: $s \times ASE(QL_h)$ is roughly equal to a constant for all values of s considered.

The QL_m estimator performs significantly worse than QL_h , with E_{10} abandonment. The ratio $ASE(QL_m)/ASE(QL_h)$ ranges from about 1.5 for $s = 100$ to about 6.5 for $s = 1000$. The RRASE of QL_m ranges from about 13% for $s = 100$ to about 10% for $s = 1000$. Figure 4 shows that QL_m is not asymptotically correct as s increases.

The least effective estimator is, yet again, the HOL estimator. The RRASE of HOL ranges from about 27% for $s = 100$ to about 25% for $s = 1000$. The difference in performance between HOL and QL_h is remarkable: $ASE(HOL)/ASE(QL_h)$ ranges from roughly 7 for $s = 100$ to roughly 33 for $s = 1000$. Figure 4 shows that $s \times ASE(HOL)$ increases linearly (and steeply) as s increases.

9.4. Results for Other Models

We consider general service-time and abandonment-time distributions in Ibrahim and Whitt (2009c). For the service-time distribution, we consider M , D , and H_2 . For the abandonment-time distribution, we consider M , H_2 , and E_{10} . We consider different combinations of service-time and abandonment-time distributions. These additional simulation results are consistent with those reported above: The QL_m estimator remains effective with M abandonment, even when the service-time distribution is not nearly exponential. With H_2 and E_{10} abandonment, QL_h outperforms QL_m , especially when the number of servers is large. The HOL estimator remains the least effective estimator, under the MSE criterion, in all models considered.

10. Conclusions

In this paper, we studied the performance of alternative delay estimators in the $M_t/GI/s$ and $M_t/GI/s+GI$ queueing models, which have a non-homogeneous Poisson process. We concentrated on the HOL estimator, which is equal to the elapsed delay of the customer at the HOL, at the time of arrival. We did so with the understanding, based on our previous work, that results for HOL should apply equally well to the delay of the LES. A main conclusion is that the performance of these delay-history-based delay estimators can degrade in face of time-varying arrivals, which often occurs in practice; that is dramatically shown in Figure 2.

As a consequence, we developed refinements of HOL, in particular, HOL_r in (7) for $M_t/GI/s$ and QL_h in (24) for $M_t/GI/s+GI$. Simulation experiments in sections 6 and 9 showed that these estimators effectively cope with both time-varying arrivals and non-exponential service-time and abandon-time distributions. We also established analytical results supporting HOL_r in section 5. We quantified the difference in performance between QL and HOL_r and found that the ratio of their respective MSE's is roughly equal to 2, especially for high values of the traffic intensity, ρ ; see (13).

However, the new refined estimators lose some of their appeal compared with the simple HOL and LES estimators, because they require information about the model, in particular, the arrival-rate function and the mean time between successive departures. Hence, in section 7 we proposed ways to estimate the required information. Even if we rely on real-time estimation of the mean time between successive departures, we showed that we can obtain suitably accurate estimates without requiring that the observation interval be too long. Table 3 shows that the HOL_r estimator remains effective even if the information is known imperfectly.

Our general strategy for creating the refined HOL estimators has been to approximate the mean conditional delay, given the observed HOL delay by (i) approximating the QL, given the observed HOL delay, and (ii) approximating the expected delay given the QL. As a consequence, direct QL-based delay estimators should be preferred if the QL is known. However, in section 1.2 we observed that there are complex service systems such as web chat and ticket queues for which the QL is not known.

Acknowledgments

This research was supported by NSF Grants DMI-0457095 and CMMI 0948190.

References

Aksin, O. Z., M. Armony, V. Mehrotra. 2007. The modern call-center: A multi-disciplinary perspective on operations management research. *Prod. Oper. Manag.* **16**(6): 665–688.

- Aldor-Noiman, S. 2006. Forecasting demand for a telephone call center: Analysis of desired versus attainable precision. Unpublished masters thesis, Technion-Israel Institute of Technology, Haifa, Israel.
- Allon, G., A. Bassambo, I. Gurvich. 2009. We will be right with you: Managing customer with vague promises. Working paper, Northwestern University, Evanston, IL.
- Armony, M., C. Maglaras. 2004. Contact centers with a call-back option and real-time delay information. *Oper. Res.* **52**: 527–545.
- Armony, M., N. Shimkin, W. Whitt. 2009. The impact of delay announcements in many-server queues with abandonments. *Oper. Res.* **57**: 66–81.
- Avramidis, A. N., A. Deslauriers, P. L'Ecuyer. 2004. Modeling daily arrivals to a telephone call center. *Manage. Sci.* **50**: 896–908.
- Barlow, R. E., F. Proschan. 1975. *Statistical Theory of Reliability and Life Testing*. Holt, Rinehart and Winston, New York.
- Brown, L., N. Gans, A. Mandelbaum, A. Sakov, H. Shen, S. Zeltyn, L. Zhao. 2005. Statistical analysis of a telephone call center: A queueing-science perspective. *J. Am. Stat. Assoc.* **100**: 36–50.
- Garnett, O., A. Mandelbaum, M. I. Reiman. 2002. Designing a call center with impatient customers. *Manuf. Serv. Oper. Manage.* **5**: 79–141.
- Green, L., P. Kolesar, W. Whitt. 2007. Coping with time-varying demand when setting staffing requirements for a service system. *Prod. Oper. Manag.* **16**: 13–39.
- Guo, P., P. Zipkin. 2007. Analysis and comparison of queues with different levels of delay information. *Manage. Sci.* **53**: 962–970.
- Heyman, D., W. Whitt. 1984. The asymptotic behavior of queues with time-varying arrival rates. *J. Appl. Probab.* **21**: 143–156.
- Ibrahim, R., W. Whitt. 2009a. Real-time delay estimation based on delay history. *Manuf. Serv. Oper. Manage.* **11**: 397–415.
- Ibrahim, R., W. Whitt. 2009b. Real-time delay estimation in overloaded multiserver queues with abandonments. *Manage. Sci.* **55**: 1729–1742.
- Ibrahim, R., W. Whitt. 2009c. Supplement to "Real-time delay estimation based on delay history in many-server service systems with time-varying arrivals," IEOR Department, Columbia University, New York, NY. Available at <http://columbia.edu/~rei2101>
- Jongbloed, G., G. Koole. 2001. Managing uncertainty in call centers using Poisson mixtures. *Appl. Stoch. Models Bus. Ind.* **17**: 307–318.
- Jouini, O., Y. Dallery, Z. Aksin. 2007. Modeling call centers with delay information. Working paper, Koc University, Turkey.
- Mandelbaum, A., A. Sakov, S. Zeltyn. 2000. Empirical analysis of a call center. Technical report, Faculty of Industrial Engineering and Management, Technion, Israel.
- Massey, W., W. Whitt. 1994. A stochastic model to capture space and time dynamics in wireless communication systems. *Probab. Eng. Inf. Sci.* **8**: 541–569.
- Nakibly, E. 2002. Predicting waiting times in telephone service systems. MS thesis, the Technion, Haifa, Israel.
- Ross, S. 1996. *Stochastic Processes*. 2nd edn. Wiley, New York.
- Shen, H., J. Huang. 2008. Interday forecasting and intraday updating of call center arrivals. *Manuf. Serv. Oper. Manage.* **10**: 601–623.
- Whitt, W. 1999a. Predicting queueing delays. *Manage. Sci.* **45**: 870–888.
- Whitt, W. 1999b. Improving service by informing customers about anticipated delays. *Manage. Sci.* **45**: 192–207.
- Whitt, W. 2005. Engineering solution of a basic call-center model. *Manage. Sci.* **51**: 221–235.
- Xu, S. H., L. Gao, J. Ou. 2007. Service performance analysis and improvement for a ticket queue with balking customers. *Manage. Sci.* **53**: 971–990.