

REAL-TIME QUEUES IN HEAVY TRAFFIC WITH EARLIEST-DEADLINE-FIRST QUEUE DISCIPLINE¹

BY BOGDAN DOYTCHINOV, JOHN LEHOCZKY AND STEVEN SHREVE

*Worcester Polytechnic Institute, Carnegie Mellon University
and Carnegie Mellon University*

This paper introduces a new aspect of queueing theory, the study of systems that service customers with specific timing requirements (e.g., due dates or deadlines). Unlike standard queueing theory in which common performance measures are customer delay, queue length and server utilization, real-time queueing theory focuses on the ability of a queue discipline to meet customer timing requirements, for example, the fraction of customers who meet their requirements and the distribution of customer lateness. It also focuses on queue control policies to reduce or minimize lateness, although these control aspects are not explicitly addressed in this paper.

To study these measures, we must keep track of the lead times (deadline minus current time) of each customer; hence, the system state is of unbounded dimension. A heavy traffic analysis is presented for the earliest-deadline-first scheduling policy. This analysis decomposes the behavior of the real-time queue into two parts: the number in the system (which converges weakly to a reflected Brownian motion with drift) and the set of lead times given the queue length. The lead-time profile has a limit that is a nonrandom function of the limit of the scaled queue length process. Hence, in heavy traffic, the system can be characterized as a diffusion evolving on a one-dimensional manifold of lead-time profiles. Simulation results are presented that indicate that this characterization is surprisingly accurate. A discussion of open research questions is also presented.

1. Introduction. This paper introduces a new aspect of queueing theory, the study of systems that service customers with individual timing requirements (e.g., due dates or deadlines). Such systems arise naturally in manufacturing in which orders have due dates. A second category of examples arises in real-time computer and communications systems. Such systems might involve the transmission of digitized voice, video or images over a network. These transmissions must reach their destination within specific deadlines to maintain the integrity of the communication (e.g., voice conversation, teleconference or movie). Real-time computer systems also control much of modern technology, for example, engines and braking systems in automobiles, all avionic systems (including air traffic control) and all aspects of modern manufacturing facilities. Computerized control systems must receive and react to state infor-

Received August 1999; revised March 2000.

¹Supported in part by the Office of Naval Research contract N00014-92-J-1524, by DARPA under contracts F30602-96-1-0160 and N66001-97-C-8527, and by the Army Research Office under contract DAAH04-95-1-0226.

AMS 2000 *subject classifications*. Primary 60K25; secondary 60G57, 60J65.

Key words and phrases. Due dates, heavy traffic, queueing, diffusion limits, random measures.

mation within a fixed, often stringent, time interval to maintain proper control over the system. Failure to meet task timing requirements in safety-critical applications can have serious consequences. Thus, the conditions for correct performance of a real-time system include both the logical correctness of each of the tasks that it executes and the timing correctness of those tasks. Over the last decade, there have been significant strides made in the development of a theory of hard real-time systems, systems in which tasks must be completed before their deadline elapses. The reader is referred to the handbook by Klein, Ralya, Pollak, Obenza and Gonzalez-Harbour [15] for a description of this theory that addresses many practical considerations encountered in computer systems such as operating system overhead, hardware architecture details, concurrency control and other sorts of blocking and task precedence relationships.

The scheduling theory described in [15] assumes an essentially deterministic environment. For example, task arrivals are modeled as the superposition of periodic arrival processes, and task service times are deterministic and given by the worst case execution of each task type. Two principal approaches have been developed to assess the design of real-time systems with periodic task arrivals: one is based on a fixed task priority structure (exemplified by *rate monotonic scheduling*) and the other is based on dynamic priorities [exemplified by the *earliest-deadline-first* (EDF) approach to scheduling]. These two scheduling algorithms were analyzed by Liu and Layland [20] and the EDF scheduling algorithm was shown to be optimal for this scheduling problem. In some systems, especially communications systems, only a small number of bits are available in each packet to represent the task's priority; thus EDF cannot be fully implemented because it can require an unlimited number of distinct priority categories. Nonetheless, we introduce real-time queueing theory in the context of the EDF scheduling algorithm, since this algorithm is optimal under some conditions. Panwar and Towsley [22] showed that EDF maximizes the fraction of customers meeting their deadlines within the class of work conserving policies that allow preemption in GI/M/1 queues where customers have general deadlines. Bounds on the performance of EDF for M/M/1 queues in which customers have exponential deadlines were developed by Hong, Tan and Towsley [10]. We also discuss the *first-in-first-out* (FIFO) queue discipline in Section 5 of this paper.

There are major limitations to any theory that requires periodic arrivals and assumes worst-case execution times. These assumptions are quite narrow and limit the range of systems that can be studied. Multimedia applications or real-time communications can exhibit substantial variability in the arrival of tasks and their work requirements. For real-time systems for which the task sets exhibit substantial variability, we would like to develop approaches based on queueing theory, a theory that was designed to model and predict stochastic system behavior with resource contention. This theory allows randomness in the task arrivals and task execution times. The difficulty with queueing theory is that it typically does not allow for explicit consideration of dynamically changing task timing requirements. Instead, it only permits priorities

that allow important tasks or tasks with initial short timing requirements to receive preferential treatment. Much of queueing theory focuses on general system performance measures, such as task delay, queue lengths and processor utilization, and these are usually computed under equilibrium assumptions. It does not model the timing requirements of each customer nor does it analyze the ability of a scheduling algorithm to meet those timing requirements. What is needed is a new theory that combines the focus on meeting task timing requirements as studied in real-time scheduling theory with the focus on stochastic task sets as studied in queueing theory. This paper represents a step in the direction of building such a theory, hence the name *real-time queueing theory*.

To study whether tasks or customers meet their timing requirements, we must keep track of the customer lead times, where the lead time is the time remaining until the deadline elapses, that is,

$$\text{lead time} = \text{deadline} - \text{current time.}$$

Customer lead times decrease linearly while a customer is in the queue. Because the lead time must be tracked for each customer, the dimension of the system state is the number of customers in the queue, which is unbounded. This causes analytic difficulties. In spite of this unbounded dimension, a heavy traffic analysis can be carried out. This analysis decomposes the behavior of the real-time queue into two parts: the number in the system, say $Q(t)$ (which is shown under the heavy-traffic scaling to converge weakly to a reflected Brownian motion with drift) and the set of lead times, $(L_1(t), \dots, L_{Q(t)}(t))$ (we refer to this as the lead-time profile), conditional on the queue length. It is convenient to think of this profile as a random counting measure on \mathbb{R} . In heavy traffic, under the earliest-deadline-first queue discipline, it will be shown that when suitably scaled, the lead-time profile converges to a nonrandom function of the limit of the scaled queue length process, the particular function being determined by the distribution of initial deadline of arriving customers. Hence, in heavy traffic, the unbounded dimension process collapses to a one-dimensional process and we can conceptualize the real-time queueing process as a diffusion evolving on a one-dimensional manifold of lead-time profiles. Simulation results, presented in Section 4, indicate that this characterization is surprisingly accurate.

This work is based on the long tradition of heavy-traffic queueing theory pioneered by Kingman [14]. This research was generalized in scope and system complexity by a number of authors; for example, see Iglehart and Whitt [12, 13], who studied the multiple server case, for a review of this early literature. The use of heavy-traffic theory in the study of the behavior of priority queues was initiated by the work of Whitt [30], Hooke [11], Harrison [5] and Kyprianou [16]. The phenomenon of state space collapse, which was originally observed by Reiman [26, 27], also occurs in our work. Specifically, the lead-time profiles have the dimension of the number of customers in the queue, which is unbounded. Nevertheless, in heavy traffic, those random profiles con-

verge to a deterministic manifold of profiles indexed by the queue length, a one-dimensional parameter.

Heavy-traffic queueing theory has evolved greatly over the last 25 years, especially for queueing networks that carry multiple customer types. A great increase in interest in this research area came with the work of Harrison and co-authors (e.g., [6–9] and Peterson [24]). The EDF queue discipline studied in this paper is related to multiclass queues, although in EDF there are an infinite number of distinct priority classes and customers change classes as they wait in the queue. Most of the work in heavy-traffic queueing networks studies the behavior of queue lengths and workloads, rather than focusing on the lead times of individual customers. We expect that our results on the convergence of the lead-time profiles for the single queue case will carry over to networks, but we do not study network behavior in this paper. Similarly, we expect that optimal control methods can be applied to real-time queues to control customers' lateness in the way that many researchers have used these methods to optimize inventory holding costs; see, for example, Harrison and Wein [9, 29].

There is some recent work on heavy-traffic approximations for systems that handle customers with due dates. Of particular importance are the papers by Van Mieghem [28], Markowitz and Wein [21], Doytchinov [3] and Lehoczký [17–19]. Van Mieghem studied a single server multiclass queueing system with k distinct customer classes. Each class has an associated convex cost of delay, $C_k(\tau)$, with derivative $c_k(\tau)$. The objective is to minimize the total delay cost incurred over a finite time horizon. The paper studied the “generalized $c\mu$ policy,” which schedules the customer having maximum value $\mu_k c_k(a_k(t))$, where μ_k is the service rate for class k and $a_k(t)$ is the age of the oldest customer in class k . Customers are served in FIFO order in each class, which is equivalent to EDF within each class. This policy is shown to be asymptotically optimal in heavy traffic. Generalizations to a countable number of customer classes and several homogeneous servers in a nonstationary, deterministic or stochastic environment are also considered.

Markowitz and Wein [21] studied the single machine scheduling problem in a manufacturing context using heavy-traffic methods. They give a unified treatment that permits setup costs, customer due dates and a mixture of standardized and customized products. The analysis assumes a cyclic policy in which different products must be produced in a fixed sequence, but the machine busy/idle policy and lot-sizing decisions are dynamic. As such, the system resembles a polling system. A heavy-traffic averaging principle such as characterized by Coffman, Puhalskii and Reiman [2] is assumed to hold and, subject to this assumption, the optimal policy is determined. Their paper gave a detailed discussion of the interactions between the setup, due date and product mix factors.

Doytchinov [3] developed a partial differential equation-based approach to the study of real-time M/M/1 queues in which the arrivals have constant deadlines. In this case, the EDF and FIFO queue disciplines are identical. His

methodology proved that the lead-time profiles converge to a uniform distribution in heavy traffic.

Lehoczky [17] gave an informal analysis for the M/M/1 queue based on representing the lead-time profile as a measure-valued Markov process and then arguing, under heavy traffic conditions, that the generator converges to that of a deterministic profile conditional on the queue length. This was done both for EDF and for processor sharing. Lehoczky [18] used these results to study the behavior of various queue control policies to reduce customer lateness. Lehoczky [19] extended the analysis to Jackson networks.

This paper is organized as follows. In Section 2, we present the basic model, assumptions and notation. Section 3 gives the major theorems that describe the heavy-traffic limiting behavior of EDF real-time single server queueing systems. Section 4 presents simulation results that illustrate the accuracy of the theory. Section 5 presents some conjectures for the extension of the theory of Section 3. Appendixes A and B are included to set notation. Appendix A collects key definitions and theorems related to weak convergence of measures on metric spaces, and Appendix B recalls classical heavy-traffic theorems.

2. The basic model, assumptions and notation. We first define the basic real-time queueing theory model. Because we shall ultimately pass to a heavy-traffic limit, we posit a sequence of queueing systems, indexed by n . The assumptions on the n th queueing system are the following:

- A1. There is a single station serving customers.
- A2. Customer interarrival times are determined by the sequence of strictly positive i.i.d. random variables $\{u_j^{(n)}\}_{j=1}^\infty$ with $E[u_j^{(n)}] = 1/\lambda^{(n)}$ and $\text{Var}[u_j^{(n)}] = (\alpha^{(n)})^2$.
- A3. Customers have service requirements that are determined by the sequence of nonnegative i.i.d. random variables $\{v_j^{(n)}\}_{j=1}^\infty$ with $E[v_j^{(n)}] = 1/\mu^{(n)}$ and $\text{Var}[v_j^{(n)}] = (\beta^{(n)})^2$.
- A4. Each customer arrives with a hard deadline (initial lead time) $L_j^{(n)}$. These initial lead times are i.i.d. with distribution given by

$$(2.1) \quad \mathbb{P}(L_j^{(n)} \leq \sqrt{n}y) = G(y),$$

where G is a right-continuous cumulative distribution function. We define

$$(2.2) \quad y^* \triangleq \min\{y \in \mathbb{R}; G(y) = 1\}$$

and assume that y^* is finite.

- A5. The sequences $\{u_j^{(n)}\}_{j=1}^\infty$, $\{v_j^{(n)}\}_{j=1}^\infty$ and $\{L_j^{(n)}\}_{j=1}^\infty$ are mutually independent.
- A6. Customers are served using the earliest-deadline-first queue discipline, that is, the server always services the customer with the shortest lead time.
- A7. Preemption is permitted (we assume preempt–resume). There is no setup, switchover or any other type of overhead.

A8. Late customers (customers with negative lead times) stay in the queue until served to completion.

A9. The queue is empty at time zero.

The interarrival times, service times, initial lead times and queue discipline completely determine the behavior of the queue. From them we can derive the *customer arrival times*

$$S_k^{(n)} \triangleq \sum_{j=1}^k u_j^{(n)},$$

with $S_0^{(n)} \triangleq 0$, and the *customer arrival process*

$$A^{(n)}(t) \triangleq \max\{k | S_k^{(n)} \leq t\}.$$

The *work arrival process*

$$V^{(n)}(t) \triangleq \sum_{j=1}^{\lfloor t \rfloor} v_j^{(n)}$$

records the amount of work that arrives with the first $\lfloor t \rfloor$ customers, and the *netput process*

$$N^{(n)}(t) \triangleq V^{(n)}(A^{(n)}(t)) - t$$

measures the work remaining in queue at time t , provided that the server is never idle up to time t . The *cumulative idleness process*

$$I^{(n)}(t) \triangleq - \min_{0 \leq s \leq t} N^{(n)}(s),$$

gives the amount of time the server is idle, and adding this to the netput process, we obtain the *workload process*

$$W^{(n)}(t) \triangleq N^{(n)}(t) + I^{(n)}(t),$$

which records the amount of work in the queue, taking server idleness into account. All the foregoing above processes are independent of the queue service discipline, provided that the server is never idle when there are customers in the queue. The *queue length process* $Q^{(n)}(t)$, which is the number of customers in the queue at time t , depends on the queue discipline. All these processes are *right-continuous with left-hand limits (RCLL)*.

To obtain heavy-traffic limits, we must scale and sometimes center the preceding processes. The real-valued processes whose limits we shall consider are

$$(2.3) \quad \widehat{A}^{(n)}(t) \triangleq \frac{1}{\sqrt{n}}(A^{(n)}(nt) - \lambda^{(n)}nt),$$

$$(2.4) \quad \widehat{V}^{(n)}(t) \triangleq \frac{1}{\sqrt{n}} \sum_{j=1}^{\lfloor nt \rfloor} \left(v_j^{(n)} - \frac{1}{\mu^{(n)}} \right),$$

$$(2.5) \quad \widehat{N}^{(n)}(t) \triangleq \frac{1}{\sqrt{n}}(V^{(n)}(A^{(n)}(nt)) - nt),$$

$$(2.6) \quad \widehat{I}^{(n)}(t) \triangleq \frac{1}{\sqrt{n}} I^{(n)}(nt),$$

$$(2.7) \quad \widehat{W}^{(n)}(t) \triangleq \frac{1}{\sqrt{n}} W^{(n)}(nt) = \widehat{N}^{(n)}(t) + \widehat{I}^{(n)}(t),$$

$$(2.8) \quad \widehat{Q}^{(n)}(t) \triangleq \frac{1}{\sqrt{n}} Q^{(n)}(nt).$$

Heavy-traffic assumptions. Define the traffic intensity $\rho^{(n)} \triangleq \lambda^{(n)}/\mu^{(n)}$. The following assumptions shall be in force throughout:

$$(2.9) \quad \lim_{n \rightarrow \infty} \sqrt{n}(1 - \rho^{(n)}) = \gamma > 0,$$

$$(2.10) \quad \lim_{n \rightarrow \infty} \lambda^{(n)} = \lambda > 0, \quad \lim_{n \rightarrow \infty} \mu^{(n)} = \lambda,$$

$$(2.11) \quad \lim_{n \rightarrow \infty} \alpha^{(n)} = \alpha, \quad \lim_{n \rightarrow \infty} \beta^{(n)} = \beta.$$

We also impose the usual Lindeberg condition on the interarrival and service times:

$$(2.12) \quad \begin{aligned} & \lim_{n \rightarrow \infty} \mathbb{E} \left[\left(u_j^{(n)} - \frac{1}{\lambda^{(n)}} \right)^2 \mathbb{1}_{\{|u_j^{(n)} - 1/\lambda^{(n)}| > c\sqrt{n}\}} \right] \\ & = \lim_{n \rightarrow \infty} \mathbb{E} \left[\left(v_j^{(n)} - \frac{1}{\mu^{(n)}} \right)^2 \mathbb{1}_{\{|v_j^{(n)} - 1/\mu^{(n)}| > c\sqrt{n}\}} \right] = 0 \quad \forall c > 0. \end{aligned}$$

It is a standard result (see Corollary B.4) that the triple $(\widehat{N}^{(n)}, \widehat{I}^{(n)}, \widehat{W}^{(n)})$ converges weakly to (N^*, I^*, W^*) , where N^* is a Brownian motion with drift and

$$I^*(t) \triangleq - \min_{0 \leq s \leq t} N^*(s),$$

$$(2.13) \quad W^*(t) \triangleq N^*(t) + I^*(t).$$

The process W^* is a reflected Brownian motion with drift, and I^* causes the reflection. Furthermore, the scaled queue length process $\widehat{Q}^{(n)}$ converges weakly to λW^* (see Corollary 3.2).

EDF-related processes. With the EDF queue discipline, customers are served in order of increasing lead times. Any two customers in the queue will maintain their relative order until they depart; however, arriving customers may preempt and move directly into service if they have a sufficiently short initial lead time. To study the behavior of the EDF queue discipline, it is useful to keep track of the lead time of the customer currently in service and the largest lead time of all customers who have ever been in service. We define the *frontier*

$$F^{(n)}(t) \triangleq \left\{ \begin{array}{l} \text{largest lead time of any customer who has ever been in} \\ \text{service, whether still present or not, or } \sqrt{n} y^* - t, \text{ if this} \\ \text{quantity is larger than the former one} \end{array} \right\}.$$

We adopt the usual convention that until the arrival of the first customer, the largest lead time of any customer who has ever been in service is $-\infty$, and hence $F^{(n)}(t) = \sqrt{n} y^* - t$ for $0 \leq t < S_1^{(n)}$. We define the *current lead time*

$$C^{(n)}(t) \triangleq \left\{ \begin{array}{l} \text{lead time of the customer in service,} \\ \text{or } F^{(n)}(t) \text{ if the queue is empty} \end{array} \right\}.$$

Under the EDF queue discipline, there is no customer in the queue with lead time smaller than $C^{(n)}(t)$ and there has never been a customer in service whose lead time, if the customer were still present, would exceed $F^{(n)}(t)$. Furthermore, $C^{(n)}(t) \leq F^{(n)}(t)$ for all $t \geq 0$. Both $F^{(n)}$ and $C^{(n)}$ are RCLL processes.

At time t all customers in the system have lead times equal to or greater than $C^{(n)}(t)$; if the queue is nonempty, $C^{(n)}(t)$ is the left support point of the random counting measure, which puts a unit point mass at the lead time of each customer in the queue at time t . In spite of this, the frontier is more important than $C^{(n)}(t)$ in the analysis of the EDF queue discipline for two reasons. The first reason is that in heavy traffic the scaled number of customers with lead times between $C^{(n)}(t)$ and $F^{(n)}(t)$ is negligible. Customers at time t with lead times in the interval $[C^{(n)}(t), F^{(n)}(t))$, if any, are part of a special type of busy period. For a nonempty queue, this busy period was initiated by a customer arrival that preempted a customer, ℓ , in service [the preempted customer with current lead time $F^{(n)}(t)$]. This busy period was possibly sustained by other arrivals, each of which had, at the time of its arrival, a lead time shorter than ℓ 's lead time. Because ℓ 's lead time decreases linearly with time, the traffic intensity associated with the customers that sustain this special busy period decreases with time. As shown subsequently in Proposition 3.6, under the heavy-traffic assumptions the scaled number of customers having lead times at time t taking values in $[C^{(n)}(t), F^{(n)}(t))$ converges to 0 as $n \rightarrow \infty$. It follows that in heavy traffic the occupancy of the queue consists essentially of customers with lead times in $[F^{(n)}(t), \infty)$.

The second feature of the frontier is that customers with lead times in $(F^{(n)}(t), \infty)$ at time t have never received any service; their lead-time profile is determined entirely by the arrival process. Because only the arrival process is involved, this profile can be determined and in heavy traffic converges to a nonrandom function of the limit of the scaled queue length.

Although this paper focuses on a heavy-traffic analysis of a single server queue using the EDF queue discipline, it is worth noting that some of the results can be expected to carry over in non-heavy-traffic conditions. For example, if the traffic intensity were not near 1, but the queue length happened to be relatively long, then most of the customers in the system would have lead times that take values in $(F^{(n)}(t), \infty)$, and their profile would be determined solely by the arrival process, not the service process. Consequently, the lead-time profiles would be the same as those predicted for that queue length under heavy-traffic conditions.

Measure-valued processes. At any instant of time, the system consists of a set of customers, each of which has a specific lead time and a remaining work requirement. We wish to characterize the instantaneous lead-time profile of the customers. It is convenient to think of this profile as a counting measure on \mathbb{R} . In this section, we define a collection of RCLL measure-valued processes that will be useful in the analysis.

Queue length measure:

$$\mathcal{Q}^{(n)}(t)(B) \triangleq \left\{ \begin{array}{l} \text{number of customers in the queue at time } t \\ \text{having lead times at time } t \text{ in } B \subset \mathbb{R} \end{array} \right\}.$$

Workload measure:

$$\mathcal{W}^{(n)}(t)(B) \triangleq \left\{ \begin{array}{l} \text{work at time } t \text{ associated with customers in} \\ \text{the queue having lead times at time } t \text{ in } B \subset \mathbb{R} \end{array} \right\}.$$

Customer arrival measure:

$$\mathcal{A}^{(n)}(t)(B) \triangleq \left\{ \begin{array}{l} \text{number of all arrivals by time } t \text{ having} \\ \text{lead times at time } t \text{ in } B \subset \mathbb{R} \end{array} \right\}.$$

Workload arrival measure:

$$\mathcal{V}^{(n)}(t)(B) \triangleq \left\{ \begin{array}{l} \text{work associated with all arrivals by time} \\ t \text{ having lead times at time } t \text{ in } B \subset \mathbb{R} \end{array} \right\}.$$

The following relationships easily follow:

$$\begin{aligned} \mathbf{Q}^{(n)}(t) &= \mathcal{Q}^{(n)}(t)(\mathbb{R}), \\ \mathbf{A}^{(n)}(t) &= \mathcal{A}^{(n)}(t)(\mathbb{R}), \\ \mathbf{V}^{(n)}(\mathbf{A}^{(n)}(t)) &= \mathcal{V}^{(n)}(t)(\mathbb{R}), \\ \mathcal{A}^{(n)}(t)(B) &= \sum_{j=1}^{A^{(n)}(t)} \mathbb{1}_{\{L_j^{(n)} - (t - S_j^{(n)}) \in B\}} \\ &= \sum_{j=1}^{\infty} \mathbb{1}_{\{S_j^{(n)} \in B + t - L_j^{(n)}, S_j^{(n)} \leq t\}}, \\ \mathcal{V}^{(n)}(t)(B) &= \sum_{j=1}^{A^{(n)}(t)} v_j^{(n)} \mathbb{1}_{\{L_j^{(n)} - (t - S_j^{(n)}) \in B\}} \\ &= \sum_{j=1}^{\infty} v_j^{(n)} \mathbb{1}_{\{S_j^{(n)} \in B + t - L_j^{(n)}, S_j^{(n)} \leq t\}}. \end{aligned}$$

Scaled EDF-related processes. For the processes just defined under the EDF queue discipline, we use the heavy-traffic scalings

$$\begin{aligned} \widehat{F}^{(n)}(t) &\triangleq \frac{1}{\sqrt{n}} F^{(n)}(nt), & \widehat{C}^{(n)}(t) &\triangleq \frac{1}{\sqrt{n}} C^{(n)}(nt), \\ \widehat{\mathcal{Q}}^{(n)}(t) &\triangleq \frac{1}{\sqrt{n}} \mathcal{Q}^{(n)}(nt)(\sqrt{n}B), & \widehat{\mathcal{W}}^{(n)}(t) &\triangleq \frac{1}{\sqrt{n}} \mathcal{W}^{(n)}(nt)(\sqrt{n}B). \end{aligned}$$

We define also

$$\begin{aligned} \widehat{\mathcal{A}}^{(n)}(t)(B) &\triangleq \frac{1}{\sqrt{n}} \mathcal{A}^{(n)}(nt)(\sqrt{n}B), \\ &= \frac{1}{\sqrt{n}} \sum_{j=1}^{A^{(n)}(nt)} \mathbf{1}_{\{L_j^{(n)} - (nt - S_j^{(n)}) \in \sqrt{n}B\}} \\ &= \frac{1}{\sqrt{n}} \sum_{j=1}^{\infty} \mathbf{1}_{\{S_j^{(n)} \in \sqrt{n}B + nt - L_j^{(n)}, S_j^{(n)} \leq nt\}}, \\ \widehat{\gamma}^{(n)}(t)(B) &\triangleq \frac{1}{\sqrt{n}} \gamma^{(n)}(nt)(\sqrt{n}B) \\ &= \frac{1}{\sqrt{n}} \sum_{j=1}^{A^{(n)}(nt)} v_j^{(n)} \mathbf{1}_{\{L_j^{(n)} - (nt - S_j^{(n)}) \in \sqrt{n}B\}} \\ &= \frac{1}{\sqrt{n}} \sum_{j=1}^{\infty} v_j^{(n)} \mathbf{1}_{\{S_j^{(n)} \in \sqrt{n}B + nt - L_j^{(n)}, S_j^{(n)} \leq nt\}}. \end{aligned}$$

3. Heavy-traffic analysis. We set

$$(3.1) \quad H(y) \triangleq \int_y^{\infty} (1 - G(\eta)) d\eta = \begin{cases} \int_y^{y^*} (1 - G(\eta)) d\eta, & \text{if } y \leq y^*, \\ 0, & \text{if } y > y^*. \end{cases}$$

The function H maps $(-\infty, y^*]$ onto $[0, \infty)$ and is strictly decreasing and Lipschitz continuous with Lipschitz constant 1 on $(-\infty, y^*]$. Therefore, there exists a continuous inverse function H^{-1} that maps $[0, \infty)$ onto $(-\infty, y^*]$. We next define what we shall ultimately show is the *limiting scaled frontier process*

$$(3.2) \quad F^*(t) \triangleq H^{-1}(W^*(t)), \quad t \geq 0,$$

where W^* is as in (2.13).

In this section, we prove weak convergence of $\widehat{\mathcal{W}}^{(n)}$ and $\widehat{\mathcal{Q}}^{(n)}$ as measure-valued processes. Weak convergence of measure-valued processes is a special case of weak convergence of metric-space-valued random objects, which is reviewed in Appendix A. We summarize its salient features here.

Denote by \mathcal{M} the set of all finite, nonnegative measures on $\mathcal{B}(\mathbb{R})$, the Borel subsets of \mathbb{R} . Under the weak topology, \mathcal{M} is separable. We can define a metric $d_{\mathcal{M}}$ on \mathcal{M} that is consistent with the weak topology on \mathcal{M} .

We now define $D_{\mathcal{M}}[0, \infty)$, the space of RCLL measure-valued functions on $[0, \infty)$. An *RCLL measure-valued process* is a random object that takes values in $D_{\mathcal{M}}[0, \infty)$, where in $D_{\mathcal{M}}[0, \infty)$ we use the Borel σ -algebra (generated by the open sets under the Skorohod topology). A sequence $\{X_n\}_{n=1}^\infty$ of RCLL measure-valued processes *converges weakly* to an RCLL measure-valued process X if the measures induced on $D_{\mathcal{M}}[0, \infty)$ by X_n converge weakly to the measure induced on $D_{\mathcal{M}}[0, \infty)$ by X , that is, for every bounded, continuous (or equivalently, uniformly continuous) function $F: D_{\mathcal{M}}[0, \infty) \rightarrow \mathbb{R}$, we have

$$\lim_{n \rightarrow \infty} \mathbb{E}_n F(X_n) = \mathbb{E}^* F(X).$$

The expectation operators \mathbb{E}_n and \mathbb{E}^* reflect the fact that each X_n may be defined on a probability space with a probability measure depending on the index n , and all these spaces may differ from the space on which X is defined. In the application of this paper, the prelimit processes are all defined on the same space and we write \mathbb{E} rather than \mathbb{E}_n .

The main result of this section is the following.

THEOREM 3.1. *Let $\widehat{\mathcal{W}}^*$ and $\widehat{\mathcal{D}}^*$ be the measure-valued processes defined by*

$$(3.3) \quad \widehat{\mathcal{W}}^*(t)(B) = \int_{B \cap [F^*(t), \infty)} (1 - G(y)) \, dy, \quad \widehat{\mathcal{D}}^*(t)(B) = \lambda \widehat{\mathcal{W}}^*(t)(B)$$

for all Borel sets $B \subset \mathbb{R}$. Then the measure-valued processes $\widehat{\mathcal{W}}^{(n)}$ and $\widehat{\mathcal{D}}^{(n)}$ converge weakly to $\widehat{\mathcal{W}}^$ and $\widehat{\mathcal{D}}^*$, respectively.*

COROLLARY 3.2. *Under the earliest-deadline-first queue discipline, the scaled queue length processes $\widehat{Q}^{(n)}$ defined by (2.8) converge weakly to λW^* .*

PROOF OF COROLLARY 3.2. We note that

$$\begin{aligned} \widehat{\mathcal{W}}^*(t)(\mathbb{R}) &= H(F^*(t)) = W^*(t), \\ \widehat{\mathcal{D}}^*(t)(\mathbb{R}) &= \lambda W^*(t), \\ \widehat{\mathcal{D}}^{(n)}(t)(\mathbb{R}) &= \widehat{Q}^{(n)}(t). \end{aligned}$$

The mapping from \mathcal{M} into \mathbb{R} , which maps each $\mu \in \mathcal{M}$ to its total mass $\mu(\mathbb{R})$, is continuous. By Theorem 3.1 and the continuous mapping theorem, Theorem A.1,

$$\widehat{\mathcal{D}}^{(n)}(\mathbb{R}) \Rightarrow \widehat{\mathcal{D}}^*(\mathbb{R})$$

or, equivalently,

$$\widehat{Q}^{(n)} \Rightarrow \lambda W^*. \quad \square$$

The proof of Theorem 3.1 is given at the end of this section. To prove this result, we first examine the convergence of the measure-valued processes $\widehat{\mathcal{V}}^{(n)}$ and $\widehat{\mathcal{A}}^{(n)}$. Recall that these processes keep track of arrived work and arrived customers, but not departures.

PROPOSITION 3.3. *Let $-\infty < y < y^*$ and $T > 0$ be given. As $n \rightarrow \infty$,*

$$(3.4) \quad \sup_{0 \leq t \leq T} |\widehat{\mathcal{Y}}^{(n)}(t)(y, \infty) + H(y + \sqrt{n}t) - H(y)| \xrightarrow{P} 0,$$

$$(3.5) \quad \sup_{0 \leq t \leq T} |\widehat{\mathcal{X}}^{(n)}(t)(y, \infty) + \lambda H(y + \sqrt{n}t) - \lambda H(y)| \xrightarrow{P} 0.$$

PROOF. For (3.4), let $\varepsilon > 0$ be given and choose a partition $y = \eta_0 < \eta_1 < \dots < \eta_M = y^*$ such that $|\eta_{m+1} - \eta_m| \leq \varepsilon$ for every $m = 0, \dots, M - 1$. Then the following inequalities hold, for each $\bar{m} = 1, \dots, M$:

$$(3.6) \quad \begin{aligned} & -\varepsilon + \sum_{m=0}^{\bar{m}-1} (1 - G(\eta_m))(\eta_{m+1} - \eta_m) \\ & \leq \int_y^{\eta_{\bar{m}}} (1 - G(\eta)) d\eta \leq \varepsilon + \sum_{m=0}^{\bar{m}-1} (1 - G(\eta_{m+1}))(\eta_{m+1} - \eta_m). \end{aligned}$$

To see why this is true, observe that for each $m = 0, \dots, M - 1$ we have

$$\begin{aligned} \int_{\eta_m}^{\eta_{m+1}} (1 - G(\eta)) d\eta & \leq (1 - G(\eta_m))(\eta_{m+1} - \eta_m) \\ & = (1 - G(\eta_{m+1}))(\eta_{m+1} - \eta_m) \\ & \quad + (G(\eta_{m+1}) - G(\eta_m))(\eta_{m+1} - \eta_m) \\ & \leq (1 - G(\eta_{m+1}))(\eta_{m+1} - \eta_m) + \varepsilon(G(\eta_{m+1}) - G(\eta_m)). \end{aligned}$$

Summing up the preceding inequality for $m = 0, \dots, \bar{m} - 1$ and cancelling the “telescoping” terms gives the right inequality in (3.6). The left inequality is obtained in a similar way.

For $0 \leq t \leq T$, we have

$$\begin{aligned} \widehat{\mathcal{Y}}^{(n)}(t)(y, \infty) & = \frac{1}{\sqrt{n}} \sum_{j=1}^{\infty} v_j^{(n)} \mathbb{1}_{\{nt + \sqrt{n}y - L_j^{(n)} < S_j^{(n)} \leq nt\}} \\ & \leq \frac{1}{\sqrt{n}} \sum_{m=1}^M \sum_{j=1}^{\infty} v_j^{(n)} \mathbb{1}_{\{\sqrt{n}\eta_{m-1} < L_j^{(n)} \leq \sqrt{n}\eta_m\}} \mathbb{1}_{\{nt - \sqrt{n}(\eta_m - y) < S_j^{(n)} \leq nt\}} \\ & = \sum_{m=1}^M \widehat{Y}_m^{(n)}(t) + \sum_{m=1}^M \widehat{U}_m^{(n)}(t), \end{aligned}$$

where

$$\begin{aligned} \widehat{Y}_m^{(n)}(t) & \triangleq \frac{1}{\sqrt{n}} \sum_{j=1}^{\infty} \left[v_j^{(n)} \mathbb{1}_{\{\sqrt{n}\eta_{m-1} < L_j^{(n)} \leq \sqrt{n}\eta_m\}} - \frac{1}{\mu^{(n)}} (G(\eta_m) - G(\eta_{m-1})) \right] \\ & \quad \times \mathbb{1}_{\{nt - \sqrt{n}(\eta_m - y) < S_j^{(n)} \leq nt\}}, \\ \widehat{U}_m^{(n)}(t) & \triangleq \frac{1}{\mu^{(n)}\sqrt{n}} \sum_{j=1}^{\infty} (G(\eta_m) - G(\eta_{m-1})) \mathbb{1}_{\{nt - \sqrt{n}(\eta_m - y) < S_j^{(n)} \leq nt\}}. \end{aligned}$$

To see that for each m , $\widehat{Y}_m^{(n)} \Rightarrow 0$ as $n \rightarrow \infty$, we define the sequence of non-negative i.i.d. random variables

$$\tilde{v}_j^{(n)} \triangleq v_j^{(n)} \mathbb{1}_{\{\sqrt{n} \eta_{m-1} < L_j^{(n)} \leq \sqrt{n} \eta_m\}}, \quad j = 1, 2, \dots,$$

and set

$$\begin{aligned} \tilde{V}_m^{(n)}(t) &\triangleq \frac{1}{\sqrt{n}} \sum_{j=1}^{\lfloor nt \rfloor} (\tilde{v}_j^{(n)} - \mathbb{E} \tilde{v}_j^{(n)}) \\ &= \frac{1}{\sqrt{n}} \sum_{j=1}^{\lfloor nt \rfloor} \left[v_j^{(n)} \mathbb{1}_{\{\sqrt{n} \eta_{m-1} < L_j^{(n)} \leq \sqrt{n} \eta_m\}} - \frac{1}{\mu^{(n)}} (G(\eta_m) - G(\eta_{m-1})) \right]. \end{aligned}$$

We may write

$$\begin{aligned} \widehat{Y}_m^{(n)}(t) &= \tilde{V}_m^{(n)} \left(\frac{1}{n} A^{(n)}(nt) \right) - \tilde{V}_m^{(n)} \left(\frac{1}{n} A^{(n)}((nt - \sqrt{n}(\eta_m - y))^+) \right) \\ &= \tilde{V}_m^{(n)} \left(\frac{1}{\sqrt{n}} \widehat{A}^{(n)}(t) + \lambda^{(n)} t \right) - \tilde{V}_m^{(n)} \left(\frac{1}{\sqrt{n}} \widehat{A}^{(n)} \left(\left(t - \frac{1}{\sqrt{n}}(\eta_m - y) \right)^+ \right) \right) \\ &\quad + \left(\lambda^{(n)} t - \frac{1}{\sqrt{n}} \lambda^{(n)}(\eta_m - y) \right)^+. \end{aligned}$$

Theorem B.1 (with $v_j^{(n)}$ replaced by $\tilde{v}_j^{(n)}$) implies that $\{\tilde{V}_m^{(n)}\}_{n=1}^\infty$ has a continuous weak limit \tilde{V}_m^* . The differencing theorem, Theorem A.3, and Theorem B.2 imply that $\widehat{Y}_m^{(n)} \Rightarrow 0$ on $[0, T]$, and hence

$$\sup_{0 \leq t \leq T} \left| \sum_{m=1}^M \widehat{Y}_m^{(n)}(t) \right| \xrightarrow{P} 0.$$

For the analysis of $\widehat{U}_m^{(n)}(t)$, we observe that

$$\begin{aligned} \widehat{U}_m^{(n)}(t) &= \frac{1}{\mu^{(n)} \sqrt{n}} (G(\eta_m) - G(\eta_{m-1})) \\ &\quad \times \left[A^{(n)}(nt) - A^{(n)}((nt - \sqrt{n}(\eta_m - y))^+) \right] \\ &= \frac{1}{\mu^{(n)}} (G(\eta_m) - G(\eta_{m-1})) \\ &\quad \times \left[\widehat{A}^{(n)}(t) - \widehat{A}^{(n)} \left(\left(t - \frac{1}{\sqrt{n}}(\eta_m - y) \right)^+ \right) \right] \\ &\quad + \lambda^{(n)} \sqrt{n} t - \lambda^{(n)} (\sqrt{n} t - (\eta_m - y))^+. \end{aligned}$$

As $n \rightarrow \infty$, the sequence of processes

$$\left\{ \widehat{A}^{(n)}(t) - \widehat{A}^{(n)}\left(\left(t - \frac{1}{\sqrt{n}}(\eta_m - y)\right)^+\right); t \geq 0 \right\}_{n=1}^{\infty}$$

converges weakly to zero. Hence, the weak limit of $\sum_{m=1}^M \widehat{U}_m^{(n)}(t)$ is the weak limit of

$$(3.7) \quad \sum_{m=1}^M (G(\eta_m) - G(\eta_{m-1}))[\sqrt{n}t - (\sqrt{n}t - (\eta_m - y))^+].$$

With n and t fixed, we define

$$\bar{m}(t) = \max \left\{ m \in \{0, 1, \dots, M\}; t \geq \frac{1}{\sqrt{n}}(\eta_m - y) \right\}.$$

Then (3.7) becomes

$$\begin{aligned} & \sum_{m=1}^{\bar{m}(t)} (G(\eta_m) - G(\eta_{m-1}))(\eta_m - y) + \sum_{m=\bar{m}(t)+1}^M (G(\eta_m) - G(\eta_{m-1}))\sqrt{n}t \\ & \leq \sum_{m=1}^{\bar{m}(t)} (1 - G(\eta_{m-1}))(\eta_m - y) - \sum_{m=1}^{\bar{m}(t)} (1 - G(\eta_m))(\eta_m - y) \\ & \quad + \mathbb{1}_{\{\bar{m}(t) \leq M-1\}}(\eta_{\bar{m}(t)+1} - y)(G(\eta_M) - G(\eta_{\bar{m}(t)})) \\ & = \sum_{m=0}^{\bar{m}(t)-1} (1 - G(\eta_m))(\eta_{m+1} - y) - \sum_{m=0}^{\bar{m}(t)-1} (1 - G(\eta_m))(\eta_m - y) \\ & \quad - (1 - G(\eta_{\bar{m}(t)}))(\eta_{\bar{m}(t)} - y) \\ & \quad + \mathbb{1}_{\{\bar{m}(t) \leq M-1\}}(\eta_{\bar{m}(t)+1} - y)(1 - G(\eta_{\bar{m}(t)})). \end{aligned}$$

If $\bar{m}(t) = M$, then $1 - G(\eta_{\bar{m}(t)}) = 0$ and we have

$$\sum_{m=0}^{\bar{m}(t)-1} (1 - G(\eta_m))(\eta_{m+1} - \eta_m) \leq \varepsilon + \int_y^{y^*} (1 - G(\eta)) d\eta,$$

where we have used (3.6). If $\bar{m}(t) < M$, we have again from (3.6) that

$$\begin{aligned} \sum_{m=0}^{\bar{m}(t)} (1 - G(\eta_m))(\eta_{m+1} - \eta_m) & \leq \varepsilon + \int_y^{\eta_{\bar{m}(t)+1}} (1 - G(\eta)) d\eta \\ & \leq 2\varepsilon + \int_y^{\eta_{\bar{m}(t)}} (1 - G(\eta)) d\eta \\ & \leq 2\varepsilon + \int_y^{y+\sqrt{n}t} (1 - G(\eta)) d\eta. \end{aligned}$$

In the former case, when $\bar{m}(t) = M$, we have $y + \sqrt{n}t \geq y^*$. Since $G(\eta) = 1$ for $\eta \geq y^*$, in both cases we have on (3.7) the upper bound

$$2\varepsilon + \int_y^{y+\sqrt{n}t} (1 - G(\eta)) d\eta = 2\varepsilon - H(y + \sqrt{n}t) + H(y).$$

We conclude that

$$\sup_{0 \leq t \leq T} \left[\sum_{m=1}^M \widehat{U}_m^{(n)}(t) + H(y + \sqrt{n}t) - H(y) - 2\varepsilon \right]^+ \xrightarrow{P} 0.$$

Since $\varepsilon > 0$ is arbitrary, we have in fact shown

$$\sup_{0 \leq t \leq T} [\widehat{\mathcal{Y}}^{(n)}(t)(y, \infty) + H(y + \sqrt{n}t) - H(y)]^+ \xrightarrow{P} 0.$$

To complete the proof of (3.4), we use the lower bound

$$\widehat{\mathcal{Y}}^{(n)}(t)(y, \infty) \geq \sum_{m=1}^M \check{Y}_m^{(n)}(t) + \sum_{m=1}^M \check{U}_m^{(n)}(t),$$

where

$$\begin{aligned} \check{Y}_m^{(n)}(t) &\triangleq \frac{1}{\sqrt{n}} \sum_{j=1}^{\infty} \left[v_j^{(n)} \mathbb{1}_{\{\sqrt{n}\eta_{m-1} < L_j^{(n)} \leq \sqrt{n}\eta_m\}} - \frac{1}{\mu^{(n)}} (G(\eta_m) - G(\eta_{m-1})) \right] \\ &\quad \times \mathbb{1}_{\{nt - \sqrt{n}(\eta_{m-1} - y) < S_j^{(n)} \leq nt\}}, \\ \check{U}_m^{(n)}(t) &\triangleq \frac{1}{\mu^{(n)} \sqrt{n}} \sum_{j=1}^{\infty} (G(\eta_m) - G(\eta_{m-1})) \mathbb{1}_{\{nt - \sqrt{n}(\eta_{m-1} - y) < S_j^{(n)} \leq nt\}}. \end{aligned}$$

By the same argument used to show that $\widehat{Y}_m^{(n)} \Rightarrow 0$, we may show that $\check{Y}_m^{(n)} \Rightarrow 0$. In place of (3.7), we have now

$$(3.8) \quad \sum_{m=1}^M (G(\eta_m) - G(\eta_{m-1})) [\sqrt{n}t - (\sqrt{n}t - (\eta_{m-1} - y))^+]$$

and we need to lower bound this quantity. With n and t fixed, we define $\bar{m}(t)$ as before and (3.8) becomes

$$\begin{aligned} &\sum_{m=1}^{M \wedge (\bar{m}(t)+1)} (G(\eta_m) - G(\eta_{m-1})) (\eta_{m-1} - y) + \sum_{m=\bar{m}(t)+2}^M (G(\eta_m) - G(\eta_{m-1})) \sqrt{n}t \\ &\geq \sum_{m=1}^{M \wedge (\bar{m}(t)+1)} (1 - G(\eta_{m-1})) (\eta_{m-1} - y) \\ &\quad - \sum_{m=1}^{M \wedge (\bar{m}(t)+1)} (1 - G(\eta_m)) (\eta_{m-1} - y) \\ &\quad + \mathbb{1}_{\{\bar{m}(t) \leq M-2\}} (\eta_{\bar{m}(t)} - y) (G(\eta_M) - G(\eta_{\bar{m}(t)+1})) \end{aligned}$$

$$\begin{aligned}
 &= \sum_{m=1}^{\bar{m}(t)} (1 - G(\eta_m))(\eta_m - y) - \sum_{m=1}^{\bar{m}(t)} (1 - G(\eta_m))(\eta_{m-1} - y) \\
 &\quad - \mathbb{1}_{\{\bar{m}(t) \leq M-1\}} (1 - G(\eta_{\bar{m}(t)+1}))(\eta_{\bar{m}(t)} - y) \\
 &\quad + \mathbb{1}_{\{\bar{m}(t) \leq M-2\}} (1 - G(\eta_{\bar{m}(t)+1}))(\eta_{\bar{m}(t)} - y) \\
 &= \sum_{m=1}^{\bar{m}(t)} (1 - G(\eta_m))(\eta_m - \eta_{m-1}) - \mathbb{1}_{\{\bar{m}(t)=M-1\}} (1 - G(\eta_M))(\eta_{M-1} - y) \\
 &\geq -\varepsilon + \int_y^{\eta_{\bar{m}(t)}} (1 - G(\eta)) d\eta \\
 &\geq -2\varepsilon + \int_y^{y+\sqrt{n}t} (1 - G(\eta)) d\eta \\
 &= -2\varepsilon - H(y + \sqrt{n}t) + H(y).
 \end{aligned}$$

It follows that

$$\sup_{0 \leq t \leq T} \left[\sum_{m=1}^M \check{U}_m^{(n)}(t) + H(y + \sqrt{n}t) - H(y) + 2\varepsilon \right]^- \xrightarrow{P} 0.$$

Since $\varepsilon > 0$ is arbitrary, we have in fact shown

$$\sup_{0 \leq t \leq T} [\widehat{\mathcal{Y}}^{(n)}(t)(y, \infty) + H(y + \sqrt{n}t) - H(y)]^- \xrightarrow{P} 0$$

and (3.4) is proved.

The proof of (3.5) is accomplished by repeating the foregoing proof, replacing $v_j^{(n)}$ and $1/\mu^{(n)} = \mathbb{E}v_j^{(n)}$ everywhere by 1. \square

Using a Glivenko–Cantelli type of argument, we can upgrade Proposition 3.3 to make the convergence uniform with respect to y on compact intervals:

PROPOSITION 3.4. *Let $-\infty < y_0 < y^*$ and $T > 0$ be given. As $n \rightarrow \infty$,*

$$(3.9) \quad \sup_{y_0 \leq y \leq y^*} \sup_{0 \leq t \leq T} |\widehat{\mathcal{Y}}^{(n)}(t)(y, \infty) + H(y + \sqrt{n}t) - H(y)| \xrightarrow{P} 0,$$

$$(3.10) \quad \sup_{y_0 \leq y \leq y^*} \sup_{0 \leq t \leq T} |\widehat{\mathcal{A}}^{(n)}(t)(y, \infty) + \lambda H(y + \sqrt{n}t) - \lambda H(y)| \xrightarrow{P} 0.$$

PROOF. Let $\varepsilon > 0$ be given. We will produce an N such that, for all $n \geq N$,

$$\mathbb{P} \left\{ \sup_{y_0 \leq y \leq y^*} \sup_{0 \leq t \leq T} |\widehat{\mathcal{Y}}^{(n)}(t)(y, \infty) + H(y + \sqrt{n}t) - H(y)| \geq \varepsilon \right\} < \varepsilon.$$

To do this, we first choose a partition $y_0 < y_1 < \dots < y_M = y^*$, such that $|y_{m+1} - y_m| < \varepsilon/2$ for $m = 0, \dots, M - 1$ and hence

$$(3.11) \quad 0 \leq H(y_m) - H(y_{m+1}) < \frac{\varepsilon}{2} \quad \text{for } m = 0, \dots, M - 1.$$

According to Proposition 3.3, we can find n_0, n_1, \dots, n_M such that, for $m = 0, 1, \dots, M$ and for $n \geq n_m$,

$$(3.12) \quad \mathbb{P} \left\{ \sup_{0 \leq t \leq T} \left| \widehat{\mathcal{Y}}^{(n)}(t)(y_m, \infty) + H(y_m + \sqrt{n}t) - H(y_m) \right| \geq \frac{\varepsilon}{2} \right\} < \frac{\varepsilon}{2(M+1)}.$$

We choose $N = \max\{n_0, n_1, \dots, n_M\}$. Using first the monotonicity of $H(y)$ and of $\widehat{\mathcal{Y}}^{(n)}(t)(y, \infty) + H(y + \sqrt{n}t)$ with respect to y , and then (3.11) and (3.12) we see that for $n \geq N$,

$$\begin{aligned} & \mathbb{P} \left\{ \sup_{y_0 \leq y \leq y^*} \sup_{0 \leq t \leq T} \left| \widehat{\mathcal{Y}}^{(n)}(t)(y, \infty) + H(y + \sqrt{n}t) - H(y) \right| \geq \varepsilon \right\} \\ & \leq \sum_{m=0}^{M-1} \mathbb{P} \left\{ \sup_{y_m \leq y \leq y_{m+1}} \sup_{0 \leq t \leq T} \left| \widehat{\mathcal{Y}}^{(n)}(t)(y, \infty) + H(y + \sqrt{n}t) - H(y) \right| \geq \varepsilon \right\} \\ & \leq \sum_{m=0}^{M-1} \mathbb{P} \left\{ \sup_{0 \leq t \leq T} \left| \widehat{\mathcal{Y}}^{(n)}(t)(y_m, \infty) + H(y_m + \sqrt{n}t) - H(y_{m+1}) \right| \geq \varepsilon \right\} \\ & \quad + \sum_{m=0}^{M-1} \mathbb{P} \left\{ \sup_{0 \leq t \leq T} \left| \widehat{\mathcal{Y}}^{(n)}(t)(y_{m+1}, \infty) + H(y_{m+1} + \sqrt{n}t) - H(y_m) \right| \geq \varepsilon \right\} \\ & \leq \sum_{m=0}^M 2\mathbb{P} \left\{ \sup_{0 \leq t \leq T} \left| \widehat{\mathcal{Y}}^{(n)}(t)(y_m, \infty) + H(y_m + \sqrt{n}t) - H(y_m) \right| \geq \frac{\varepsilon}{2} \right\} \\ & < 2(M+1) \frac{\varepsilon}{2(M+1)} = \varepsilon. \end{aligned}$$

This proves (3.9); the proof of (3.10) is analogous. \square

COROLLARY 3.5. *Let $-\infty < y_0 < y^*$ and $T > 0$ be given. As $n \rightarrow \infty$,*

$$(3.13) \quad \sup_{y_0 \leq y \leq y^*} \sup_{0 \leq t \leq T} \widehat{\mathcal{Y}}^{(n)}(t)\{y\} \xrightarrow{P} 0,$$

$$(3.14) \quad \sup_{y_0 \leq y \leq y^*} \sup_{0 \leq t \leq T} \widehat{\mathcal{A}}^{(n)}(t)\{y\} \xrightarrow{P} 0.$$

PROOF. We give the proof of (3.13); the proof of (3.14) is analogous. For $\delta \in (0, 1)$, $y_0 \leq y \leq y^*$ and $0 \leq t \leq T$, the monotonicity and Lipschitz continuity

with constant 1 of the function H imply

$$\begin{aligned} & \widehat{\mathcal{V}}^{(n)}(t)(y, \infty) + H(y + \sqrt{n}t) - H(y) \\ & \leq \widehat{\mathcal{V}}^{(n)}(t)[y, \infty) + H(y + \sqrt{n}t) - H(y) \\ & \leq \widehat{\mathcal{V}}^{(n)}(t)(y - \delta, \infty) + H(y - \delta + \sqrt{n}t) - H(y - \delta) + \delta \\ & \leq \delta + \sup_{y_0 - 1 \leq \eta \leq y^*} \sup_{0 \leq s \leq T} |\widehat{\mathcal{V}}^{(n)}(s)(\eta, \infty) + H(\eta + \sqrt{n}s) - H(\eta)|. \end{aligned}$$

It follows that

$$\begin{aligned} \widehat{\mathcal{V}}^{(n)}(t)\{y\} & \leq \delta + \sup_{y_0 - 1 \leq \eta \leq y^*} \sup_{0 \leq s \leq T} |\widehat{\mathcal{V}}^{(n)}(s)(\eta, \infty) + H(\eta + \sqrt{n}s) - H(\eta)| \\ & \quad - [\widehat{\mathcal{V}}^{(n)}(t)(y, \infty) + H(y + \sqrt{n}t) - H(y)]. \end{aligned}$$

Proposition 3.4 implies that the limit of the right-hand side, in probability, is δ . Since δ is arbitrary, we have (3.13). \square

The heavy-traffic analysis of the queueing system with due dates depends critically on the following proposition, which asserts that the number of customers whose lead times lie between the current lead time $C^{(n)}(t)$ and the frontier $F^{(n)}(t)$ and the work associated with these customers is negligible.

PROPOSITION 3.6. *The processes $\widehat{\mathcal{D}}^{(n)}[\widehat{C}^{(n)}, \widehat{F}^{(n)}]$ and $\widehat{\mathcal{W}}^{(n)}[\widehat{C}^{(n)}, \widehat{F}^{(n)}]$ converge weakly to zero as $n \rightarrow \infty$.*

PROOF. We fix $T > 0$ and establish the convergence on $[0, T]$. For this we follow ideas of Peterson [24].

Let $y \leq y^*$ be given. For $t \geq 0$, we set

$$\widehat{T}^{(n)}(t) \triangleq \frac{1}{\sqrt{n}} \sum_{j=1}^{\lfloor nt \rfloor} \left(v_j^{(n)} \mathbb{1}_{\{L_j^{(n)} \leq \sqrt{n}y\}} - \frac{1}{\mu^{(n)}} G(y) \right).$$

According to Theorem B.1, $\widehat{T}^{(n)}$ converges weakly to a Brownian motion.

Next define

$$\tau^{(n)}(t) \triangleq \sup\{s \in [0, t]; \widehat{C}^{(n)}(s) = \widehat{F}^{(n)}(s)\}.$$

By assumption, $\widehat{C}^{(n)}(0) = \widehat{F}^{(n)}(0) = y^*$ and so $\tau^{(n)}(t) \leq t$, that is, the supremum is not over the empty set. We first show that

$$(3.15) \quad t - \tau^{(n)}(t) \Rightarrow 0$$

and subsequently show that

$$(3.16) \quad \sqrt{n}(t - \tau^{(n)}(t)) \Rightarrow 0,$$

where the convergence in (3.15) and (3.16) is for processes on $[0, T]$.

Observing that $\widehat{\mathcal{W}}^{(n)}[\widehat{C}^{(n)}, \widehat{F}^{(n)}]$ is RCLL, we note that

$$(3.17) \quad \begin{aligned} & \widehat{\mathcal{W}}^{(n)}(\tau^{(n)}(t)-)[\widehat{C}^{(n)}(\tau^{(n)}(t)-), \widehat{F}^{(n)}(\tau^{(n)}(t)-)] \\ &= \widehat{\mathcal{W}}^{(n)}(\tau^{(n)}(t)-)(\emptyset) = 0, \end{aligned}$$

$$(3.18) \quad \begin{aligned} & \widehat{\mathcal{W}}^{(n)}(\tau^{(n)}(t))[\widehat{C}^{(n)}(\tau^{(n)}(t)), \widehat{F}^{(n)}(\tau^{(n)}(t))] \\ & \leq \frac{1}{\sqrt{n}} \max_{1 \leq j \leq A^{(n)}(nt)} v_j^{(n)} \leq \max_{0 \leq s \leq T} [\widehat{N}^{(n)}(s) - \widehat{N}^{(n)}(s-)]. \end{aligned}$$

As long as there are customers with lead times in the unscaled interval $[C^{(n)}, F^{(n)}]$, the unscaled frontier $F^{(n)}$ decreases at rate 1 per unit time. Therefore, for $s \in (n\tau^{(n)}(t), nt]$,

$$(3.19) \quad F^{(n)}(s) = F^{(n)}(n\tau^{(n)}(t)) - (s - n\tau^{(n)}(t)).$$

For what follows, it will be helpful to introduce some notation. We define

$$(3.20) \quad \begin{aligned} D^{(n)}(t) &= \sum_{j=1}^{\infty} v_j^{(n)} \mathbb{1}_{\{n\tau^{(n)}(t) < S_j^{(n)} \leq nt\}} \\ & \quad \times \mathbb{1}_{\{L_j^{(n)} - (nt - S_j^{(n)}) < F^{(n)}(n\tau^{(n)}(t)) - n(t - \tau^{(n)}(t))\}}. \end{aligned}$$

Observe that because of (3.19), whenever $n\tau^{(n)}(t) < S_j^{(n)} \leq nt$, the condition

$$L_j^{(n)} - (nt - S_j^{(n)}) < F^{(n)}(n\tau^{(n)}(t)) - n(t - \tau^{(n)}(t))$$

is equivalent to $L_j^{(n)} < F^{(n)}(S_j^{(n)})$. In other words, $D^{(n)}(t)$ is the work associated with customers arriving within the time interval $(n\tau^{(n)}(t), nt]$ whose lead times upon arrival are to the left of the frontier.

We now note that on the time interval $(n\tau^{(n)}(t), nt]$ the server is never idle, which means that the workload is being decreased by the server at a constant rate 1. This gives us the estimate

$$\begin{aligned} 0 &\leq \mathcal{W}^{(n)}(nt)[C^{(n)}(nt), F^{(n)}(nt)] \\ &= \mathcal{W}^{(n)}(n\tau^{(n)}(t))[C^{(n)}(n\tau^{(n)}(t)), F^{(n)}(n\tau^{(n)}(t))] + D^{(n)}(t) - n(t - \tau^{(n)}(t)) \end{aligned}$$

or, after scaling,

$$(3.21) \quad \begin{aligned} 0 &\leq \widehat{\mathcal{W}}^{(n)}(t)[\widehat{C}^{(n)}(t), \widehat{F}^{(n)}(t)] \\ &= \widehat{\mathcal{W}}^{(n)}(\tau^{(n)}(t))[\widehat{C}^{(n)}(\tau^{(n)}(t)), \widehat{F}^{(n)}(\tau^{(n)}(t))] \\ & \quad + \frac{1}{\sqrt{n}} D^{(n)}(t) - \sqrt{n}(t - \tau^{(n)}(t)). \end{aligned}$$

Next, we estimate the term $(1/\sqrt{n})D^{(n)}(t)$. For $y \leq y^*$ we have either

$$(3.22) \quad t - \tau^{(n)}(t) < \frac{1}{\sqrt{n}}(y^* - y)$$

or else

$$(3.23) \quad n\tau^{(n)}(t) + \sqrt{n}(y^* - y) \leq nt.$$

In the latter case,

$$\begin{aligned}
 \frac{1}{\sqrt{n}}D^{(n)}(t) &\leq \frac{1}{\sqrt{n}} \sum_{j=1}^{\infty} v_j^{(n)} \mathbb{1}_{\{n\tau^{(n)}(t) < S_j^{(n)} \leq n\tau^{(n)}(t) + \sqrt{n}(y^* - y)\}} \\
 &\quad + \frac{1}{\sqrt{n}} \sum_{j=1}^{\infty} v_j^{(n)} \mathbb{1}_{\{n\tau^{(n)}(t) + \sqrt{n}(y^* - y) < S_j^{(n)} \leq nt\}} \mathbb{1}_{\{L_j^{(n)} \leq \sqrt{n}y\}} \\
 &= \frac{1}{\sqrt{n}}V^{(n)}(A^{(n)}(n\tau^{(n)}(t) + \sqrt{n}(y^* - y))) \\
 &\quad - \frac{1}{\sqrt{n}}V^{(n)}(A^{(n)}(n\tau^{(n)}(t))) + \widehat{T}^{(n)}\left(\frac{1}{n}A^{(n)}(nt)\right) \\
 &\quad - \widehat{T}^{(n)}\left(\frac{1}{n}A^{(n)}(n\tau^{(n)}(t) + \sqrt{n}(y^* - y))\right) \\
 &\quad + \frac{G(y)}{\mu^{(n)}\sqrt{n}}[A^{(n)}(nt) - A^{(n)}(n\tau^{(n)}(t) + \sqrt{n}(y^* - y))] \\
 &= \widehat{N}^{(n)}\left(\tau^{(n)}(t) + \frac{1}{\sqrt{n}}(y^* - y)\right) - \widehat{N}^{(n)}(\tau^{(n)}(t)) + y^* - y \\
 &\quad + \widehat{T}^{(n)}\left(\frac{1}{\sqrt{n}}\widehat{A}^{(n)}(t) + \lambda^{(n)}t\right) \\
 (3.24) \quad &\quad - \widehat{T}^{(n)}\left(\frac{1}{\sqrt{n}}\widehat{A}^{(n)}\left(\tau^{(n)}(t) + \frac{1}{\sqrt{n}}(y^* - y)\right)\right. \\
 &\quad \quad \left. + \lambda^{(n)}\tau^{(n)}(t) + \frac{\lambda^{(n)}}{\sqrt{n}}(y^* - y)\right) \\
 &\quad + \frac{G(y)}{\mu^{(n)}}\left[\widehat{A}^{(n)}(t) - \widehat{A}^{(n)}\left(\tau^{(n)}(t) + \frac{1}{\sqrt{n}}(y^* - y)\right)\right] \\
 &\quad + G(y)\rho^{(n)}\sqrt{n}(t - \tau^{(n)}(t)) - G(y)\rho^{(n)}(y^* - y) \\
 &\leq \left[\widehat{N}^{(n)}\left(\tau^{(n)}(t) + \frac{1}{\sqrt{n}}(y^* - y)\right) - \widehat{N}^{(n)}(\tau^{(n)}(t))\right] \\
 &\quad + \left[\widehat{T}^{(n)}\left(\frac{1}{\sqrt{n}}\widehat{A}^{(n)}(t) + \lambda^{(n)}t\right)\right. \\
 &\quad \quad \left. - \widehat{T}^{(n)}\left(\frac{1}{\sqrt{n}}\widehat{A}^{(n)}\left(\tau^{(n)}(t) + \frac{1}{\sqrt{n}}(y^* - y)\right)\right)\right. \\
 &\quad \quad \left. + \lambda^{(n)}\tau^{(n)}(t) + \frac{\lambda^{(n)}}{\sqrt{n}}(y^* - y)\right]
 \end{aligned}$$

$$\begin{aligned}
& + \frac{G(y)}{\mu^{(n)}} \left[\widehat{A}^{(n)}(t) - \widehat{A}^{(n)}\left(\tau^{(n)}(t) + \frac{1}{\sqrt{n}}(y^* - y)\right) \right] \\
& + G(y)\rho^{(n)}\sqrt{n}(t - \tau^{(n)}(t)) + (1 - G(y)\rho^{(n)})(y^* - y).
\end{aligned}$$

Continuing this inequality we may write

$$\begin{aligned}
\frac{1}{\sqrt{n}}D^{(n)}(t) & \leq \max_{0 \leq s \leq T} \left| \widehat{N}^{(n)}\left(s + \frac{1}{\sqrt{n}}(y^* - y)\right) - \widehat{N}^{(n)}(s) \right| \\
& + 2 \max_{0 \leq s \leq T} \left| \widehat{T}^{(n)}\left(\frac{1}{\sqrt{n}}\widehat{A}^{(n)}(s) + \lambda^{(n)}s\right) \right| \\
& + \frac{2G(y)}{\mu^{(n)}} \max_{0 \leq s \leq T} \left| \widehat{A}^{(n)}(s) \right| + G(y)\rho^{(n)}\sqrt{n}(t - \tau^{(n)}(t)) \\
& + (1 - G(y)\rho^{(n)})(y^* - y).
\end{aligned}$$

Assume now that $y < y^*$. Substituting the preceding inequality into (3.21), using (3.18) and dividing by $(1 - G(y)\rho^{(n)})\sqrt{n}$, we obtain

$$\begin{aligned}
(3.25) \quad 0 \leq t - \tau^{(n)}(t) & \leq \frac{1}{(1 - G(y)\rho^{(n)})\sqrt{n}} \\
& \times \left\{ \max_{0 \leq s \leq T} [\widehat{N}^{(n)}(s) - \widehat{N}^{(n)}(s-)] \right. \\
& + \max_{0 \leq s \leq T} \left[\widehat{N}^{(n)}\left(s + \frac{1}{\sqrt{n}}(y^* - y)\right) - \widehat{N}^{(n)}(s) \right] \\
& + 2 \max_{0 \leq s \leq T} \left| \widehat{T}^{(n)}\left(\frac{1}{\sqrt{n}}\widehat{A}^{(n)}(s) + \lambda^{(n)}s\right) \right| \\
& \left. + \frac{2G(y)}{\mu^{(n)}} \max_{0 \leq s \leq T} \left| \widehat{A}^{(n)}(s) \right| \right\} + \frac{1}{\sqrt{n}}(y^* - y).
\end{aligned}$$

Of course, if (3.22) holds, then (3.25) does as well. From (3.25) we have (3.15).

Next we substitute inequality (3.24) into (3.21), using (3.18), and dividing by $(1 - G(y)\rho^{(n)})$ to get

$$\begin{aligned}
0 \leq \sqrt{n}(t - \tau^{(n)}(t)) & \leq \frac{1}{(1 - G(y)\rho^{(n)})} \\
& \times \left\{ \max_{0 \leq s \leq T} [\widehat{N}^{(n)}(s) - \widehat{N}^{(n)}(s-)] \right. \\
& + \left[\widehat{N}^{(n)}\left(\tau^{(n)}(t) + \frac{1}{\sqrt{n}}(y^* - y)\right) - \widehat{N}^{(n)}(\tau^{(n)}(t)) \right] \\
& + \left[\widehat{T}^{(n)}\left(\frac{1}{\sqrt{n}}\widehat{A}^{(n)}(t) + \lambda^{(n)}t\right) \right. \\
& \quad \left. - \widehat{T}^{(n)}\left(\frac{1}{\sqrt{n}}\widehat{A}^{(n)}\left(\tau^{(n)}(t) + \frac{1}{\sqrt{n}}(y^* - y)\right)\right) \right]
\end{aligned}$$

$$\begin{aligned}
 & + \lambda^{(n)} \tau^{(n)}(t) + \frac{\lambda^{(n)}}{\sqrt{n}}(y^* - y) \Big] \\
 & + \frac{G(y)}{\mu^{(n)}} \left[\widehat{A}^{(n)}(t) - \widehat{A}^{(n)} \left(\tau^{(n)}(t) + \frac{1}{\sqrt{n}}(y^* - y) \right) \right]^+ + (y^* - y),
 \end{aligned}$$

where the positive part on the terms in $\{\dots\}$ ensures that this inequality holds even if (3.23) fails [and hence (3.22) holds]. As $n \rightarrow \infty$, (3.15) and the time change and differencing theorems, Theorems A.2 and A.3, imply that the right-hand side has limit $y^* - y$, that is,

$$\left[\sqrt{n}(t - \tau^{(n)}(t)) - (y^* - y) \right]^+ \Rightarrow 0.$$

Since $y < y^*$ is arbitrary, we must have (3.16).

Using (3.18), (3.21) and inequality (3.24) with $y = y^*$, we obtain

$$\begin{aligned}
 0 \leq \widehat{\mathcal{W}}^{(n)}(t) [\widehat{C}^{(n)}(t), \widehat{F}^{(n)}(t)] & \leq \max_{0 \leq s \leq T} [\widehat{N}^{(n)}(s) - \widehat{N}^{(n)}(s-)] \\
 & + \left[\widehat{T}^{(n)} \left(\frac{1}{\sqrt{n}} \widehat{A}^{(n)}(t) + \lambda^{(n)} t \right) - \widehat{T}^{(n)} \left(\frac{1}{\sqrt{n}} \widehat{A}^{(n)} \left(\tau^{(n)}(t) \right) + \lambda^{(n)} \tau^{(n)}(t) \right) \right] \\
 & + \frac{1}{\mu^{(n)}} \left[\widehat{A}^{(n)}(t) - \widehat{A}^{(n)}(\tau^{(n)}(t)) \right] - (1 - \rho^{(n)}) \sqrt{n}(t - \tau^{(n)}(t)).
 \end{aligned}$$

Once again the time change and differencing theorems, Theorems A.2 and A.3, show that the right-hand side has limit zero. This implies

$$\widehat{\mathcal{W}}^{(n)}[\widehat{C}^{(n)}, \widehat{F}^{(n)}] \Rightarrow 0.$$

Similarly,

$$\begin{aligned}
 0 \leq \widehat{\mathcal{D}}^{(n)}(t) [\widehat{C}^{(n)}(t), \widehat{F}^{(n)}(t)] & \leq \frac{1}{\sqrt{n}} [1 + A^{(n)}(nt) - A^{(n)}(n\tau^{(n)}(t))] \\
 & = \frac{1}{\sqrt{n}} + \widehat{A}^{(n)}(t) - \widehat{A}^{(n)}(\tau^{(n)}(t)) + \lambda^{(n)} \sqrt{n}(t - \tau^{(n)}(t)).
 \end{aligned}$$

The time change and differencing theorems, Theorems A.2 and A.3, and the convergence (3.16) imply that the right-hand side has limit zero. This implies

$$\widehat{\mathcal{D}}^{(n)}[\widehat{C}^{(n)}, \widehat{F}^{(n)}] \Rightarrow 0. \quad \square$$

We next examine the limit of the scaled frontier process $\widehat{F}^{(n)}$. Since $\sqrt{n} y^* - nt \leq F^{(n)}(nt) \leq \sqrt{n} y^*$ at all times, we have the bounds

$$(3.26) \quad y^* - \sqrt{n} t \leq \widehat{F}^{(n)}(t) \leq y^*, \quad t \geq 0.$$

The following lemma provides a tightness bound from below.

LEMMA 3.7. *For every $T > 0$ and $\varepsilon > 0$, there exists $y \in (-\infty, y^*)$ such that for all n ,*

$$\mathbb{P}\left\{\inf_{0 \leq t \leq T} \widehat{F}^{(n)}(t) < y\right\} < \varepsilon.$$

PROOF. By definition,

$$W^{(n)}(t) = \mathscr{W}^{(n)}(t)(\mathbb{R}) = \mathscr{W}^{(n)}(t)[C^{(n)}(t), F^{(n)}(t)] + \mathscr{W}^{(n)}(t)[F^{(n)}(t), \infty).$$

Scaling this equation, we obtain

$$\widehat{W}^{(n)}(t) = \widehat{\mathscr{W}}^{(n)}(t)[\widehat{C}^{(n)}(t), \widehat{F}^{(n)}(t)] + \widehat{\mathscr{W}}^{(n)}(t)[\widehat{F}^{(n)}(t), \infty).$$

Corollary B.4 implies $\widehat{W}^{(n)} \Rightarrow W^*$, where W^* is a reflected Brownian motion with drift, and Proposition 3.6 shows that $\widehat{\mathscr{W}}^{(n)}[\widehat{C}^{(n)}, \widehat{F}^{(n)}] \Rightarrow 0$. Therefore,

$$(3.27) \quad \widehat{\mathscr{W}}^{(n)}[\widehat{F}^{(n)}, \infty] \Rightarrow W^*.$$

Fix $T > 0$. The continuous mapping theorem, Theorem A.1, applied to (3.27) yields

$$(3.28) \quad \max_{0 \leq t \leq T} \widehat{W}^{(n)}(t)[\widehat{F}^{(n)}(t), \infty] \Rightarrow \max_{0 \leq t \leq T} W^*(t).$$

At time t , no customer with lead time in $(F^{(n)}(t), \infty)$ has ever been in service, so $0 \leq \widehat{\mathscr{Y}}^{(n)}(t)(\widehat{F}^{(n)}(t), \infty) \leq \widehat{\mathscr{W}}^{(n)}(t)[\widehat{F}^{(n)}(t), \infty]$. From (3.28) we see then that the sequence of random variables $\{\max_{0 \leq t \leq T} \widehat{\mathscr{Y}}^{(n)}(t)(\widehat{F}^{(n)}(t), \infty)\}_{n=1}^{\infty}$ is tight. Let $\varepsilon > 0$ be given. Because $\lim_{y \rightarrow -\infty} \bar{H}(y) = \infty$, we may choose $y < y^*$ so that for each n ,

$$\max_{0 \leq t \leq T} \widehat{\mathscr{Y}}^{(n)}(t)(\widehat{F}^{(n)}(t), \infty) \leq \sqrt{H(y)} \quad \text{on } A_n,$$

where the event A_n satisfies $\mathbb{P}(A_n) \geq 1 - \varepsilon/3$. Proposition 3.3 and the continuous mapping theorem, Theorem A.1, imply the existence of N such that for every $n \geq N$,

$$\min_{0 \leq t \leq T} [\widehat{\mathscr{Y}}^{(n)}(t)(y, \infty) + H(y + \sqrt{n}t)] \geq \frac{1}{2}H(y) \quad \text{on } B_n,$$

where the event B_n satisfies $\mathbb{P}(B_n) \geq 1 - \varepsilon/3$.

Now $\mathbb{P}(A_n \cap B_n) \geq 1 - 2\varepsilon/3$ and on $A_n \cap B_n$,

$$\begin{aligned} \sqrt{H(y)} &\geq \max_{0 \leq t \leq T} \widehat{\mathscr{Y}}^{(n)}(t)(\widehat{F}^{(n)}(t), \infty) \\ &\geq \max_{0 \leq t \leq T} \widehat{\mathscr{Y}}^{(n)}(t)(y, \infty) \mathbb{1}_{\{\widehat{F}^{(n)}(t) < y\}} \\ &\geq \max_{0 \leq t \leq T} [\widehat{\mathscr{Y}}^{(n)}(t)(y, \infty) + H(y + \sqrt{n}t)] \mathbb{1}_{\{\widehat{F}^{(n)}(t) < y\}}, \end{aligned}$$

because $y + \sqrt{n}t \geq y^*$ on $\{\widehat{F}^{(n)}(t) < y\}$ [see (3.26)] and $H(y + \sqrt{n}t) = 0$. Continuing, we have on $A_n \cap B_n$ that

$$\sqrt{H(y)} \geq \frac{1}{2}H(y) \max_{0 \leq t \leq T} \mathbb{1}_{\{\widehat{F}^{(n)}(t) < y\}} = \frac{1}{2}H(y) \mathbb{1}_{\{\inf_{0 \leq t \leq T} \widehat{F}^{(n)}(t) < y\}},$$

which implies

$$\begin{aligned} \frac{2}{\sqrt{H(y)}} &\geq \mathbb{E} \left[\mathbb{1}_{\{\inf_{0 \leq t \leq T} \widehat{F}^{(n)}(t) < y\}} \mathbb{1}_{A_n \cap B_n} \right] \\ &= \mathbb{E} \left[\mathbb{1}_{\{\inf_{0 \leq t \leq T} \widehat{F}^{(n)}(t) < y\}} (1 - \mathbb{1}_{(A_n \cap B_n)^c}) \right] \\ &\geq \mathbb{P} \left(\inf_{0 \leq t \leq T} \widehat{F}^{(n)}(t) < y \right) - \mathbb{P}((A_n \cap B_n)^c) \\ &\geq \mathbb{P} \left(\inf_{0 \leq t \leq T} \widehat{F}^{(n)}(t) < y \right) - \frac{2\varepsilon}{3}. \end{aligned}$$

In other words,

$$\mathbb{P} \left(\inf_{0 \leq t \leq T} \widehat{F}^{(n)}(t) < y \right) \leq \frac{2\varepsilon}{3} + \frac{2}{\sqrt{H(y)}}$$

and by choosing $|y|$ larger if necessary, we may ensure that $2/\sqrt{H(y)} < \varepsilon/3$. \square

COROLLARY 3.8. *The processes $\widehat{\mathcal{D}}^{(n)}[\widehat{C}^{(n)}, \widehat{F}^{(n)}]$ and $\widehat{\mathcal{W}}^{(n)}[\widehat{C}^{(n)}, \widehat{F}^{(n)}]$ converge weakly to zero as $n \rightarrow \infty$. In particular,*

$$(3.29) \quad \widehat{\mathcal{V}}^{(n)}(\widehat{F}^{(n)}, \infty) \Rightarrow W^*.$$

PROOF. In light of Proposition 3.6, it suffices to show that $\widehat{\mathcal{D}}^{(n)}(t)\{\widehat{F}^{(n)}(t)\}$ and $\widehat{\mathcal{W}}^{(n)}\{\widehat{F}^{(n)}(t)\}$ converge weakly to zero. We first choose $T > 0$, $\varepsilon > 0$ and y_0 so that $\mathbb{P}\{\inf_{0 \leq t \leq T} \widehat{F}^{(n)}(t) < y_0\} < \varepsilon$. On the complementary set $\{\inf_{0 \leq t \leq T} \widehat{F}^{(n)}(t) \geq y_0\}$, we have

$$\begin{aligned} \sup_{0 \leq t \leq T} \widehat{W}^{(n)}\{\widehat{F}^{(n)}(t)\} &\leq \sup_{0 \leq t \leq T} \widehat{\mathcal{V}}^{(n)}(t)\{\widehat{F}^{(n)}(t)\} \\ &\leq \sup_{y_0 \leq y \leq y^*} \sup_{0 \leq t \leq T} \widehat{\mathcal{V}}^{(n)}(t)\{y\}. \end{aligned}$$

Corollary 3.5 implies that the last expression converges in probability to zero. The proof for $\widehat{\mathcal{D}}^{(n)}(t)\{\widehat{F}^{(n)}(t)\}$ is analogous.

Finally, because at time t no customer with lead time in $(F^{(n)}(t), \infty)$ has ever been in service, we have

$$0 \leq \widehat{W}^{(n)}(t)[\widehat{F}^{(n)}(t), \infty) - \widehat{\mathcal{V}}^{(n)}(t)(\widehat{F}^{(n)}(t), \infty) = \widehat{W}^{(n)}\{\widehat{F}^{(n)}(t)\}.$$

From (3.27), we conclude that (3.29) holds. \square

PROPOSITION 3.9. *Let $-\infty < y_0 < y^*$ and $T > 0$ be given. As $n \rightarrow \infty$,*

$$(3.30) \quad \sup_{y_0 \leq y \leq y^*} \sup_{0 \leq t \leq T} |\widehat{\mathcal{Y}}^{(n)}(t)(\widehat{F}^{(n)}(t) \vee y, \infty) - H(\widehat{F}^{(n)}(t) \vee y)| \xrightarrow{P} 0,$$

$$(3.31) \quad \sup_{y_0 \leq y \leq y^*} \sup_{0 \leq t \leq T} |\widehat{\mathcal{Z}}^{(n)}(t)(\widehat{F}^{(n)}(t) \vee y, \infty) - \lambda H(\widehat{F}^{(n)}(t) \vee y)| \xrightarrow{P} 0.$$

PROOF. Because of the inequality $\widehat{F}^{(n)}(t) + \sqrt{n}t \geq y^*$ and the fact that $H(y) = 0$ for $y \geq y^*$, the quantities on the left-hand sides of relations (3.9) and (3.10) in Proposition 3.4 dominate those on the left-hand sides of (3.30) and (3.31), respectively. \square

PROPOSITION 3.10. $\widehat{F}^{(n)} \Rightarrow F^* \stackrel{\Delta}{=} H^{-1}(W^*)$.

PROOF. Let $T > 0$ and $\varepsilon > 0$ be given. Using Lemma 3.7, we may choose $y_0 < y^*$ so that $A_n \stackrel{\Delta}{=} \{\inf_{0 \leq t \leq T} \widehat{F}^{(n)}(t) > y_0\}$ satisfies $\mathbb{P}(A_n) \geq 1 - \varepsilon$ for every n . According to Proposition 3.9, there is an N such that for each $n \geq N$, there is an event B_n with $\mathbb{P}(B_n) \geq 1 - \varepsilon$ and on the event B_n , we have

$$\sup_{y_0 \leq y \leq y^*} \sup_{0 \leq t \leq T} |\widehat{\mathcal{Y}}^{(n)}(t)(\widehat{F}^{(n)}(t) \vee y, \infty) - H(\widehat{F}^{(n)}(t) \vee y)| < \varepsilon.$$

On the intersection $A_n \cap B_n$, we have in particular

$$\begin{aligned} & \sup_{0 \leq t \leq T} |\widehat{\mathcal{Y}}^{(n)}(t)(\widehat{F}^{(n)}(t), \infty) - H(\widehat{F}^{(n)}(t))| \\ &= \sup_{0 \leq t \leq T} |\widehat{\mathcal{Y}}^{(n)}(t)(\widehat{F}^{(n)}(t) \vee y_0, \infty) - H(\widehat{F}^{(n)}(t) \vee y_0)| < \varepsilon. \end{aligned}$$

It follows that

$$\sup_{0 \leq t \leq T} |\widehat{\mathcal{Y}}^{(n)}(t)(\widehat{F}^{(n)}(t), \infty) - H(\widehat{F}^{(n)}(t))| \xrightarrow{P} 0.$$

Relation (3.29) shows that

$$(3.32) \quad H(\widehat{F}^{(n)}) \Rightarrow W^*.$$

Applying the continuous function H^{-1} to both sides of (3.32), we obtain the desired result $\widehat{F}^{(n)} \Rightarrow H^{-1}(W^*)$ from the continuous mapping theorem, Theorem A.1. \square

PROPOSITION 3.11. *Let $T > 0$ be given. As $n \rightarrow \infty$,*

$$(3.33) \quad \sup_{y \in \mathbb{R}} \sup_{0 \leq t \leq T} |\widehat{\mathcal{Y}}^{(n)}(t)(y, \infty) - H(\widehat{F}^{(n)}(t) \vee y)| \xrightarrow{P} 0,$$

$$(3.34) \quad \sup_{y \in \mathbb{R}} \sup_{0 \leq t \leq T} |\widehat{\mathcal{Z}}^{(n)}(t)(y, \infty) - \lambda H(\widehat{F}^{(n)}(t) \vee y)| \xrightarrow{P} 0.$$

PROOF. For $y \geq y^*$,

$$\widehat{\mathcal{W}}^{(n)}(t)(y, \infty) = H(\widehat{F}^{(n)}(t) \vee y) = 0.$$

For $y < y^*$, $0 \leq t \leq T$,

$$\begin{aligned} & |\widehat{\mathcal{W}}^{(n)}(t)(y, \infty) - H(\widehat{F}^{(n)}(t) \vee y)| \\ & \leq |\widehat{\mathcal{W}}^{(n)}(t)(\widehat{F}^{(n)}(t) \vee y, \infty) - H(\widehat{F}^{(n)}(t) \vee y)| + \widehat{\mathcal{W}}^{(n)}(t)[\widehat{C}^{(n)}(t), \widehat{F}^{(n)}(t)] \\ & = |\widehat{\mathcal{Y}}^{(n)}(t)(\widehat{F}^{(n)}(t) \vee y, \infty) - H(\widehat{F}^{(n)}(t) \vee y)| + \widehat{\mathcal{W}}^{(n)}(t)[\widehat{C}^{(n)}(t), \widehat{F}^{(n)}(t)]. \end{aligned}$$

Corollary 3.8 implies that the second term on the right-hand side has limit zero. The first also has limit zero in probability, uniformly over $y \in [y_0, y^*]$, $t \in [0, T]$, for each fixed y_0 and T , because of Proposition 3.9. However, Lemma 3.7 permits us to extend this result to uniform convergence over $y \leq y^*$ and $t \in [0, T]$. This establishes (3.33); the proof of (3.34) is similar. \square

We are now prepared to prove Theorem 3.1.

PROOF OF THEOREM 3.1. We define a mapping $\psi: \mathbb{R} \rightarrow \mathcal{A}$ by the formula

$$\psi(x)(B) \triangleq \int_{B \cap [x, \infty)} (1 - G(\eta)) d\eta, \quad \text{for } x \in \mathbb{R}, B \in \mathcal{B}(\mathbb{R}).$$

Observe that, for $x_1, x_2 \in \mathbb{R}$,

$$\sup_{B \in \mathcal{B}(\mathbb{R})} |\psi(x_1)(B) - \psi(x_2)(B)| \leq \int_{x_1 \wedge x_2}^{x_1 \vee x_2} (1 - G(\eta)) d\eta \leq |x_2 - x_1|,$$

which shows that the mapping ψ is continuous. According to Proposition 3.10,

$$\widehat{F}^{(n)} \Rightarrow F^*.$$

By the continuous mapping theorem, Theorem A.1,

$$(3.35) \quad \psi(\widehat{F}^{(n)}) \Rightarrow \psi(F^*) = \widehat{\mathcal{W}}^*.$$

On the other hand, according to Proposition 3.11,

$$(3.36) \quad \sup_{y \in \mathbb{R}} \sup_{0 \leq t \leq T} |\widehat{\mathcal{W}}^{(n)}(t)(y, \infty) - \psi(\widehat{F}^{(n)}(t))(y, \infty)| \xrightarrow{P} 0$$

[this is a rewriting of (3.33)]. Combining (3.35) and (3.36), we see that $\widehat{\mathcal{W}}^{(n)} \Rightarrow \widehat{\mathcal{W}}^*$. The proof of $\widehat{\mathcal{Q}}^{(n)} \Rightarrow \widehat{\mathcal{Q}}^*$ is analogous. \square

4. Simulation results. In this section, we use simulation to verify the predictive value of the theory of the previous sections. In the previous sections, we actually considered a sequence of queueing systems, indexed by n , whereas here we want to consider a single queueing system. We imagine that this single system is a member of the sequence of the previous sections that corresponds to a large value of n . We first recast the definitions of the previous sections in such a way that this parameter n does not appear.

Suppressing the time variable t , we recall the definitions of Section 2. We denoted the queue length in the n th system by $Q^{(n)}$ and the scaled queue length by $\widehat{Q}^{(n)} = (1/\sqrt{n})Q^{(n)}$, which, for large values of n , is approximately equal to $Q^* = \lambda W^*$ (Corollary 3.2). The workload and scaled workload, respectively, are $W^{(n)}$ and $\widehat{W}^{(n)} = (1/\sqrt{n})W^{(n)}$. The “frontier” (see Section 2 for the definition) is $F^{(n)}$ and the scaled frontier is $\widehat{F}^{(n)} = (1/\sqrt{n})F^{(n)}$. Finally, there are the measure-valued processes $\mathcal{Q}^{(n)}$ and $\widehat{\mathcal{Q}}^{(n)}$. We shall be interested particularly in $\mathcal{Q}^{(n)}(x, \infty)$, which tells us the number of customers whose lead times exceed x , and in $\widehat{\mathcal{Q}}^{(n)}(y, \infty) = (1/\sqrt{n})\mathcal{Q}^{(n)}(\sqrt{n}y, \infty)$, where we continue to suppress the time variable t .

Recall that customers arrive with lead-time distribution given by (2.1):

$$\mathbb{P}(L_j^{(n)} \leq \sqrt{n}y) = G(y).$$

We define $G_n(x) = G(x/\sqrt{n})$, so that

$$(4.1) \quad \mathbb{P}(L_j^{(n)} \leq x) = G_n(x)$$

is the cumulative distribution function of the lead times in the n th queueing system. The limit of $\widehat{\mathcal{Q}}^{(n)}$ is characterized in terms of the function H of (3.1):

$$H(y) \triangleq \int_y^\infty (1 - G(\eta)) d\eta.$$

In this section, we will need the function

$$(4.2) \quad H_n(x) = \sqrt{n} H\left(\frac{x}{\sqrt{n}}\right) = \int_x^\infty (1 - G_n(\xi)) d\xi,$$

whose inverse is $H_n^{-1}(y) = \sqrt{n}H^{-1}(y/\sqrt{n})$.

According to Theorem 3.1, for large values of n ,

$$(4.3) \quad \widehat{\mathcal{Q}}(y, \infty) \approx \lambda H(y \vee F^*).$$

Moreover, $F^* = H^{-1}(W^*) = H^{-1}(Q^*/\lambda)$. Multiplying (4.3) by \sqrt{n} and replacing y by x/\sqrt{n} , we obtain

$$(4.4) \quad \mathcal{Q}^{(n)}(x, \infty) \approx \lambda H_n(x \vee \sqrt{n}F^*).$$

Because $H_n(\sqrt{n}F^*) = \sqrt{n}H(F^*) = (\sqrt{n}/\lambda)Q^* \approx (1/\lambda)Q^{(n)}$, we also obtain

$$(4.5) \quad \sqrt{n}F^* \approx H_n^{-1}\left(\frac{1}{\lambda}Q^{(n)}\right).$$

We define

$$(4.6) \quad F_n \triangleq H_n^{-1}\left(\frac{1}{\lambda}Q^{(n)}\right),$$

so that (4.5) becomes $\sqrt{n}F^* \approx F_n$ and (4.4) becomes

$$(4.7) \quad \mathcal{D}^{(n)}(x, \infty) \approx \lambda H_n(x \vee F_n), \quad x \geq 0.$$

Note that F_n is not the frontier $F^{(n)}$ defined in Section 2. However,

$$\begin{aligned} \frac{1}{\sqrt{n}}(F_n - F^{(n)}) &= H^{-1}\left(\frac{1}{\lambda}\widehat{Q}^{(n)}\right) - \widehat{F}^{(n)} \\ &\Rightarrow H^{-1}(W^*) - F^* = 0. \end{aligned}$$

Relations (4.6) and (4.7) connect the unscaled queue length $Q^{(n)}$ with the number of customers whose unscaled lead times exceed x , and the function H_n appearing in these relations can be computed from the cumulative distribution function G_n of the unscaled lead-time distribution. These relations can be verified by simulation without knowledge of the parameter n .

The function of x appearing on the right-hand side of (4.7) is nonincreasing, with limit $Q^{(n)}$ at $-\infty$ and limit zero at ∞ . Therefore,

$$(4.8) \quad F_{\text{thy}}(x) = 1 - \frac{\lambda}{Q^{(n)}}H_n(x \vee F_n), \quad x \geq 0,$$

is a cumulative distribution function. According to (4.7), $F_{\text{thy}}(x)$ should approximate the fraction of customers in queue whose lead times are less than or equal to x . Since the parameter n is irrelevant, we henceforth omit it in our discussion of (4.8).

We present simulation results that illustrate the accuracy this approximation. In the various experiments, we simulate an M/M/1 queue using the EDF queue discipline, usually with $\lambda = 0.95$ or 0.99 , $\mu = 1.0$, $\rho = 0.95$ or 0.99 . (The theory applies to GI/G/1 queues, but only M/M/1 systems are simulated in this section. Limited experience with the simulation of GI/G/1 queues suggests that the accuracy of the approximations in this section are representative of more general systems.) According to the theory developed in Section 3, if we were randomly to stop the simulation at any point, observe the current number in the queue, Q , and find that Q is sufficiently large, then the corresponding instantaneous lead-time profile, expressed as an empirical cdf, should be given approximately by (4.8).

For real-time queueing theory to be useful in practice, it is important that it can be applied in cases in which the queue length Q is moderate in size. However, when Q is moderate, we expect the lead-time profiles to exhibit substantial variability, and it is not at all clear that the asymptotic form given by (4.8) is appropriate. The simulations presented in this section are designed to address this issue.

For each simulation run, a particular deadline distribution, $G = G_n$, and queue length, $Q = Q^{(n)}$, is chosen. The run is initiated with an empty queue

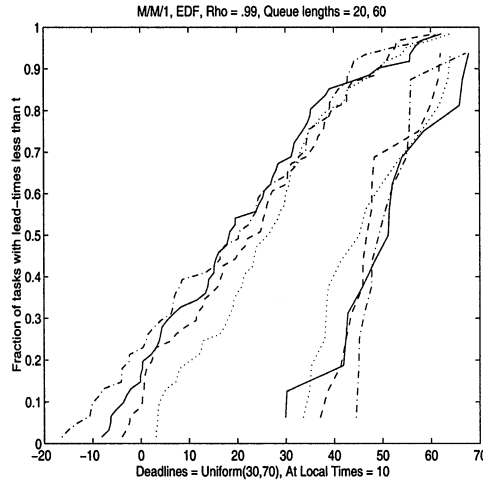


FIG. 1. Sample lead-time profiles, $Q = 20$, $Q = 60$.

and continues until the instant the local time at level Q reaches a prespecified value, 10 for the results presented in this paper. At that instant, the lead-time profile is recorded. This same experiment is repeated a total of N times; hence, the N profiles can be thought of as independent random objects. We wish to assess how close they are to their predicted form given by (4.8).

The deadline distribution G used in Figures 1–5 is a $\text{Uniform}(30, 70)$ distribution. Figure 1 illustrates the variability in the lead-time profiles for small to moderate values of Q . This figure shows the first four lead-time profiles recorded when $Q = 20$ (then when $Q = 60$) and the accumulated local time is 10. The lead-time profiles are actually 20- (or 60-) dimensional vectors, but are plotted here as line segments for visual convenience. Notice that although the profiles have similar shapes, they exhibit substantial variability.

4.1. *Uniform deadlines.* The first deadline distribution considered was the $\text{Uniform}(A, B)$. For this distribution, the frontier $F = F_n$ is given by

$$F = \begin{cases} B - \sqrt{2W(B - A)}, & \text{if } W \leq \frac{B - A}{2}, \\ \frac{B + A}{2} - W, & \text{if } W > \frac{B - A}{2}, \end{cases}$$

where throughout we use the notation $W = Q/\lambda$. The theoretical cdf defined by (4.8) for this $\text{Uniform}(A, B)$ deadline case is given by one of two forms depending on the magnitude of W . If $W \geq (B - A)/2$, then

$$F_{\text{thy}}(x) = \begin{cases} 0, & \text{if } x \leq F, \\ 1 - \frac{1}{W} \left(\frac{A + B}{2} - x \right), & \text{if } F < x < A, \\ 1 - \frac{(B - x)^2}{2W(B - A)}, & \text{if } A \leq x \leq B, \end{cases}$$

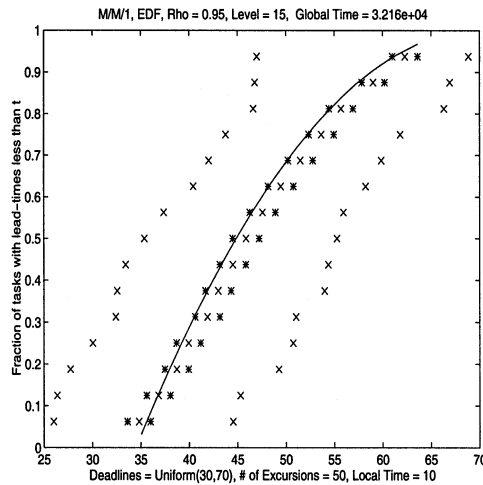


FIG. 2. Profiles: mean, max, min and theory, $Q = 15, N = 50$.

whereas if $W < (B - A)/2$, then

$$F_{\text{thy}}(x) = \begin{cases} 0, & \text{if } x \leq F, \\ 1 - \left(\frac{B - x}{B - F}\right)^2, & \text{if } F < x \leq B, \\ 1, & \text{if } B \leq x. \end{cases}$$

If we substitute the simulated customer lead times into F_{thy} , the resulting empirical cdf should correspond to a Uniform(0, 1). In addition, from F_{thy} the quantiles of the lead-time distribution can be determined by solving $F_{\text{thy}}^{-1}(p) = x_p$ for $0 < p < 1$. For $W \geq (B - A)/2$ this gives

$$x_p = \begin{cases} F + pW, & \text{if } 0 < p < 1 - \frac{B - A}{2W}, \\ B - \sqrt{2(1 - p)W(B - A)}, & \text{if } 1 - \frac{B - A}{2W} < p < 1, \end{cases}$$

whereas for $W < (B - A)/2$ we have

$$x_p = B - (B - F)\sqrt{1 - p}, \quad 0 < p < 1.$$

The theoretical profiles for any particular Q and G in all the following figures are obtained by connecting the points $\{(x_p, p), p = 1/(2Q + 1), \dots, 2Q/(2Q + 1)\}$.

Interestingly, in spite of the substantial variability in each profile, if we average those profiles (by averaging each of the components of the N distinct Q -dimensional vectors), the result is a very smooth profile that is nearly identical to the theoretical cdf given by (4.8). Figures 2–5 show the mean profile, the componentwise minimum profile and the componentwise maximum profile for the cases $Q = 15, 20, 40$ and 60 , each component of which is denoted by an “x”. The componentwise minimum and maximum form an envelop for all N

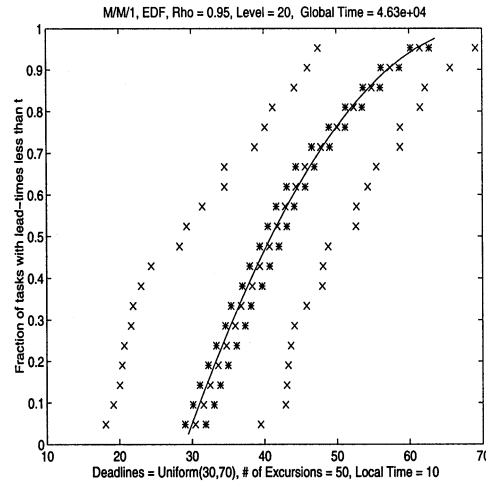


FIG. 3. Profiles: mean, max, min and theory, $Q = 20$, $N = 50$.

profiles generated by the particular simulation run. In Figures 2 and 3, 95% confidence intervals for the mean are determined for each of the Q quantiles. These confidence limits are denoted by “*” and have been constructed independently for each of the Q components. In subsequent figures, the confidence interval indications are omitted because they become visually distracting. In all cases, the traffic intensity is 0.95 and the profile is recorded when exactly 10 units of local time at the indicated queue length have elapsed. The global time at which the last lead-time profile is taken is recorded at the top of each figure.

The profile is approximately correct for $Q = 15$, but there are systematic departures evident. Figures 3–5 suggest that the profile form given by (4.8) is nearly exact as a mean value for $Q \geq 20$, because the theoretical curve is always within the 95% confidence limits.

In considering the behavior of a real-time queueing system, it is important to put bounds or confidence sets around the profiles described by (4.8), which capture a large fraction of profiles. Such bounds could be used to determine access control policies that would prevent customer lateness (at the expense of losing customers through admission denial). The minimum and maximum curves offer some idea of how wide such profiles must be and how wide the confidence regions must be. Presumably, the bounds can be constructed using large deviation theory; however, this is not studied any further in this paper. Although the minimum and maximum limits seem fairly wide, we expect that nearly all of the empirical profiles will be within an $O(1)$ distance from the theoretical profile in which lead times are $O(\sqrt{n})$.

These plots are representative of many such plots for a variety of bounded deadline distributions and traffic intensities. The greater the variability in G , the larger Q must be for the average profiles to agree with (4.8).

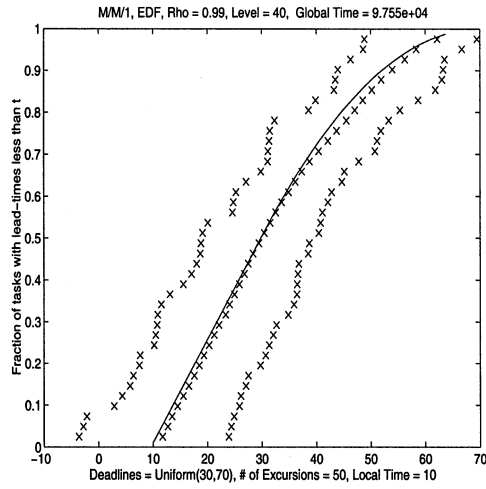


FIG. 4. Profiles: mean, max, min and theory, $Q = 40$, $N = 50$.

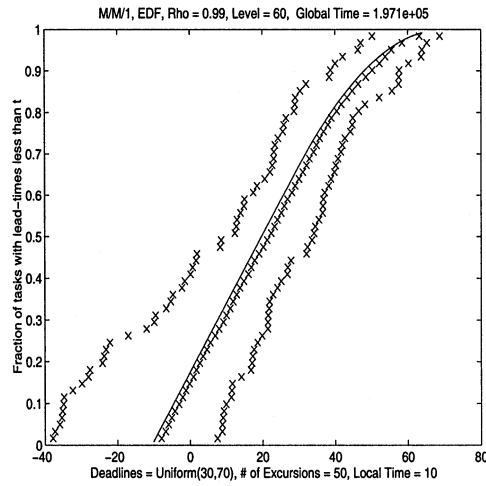


FIG. 5. Profiles: mean, max, min and theory, $Q = 60$, $N = 50$.

The accumulated local time at which the profiles are recorded also has important consequences. These issues are both addressed in the next section

5. Future research. In this section, we introduce additional simulation results that illustrate potential generalizations of the real-time queueing theory developed in this paper and some additional issues.

5.1. More general deadline distributions. The theory developed in Sections 3 and 4 used the assumption that the deadline distribution was bounded above by some finite constant. Here, we illustrate that this assumption appears to be unnecessary. Two distributions that are not bounded above are considered: the exponential(α) distribution with mean $1/\alpha$ and the Pareto(α, B) distribution.

We begin with the exponential(α) distribution with mean $1/\alpha$. All moments of this distribution are finite, but it is not bounded above as in the uniform case. For this distribution, the frontier is given by

$$F = \begin{cases} -\frac{1}{\alpha} \log(\alpha W), & \text{if } W \leq \frac{1}{\alpha}, \\ \frac{1}{\alpha} - W, & \text{if } W > \frac{1}{\alpha}. \end{cases}$$

The theoretical cdf for the lead-time profiles for this exponential(α) case takes on one of two forms depending upon whether $W \geq 1/\alpha$ or not. For the case in which $W \geq 1/\alpha$, we have

$$F_{\text{thy}}(x) = \begin{cases} 0, & \text{if } x \leq F, \\ 1 - \frac{1}{\alpha W} (1 - \alpha x), & \text{if } F < x < 1, \\ 1 - \frac{1}{\alpha W} \exp(-\alpha x), & \text{if } x \geq 1, \end{cases}$$

whereas for $W < 1/\alpha$ we have

$$F_{\text{thy}}(x) = \begin{cases} 0, & \text{if } x \leq F, \\ 1 - \frac{1}{\alpha W} \exp(-\alpha x), & \text{if } x > F. \end{cases}$$

For $W \geq 1/\alpha$, the quantiles are given by

$$x_p = \begin{cases} \frac{1}{\alpha} - W(1 - p), & \text{if } 0 < p < 1 - \frac{1}{\alpha W}, \\ -\frac{1}{\alpha} \log(\alpha W(1 - p)), & \text{if } 1 - \frac{1}{\alpha W} \leq p < 1, \end{cases}$$

whereas for $W < 1/\alpha$ they are given by

$$x_p = -\frac{1}{\alpha} \log((1 - p)\alpha W), \quad 0 < p < 1.$$

For the simulation, we chose $\alpha = 0.02$, giving a mean of 50. We again simulate 50 independent profiles taken when 10 units of local time have been reached. Figure 6 shows the mean profile for $Q = 20$. Again, the shape is generally correct, but systematic departures are evident. Because this deadline distribution will result in a few customers in the queue having very large lead times, it is more informative to use Q-Q plots to judge the agreement

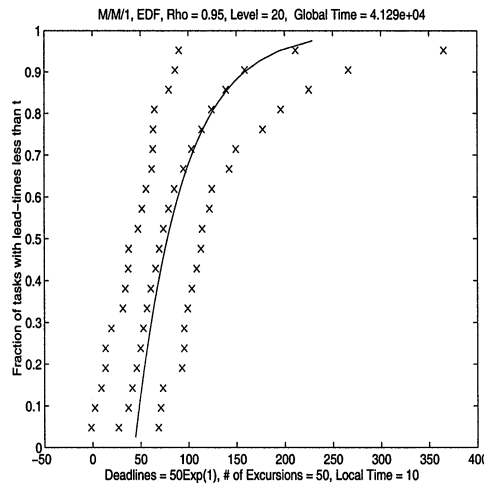


FIG. 6. Profiles: mean, max, min and theory, $Q = 20, N = 50$.

between the empirical and the theoretical cdf. A Q-Q plot is obtained by plotting $\{(F_{thy}(L_i), i/(Q + 1)), 1 \leq i \leq Q\}$. If the lead times (L_1, \dots, L_Q) are a random sample from F_{thy} , then these points should lie close to a 45° line connecting $(0, 0)$ with $(1, 1)$. Figure 7 gives the Q-Q plot that corresponds to Figure 6. Figure 8 presents the Q-Q plot for the same exponential deadline case when $Q = 40$. When $Q = 20$, systematic departures between the average empirical cdf and the theoretical cdf are evident, especially in the left-hand tail. Nevertheless, when $Q = 40$, the agreement is nearly exact.

We next consider a Pareto(α, B) deadline distribution. This distribution is characterized by the cdf

$$G(x) = \begin{cases} 0, & \text{if } \frac{x}{B} < 1, \\ 1 - \left(\frac{B}{x}\right)^{(\alpha-1)}, & \text{if } \frac{x}{B} \geq 1, \end{cases}$$

for $B > 0$ and $\alpha > 1$. This distribution has no moments of order $\alpha - 1$ or higher; hence, it has a very heavy right-hand tail. Indeed, for $1 < \alpha \leq 2$, the function $H(y) = \infty$ for all finite y , and the proposed lead-time profile given by (4.8) does not exist. Nevertheless, we show simulation results for $\alpha = 3$ and $\alpha = 6$ which demonstrate that there is a stable lead-time profile associated with this family of distributions for $\alpha > 2$. For $\alpha > 2$, the frontier is given by

$$F = \begin{cases} B\left(\frac{B}{(\alpha - 2)W}\right)^{1/(\alpha-2)}, & \text{if } W \leq \frac{B}{\alpha - 2}, \\ \frac{\alpha - 1}{\alpha - 2}B - W, & \text{if } W > \frac{B}{\alpha - 2}. \end{cases}$$

The theoretical cdf for this Pareto(α, B) case is given by one of two forms depending upon whether $W \geq B/(\alpha - 2)$ or not. For the case in which

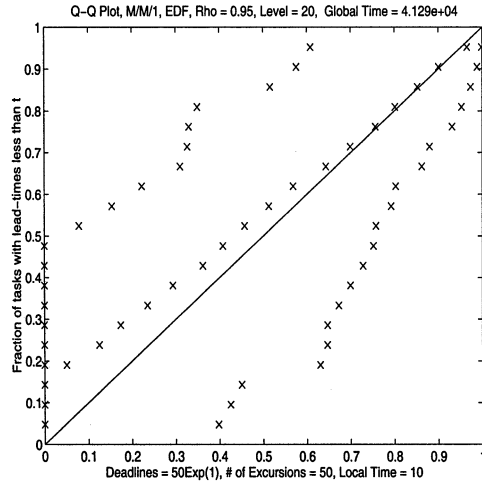


FIG. 7. *Q-Q profiles: mean, max, min and theory, Q = 20, N = 50.*

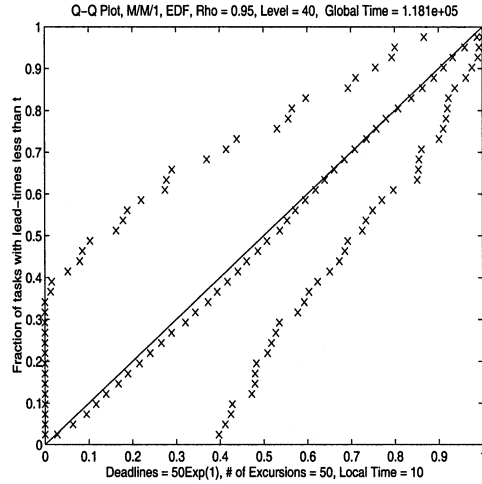


FIG. 8. *Q-Q profiles: mean, max, min and theory, Q = 40, N = 50.*

$W \geq B/(\alpha - 2)$, we have

$$F_{\text{thy}}(x) = \begin{cases} 0, & \text{if } x \leq F, \\ 1 - \frac{1}{W} \left(\frac{\alpha - 1}{\alpha - 2} B - x \right), & \text{if } F \leq x < B, \\ 1 - \frac{B}{W(\alpha - 2)} \left(\frac{B}{x} \right)^{\alpha - 2}, & \text{if } x \geq B, \end{cases}$$

whereas for $W < B/(\alpha - 2)$ we have

$$F_{\text{thy}}(x) = \begin{cases} 0, & \text{if } x \leq F, \\ 1 - \frac{B}{W(\alpha - 2)} \left(\frac{B}{x} \right)^{\alpha - 2}, & \text{if } x \geq F. \end{cases}$$

For $W \geq B/(\alpha - 2)$ the quantiles are given by

$$x_p = \begin{cases} \frac{\alpha - 1}{\alpha - 2} B - (1 - p)W, & \text{if } 0 < p < 1 - \frac{B}{(\alpha - 2)W}, \\ B \left(\frac{B}{(\alpha - 2)W(1 - p)} \right)^{1/(\alpha - 2)}, & \text{if } 1 - \frac{B}{(\alpha - 2)W} < p < 1, \end{cases}$$

whereas for $W < B/(\alpha - 2)$ they are given by

$$x_p = B \left(\frac{B}{(\alpha - 2)W(1 - p)} \right)^{1/(\alpha - 2)}, \quad 0 < p < 1.$$

Figures 9 and 10 give Q-Q plots for the Pareto(6, 40) distribution, a distribution with mean 50 and a relatively heavy tail. The traffic intensity was increased to 0.99 and the number of profiles was increased from 50 to 100. Figure 9 corresponds to $Q = 20$, while Figure 10 corresponds to $Q = 40$. The variability in the Q-Q plots is quite large; however, there is good agreement between the average profile and the theoretical distributions. For $Q = 20$, there are some systematic departures in the tails, but for $Q = 40$, the agreement is very good except in the upper tail.

Figures 11 and 12 present results for the more extreme Pareto(3, 25), a distribution with mean 50, but infinite variance. This distribution has an extremely heavy right-hand tail, and very long lead times will occur as Q increases. For this case, we considered longer queue lengths, with Figure 11 corresponding to $Q = 40$ and Figure 12 corresponding to $Q = 60$. Figure 11 shows systematic departures in the left-hand tail, while Figure 12 shows near exact agreement between the average profile and the theoretical distribution.

5.2. FIFO queue discipline. This paper studied the EDF queue discipline; however, the results can be used heuristically to determine lead-time profiles for the behavior of FIFO queues. Suppose that arriving customers have deadlines given by distribution G , but are serviced in FIFO order. This queue discipline does not require knowledge of the customer deadlines, and a customer's instantaneous lead time is equal to its initial deadline minus its time in the queue. The time in the queue can be determined using (4.8) first by assuming all customers have deadlines 0. If all customers have deadline 0, then

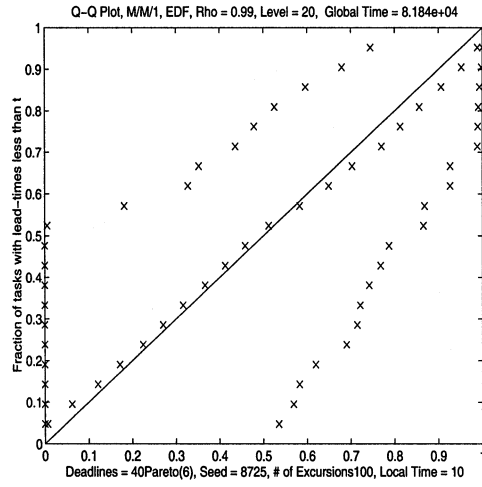


FIG. 9. Q-Q profiles: mean, max, min and theory $Q = 20$, $N = 100$.

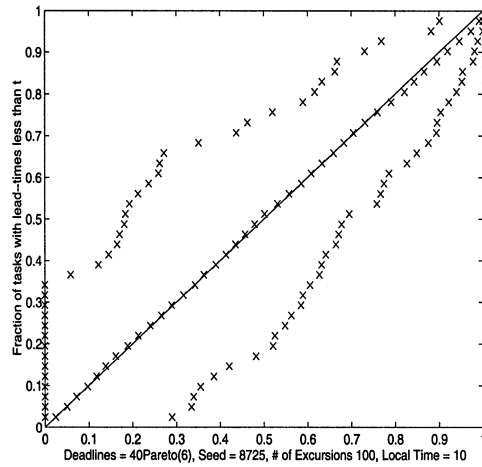


FIG. 10. Profiles: mean, max, min and theory, $Q = 40$, $N = 100$.

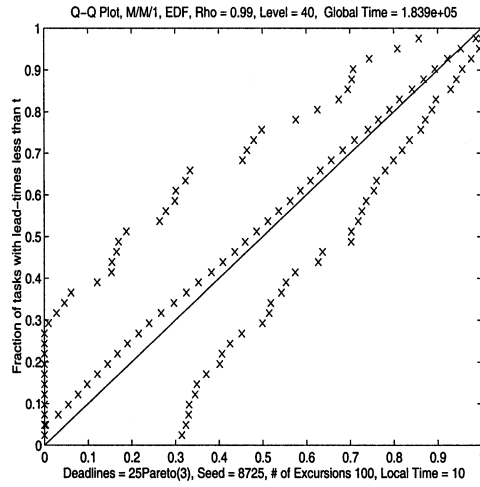


FIG. 11. Profiles: mean, max, min and theory, $Q = 40$, $N = 100$.

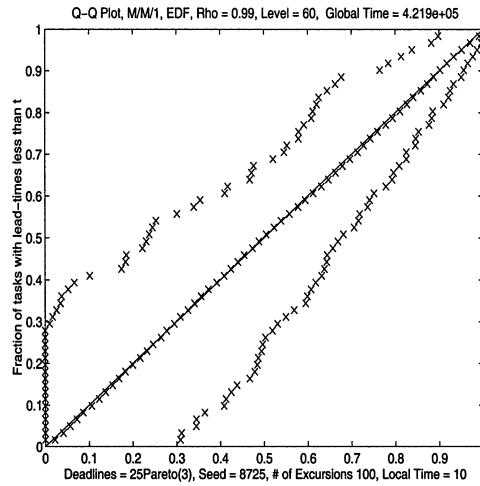


FIG. 12. Profiles: mean, max, min and theory, $Q = 60$, $N = 100$.

the EDF queue discipline is equivalent to FIFO, and any customer’s instantaneous lead time is equal to the negative of its time in the queue. Profiles of customers’ times in the queue can be approximated by using (4.8) with deadline distribution corresponding to point mass at 0. In this case, the resulting distribution will be $\text{Uniform}(-W, 0)$. By adding back their actual deadlines to their time in queue, we can recover the customer lead times. Consequently, if we were to order the customers in a FIFO queue by lead time (the FIFO ordering is by time in the queue), then the resulting lead-time profile should be the convolution of G with a $\text{Uniform}(-W, 0)$ distribution.

Figures 13 and 15 illustrate the lead-time profiles for a FIFO queue assuming G is exponential(1/50), $Q = 20$ and $Q = 40$. Figures 14 and 16 present the corresponding Q-Q plots for more accurate assessment of agreement. In this case, the frontier is given by $F = -W = -Q/\lambda$. The theoretical cdf is given by

$$F_{\text{thy}}(x) = \begin{cases} 0, & \text{if } x < -W, \\ 1 + \frac{x}{W} + \frac{1}{\alpha W}(e^{-\alpha(x+W)} - 1), & \text{if } -W < x < 0, \\ 1 - \frac{1 - e^{-\alpha W}}{\alpha W}e^{-\alpha x}, & \text{if } 0 \leq x. \end{cases}$$

There are two different cases associated with finding the quantiles for this distribution. First, when $0 < p < 1 - (1 - e^{-\alpha W})/\alpha W$, then $x_p = y/\alpha - W$, where y is the solution of the equation

$$y + e^{-y} = 1 + \alpha W p.$$

For $1 - (1 - e^{-\alpha W})/\alpha W \leq p < 1$,

$$x_p = -\frac{1}{\alpha} \log\left(\frac{(1 - p)\alpha W}{1 - e^{-\alpha W}}\right).$$

Again, the agreement is reasonable for $Q = 20$ and excellent for $Q = 40$.

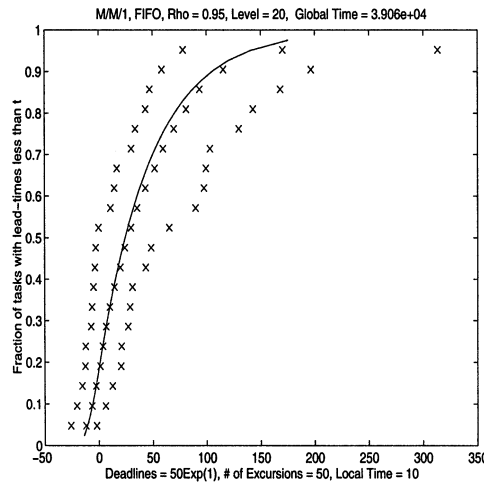


FIG. 13. Profiles: mean, max, min and theory, $Q = 20$, $N = 50$.

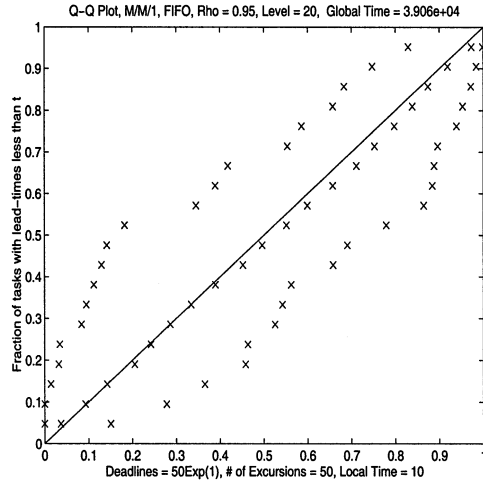


FIG. 14. *Q-Q profiles: mean, max, min and theory, $Q = 20, N = 50$.*

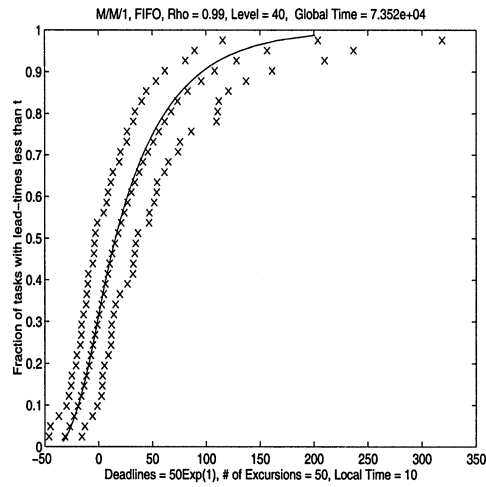


FIG. 15. *Profiles: mean, max, min and theory, $Q = 40, N = 50$.*

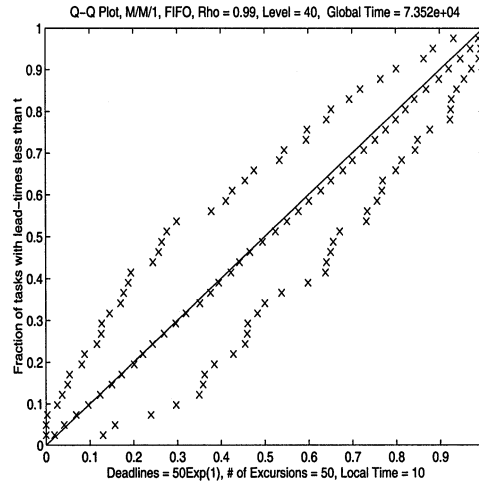


FIG. 16. *Q-Q profiles: mean, max, min and theory, $Q = 40$, $N = 50$.*

5.3. *Profiles at hitting times.* Real-time queueing theory should be useful for developing and analyzing control policies to reduce or eliminate customer lateness. One simple policy would be a threshold policy in which arrivals would be denied admission and be lost if Q reached some specified level. In principle, the threshold could be chosen based on the predicted lead-time profiles, for example, by choosing the threshold so that the frontier is bounded away from 0 by some confidence margin. This analysis would use the profile when Q first hit the threshold. Interestingly, the profile obtained at the time at which a level is first hit can be systematically different from the profile predicted by (4.8). Figure 17 gives a representative example. The simulation parameters are identical to those presented in Figure 5 except that the profile is recorded at the instant that $Q = 60$ rather than after 10 units of local time have been accumulated at level 60. The mean empirical cdf curve is shifted away from the theoretical profile. The introduction of the hitting time, which guarantees that the queue length has never exceeded the level in question, creates systematic distortions in the profiles. It will be important to develop corrections to (4.8) that incorporate this stopping time bias.

5.4. *The non-heavy-traffic case.* The theory presented in Section 3 was developed assuming the traffic intensity approaches 1 as $n \rightarrow \infty$ [see (2.9)]. Interestingly, suppose we do not assume this heavy-traffic condition, but simulate an EDF system with moderate traffic intensity. If the simulation is continued until a suitable amount of local time (say 10 units) at a large enough queue level (say $Q = 40$) is obtained and the lead-time profile is recorded, then that profile will be essentially indistinguishable from those recorded under heavy-traffic conditions. Hence, it appears that for the EDF queue discipline,

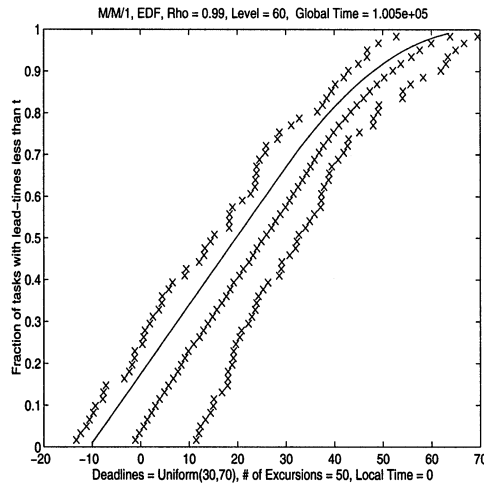


FIG. 17. Q - Q profiles: mean, max, min and theory, $Q = 60$, $N = 50$.

heavy-traffic theory calculations can be used to make accurate lead-time profile predictions under moderate traffic conditions. Of course, if the traffic intensity is moderate, long queue lengths are relatively infrequent compared with heavy-traffic conditions. Much more research is needed to assess the accuracy of the conjectures presented in this section. Nevertheless, it appears that the heavy-traffic approximations developed for real-time queues will have many important applications under nonextreme traffic conditions.

APPENDIX A

Weak convergence. The following standard results can be found in or are easily derived from assertions found in Parthasarathy ([23], Chapter II) and Billingsley ([1], Section 17). In this Appendix, we state a version of these results needed for this paper.

Let S be a separable metric space and let $\mathcal{M}(S)$ be the set of finite measures defined on the σ -algebra of Borel subsets of S . We endow $\mathcal{M}(S)$ with the *weak topology*, whereby a sequence of finite measures $\{\mu_n\}_{n=1}^\infty$ converges to a finite measure μ if and only if $\lim_{n \rightarrow \infty} \int_S g d\mu_n = \int_S g d\mu$ for every bounded, continuous function g mapping S into \mathbb{R} . The weak topology on $\mathcal{M}(S)$ is metrizable and $\mathcal{M}(S)$ is locally compact.

Now let $\{X_n\}_{n=1}^\infty$ be a sequence of S -valued random objects defined on respective probability spaces $(\Omega_n, \mathcal{F}_n, \mathbb{P}_n)$, which may depend on n , and let X be an S -valued random object defined on a probability space $(\Omega^*, \mathcal{F}^*, \mathbb{P}^*)$. We say X_n *converges weakly* to X and we write $X_n \Rightarrow X$ if the sequence of probability measures μ_n induced on S by X_n converges weakly to the probability measure induced on S by X .

THEOREM A.1 (Continuous mapping theorem). *Let $\{X_n\}_{n=1}^\infty$ be a sequence of S -valued random objects converging weakly to another S -valued random object X . Let $f: S \rightarrow U$ be a measurable function from S to another metric space U and assume f is continuous on the support of X . Then $f(X_n) \Rightarrow f(X)$.*

Let (S, ρ) be a locally compact separable metric space and let $T > 0$ be given. A separable metric space that shall concern us is $D_S[0, T]$, the space of right-continuous functions with left-hand limits (hereafter called RCLL functions) from $[0, T]$ to S , equipped with the Skorohod metric

$$d_T(x, y) = \inf_{\lambda} \left\{ \sup_{0 \leq t \leq T} \rho(x(\lambda(t)), y(t)) + \sup_{0 \leq t \leq T} |\lambda(t) - t| \right\}, \quad x, y \in D_S[0, T],$$

where the infimum is over all strictly increasing functions λ mapping $[0, T]$ onto itself.

In this paper, most processes are in fact defined on $[0, \infty)$. The space $D_S[0, \infty)$ of RCLL S -valued functions defined on $[0, \infty)$ has a metric d_∞ with the property that whenever x and y are in $D_S[0, T]$ and their restrictions $x|_{[0, T]}$ and $y|_{[0, T]}$ to $[0, T]$ agree, then $d_\infty(x, y) \leq e^{-T}$ ([4], Chapter 3, Section 5). If $\{x_n\}_{n=1}^\infty$ is a sequence in $D_S[0, \infty)$, x is another function in $D_S[0, \infty)$, and for every $T > 0$ the sequence of restrictions $\{x_n|_{[0, T]}\}_{n=1}^\infty$ converges in $D_S[0, T]$ to $x|_{[0, T]}$, then $\{x_n\}_{n=1}^\infty$ converges in $D_S[0, \infty)$ to x . The converse holds if x is continuous.

Now let $\{X_n(t); 0 \leq t \leq T\}_{n=1}^\infty$ be a sequence of RCLL S -valued processes defined on $[0, T]$. These induce a sequence of measures on $D_S[0, T]$. If this sequence converges weakly to the measure induced by another RCLL S -valued process $\{X(t); 0 \leq t \leq T\}$, then we say that the sequence of processes $\{X_n\}_{n=1}^\infty$ converges weakly to the process X and write $X_n \Rightarrow X$. The definition of weak convergence of a sequence of RCLL S -valued processes on $[0, \infty)$ is similar. Such a sequence converges weakly to a continuous process $\{X(t); 0 \leq t < \infty\}$ if and only if, for every $T > 0$, the sequence of restricted processes $\{X_n(t); 0 \leq t \leq T\}$ converges weakly to the restricted process $\{X(t); 0 \leq t \leq T\}$.

THEOREM A.2 (Time change theorem). *Suppose the sequence of RCLL S -valued processes $\{X_n(t); 0 \leq t < \infty\}_{n=1}^\infty$ converges weakly to a continuous, S -valued process $\{X(t); 0 \leq t < \infty\}$. Suppose further that the sequence of RCLL $[0, \infty)$ -valued processes $\{\Phi_n(t); 0 \leq t < \infty\}_{n=1}^\infty$ converges weakly to a nonrandom continuous $[0, \infty)$ -valued process $\{\Phi(t); 0 \leq t < \infty\}$. Then*

$$X_n \circ \Phi_n \Rightarrow X \circ \Phi.$$

THEOREM A.3 (Differencing theorem). *Suppose the sequence of RCLL S -valued processes $\{X_n(t); 0 \leq t < \infty\}_{n=1}^\infty$ converges weakly to a continuous, S -valued process $\{X(t); 0 \leq t < \infty\}$. Suppose further that the sequences of RCLL $[0, \infty)$ -valued processes $\{\Phi_n(t); 0 \leq t < \infty\}_{n=1}^\infty$ and $\{\Psi_n(t); 0 \leq t < \infty\}_{n=1}^\infty$ converge weakly to the identically zero process. Then the sequence*

of processes

$$Y_n(t) \triangleq \rho(X_n(t + \Phi_n(t)), X_n(t + \Psi_n(t)))$$

converges weakly to the identically zero process.

APPENDIX B

Functional central limit theorem. This Appendix summarizes classical heavy-traffic limit results for a sequence of queues. It is included here primarily to establish notation for the main body of the paper. Recall the definitions $S_0^{(n)} \triangleq 0$ and for $k \geq 1$, $S_k^{(n)} \triangleq \sum_{j=1}^k u_j^{(n)}$, where for each n , $\{u_j^{(n)}\}_{j=1}^\infty$ is a sequence of independent, identically distributed strictly positive random variables with mean $1/\lambda^{(n)}$ and standard deviation $\alpha^{(n)}$. In the n th queue, $S_k^{(n)}$ is the arrival time of the k th customer. The number of customers arrived by time t is $A^{(n)}(t) \triangleq \max\{k \geq 0; S_k^{(n)} \leq t\}$. We define the *centered and scaled arrival process*

$$\widehat{A}^{(n)}(t) \triangleq \frac{1}{\sqrt{n}}[A^{(n)}(nt) - \lambda^{(n)}nt], \quad t \geq 0.$$

Recall also the definition of the *centered and scaled work arrival process*

$$\widehat{V}^{(n)}(t) \triangleq \frac{1}{\sqrt{n}} \sum_{j=1}^{\lfloor nt \rfloor} \left(v_j^{(n)} - \frac{1}{\mu^{(n)}} \right),$$

where for each n , $\{v_j^{(n)}\}_{j=1}^\infty$ is a sequence of independent, identically distributed random variables with mean $\mu^{(n)}$ and variance $\beta^{(n)}$. The work arrival process is

$$V^{(n)}(t) \triangleq \sum_{j=1}^{\lfloor nt \rfloor} v_j^{(n)},$$

and the centered and scaled netput process is

$$\widehat{N}^{(n)}(t) \triangleq \frac{1}{\sqrt{n}}[V^{(n)}(A^{(n)}(nt)) - nt].$$

We impose the heavy-traffic assumptions (2.9)–(2.11), which are in force throughout.

Suppose B is a standard Brownian motion, and μ and σ are constants. Then $B^*(t) = \mu t + \sigma B(t)$ is a Brownian motion with drift μ and variance σ^2 per unit time. We denote this by writing $B^* \sim BM(\mu, \sigma^2)$. The following theorems are consequences of Prohorov ([25], Theorem 3.1), used to extend Billingsley ([1], Section 17.3).

THEOREM B.1. *The sequence of processes $\{\widehat{V}^{(n)}\}_{n=1}^{\infty}$ converges weakly to a process $V^* \sim BM(0, \beta^2)$.*

THEOREM B.2. *The sequence of processes $\{\widehat{A}^{(n)}\}_{n=1}^{\infty}$ converges weakly to a process $A^* \sim BM(0, \alpha^2 \lambda^3)$.*

THEOREM B.3. *The sequence $\widehat{N}^{(n)}$ converges weakly to $(1/\lambda)A^* + V^* \circ \lambda e - \gamma t$, where $A^* \sim B(0, \alpha^2 \lambda^3)$, $V^* \sim B(0, \beta^2)$, A^* and V^* are independent, and e is the identity function $e(t) = t$ for all $t \in [0, 1]$.*

COROLLARY B.4. *Let $N^* = (1/\lambda)A^* + V^* \circ \lambda e - \gamma t$ be the Brownian motion with drift in Theorem B.3, and define*

$$I^*(t) \triangleq - \min_{0 \leq s \leq t} N^*(s),$$

$$W^*(t) \triangleq N^*(t) + I^*(t).$$

Then

$$(\widehat{N}^{(n)}, \widehat{I}^{(n)}, \widehat{W}^{(n)}) \Rightarrow (N^*, I^*, W^*).$$

Acknowledgments. We wish to acknowledge an initial discussion in 1995 with Frank Kelly of Cambridge University and David Aldous of University of California, Berkeley. Frank Kelly first produced the formula (4.8), and David Aldous provided an equilibrium analysis proof in an unpublished manuscript. Also, in 1997, Martin Reiman of Lucent Technologies provided insights into how to use heavy-traffic methodology (as opposed to Markov methods) to analyze real-time queues that led directly to this paper. Finally, we acknowledge an anonymous referee, who pointed out an error in an earlier version of this paper and suggested a simplification of Theorem 3.1.

REFERENCES

- [1] BILLINGSLEY, P. (1968). *Convergence of Probability Measures*. Wiley, New York.
- [2] COFFMAN, E. G., JR., PUHALSKII, A. A. and REIMAN, M. I. (1995). Polling systems with zero switchover times: A heavy traffic-traffic averaging principle. *Ann. Appl. Probab.* **5** 681–719.
- [3] DOYTCHINOV, B. (1997). Heavy traffic limits of queues with due dates. Ph.D. dissertation, Dept. Mathematical Sciences, Carnegie Mellon Univ.
- [4] ETHIER, S. N. and KURTZ, T. G. (1986). *Markov Processes. Characterization and Convergence*. Wiley, New York.
- [5] HARRISON, J. M. (1973). A limit theorem for priority queues in heavy traffic. *J. Appl. Probab.* **10** 907–912.
- [6] HARRISON, J. M. (1985). *Brownian Motion and Stochastic Flow Systems*. Wiley, New York.
- [7] HARRISON, J. M. (1988). Brownian models of queueing networks with heterogeneous customer populations. In *Proceedings of the IMA Workshop on Stochastic Differential Systems* 147–186. Springer, New York.

- [8] HARRISON, J. M. and NGUYEN, V. (1993). Brownian models of multiclass queueing networks: Current status and open problems. *Queueing Systems Theory Appl.* **13** 5–40.
- [9] HARRISON, J. M. and WEIN, L. J. (1990). Scheduling networks of queues: Heavy traffic analysis of a two-station closed network. *Oper. Res.* **38** 1052–1064.
- [10] HONG, J., TAN, X. and TOWSLEY, D. (1989). A performance analysis of minimum laxity and earliest deadline scheduling in a real-time system. *IEEE Trans. Comput.* **38** 1736–1744.
- [11] HOOKE, J. (1970). On some limit theorems for the GI/G/1 queue. *J. Appl. Probab.* **7** 634–640.
- [12] IGLEHART, D. and WHITT, W. (1970). Multiple channel queues in heavy traffic I. *Adv. in Appl. Probab.* **2** 150–177.
- [13] IGLEHART, D. and WHITT, W. (1970). Multiple channel queues in heavy traffic II. *Adv. in Appl. Probab.* **2** 355–364.
- [14] KINGMAN, J. F. C. (1961). The single server queue in heavy traffic. *Proceedings of the Cambridge Philosophical Society* **48** 277–289.
- [15] KLEIN, M., RALYA, T., POLLAK, B., OBENZA, R. and GONZALEZ-HARBOUR, M. (1993). *A Practitioner's Handbook for Real-Time Analysis*. Kluwer, Dordrecht.
- [16] KYPRIANOU, E. (1971). The virtual waiting time of the GI/G/1 queue in heavy traffic. *Adv. in Appl. Probab.* **3** 249–268.
- [17] LEHOCZKY, J. P. (1996). Real-time queueing theory. In *Proceedings of the IEEE Real-Time Systems Symposium* 186–195. IEEE, New York.
- [18] LEHOCZKY, J. P. (1997). Using real-time queueing theory to control lateness in real-time systems. *Performance Evaluation Review* **25** 158–168.
- [19] LEHOCZKY, J. P. (1997). Real-time queueing network theory. In *Proceedings of the IEEE Real-Time Systems Symposium* 58–67. IEEE, New York.
- [20] LIU, C. L. and LAYLAND, J. W. (1973). Scheduling algorithms for multiprogramming in a hard real-time environment. *J. Assoc. Comput. Mach.* **20**(1) 40–61.
- [21] MARKOWITZ, D. M. and WEIN, L. M. (1996). Heavy traffic analysis of dynamic cyclic policies: A unified treatment of the single machine scheduling problem. Preprint, Sloan School of Management, MIT.
- [22] PANWAR, S. and TOWSLEY, D. (1988). On the optimality of the STE rule for multiple server queues that serve customers with deadlines. Technical Report 88-81, Dept. Computer and Information Science, Univ. Massachusetts, Amherst.
- [23] PARTHASARATHY, K. R. (1967). *Probability Measures on Metric Spaces*, Academic Press, New York.
- [24] PETERSON, W. P. (1991). A heavy traffic limit theorem for networks of queues with multiple customer types. *Math. Oper. Res.* **16** 90–118.
- [25] PROKHOROV, YU. (1956). Convergence of random processes and limit theorems in probability theory. *Theory Probab. Appl.* **1** 157–214.
- [26] REIMAN, M. (1983). Some diffusion approximations with state space collapse. In *Proceedings of the International Seminar on Modeling and Performance Evaluation Methodology*. Springer, Berlin.
- [27] REIMAN, M. (1988). A multiclass feedback queue in heavy traffic. *Adv. in Appl. Probab.* **20** 179–207.
- [28] VAN MIEGHEM, J. A. (1995). Dynamic scheduling with convex delay costs: The generalized $c\mu$ rule. *Ann. Appl. Probab.* **5** 809–833.
- [29] WEIN, L. J. (1990). Scheduling networks of queues: Heavy traffic analysis of a two-station network with controllable inputs. *Oper. Res.* **38** 1065–1078.
- [30] WHITT, W. (1971). Weak convergence theorems for priority queues: Preemptive-resume discipline. *J. Appl. Probab.* **8** 74–94.

B. DOYTCHINOV
DEPARTMENT OF MATHEMATICAL SCIENCES
WORCESTER POLYTECHNIC INSTITUTE
WORCESTER, MASSACHUSETTS 01609
E-MAIL: bogdand@wpi.edu

J. LEHOCZKY
DEPARTMENT OF STATISTICS
CARNEGIE MELLON UNIVERSITY
PITTSBURGH, PENNSYLVANIA 15213
E-MAIL: jpl@stat.cmu.edu

S. SHREVE
DEPARTMENT OF MATHEMATICAL SCIENCES
CARNEGIE MELLON UNIVERSITY
PITTSBURGH, PENNSYLVANIA 15213
E-MAIL: shreve@cmu.edu