

## MULTIPLE-INPUT HEAVY-TRAFFIC REAL-TIME QUEUES

BY LUKASZ KRUK,<sup>1</sup> JOHN LEHOCZKY, STEVEN SHREVE<sup>2</sup> AND  
SHU-NGAI YEUNG

*Maria Curie-Skłodowska University, Carnegie Mellon University,  
Carnegie Mellon University and AT&T Labs*

A single queueing station that serves  $K$  input streams is considered. Each stream is an independent renewal process, with customers having random lead times. Customers are served by processor sharing across streams. Within each stream, two disciplines are considered—earliest deadline first and first-in, first-out. The set of current lead times of the  $K$  streams is modeled as a  $K$ -dimensional vector of random counting measures on  $\mathbb{R}$ , and the limit of this vector of measure-valued processes is obtained under heavy traffic conditions.

**1. Introduction.** Over the last 10 years, communication technology has become dramatically more sophisticated. There are now many different types of communication services available and an ever-increasing demand for those services. Real-time communication services, while currently a small portion of the total demand, are of growing importance. In view of the increase in demand for real-time services (such as videoconferencing in which voice and video must be delivered in a timely fashion to ensure continuity of the image and sound), it is becoming important to develop models and analysis techniques first to understand and then to better control these networks. For real-time traffic, in addition to standard measures of stability and delay, one must also be concerned whether individual application packets are meeting their timing requirements. Thus, measuring average packet delay is not sufficient to determine whether a particular queue scheduling policy can satisfy real-time traffic requirements.

There has been important recent work that does develop analytic methods to determine stability and calculate delay. An excellent introduction to the scheduling of computer networks is given in the textbook by Keshav [12], Chapter 9. In addition, the paper by Zhang [21] surveys scheduling algorithms and associated performance bounds. More specifically, Cruz developed a calculus for network delay in the single-node [6] and the multiple-node [7] cases. The theory he developed is different from standard queueing theoretic models, because he used nonprobabilistic data stream arrival models that satisfy burstiness constraints.

---

Received November 2000; revised July 2001.

<sup>1</sup>Supported by NSF Grant DMS-98-03791.

<sup>2</sup>Supported by NSF Grants DMS-98-02464 and DMS-01-03814.

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

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

Using this assumption, he was able to derive bounds on delay and buffering requirements for network elements. Parekh and Gallager [14, 15] studied a new network resource allocation policy called the generalized (or weighted) processor sharing (GPS) algorithm or weighted fair queueing (WFQ) and combined this with a leaky bucket admission control policy. Demers, Keshav and Shenker [8] presented an analysis of these approaches. For these scheduling and control mechanisms, the authors developed worst-case packet delay and worst-case session backlogs for both the single-node and the multiple-node cases. Related research was published by Chang [4]. He also addressed the stability problem by determining conditions on queueing networks that result in bounded queue lengths and bounded delays for customers and gave conditions for the queue length distribution to have an exponential tail.

Most of these papers present a worst-case analysis, leading to bounds for the system that may be quite pessimistic. Thus, it is important to find methods that offer more realistic answers. An additional drawback is that the methods do not apply directly to real-time traffic, where performance must be measured at the packet level, rather than at an aggregated level. This paper introduces methods for some common queue scheduling policies to assess whether those packet timing requirements can be met.

In this paper, we consider a single-node system with  $K$  independent arrival streams (sessions). Each stream creates an arrival process of packets, and each packet has an individual service time and deadline, the time by which its transmission must be completed. Packets from individual streams are queued in separate buffers. The  $K$  streams are jointly serviced by a single server that allocates its total service capacity across the streams. We consider the *head-of-the-line processor sharing* (HOL-PS) service policy. This policy assumes that the total service capacity is evenly allocated to each stream having packets ready for transmission. Thus, if  $k$  of the  $K$  streams have one or more packets available for service, then each receives a fraction  $1/k$  of the server's capacity. We also introduce a weighting scheme to allow different streams to receive different relative amounts of service. We define a set of weights,  $\{w_i, 1 \leq i \leq K\}$ ,  $w_i \geq 0$ ,  $\sum_{i=1}^K w_i = 1$ . If two streams  $i$  and  $j$  have packets ready for transmission, then the relative service rate applied to streams  $i$  and  $j$  is the ratio  $w_i/w_j$ .

We assume that the packets from each stream are queued in different buffers, either in deadline order or in FIFO (first-in, first-out) order. Ordinarily, packets from a common stream will have a common, deterministic deadline; however, our approach permits us to consider the more general case in which packets from a single stream have random deadlines. For real-time systems, the service policy should take into account the packet timing requirements, and the earliest-deadline-first (EDF) algorithm for a single queue is generally optimal. The FIFO queueing discipline arises as a special case of EDF when the packets have constant deadlines.

Our approach follows the methods of Doytchinov, Lehoczky and Shreve [9]. (We note that the single-station results of [9] have recently been extended to feed-forward networks [20] and acyclic networks [13].) We study the heavy-traffic case in which the total traffic intensity on the server approaches 1. The HOL-PS queue discipline is work conserving, so under reasonable conditions on the arrival and service processes, the total workload on the system will converge to a drifted, reflected Brownian motion process. It is, however, more complicated to determine how that workload will be distributed across the  $K$  different queues and whether the packets in those queues will meet their timing requirements. To determine this, we define a lead time for each packet (the lead time is the time until the packet deadline) and a *lead time profile*, a counting measure on  $\mathbb{R}$  putting unit mass at the lead time of each customer in queue, for each of the  $K$  queues. Under heavy-traffic conditions, we prove that these profiles, when conditioned on the workload, converge to a deterministic form (including a null profile when the scaled queue length process converges to 0), which depends on the task deadlines, the queue discipline and the traffic intensities for each of the queues. The ability of this service discipline to meet packet timing requirements can be inferred from these profiles.

This paper is organized as follows. Section 2 presents the model, the assumptions and the notation. Section 3 introduces the heavy-traffic assumptions. Section 4 develops the measure-valued processes that will be needed to study the lead time profiles. Sections 5 and 6 present limit theorems for those measure-valued processes, assuming individual packet streams are queued using either EDF or FIFO, respectively. Section 7 presents simulation results illustrating the accuracy of the theory. Appendix A presents an analysis when the traffic streams are not balanced and one of the streams is dominant, and Appendix B provides a technical result needed for Appendix A.

**2. The model.** We have a sequence of single-station queueing systems, each with  $K$  arrival processes. The queueing systems are indexed by superscript  $(n)$ , and the arrival processes within each queueing system are indexed by  $k$ .

The interarrival times for the  $k$ th arrival process are  $\{u_{k,j}^{(n)}\}_{j=1}^{\infty}$ , a sequence of positive, independent, identically distributed random variables with common mean  $1/\lambda_k^{(n)}$  and standard deviation  $\alpha_k^{(n)}$ . The service times are  $\{v_{k,j}^{(n)}\}_{j=1}^{\infty}$ , another sequence of positive, independent, identically distributed random variables with common mean  $1/\mu_k^{(n)}$  and standard deviation  $\beta_k^{(n)}$ . We assume that each queue is empty at time 0.

For arrival stream  $k$ , we define the *customer arrival times*

$$(2.1) \quad S_{k,0}^{(n)} \triangleq 0, \quad S_{k,j}^{(n)} \triangleq \sum_{i=1}^j u_{k,i}^{(n)}, \quad j \geq 1,$$

the *customer arrival process*

$$(2.2) \quad A_k^{(n)}(t) \triangleq \max\{j; S_{k,j}^{(n)} \leq t\}, \quad t \geq 0,$$

and the *work arrival process*

$$(2.3) \quad V_k^{(n)}(t) \triangleq \sum_{j=1}^{\lfloor t \rfloor} v_{k,j}^{(n)}, \quad t \geq 0.$$

The work that has arrived at queue  $k$  by time  $t$  is then  $V_k^{(n)}(A_k^{(n)}(t))$ .

Each customer arrives with an initial lead time  $L_{k,j}^{(n)}$ , the time between the arrival time and the deadline for completion of service for that customer. These initial lead times are independent and identically distributed with distribution given by

$$(2.4) \quad \mathbb{P}\{L_{k,j}^{(n)} \leq \sqrt{n}y\} = G_k(y),$$

where  $G_k$  is a right-continuous cumulative distribution function for each  $k = 1, \dots, K$ . We define

$$(2.5) \quad y_k^* \triangleq \min\{y \in \mathbb{R}; G_k(y) = 1\}$$

and assume that  $y_k^*$  is finite for every  $k = 1, \dots, K$ . We assume that, for every  $n$ , the sequences  $\{u_{k,j}^{(n)}\}$ ,  $\{v_{k,j}^{(n)}\}$  and  $\{L_{k,j}^{(n)}\}$  are mutually independent over  $k = 1, \dots, K$  and  $j = 1, 2, \dots$ .

We now describe the allocation of service capacity to each of the  $K$  arrival streams in the  $n$ th queueing system. To this end, we introduce nonnegative random *weight* processes  $w_k^{(n)}(t)$ ,  $k = 1, \dots, K$ , such that  $\sum_{k=1}^K w_k^{(n)}(t) = 1$  for every  $t$  almost surely. The process  $w_k^{(n)}(t)$  represents the proportion of the server's capacity that would be assigned to queue  $k$  at time  $t$  if all the queues are nonempty at that time. If one or more queues is empty but queue  $k$  is not, then  $w_k^{(n)}(t)$  divided by the sum of the weights associated with the nonempty queues is the proportion of the server's capacity assigned to queue  $k$ . We allow the weight processes to take the value 0, but we make the following assumption, which guarantees that the proportion just described is defined.

**ASSUMPTION 2.1.** *At every time  $t$  when there is at least one customer present in the  $n$ th queueing system, there is at least one nonempty queue  $\ell$  for which  $w_\ell^{(n)}(t) > 0$ .*

An important consequence of Assumption 2.1 is that the server is always fully utilized whenever there are customers anywhere in the system.

Throughout Sections 2–5, we assume that customers within each of the  $K$  queues are served using the earliest-deadline-first (EDF) queue discipline; that is, the server always serves the customer with the shortest lead time. Preemption is permitted (we assume preempt–resume). There is no setup, switchover or other type of overhead. Late customers (customers with negative lead times) stay in queue until served to completion. In Section 6, customers are served using the first-in, first-out (FIFO) discipline within each queue.

Denoting by  $W_k^{(n)}(t)$  the work remaining in queue  $k$  at time  $t$ , we define the *idleness rate* for queue  $k$  to be

$$(2.6) \quad \dot{i}_k^{(n)}(t) \triangleq \begin{cases} 1, & \text{if } W_k^{(n)}(t) = 0, \\ 0, & \text{if } W_k^{(n)}(t) > 0. \end{cases}$$

The *entitlement rate* for queue  $k$  is then defined as

$$(2.7) \quad \dot{T}_k^{(n)}(t) \triangleq \frac{w_k^{(n)}(t)}{w_k^{(n)}(t) + \sum_{\ell \neq k} w_\ell^{(n)}(t)(1 - \dot{i}_\ell^{(n)}(t))}.$$

Because of Assumption 2.1, the denominator in (2.7) is not 0 as long as there is at least one customer present in the system. If this is not the case and if  $w_k^{(n)}(t) = 0$ , we interpret  $\frac{0}{0}$  in (2.7) as 0. The entitlement rate for queue  $k$  is the fraction of the server capacity queue  $k$  is entitled to receive at time  $t$  and would receive if there is work in queue  $k$  at that time.

We may alternatively represent the server entitlement rate for a queue in terms of joint idleness rates. By a *multi-index*  $\alpha = (\ell_1, \dots, \ell_{|\alpha|})$ , we shall mean a finite set of distinct indices  $\{\ell_1, \dots, \ell_{|\alpha|}\}$  between 1 and  $K$ . We denote by  $|\alpha|$  the cardinality of the set  $\alpha$ . Given a multi-index  $\alpha$ , we define the *joint idleness rate*

$$(2.8) \quad \dot{j}_\alpha^{(n)}(t) \triangleq \begin{cases} 1, & \text{if } W_\ell^{(n)}(t) = 0 \ \forall \ell \in \alpha, \ W_\ell^{(n)}(t) > 0 \ \forall \ell \notin \alpha, \\ 0, & \text{otherwise.} \end{cases}$$

By convention,

$$(2.9) \quad \dot{j}_\emptyset^{(n)}(t) = \begin{cases} 1, & \text{if } W_\ell^{(n)}(t) > 0 \ \forall \ell, \\ 0, & \text{otherwise.} \end{cases}$$

Finally, we define  $A_k(m)$  to be the set of all multi-indexes of cardinality  $m$  in which the index  $k$  appears and  $A_k^c(m)$  to be the set of all multi-indexes of cardinality  $m$  in which the index  $k$  does not appear. By convention,  $A_k(0) = \emptyset$  and  $A_k^c(0) = \{\emptyset\}$ . Then the idleness rate for queue  $k$  is

$$(2.10) \quad \dot{i}_k^{(n)}(t) = \sum_{m=1}^K \sum_{\alpha \in A_k(m)} \dot{j}_\alpha^{(n)}(t),$$

and the entitlement rate  $\dot{T}_k^{(n)}(t)$  can be written as

$$(2.11) \quad \begin{aligned} \dot{T}_k^{(n)}(t) &= \sum_{m=0}^{K-1} \sum_{\alpha \in A_k^c(m)} \frac{w_k^{(n)}(t)}{1 - \sum_{\ell \in \alpha} w_\ell^{(n)}(t)} \dot{j}_\alpha^{(n)}(t) \\ &+ \sum_{m=1}^K \sum_{\alpha \in A_k(m)} \frac{w_k^{(n)}(t)}{w_k^{(n)}(t) + 1 - \sum_{\ell \in \alpha} w_\ell^{(n)}(t)} \dot{j}_\alpha^{(n)}(t). \end{aligned}$$

It can happen in the above expressions that the denominator  $1 - \sum_{\ell \in \alpha} w_\ell^{(n)}(t)$  or  $w_k^{(n)}(t) + 1 - \sum_{\ell \in \alpha} w_\ell^{(n)}(t)$  is 0, but this can only occur if the numerator  $w_k^{(n)}(t)$  is also 0. In such a case, we interpret  $\frac{0}{0} \cdot 0$  as 0. At each time, at most one positive term appears on the right-hand side of (2.10) and (2.11). We define *idleness*, *joint idleness* and *entitlement* by integrating rates:

$$(2.12) \quad \begin{aligned} I_k^{(n)}(t) &= \int_0^t \dot{I}_k^{(n)}(s) ds, & J_\alpha^{(n)}(t) &= \int_0^t \dot{J}_\alpha^{(n)}(s) ds, \\ T_k^{(n)}(t) &= \int_0^t \dot{T}_k^{(n)}(s) ds. \end{aligned}$$

In particular,

$$(2.13) \quad I_k^{(n)}(t) = \sum_{m=1}^K \sum_{\alpha \in A_k(m)} J_\alpha^{(n)}(t),$$

$$(2.14) \quad \begin{aligned} T_k^{(n)}(t) &= \sum_{m=0}^{K-1} \sum_{\alpha \in A_k^c(m)} \int_0^t \frac{w_k^{(n)}(s)}{1 - \sum_{\ell \in \alpha} w_\ell^{(n)}(s)} \dot{J}_\alpha^{(n)}(s) ds \\ &+ \sum_{m=1}^K \sum_{\alpha \in A_k(m)} \int_0^t \frac{w_k^{(n)}(s)}{w_k^{(n)}(s) + 1 - \sum_{\ell \in \alpha} w_\ell^{(n)}(s)} \dot{J}_\alpha^{(n)}(s) ds. \end{aligned}$$

The *netput* for queue  $k$  is

$$(2.15) \quad N_k^{(n)}(t) \triangleq V_k^{(n)}(A_k^{(n)}(t)) - T_k^{(n)}(t),$$

which is the amount of work in the queue at time  $t$  if it has never been empty prior to time  $t$ . Because queue  $k$  might be empty, it can have *unused entitlement*, defined to be

$$(2.16) \quad U_k^{(n)}(t) \triangleq - \min_{0 \leq s \leq t} N_k^{(n)}(s).$$

The actual service received by queue  $k$  up to time  $t$  is

$$(2.17) \quad R_k^{(n)}(t) \triangleq T_k^{(n)}(t) - U_k^{(n)}(t),$$

and the workload at time  $t$  is

$$(2.18) \quad W_k^{(n)}(t) = N_k^{(n)}(t) + U_k^{(n)}(t) = V_k^{(n)}(A_k^{(n)}(t)) - R_k^{(n)}(t).$$

The unused entitlement increases at time  $t$  if and only if queue  $k$  is empty at time  $t$  and  $w_k^{(n)}(t) > 0$ ; in this case,  $\dot{U}_k^{(n)}(t) = \dot{T}_k^{(n)}(t)$ . This implies that

$$(2.19) \quad U_k^{(n)}(t) = \sum_{m=1}^K \sum_{\alpha \in A_k(m)} \int_0^t \frac{w_k^{(n)}(s)}{w_k^{(n)}(s) + 1 - \sum_{\ell \in \alpha} w_\ell^{(n)}(s)} \dot{J}_\alpha^{(n)}(s) ds.$$

Indeed, this is precisely the part of (2.14) corresponding to queue  $k$  being empty.

All the above processes are independent of the service discipline within each queue, provided that the processes  $w_k^{(n)}(t)$  remain unchanged. However, the *queue length processes*  $Q_k^{(n)}(t)$ , which are the number of customers in each queue  $k = 1, \dots, K$  at time  $t$ , depend on the queue discipline.

**3. Heavy-traffic assumptions.** We assume that, for each  $k$ , the following limits exist and are all positive:

$$(3.1) \quad \begin{aligned} \lambda_k &= \lim_{n \rightarrow \infty} \lambda_k^{(n)}, & \mu_k &= \lim_{n \rightarrow \infty} \mu_k^{(n)}, \\ \alpha_k &= \lim_{n \rightarrow \infty} \alpha_k^{(n)}, & \beta_k &= \lim_{n \rightarrow \infty} \beta_k^{(n)}. \end{aligned}$$

The *traffic intensity* for queue  $k$  in queueing system  $n$  is  $\rho_k^{(n)} = \lambda_k^{(n)} / \mu_k^{(n)}$ , and the *limiting traffic intensity* for queue  $k$  is  $\rho_k = \lambda_k / \mu_k$ . We assume

$$(3.2) \quad \sum_{k=1}^K \rho_k = 1$$

and make the *heavy-traffic assumption*

$$(3.3) \quad \lim_{n \rightarrow \infty} \sqrt{n}(\rho_k - \rho_k^{(n)}) = \gamma_k > 0$$

for each  $k = 1, \dots, K$ . In particular, there is a finite constant  $c$  such that

$$(3.4) \quad \sqrt{n}|\rho_k^{(n)} - \rho_k| \leq c$$

for all  $n$  and  $k$ . We also impose the usual Lindeberg condition on the interarrival and service times: for  $k = 1, \dots, K$ ,

$$(3.5) \quad \begin{aligned} &\lim_{n \rightarrow \infty} \mathbb{E} \left[ \left( u_{k,j}^{(n)} - (\lambda_k^{(n)})^{-1} \right)^2 \mathbb{I}_{\{|u_{k,j}^{(n)} - (\lambda_k^{(n)})^{-1}| > c\sqrt{n}\}} \right] \\ &= \lim_{n \rightarrow \infty} \mathbb{E} \left[ \left( v_{k,j}^{(n)} - (\mu_k^{(n)})^{-1} \right)^2 \mathbb{I}_{\{|v_{k,j}^{(n)} - (\mu_k^{(n)})^{-1}| > c\sqrt{n}\}} \right] = 0 \quad \forall c > 0. \end{aligned}$$

We introduce the *heavy-traffic scaling* for the workload and queue length processes

$$(3.6) \quad \widehat{W}_k^{(n)}(t) = \frac{1}{\sqrt{n}} W_k^{(n)}(nt), \quad \widehat{Q}_k^{(n)}(t) = \frac{1}{\sqrt{n}} Q_k^{(n)}(nt),$$

and the *centered heavy-traffic scaling* for the arrival process

$$(3.7) \quad \widehat{A}_k^{(n)}(t) = \frac{1}{\sqrt{n}} \left[ A_k^{(n)}(nt) - \lambda_k^{(n)} nt \right].$$

For  $k = 1, \dots, K$  and  $y \leq y_k^*$ , we define the respective *partial-work arrival process* and *centered scaled partial-work arrival process*

$$(3.8) \quad V_{k,y}^{(n)}(t) = \sum_{j=1}^{\lfloor t \rfloor} v_{k,j}^{(n)} \mathbb{I}_{\{L_{k,j}^{(n)} \leq \sqrt{ny}\}},$$

$$(3.9) \quad \widehat{V}_{k,y}^{(n)}(t) = \frac{1}{\sqrt{n}} \sum_{j=1}^{\lfloor nt \rfloor} \left( v_{k,j}^{(n)} \mathbb{I}_{\{L_{k,j}^{(n)} \leq \sqrt{ny}\}} - (\mu_k^{(n)})^{-1} G_k(y) \right),$$

and we also define

$$(3.10) \quad \widehat{Z}_{k,y}^{(n)}(t) = \frac{1}{\sqrt{n}} \left[ V_{k,y}^{(n)}(A_k^{(n)}(nt)) - G_k(y) \rho_k^{(n)} nt \right].$$

Note that the term  $(\mu_k^{(n)})^{-1} G_k(y)$  appearing in (3.9) is the expectation of  $v_{k,j}^{(n)} \mathbb{I}_{\{L_{k,j}^{(n)} \leq \sqrt{ny}\}}$ . The variance of this random variable is

$$(3.11) \quad (\beta_{k,y}^{(n)})^2 \triangleq G_k(y) (\beta_k^{(n)})^2 + (\mu_k^{(n)})^{-2} G_k(y) (1 - G_k(y)),$$

which has limit

$$(3.12) \quad (\beta_{k,y})^2 = G_k(y) (\beta_k)^2 + (\mu_k)^{-2} G_k(y) (1 - G_k(y)).$$

We use the notation  $\widehat{V}_k^{(n)} = \widehat{V}_{k,y_k^*}^{(n)}$  and  $\widehat{Z}_k^{(n)} = \widehat{Z}_{k,y_k^*}^{(n)}$ . Note that  $V_k^{(n)} = V_{k,y_k^*}^{(n)}$  and  $\beta_k^{(n)} = \beta_{k,y_k^*}^{(n)}$ .

The following theorem is due to Prokhorov [16], Theorem 3.1 (see also [2], Theorem 14.6). Here and elsewhere, the symbol  $\Rightarrow$  denotes weak convergence of measures on the space  $D[0, \infty)$  of functions from  $[0, \infty)$  to  $\mathbb{R}$  which are right-continuous with left limits. The topology of this space is a generalization of the topology introduced by Skorohod for  $D[0, 1]$ . See [2] for details.

**THEOREM 3.1.** *For every  $k = 1, \dots, K$  and every  $y \leq y_k^*$ ,*

$$(3.13) \quad (\widehat{A}_k^{(n)}, \widehat{V}_{k,y}^{(n)}) \Rightarrow (A_k^*, V_{k,y}^*),$$

where  $A_k^*$  and  $V_{k,y}^*$  are independent Brownian motions with mean 0 and variances  $\alpha_k^2 \lambda_k^3$  and  $\beta_{k,y}^2$ , respectively.

This theorem has the following corollaries.

**COROLLARY 3.2.** *For every  $k = 1, \dots, K$  and  $y \leq y_k^*$ ,*

$$(3.14) \quad \widehat{Z}_{k,y}^{(n)} \Rightarrow V_{k,y}^* \circ \lambda_k e + \frac{G_k(y)}{\mu_k} A_k^*,$$

where  $e$  is the identity function  $e(t) = t$  for all  $t \geq 0$ .



PROOF. We compute

$$\begin{aligned}
\widehat{Z}_{k,y}^{(n)}(t) &= \frac{1}{\sqrt{n}} \left[ V_{k,y}^{(n)}(A_k^{(n)}(nt)) - G_k(y) \rho_k^{(n)} nt \right] \\
(3.15) \quad &= \frac{1}{\sqrt{n}} \sum_{j=1}^{A_k^{(n)}(nt)} \left[ v_{k,j}^{(n)} \mathbb{1}_{\{L_{k,j}^{(n)} \leq \sqrt{n}y\}} - \frac{G_k(y)}{\mu_k^{(n)}} \right] + \frac{G_k(y)}{\mu_k^{(n)} \sqrt{n}} \left[ A_k^{(n)}(nt) - \lambda_k^{(n)} nt \right] \\
&= \left[ \widehat{V}_{k,y}^{(n)} \left( \frac{1}{n} A_k^{(n)}(nt) \right) - \widehat{V}_{k,y}^{(n)}(\lambda_k t) \right] + \widehat{V}_{k,y}^{(n)}(\lambda_k t) + \frac{G_k(y)}{\mu_k^{(n)}} \widehat{A}_k^{(n)}(t).
\end{aligned}$$

The term in square brackets converges to 0 because  $\frac{1}{n} A_k^{(n)}(nt) \Rightarrow \lambda_k t$  (see, e.g., Theorem A.3 of [9]).  $\square$

COROLLARY 3.3. *The process  $\sum_{k=1}^K \widehat{W}_k^{(n)}$  converges weakly to a reflected Brownian motion with drift.*

PROOF. Since  $\sum_{k=1}^K \widehat{W}_k^{(n)}$  is the total workload, a quantity that is invariant under work-conserving disciplines, this is a classical result due to Iglehart and Whitt [11]. One can also obtain it immediately from Corollary 3.2 and the equation

$$\begin{aligned}
(3.16) \quad \sum_{k=1}^K \widehat{W}_k^{(n)}(t) &= \sum_{k=1}^K \widehat{Z}_k^{(n)}(t) - \sqrt{n} \left( 1 - \sum_{k=1}^K \rho_k^{(n)} \right) t \\
&\quad - \min_{0 \leq s \leq t} \left[ \sum_{k=1}^K \widehat{Z}_k^{(n)}(s) - \sqrt{n} \left( 1 - \sum_{k=1}^K \rho_k^{(n)} \right) s \right]. \quad \square
\end{aligned}$$

For the remainder of this paper, we make the following assumption.

ASSUMPTION 3.4. *We have*

$$(3.17) \quad \widehat{W}^{(n)} \triangleq (\widehat{W}_1^{(n)}, \dots, \widehat{W}_K^{(n)}) \Rightarrow W^*,$$

where  $W^* = (W_1^*, \dots, W_K^*)$  is an RCLL (right-continuous with left-hand limits) process on  $\mathbb{R}^K$ .

In Appendix A, we will show that this assumption holds when the weight processes  $w_k^{(n)}(t)$  converge sufficiently fast to positive constants  $w_k$  and the traffic intensities are *unbalanced*, that is,

$$(3.18) \quad 0 < \frac{\rho_1}{w_1} \leq \frac{\rho_2}{w_2} \leq \dots \leq \frac{\rho_{K-1}}{w_{K-1}} < \frac{\rho_K}{w_K}.$$

Actually, in this case,  $W_k^* \equiv 0$  for all  $k < K$ . In the case of *balanced traffic intensities*,

$$(3.19) \quad \frac{\rho_1}{w_1} = \frac{\rho_2}{w_2} = \dots = \frac{\rho_{K-1}}{w_{K-1}} = \frac{\rho_K}{w_K},$$

Assumption 3.4 is verified by Ramanan and Reiman [17].

**COROLLARY 3.5.** *The processes  $W_k^*$ ,  $k = 1, \dots, K$ , are continuous.*

**PROOF.** For simplicity, we prove continuity of  $W_k^*$  on  $[0, 1]$ . For  $x \in D[0, 1]$ , let  $j(x) = \sup_{0 \leq t \leq 1} |x(t) - x(t-)|$ . By the continuous mapping theorem,  $j(\widehat{Z}_k^{(n)}) \Rightarrow 0$ . But  $j(\widehat{W}_k^{(n)}) = j(\widehat{Z}_k^{(n)})$ , and so  $j(\widehat{W}_k^{(n)}) \Rightarrow 0$ . According to Theorem 13.4 of [2], this implies continuity of the limiting process  $W_k^*$ .  $\square$

**4. Measure-valued processes.** To study whether tasks or customers meet their timing requirements, one must keep track of customer lead times, where the lead time is the time remaining until the deadline elapses, that is,

$$\text{Lead time} = \text{Deadline} - \text{Current time.}$$

In this section, we define a collection of measure-valued processes that will be useful in the analysis of the instantaneous lead time profile of the customers.

*Queue length measures:*

$$(4.1) \quad \mathcal{Q}_k^{(n)}(t)(B) \triangleq \left\{ \begin{array}{l} \text{Number of customers in queue } k \text{ at time } t \\ \text{having lead times at time } t \text{ in } B \subset \mathbb{R} \end{array} \right\}.$$

*Workload measures:*

$$(4.2) \quad \mathcal{W}_k^{(n)}(t)(B) \triangleq \left\{ \begin{array}{l} \text{Work in queue } k \text{ at time } t \text{ associated with customers} \\ \text{in this queue having lead times at time } t \text{ in } B \subset \mathbb{R} \end{array} \right\}.$$

*Customer arrival measures:*

$$(4.3) \quad \mathcal{A}_k^{(n)}(t)(B) \triangleq \left\{ \begin{array}{l} \text{Number of all arrivals in queue } k \text{ by time } t \\ \text{having lead times at time } t \text{ in } B \subset \mathbb{R} \end{array} \right\}.$$

*Workload arrival measures:*

$$(4.4) \quad \mathcal{V}_k^{(n)}(t)(B) \triangleq \left\{ \begin{array}{l} \text{Work in queue } k \text{ associated with all arrivals by} \\ \text{time } t \text{ having lead times at time } t \text{ in } B \subset \mathbb{R} \end{array} \right\}.$$

The following relationships easily follow:

$$(4.5) \quad \mathcal{Q}_k^{(n)}(t) = \mathcal{Q}_k^{(n)}(t)(\mathbb{R}), \quad \mathcal{W}_k^{(n)}(t) = \mathcal{W}_k^{(n)}(t)(\mathbb{R}),$$

$$(4.6) \quad \mathcal{A}_k^{(n)}(t) = \mathcal{A}_k^{(n)}(t)(\mathbb{R}), \quad \mathcal{V}_k^{(n)}(t) = \mathcal{V}_k^{(n)}(t)(\mathbb{R}),$$

$$(4.7) \quad \mathcal{A}_k^{(n)}(t)(B) = \sum_{j=1}^{A_k^{(n)}(t)} \mathbb{I}_{\{L_{k,j}^{(n)} - (t - S_{k,j}^{(n)}) \in B\}} = \sum_{j=1}^{\infty} \mathbb{I}_{\{S_{k,j}^{(n)} \in B + t - L_{k,j}^{(n)}, S_{k,j}^{(n)} \leq t\}},$$

$$(4.8) \quad \mathcal{V}_k^{(n)}(t)(B) = \sum_{j=1}^{A_k^{(n)}(t)} v_{k,j}^{(n)} \mathbb{I}_{\{L_{k,j}^{(n)} - (t - S_{k,j}^{(n)}) \in B\}} = \sum_{j=1}^{\infty} v_{k,j}^{(n)} \mathbb{I}_{\{S_{k,j}^{(n)} \in B + t - L_{k,j}^{(n)}, S_{k,j}^{(n)} \leq t\}}.$$

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 in each queue  $k$  and the largest lead time of all customers from queue  $k$ , whether present or departed, who have ever been in service. For  $k = 1, \dots, K$ , we define the *frontier* for queue  $k$

$$(4.9) \quad F_k^{(n)}(t) \triangleq \left\{ \begin{array}{l} \text{Largest lead time of all customers from queue } k \\ \text{who have ever been in service, whether still} \\ \text{present or not, or } \sqrt{n}y_k^* - t, \text{ if this quantity is} \\ \text{larger than the former one} \end{array} \right\},$$

and the *current lead time*:

$$(4.10) \quad C_k^{(n)}(t) \triangleq \left\{ \begin{array}{l} \text{Lead time of the customer in service in} \\ \text{queue } k \text{ or } F_k^{(n)}(t) \text{ if the queue is empty} \end{array} \right\}.$$

Prior to arrival of the first customer,  $F_k^{(n)}(t) = \sqrt{n}y_k^* - t$ . Under the EDF queue discipline, there is no customer in queue  $k$  with lead time smaller than  $C_k^{(n)}(t)$ , and there has never been a customer in service in this queue whose lead time, if the customer were still present, would exceed  $F_k^{(n)}(t)$ . Furthermore,  $C_k^{(n)}(t) \leq F_k^{(n)}(t)$  for all  $t \geq 0, k = 1, \dots, K$ . Both  $F_k^{(n)}$  and  $C_k^{(n)}$  are RCLL processes.

For the processes just defined under the EDF queue discipline, we use the following heavy-traffic scalings:

$$(4.11) \quad \widehat{F}_k^{(n)}(t) \triangleq \frac{1}{\sqrt{n}} F_k^{(n)}(nt), \quad \widehat{C}_k^{(n)}(t) \triangleq \frac{1}{\sqrt{n}} C_k^{(n)}(nt),$$

$$(4.12) \quad \widehat{\mathcal{Q}}_k^{(n)}(t) \triangleq \frac{1}{\sqrt{n}} \mathcal{Q}_k^{(n)}(nt)(\sqrt{n}B), \quad \widehat{\mathcal{W}}_k^{(n)}(t) \triangleq \frac{1}{\sqrt{n}} \mathcal{W}_k^{(n)}(nt)(\sqrt{n}B).$$

We also define

$$(4.13) \quad \begin{aligned} \widehat{\mathcal{A}}_k^{(n)}(t)(B) &\triangleq \frac{1}{\sqrt{n}} A_k^{(n)}(nt)(\sqrt{n}B) \\ &= \frac{1}{\sqrt{n}} \sum_{j=1}^{A_k^{(n)}(nt)} \mathbb{I}_{\{L_{k,j}^{(n)} - (nt - S_{k,j}^{(n)}) \in \sqrt{n}B\}} \\ &= \frac{1}{\sqrt{n}} \sum_{j=1}^{\infty} \mathbb{I}_{\{S_{k,j}^{(n)} \in \sqrt{n}B + nt - L_{k,j}^{(n)}, S_{k,j}^{(n)} \leq nt\}}, \end{aligned}$$

$$\begin{aligned}
 \widehat{\mathbf{v}}_k^{(n)}(t)(B) &\triangleq \frac{1}{\sqrt{n}} \mathbf{v}_k^{(n)}(nt)(\sqrt{n}B) \\
 (4.14) \qquad &= \frac{1}{\sqrt{n}} \sum_{j=1}^{A_k^{(n)}(nt)} v_{k,j}^{(n)} \mathbb{I}_{\{L_{k,j}^{(n)} - (nt - S_{k,j}^{(n)}) \in \sqrt{n}B\}} \\
 &= \frac{1}{\sqrt{n}} \sum_{j=1}^{\infty} v_{k,j}^{(n)} \mathbb{I}_{\{S_{k,j}^{(n)} \in \sqrt{n}B + nt - L_{k,j}^{(n)}, S_{k,j}^{(n)} \leq nt\}}.
 \end{aligned}$$

**5. Earliest-deadline-first limits.** We begin with an informal discussion of the limiting lead time profile when customers are served according to the earliest-deadline-first discipline. At each time  $t$ , there are potentially two types of customers in queue  $k$ : those whose lead times are in  $[C_k^{(n)}(t), F_k^{(n)}(t))$  and those whose lead times are in  $[F_k^{(n)}(t), \infty)$ . Customers of the former type are present only under unusual circumstances. Either there must be a customer present with lead time less than  $\sqrt{n}y_k^* - t$ , an unlikely event for  $t$  of order  $n$ , or there must be a customer in service, the customer we call  $\mathcal{C}$ , whose lead time is  $F_k^{(n)}(t)$  at time  $t$  and who must have been preempted by a customer with smaller lead time. When  $\mathcal{C}$  is preempted, the preempting customer is the only customer with lead time in  $[C_k^{(n)}(t), F_k^{(n)}(t))$ , but additional customers may later arrive, sustaining a period of time before service of  $\mathcal{C}$  is resumed. These later customers must have lead times less than or equal to  $F_k^{(n)}(t)$ , which, in turn, is strictly less than  $\sqrt{n}y_k^*$ . Hence, the arrival intensity for these customers is strictly less than  $\lambda_k^{(n)}$ , although they have the full attention of the server. This causes the period of time before service of  $\mathcal{C}$  is resumed to be brief and suggests that few customers ever have lead times in  $[C_k^{(n)}(t), F_k^{(n)}(t))$ . This intuition is made precise by Lemma 5.4.

Ignoring customers with lead times in  $[C_k^{(n)}(t), F_k^{(n)}(t))$ , we work with the lead time distribution of customers in queue  $k$  by examining those with lead times in  $[F_k^{(n)}(t), \infty)$ . Except possibly for a customer with lead time  $F_k^{(n)}(t)$ , these customers have never been in service. In order for one of these customers to have lead time  $y > F_k^{(n)}(t)$ , the customers must arrive at some time  $t - s$  prior to  $t$  and be assigned lead time  $y + s$  upon arrival. If  $G_k$  has a density, then the density of the assigned lead time distribution is  $(1/\sqrt{n})G_k'((y + s)/\sqrt{n})$  [see (2.4)], and multiplying by the arrival rate  $\lambda_k^{(n)}$ , we obtain the density of customers with lead time  $y$ :

$$\frac{\lambda_k^{(n)}}{\sqrt{n}} \int_0^t G_k' \left( \frac{y + s}{\sqrt{n}} \right) ds = \lambda_k^{(n)} \left[ G_k \left( \frac{y + t}{\sqrt{n}} \right) - G_k \left( \frac{y}{\sqrt{n}} \right) \right].$$

The heavy-traffic scaling considers the density of  $1/\sqrt{n}$  times the actual number of customers whose lead times are

$$\hat{y} = \frac{y}{\sqrt{n}} > \frac{1}{\sqrt{n}} F_k^{(n)}(t) = \widehat{F}_k^{(n)}(\hat{t})$$

at scaled time  $\hat{t} = \frac{t}{n}$ . This density is

$$\frac{\lambda_k^{(n)}}{\sqrt{n}} [G_k(\hat{y} + \sqrt{n}\hat{t}) - G_k(\hat{y})] \frac{dy}{d\hat{y}} = \lambda_k^{(n)} [G_k(\hat{y} + \sqrt{n}\hat{t}) - G_k(\hat{y})],$$

and as  $n \rightarrow \infty$ , it converges to  $\lambda_k[1 - G_k(\hat{y})]$ . The scaled number of customers in queue  $k$  at time  $t$  (scaled time  $\hat{t}$ ) is thus approximately equal to

$$\lambda_k \int_{\widehat{F}_k^{(n)}(\hat{t})}^{\infty} (1 - G_k(\eta)) d\eta,$$

and since the average work per customer is  $1/\mu_k^{(n)}$ , the scaled work in the queue can be approximated by

$$\rho_k \int_{\widehat{F}_k^{(n)}(\hat{t})}^{\infty} (1 - G_k(\eta)) d\eta.$$

As  $n \rightarrow \infty$ , this should converge to  $W_k^*(t)$  of Assumption 3.4, and this gives us an equation that characterizes  $F_k^* = \lim_{n \rightarrow \infty} \widehat{F}_k^{(n)}$ :

$$W_k^* = \rho_k \int_{F_k^*}^{\infty} (1 - G_k(\eta)) d\eta.$$

These considerations prompt the following definitions.

For  $k = 1, \dots, K$ , we set

$$(5.1) \quad H_k(y) \triangleq \int_y^{\infty} (1 - G_k(\eta)) d\eta = \begin{cases} \int_y^{y_k^*} (1 - G_k(\eta)) d\eta, & \text{if } y \leq y_k^*, \\ 0, & \text{if } y > y_k^*. \end{cases}$$

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

$$(5.2) \quad F_k^*(t) \triangleq H_k^{-1} \left( \frac{W_k^*(t)}{\rho_k} \right), \quad t \geq 0,$$

where  $W_k^*$  is as in Assumption 3.4. Let  $\widehat{F}^{(n)} = (\widehat{F}_1^{(n)}, \dots, \widehat{F}_K^{(n)})$  and  $F^* = (F_1^*, \dots, F_K^*)$ . Denote by  $\mathcal{M}$  the set of all finite, nonnegative measures on  $\mathcal{B}(\mathbb{R})$ , the Borel subsets of  $\mathbb{R}$ , and denote by  $\mathcal{M}^K$  the  $K$ -fold product of  $\mathcal{M}$ . Under the weak topology,  $\mathcal{M}$  is a separable, metrizable topological space, and so is  $\mathcal{M}^K$ .

Under Assumptions 2.1 and 3.4, we have the following generalizations of Proposition 3.10, Theorem 3.1 and Corollary 3.2 of [9]. The proofs are straightforward generalizations of those given in [9], with the exception of Lemma 5.4.

PROPOSITION 5.1. *Under the earliest-deadline-first queue discipline, we have  $\widehat{F}^{(n)} \Rightarrow F^*$ .*

THEOREM 5.2. *For  $k = 1, \dots, K$ , define measure-valued processes  $\widehat{W}_k^*$  and  $\widehat{Q}_k^*$  by*

$$(5.3) \quad \widehat{W}_k^*(t)(B) \triangleq \rho_k \int_{B \cap [F_k^*(t), \infty)} (1 - G_k(y)) dy, \quad \widehat{Q}_k^*(t)(B) \triangleq \mu_k \widehat{W}_k^*(t)$$

for all Borel sets  $B \subset \mathbb{R}$ . Under the earliest-deadline-first queue discipline, the processes  $\widehat{W}^{(n)} = (\widehat{W}_1^{(n)}, \dots, \widehat{W}_K^{(n)})$  and  $\widehat{Q}^{(n)} = (\widehat{Q}_1^{(n)}, \dots, \widehat{Q}_K^{(n)})$  converge weakly in  $D([0, \infty), \mathcal{M}^K)$  to  $\widehat{W}^* = (\widehat{W}_1^*, \dots, \widehat{W}_K^*)$  and  $\widehat{Q}^* = (\widehat{Q}_1^*, \dots, \widehat{Q}_K^*)$ , respectively.

COROLLARY 5.3. *Under the earliest-deadline-first queue discipline, the  $K$ -dimensional scaled queue length process  $\widehat{Q}^{(n)} = (\widehat{Q}_1^{(n)}, \dots, \widehat{Q}_K^{(n)})$  converges weakly to the process  $(\mu_1 W_1^*, \dots, \mu_K W_K^*)$ .*

To extend the arguments of [9] to prove the above results, only the generalization of Proposition 3.6 of [9] needs more attention. We state and prove the relevant result.

LEMMA 5.4. *Under the earliest-deadline-first queue discipline, we have, for  $k = 1, \dots, K$ ,*

$$(5.4) \quad \widehat{W}_k^{(n)}[\widehat{C}_k^{(n)}, \widehat{F}_k^{(n)}] \Rightarrow 0, \quad \widehat{Q}_k^{(n)}[\widehat{C}_k^{(n)}, \widehat{F}_k^{(n)}] \Rightarrow 0.$$

PROOF. We fix  $k$  and prove convergence on  $[0, 1]$ . For  $0 \leq t \leq 1$ , we define

$$(5.5) \quad \tau_k^{(n)}(t) \triangleq \sup\{s \in [0, t]; \widehat{C}_k^{(n)}(s) = \widehat{F}_k^{(n)}(s)\}.$$

We further define

$$(5.6) \quad \begin{aligned} & \widehat{D}_k^{(n)}(t) \\ &= \frac{1}{\sqrt{n}} \sum_{j=1}^{\infty} v_{k,j}^{(n)} \mathbb{I}_{\{n\tau_k^{(n)}(t) < S_{k,j}^{(n)} \leq nt\}} \mathbb{I}_{\{L_{k,j}^{(n)} - (nt - S_{k,j}^{(n)}) < F_k^{(n)}(n\tau_k^{(n)}(t)) - n(t - \tau_k^{(n)}(t))\}}. \end{aligned}$$

This is the scaled work brought by customers who arrive in queue  $k$  within the time interval  $(n\tau_k^{(n)}(t), nt]$  with lead times at arrival smaller than the current frontier. Because there is work with lead time in  $[C_k^{(n)}(s), F_k^{(n)}(s))$  for all  $s \in [n\tau_k^{(n)}(t), nt]$ , the effort expended by the server on customers with lead times smaller than the frontier during the period  $n\tau_k^{(n)}(t)$  to  $nt$  is  $T_k^{(n)}(nt) - T_k^{(n)}(n\tau_k^{(n)}(t))$ , the full entitlement of queue  $k$  during this period. The work at time  $nt$  with lead time smaller than the frontier is thus the work at time  $n\tau_k^{(n)}(t)$  with lead time smaller

than the frontier plus the work that arrives with lead time smaller than the frontier minus  $T_k^{(n)}(nt) - T_k^{(n)}(n\tau_k^{(n)})$ . It follows that

$$\begin{aligned}
 (5.7) \quad & 0 \leq \widehat{\mathcal{W}}_k^{(n)}(t) \left[ \widehat{\mathcal{C}}_k^{(n)}(t), \widehat{\mathcal{F}}_k^{(n)}(t) \right] \\
 & = \widehat{\mathcal{W}}_k^{(n)}(\tau_k^{(n)}(t)) \left[ \widehat{\mathcal{C}}_k^{(n)}(\tau_k^{(n)}(t)), \widehat{\mathcal{F}}_k^{(n)}(\tau_k^{(n)}(t)) \right] + \widehat{D}_k^{(n)}(t) \\
 & \quad - \frac{1}{\sqrt{n}} \left[ T_k^{(n)}(nt) - T_k^{(n)}(n\tau_k^{(n)}(t)) \right].
 \end{aligned}$$

We examine each of the three terms on the right-hand side of (5.7).

By definition,

$$\widehat{\mathcal{W}}_k^{(n)}(\tau_k^{(n)}(t)-) \left[ \widehat{\mathcal{C}}_k^{(n)}(\tau_k^{(n)}(t)-), \widehat{\mathcal{F}}_k^{(n)}(\tau_k^{(n)}(t)-) \right] = 0,$$

and so

$$(5.8) \quad \widehat{\mathcal{W}}_k^{(n)}(\tau_k^{(n)}(t)) \left[ \widehat{\mathcal{C}}_k^{(n)}(\tau_k^{(n)}(t)), \widehat{\mathcal{F}}_k^{(n)}(\tau_k^{(n)}(t)) \right] \leq j(\widehat{\mathcal{W}}_k^{(n)}),$$

where  $j(x) = \sup_{0 \leq s \leq 1} |x(s) - x(s-)|$  for  $x \in D[0, 1]$ . In light of Corollary 3.5, we conclude that

$$(5.9) \quad \widehat{\mathcal{W}}_k^{(n)}(\tau_k^{(n)}(t)) \left[ \widehat{\mathcal{C}}_k^{(n)}(\tau_k^{(n)}(t)), \widehat{\mathcal{F}}_k^{(n)}(\tau_k^{(n)}(t)) \right] \Rightarrow 0,$$

where the convergence is that of RCLL processes defined on  $[0, 1]$ .

In contrast to the proof of Proposition 3.6 in [9], we cannot assert that the workload in queue  $k$  is being decreased at rate 1 on the time interval  $(n\tau_k^{(n)}(t), nt]$ . Indeed,  $\dot{T}_k^{(n)}$  is, in general, not equal to 1. However, on the interval under consideration, the  $k$ th queue is never empty, so

$$(5.10) \quad U_k^{(n)}(nt) = U_k^{(n)}(n\tau_k^{(n)}(t)).$$

By (2.15) and (2.18), we have

$$(5.11) \quad T_k^{(n)}(s) = V_k^{(n)}(A_k^{(n)}(s)) - W_k^{(n)}(s) + U_k^{(n)}(s), \quad s \geq 0.$$

Thus the amount by which the workload in queue  $k$  has been decreased on  $(n\tau_k^{(n)}(t), nt]$  due to the activity of the server, divided by the scaling factor  $\sqrt{n}$ , is

$$\begin{aligned}
 (5.12) \quad & \frac{1}{\sqrt{n}} \left[ T_k^{(n)}(nt) - T_k^{(n)}(n\tau_k^{(n)}(t)) \right] \\
 & = \frac{1}{\sqrt{n}} \left[ V_k^{(n)}(A_k^{(n)}(nt)) - V_k^{(n)}(A_k^{(n)}(n\tau_k^{(n)}(t))) \right] \\
 & \quad - \frac{1}{\sqrt{n}} \left[ W_k^{(n)}(nt) - W_k^{(n)}(n\tau_k^{(n)}(t)) \right] \\
 & = \widehat{Z}_k^{(n)}(t) - \widehat{Z}_k^{(n)}(\tau_k^{(n)}(t)) - \left[ \widehat{W}_k^{(n)}(t) - \widehat{W}_k^{(n)}(\tau_k^{(n)}(t)) \right] \\
 & \quad + \rho_k^{(n)} \sqrt{n} (t - \tau_k^{(n)}(t)).
 \end{aligned}$$

Because  $0 \leq \tau_k^{(n)}(t) \leq t$ , we have

$$(5.13) \quad \min_{0 \leq s \leq t} \widehat{Z}_k^{(n)}(s) \leq \widehat{Z}_k^{(n)}(\tau_k^{(n)}(t)) \leq \max_{0 \leq s \leq t} \widehat{Z}_k^{(n)}(s),$$

$$(5.14) \quad \min_{0 \leq s \leq t} \widehat{Z}_{k,y}^{(n)}(s) \leq \widehat{Z}_{k,y}^{(n)}(\tau_k^{(n)}(t)) \leq \max_{0 \leq s \leq t} \widehat{Z}_{k,y}^{(n)}(s),$$

$$(5.15) \quad \min_{0 \leq s \leq t} \widehat{W}_k^{(n)}(s) \leq \widehat{W}_k^{(n)}(\tau_k^{(n)}(t)) \leq \max_{0 \leq s \leq t} \widehat{W}_k^{(n)}(s).$$

We conclude from Corollary 3.2, Assumption 3.4 and the continuity theorem that the upper and lower bounds in (5.13), (5.14) and (5.15) have weak limits, and hence

$$(5.16) \quad \frac{1}{\sqrt{n}} \left[ T_k^{(n)}(nt) - T_k^{(n)}(n\tau_k^{(n)}(t)) \right] = \rho_k^{(n)} \sqrt{n}(t - \tau_k^{(n)}(t)) + O(1).$$

Finally, to study  $\widehat{D}_k^{(n)}(t)$ , we choose  $y \leq y_k^*$  and consider two cases.

*Case I.*  $n\tau_k^{(n)}(t) + \sqrt{n}(y_k^* - y) \leq nt$ . In this case, on the event  $\{n\tau_k^{(n)}(t) + \sqrt{n}(y_k^* - y) < S_{k,j}^{(n)} \leq nt\}$ , the inequality  $L_{k,j}^{(n)} - (nt - S_{k,j}^{(n)}) < F_k^{(n)}(n\tau_k^{(n)}(t)) - n(t - \tau_k^{(n)}(t))$  implies the inequality  $L_{k,j}^{(n)} \leq \sqrt{n}y$ . Therefore,

$$(5.17) \quad \begin{aligned} \widehat{D}_k^{(n)}(t) &\leq \frac{1}{\sqrt{n}} \sum_{j=1}^{\infty} v_{k,j}^{(n)} \mathbb{I}_{\{n\tau_k^{(n)}(t) < S_{k,j}^{(n)} \leq n\tau_k^{(n)}(t) + \sqrt{n}(y_k^* - y)\}} \\ &\quad + \frac{1}{\sqrt{n}} \sum_{j=1}^{\infty} v_{k,j}^{(n)} \mathbb{I}_{\{n\tau_k^{(n)}(t) + \sqrt{n}(y_k^* - y) < S_{k,j}^{(n)} \leq nt\}} \mathbb{I}_{\{L_{k,j}^{(n)} \leq \sqrt{n}y\}} \\ &= \frac{1}{\sqrt{n}} \left[ V_k^{(n)}(A_k^{(n)}(n\tau_k^{(n)}(t) + \sqrt{n}(y_k^* - y))) - V_k^{(n)}(A_k^{(n)}(n\tau_k^{(n)}(t))) \right] \\ &\quad + \frac{1}{\sqrt{n}} \left[ V_{k,y}^{(n)}(A_k^{(n)}(nt)) - V_{k,y}^{(n)}(A_k^{(n)}(n\tau_k^{(n)}(t) + \sqrt{n}(y_k^* - y))) \right] \\ &= \widehat{Z}_k^{(n)}\left(\tau_k^{(n)}(t) + \frac{1}{\sqrt{n}}(y_k^* - y)\right) - \widehat{Z}_k^{(n)}(\tau_k^{(n)}(t)) + \rho_k^{(n)}(y_k^* - y) \\ &\quad + \widehat{Z}_{k,y}^{(n)}(t) - \widehat{Z}_{k,y}^{(n)}\left(\tau_k^{(n)}(t) + \frac{1}{\sqrt{n}}(y_k^* - y)\right) \\ &\quad + G_k(y)\rho_k^{(n)}\sqrt{n}(t - \tau_k^{(n)}(t)) - G_k(y)\rho_k^{(n)}(y_k^* - y). \end{aligned}$$

Thus, by (5.13) and (5.14),

$$(5.18) \quad \widehat{D}_k^{(n)}(t) \leq G_k(y)\rho_k^{(n)}\sqrt{n}(t - \tau_k^{(n)}(t)) + O(1).$$



We now use (5.18) and (5.16) to conclude from (5.7) that

$$(5.19) \quad 0 \leq \widehat{W}_k^{(n)}(\tau_k^{(n)}(t)) \left[ \widehat{C}_k^{(n)}(\tau_k^{(n)}(t)), \widehat{F}_k^{(n)}(\tau_k^{(n)}(t)) \right] \\ - (1 - G_k(y)) \rho_k^{(n)} \sqrt{n} (t - \tau_k^{(n)}(t)) + O(1),$$

which implies, because of (5.9), that

$$(5.20) \quad (1 - G_k(y)) \rho_k^{(n)} (t - \tau_k^{(n)}(t)) \leq O\left(\frac{1}{\sqrt{n}}\right).$$

*Case II.*  $n\tau_k^{(n)}(t) + \sqrt{n}(y_k^* - y) > nt$ . In this case,  $t - \tau_k^{(n)}(t) \leq (1/\sqrt{n}) \times (y_k^* - y)$ , and again (5.20) holds.

We choose  $y < y_k^*$  so that  $1 - G_k(y) > 0$  and conclude that

$$(5.21) \quad \tau_k^{(n)}(t) \Rightarrow t,$$

where the convergence is that of RCLL processes on  $[0, 1]$ .

We now use (5.8), (5.12) and (5.17) in (5.7) to obtain the more precise estimate

$$(5.22) \quad (1 - G_k(y)) \rho_k^{(n)} \sqrt{n} (t - \tau_k^{(n)}(t)) \\ \leq j(\widehat{W}_k^{(n)}) + \left[ \widehat{Z}_k^{(n)}\left(\tau_k^{(n)}(t) + \frac{1}{\sqrt{n}}(y_k^* - y)\right) - \widehat{Z}_k^{(n)}(\tau_k^{(n)}(t)) \right] \\ + \left[ \widehat{Z}_{k,y}^{(n)}(t) - \widehat{Z}_{k,y}^{(n)}\left(\tau_k^{(n)}(t) + \frac{1}{\sqrt{n}}(y_k^* - y)\right) \right] \\ + \left[ \widehat{W}_k^{(n)}(t) - \widehat{W}_k^{(n)}(\tau_k^{(n)}(t)) \right] - \left[ \widehat{Z}_k^{(n)}(t) - \widehat{Z}_k^{(n)}(\tau_k^{(n)}(t)) \right] \\ + (1 - G_k(y)) \rho_k^{(n)} (y_k^* - y)$$

in Case I, and the companion estimate  $\sqrt{n}(t - \tau_k^{(n)}(t)) \leq (y_k^* - y)$  in Case II. According to the differencing theorem (e.g., Theorem A.3 of [9]), every term on the right-hand side of (5.22), save the last, converges weakly to 0. Consequently,

$$(5.23) \quad (\sqrt{n}(t - \tau_k^{(n)}(t)) - (y_k^* - y))^+ \Rightarrow 0.$$

Because  $y < y_k^*$  is arbitrary and  $t - \tau_k^{(n)}(t) \geq 0$ , we conclude

$$(5.24) \quad \sqrt{n}(t - \tau_k^{(n)}(t)) \Rightarrow 0.$$

Armed with (5.24), we return once again to (5.17), this time taking  $y = y_k^*$  so that Case II is vacuous and concluding that  $\widehat{D}_k^{(n)} \Rightarrow 0$ . From (5.7), using (5.8), we see that

$$(5.25) \quad \widehat{W}_k^{(n)}(t) \left[ \widehat{C}_k^{(n)}(t), \widehat{F}_k^{(n)}(t) \right] \leq j(W_k^{(n)}) + \widehat{D}_k^{(n)}(t),$$

and hence  $\widehat{W}_k^{(n)}[\widehat{C}_k^{(n)}, \widehat{F}_k^{(n)}] \Rightarrow 0$ . For the second part of (5.4), we observe that

$$\begin{aligned}
 (5.26) \quad & \widehat{Q}_k^{(n)}(t) \left[ \widehat{C}_k^{(n)}(t), \widehat{F}_k^{(n)}(t) \right] \\
 & \leq \frac{1}{\sqrt{n}} \left[ 1 + A_k^{(n)}(nt) - A_k^{(n)}(n\tau_k^{(n)}(t)) \right] \\
 & = \frac{1}{\sqrt{n}} + \left[ \widehat{A}_k^{(n)}(t) - \widehat{A}_k^{(n)}(\tau_k^{(n)}(t)) \right] + \lambda_k^{(n)} \sqrt{n}(t - \tau_k^{(n)}(t)),
 \end{aligned}$$

and all terms on the right-hand side have limit 0.  $\square$

We may now follow the proofs of Section 3 of [9], obtaining

$$(5.27) \quad \left( \widehat{W}_1^{(n)}[\widehat{F}_1^{(n)}, \infty), \dots, \widehat{W}_K^{(n)}[\widehat{F}_K^{(n)}, \infty) \right) \Rightarrow (W_1^*, \dots, W_K^*)$$

in place of (3.27) of [9], using the continuity of the mapping

$$(5.28) \quad H^{-1} : [0, \infty)^K \xrightarrow{\text{onto}} \prod_{k=1}^K (-\infty, y_k^*]$$

defined by

$$(5.29) \quad H^{-1}(x_1, \dots, x_K) \triangleq (H_1^{-1}(x_1), \dots, H_K^{-1}(x_K))$$

to modify the proof of Proposition 3.10 of [9], and replacing  $\psi : \mathbb{R} \rightarrow \mathcal{M}$  in the proof of Theorem 3.1 of [9] by  $\psi : \mathbb{R}^K \rightarrow \mathcal{M}^K$  given by

$$\begin{aligned}
 (5.30) \quad & \psi(x_1, \dots, x_K)(B_1, \dots, B_K) \\
 & \triangleq \left( \int_{B_1 \cap [x_1, \infty)} (1 - G_1(\eta)) d\eta, \dots, \int_{B_K \cap [x_K, \infty)} (1 - G_K(\eta)) d\eta \right).
 \end{aligned}$$

**6. First-in, first-out limits.** In this section, we assume that customers within each of the  $K$  queues are served using the first-in, first-out (FIFO) queue discipline; that is, the server always services the customer with the longest time in the queue. This discipline is equivalent to EDF if all customers in each queue have initial lead times 0.

As we remarked at the end of Section 2, all the processes introduced in that section, with the exception of the queue lengths, are independent of the service discipline within each queue, provided that the server always allocates its total capacity in the same way among all queues that have work to do (i.e., the processes  $w_k^{(n)}$  defined in Section 2 remain unchanged). Thus, in this section, we adopt all the notation and the assumptions of Sections 2 and 3. We still assume that initial lead times are distributed as in (2.4), although we shall sometimes use the results of the previous section in the case that all initial lead times are 0 as a way of justifying assertions about the limiting process under FIFO. Assumptions 2.1 and 3.4 are still in force.

The aim of this section is to give a counterpart of Theorem 5.2 for the FIFO scheduling policy. We will start with the case of a single queue ( $K = 1$ ) in which the notation somewhat simplifies. In particular, we drop all the subscripts  $k$  in  $G_k$ ,  $y_k^*$ ,  $v_{k,j}^{(n)}$ , and so on, when a single queue is considered. It turns out that the ideas developed for this case can be easily generalized to  $K \geq 2$  arrival streams.

Denote by  $W^*$  the limiting (real-valued) workload in the queue, that is, the limit of  $\widehat{W}^{(n)}$ .  $W^*$  is the Brownian motion with drift described in Corollary 3.3 (with  $K = 1$ ).

The lead time of a customer in queue is the sum of two independent random variables: the lead time assigned upon arrival and the negative of the time since arrival. If the former is always 0, then the FIFO queue discipline coincides with EDF, and Theorem 5.2 implies that

$$\widehat{\mathcal{W}}^*(t)(B) = \int_{B \cap [F^*(t), \infty)} (1 - \mathbb{I}_{[0, \infty)}(y)) dy = m(B \cap [F^*(t), 0]),$$

where  $m$  denotes Lebesgue measure. In particular,  $W^*(t) = \widehat{\mathcal{W}}^*(t)(\mathbb{R}) = -F^*(t)$ . In other words,  $\widehat{\mathcal{W}}^*(t)$  is Lebesgue measure restricted to  $[-W^*(t), 0]$ .

When the lead times assigned upon arrival are nontrivial, it is reasonable to expect  $\widehat{\mathcal{W}}^*(t)$  to be the convolution of Lebesgue measure on  $[-W^*(t), 0]$  with  $dG$ , the distribution of scaled lead times. The following theorem confirms this conjecture.

**THEOREM 6.1.** *Let  $\widehat{\mathcal{W}}^*$  be the measure-valued process defined by*

$$\widehat{\mathcal{W}}^*(t)(B) \triangleq \int_B (G(\eta + W^*(t)) - G(\eta)) d\eta,$$

*that is, the process whose value is the convolution of  $dG$  with the uniform distribution on  $[-W^*(t), 0]$ . Let also  $\widehat{\mathcal{Q}}^* \triangleq \mu \widehat{\mathcal{W}}^*$ . In the case of a single queue, under the first-in, first-out queue discipline, we have*

$$(6.1) \quad \widehat{\mathcal{W}}^{(n)} \Rightarrow \widehat{\mathcal{W}}^*, \quad \widehat{\mathcal{Q}}^{(n)} \Rightarrow \widehat{\mathcal{Q}}^*.$$

Note that this result was conjectured in [9].

**PROOF.** For  $y \leq 0$ ,  $l \in \mathbb{R}$ , define

$$\mathcal{M}^{(n)}(t)([y, 0] \times (-\infty, l]) \triangleq \sum_{j=1}^{A^{(n)}(t)} v_j^{(n)} \mathbb{I}_{\{y \leq S_j^{(n)} - t\}} \mathbb{I}_{\{L_j^{(n)} \leq l\}}.$$

This is the arrived work (including departed work) whose time in queue is less than or equal to  $-y$  and whose lead time upon arrival was less than or equal to  $l$ . Similarly, for  $y \leq 0$ ,  $l \in \mathbb{R}$ , let  $\mathcal{N}^{(n)}(t)([y, 0] \times (-\infty, l])$  denote the work still present in queue whose time in queue is less than or equal to  $-y$  and whose lead

time upon arrival was less than or equal to  $l$ . The scaled versions of these two processes are

$$\widehat{\mathcal{N}}^{(n)}(t)([y, 0] \times (-\infty, l]) \triangleq \frac{1}{\sqrt{n}} \mathcal{N}^{(n)}(nt) \left( [\sqrt{n}y, 0] \times (-\infty, \sqrt{n}l] \right)$$

and

$$\begin{aligned} \widehat{\mathcal{M}}^{(n)}(t)([y, 0] \times (-\infty, l]) &\triangleq \frac{1}{\sqrt{n}} \mathcal{M}^{(n)}(nt) \left( [\sqrt{n}y, 0] \times (-\infty, \sqrt{n}l] \right) \\ &= \frac{1}{\sqrt{n}} \sum_{j=1}^{\infty} v_j^{(n)} \mathbb{I}_{\{\sqrt{n}y+nt \leq S_j^{(n)} \leq nt\}} \mathbb{I}_{\{L_j^{(n)} \leq \sqrt{n}l\}}. \end{aligned}$$

We have

$$\widehat{\mathcal{M}}^{(n)}(t)([y, 0] \times (-\infty, l]) = \widehat{Y}_m^{(n)}(t) + \widehat{U}_m^{(n)}(t),$$

where

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

Recall that, by Theorem 3.1, the process  $\widehat{V}_l^{(n)}$  defined by (3.9) (recall that we drop the subscript  $k = 1$ ) converges weakly to a Brownian motion  $V_l^*$ . Furthermore,

$$\widehat{Y}^{(n)}(t) = \widehat{V}_l^{(n)} \left( \frac{1}{n} A^{(n)}(nt) \right) - \widehat{V}_l^{(n)} \left( \frac{1}{n} A^{(n)}((nt + \sqrt{n}y)^+ -) \right).$$

But

$$\widehat{A}^{(n)}(t) = \frac{1}{\sqrt{n}} [A^{(n)}(nt) - \lambda^{(n)}nt],$$

so

$$\begin{aligned} \frac{1}{n} A^{(n)}(nt) &= \frac{1}{\sqrt{n}} \widehat{A}^{(n)}(t) + \lambda^{(n)}t \Rightarrow \lambda t, \\ \frac{1}{n} A^{(n)}((nt + \sqrt{n}y)^+ -) &= \frac{1}{\sqrt{n}} \widehat{A}^{(n)} \left( \left( t + \frac{y}{\sqrt{n}} \right)^+ - \right) + \lambda^{(n)} \left( t + \frac{y}{\sqrt{n}} \right)^+ \Rightarrow \lambda t. \end{aligned}$$

The differencing theorem implies that  $\widehat{Y}^{(n)} \Rightarrow 0$ .

On the other hand,

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

The first term on the right-hand side of (6.2) has limit 0. Indeed,

$$\widehat{A}^{(n)}(t) - \widehat{A}^{(n)}\left(\left(t + \frac{y}{\sqrt{n}}\right)^+ -\right) \Rightarrow 0$$

by the differencing theorem and

$$\widehat{A}^{(n)}\left(\left(t + \frac{y}{\sqrt{n}}\right)^+ -\right) - \widehat{A}^{(n)}\left(\left(t + \frac{y}{\sqrt{n}}\right)^+ -\right) \Rightarrow 0$$

because  $\widehat{A}^{(n)}$  converges to a continuous limit and hence the maximal jump of  $\widehat{A}^{(n)}$  on any finite time horizon converges to 0 (compare the proof of Corollary 3.5). The second term on the right-hand side of (6.2) has limit 0 because  $\lambda^{(n)}/\mu^{(n)} \rightarrow 1$  and  $0 \leq \sqrt{nt} - (\sqrt{nt} + y)^+ \leq -y$ .

Thus, we have proved that, for every  $T > 0$ ,  $y \leq 0$  and  $l \in \mathbb{R}$ ,

$$(6.3) \quad \sup_{0 \leq t \leq T} \left| \widehat{\mathcal{M}}^{(n)}(t)([y, 0] \times (-\infty, l]) + ((\sqrt{nt} + y)^+ - \sqrt{nt})G(l) \right| \xrightarrow{P} 0.$$

Clearly,  $(\sqrt{nt} + y)^+ - \sqrt{nt}$  is Lipschitz-continuous in  $y$  (with the Lipschitz constant 1). Moreover, both  $\widehat{\mathcal{M}}^{(n)}(t)([y, 0] \times (-\infty, l])$  and  $((\sqrt{nt} + y)^+ - \sqrt{nt})G(l)$  are monotone in  $y$ . Thus, by the same argument as in Proposition 3.4 of [9], we can upgrade (6.3) to

$$(6.4) \quad \sup_{y_0 \leq y \leq 0} \sup_{0 \leq t \leq T} \left| \widehat{\mathcal{M}}^{(n)}(t)([y, 0] \times (-\infty, l]) + ((\sqrt{nt} + y)^+ - \sqrt{nt})G(l) \right| \xrightarrow{P} 0$$

for every  $T > 0$ ,  $y_0 \leq 0$  and  $l \in \mathbb{R}$ . In particular, taking  $l = y^*$  in (6.4) we get, by (2.4) and (2.5),

$$(6.5) \quad \sup_{y_0 \leq y \leq 0} \sup_{0 \leq t \leq T} \left| \widehat{\mathcal{M}}^{(n)}(t)([y, 0] \times \mathbb{R}) + (\sqrt{nt} + y)^+ - \sqrt{nt} \right| \xrightarrow{P} 0.$$

A slight modification of the above argument yields

$$(6.6) \quad \sup_{y_0 \leq y \leq 0} \sup_{0 \leq t \leq T} |\widehat{\mathcal{M}}^{(n)}(t)([y, 0] \times (-\infty, l)) + ((\sqrt{nt} + y)^+ - \sqrt{nt})G(l-)| \xrightarrow{P} 0$$

for every  $T > 0$ ,  $y_0 \leq 0$  and  $l \in \mathbb{R}$ . Using (6.6) and proceeding as in the proof of the Glivenko–Cantelli theorem (see, e.g., Theorem 20.6 in [1]), we can finally upgrade (6.4) to

$$(6.7) \quad \sup_{l \in \mathbb{R}} \sup_{y_0 \leq y \leq 0} \sup_{0 \leq t \leq T} |\widehat{\mathcal{M}}^{(n)}(t)([y, 0] \times (-\infty, l)) + ((\sqrt{nt} + y)^+ - \sqrt{nt})G(l)| \xrightarrow{P} 0$$

for every  $T > 0$ ,  $y_0 \leq 0$ . Using (6.5) and proceeding as in the proof of Corollary 3.5 of [9], we get

$$\sup_{y_0 \leq y \leq 0} \sup_{0 \leq t \leq T} \widehat{\mathcal{M}}^{(n)}(t)(\{y\} \times \mathbb{R}) \xrightarrow{P} 0$$

for every  $T > 0$ ,  $y_0 \leq 0$ .

Let us introduce the notation

$$F^{(n)}(t) \triangleq C^{(n)}(t) \triangleq \begin{cases} \text{The negative of the time in queue} \\ \text{of the customer currently in service} \\ \text{or } S_{A^{(n)}(t)}^{(n)} - t \text{ if the queue is empty} \end{cases}.$$

Notice that  $C^{(n)}$  and  $F^{(n)}$  defined above would coincide with those introduced in Section 4 (and used to study the EDF queue discipline) if all the initial lead times of the customers in queue were 0. Denote also by  $\widehat{C}^{(n)}$  and  $\widehat{F}^{(n)}$  the rescaled versions of  $C^{(n)}$  and  $F^{(n)}$ :

$$\widehat{C}^{(n)}(t) \triangleq \widehat{F}^{(n)}(t) \triangleq \frac{1}{\sqrt{n}} F^{(n)}(nt).$$

By Corollary 3.8 in [9] (applied to the case of zero initial lead times), we have

$$\widehat{\mathcal{N}}^{(n)}(\{\widehat{F}^{(n)}\} \times \mathbb{R}) = \widehat{\mathcal{N}}^{(n)}([\widehat{C}^{(n)}, \widehat{F}^{(n)}] \times \mathbb{R}) \Rightarrow 0.$$

Also, Proposition 3.10 in [9] (or Proposition 5.1 of the present paper) yields

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

By a slight modification of the proof of Proposition 3.11 in [9], we can show that

$$\sup_{l \in \mathbb{R}} \sup_{y \leq 0} \sup_{0 \leq t \leq T} |\widehat{\mathcal{N}}^{(n)}(t)([y, 0] \times (-\infty, l)) + (y \vee \widehat{F}^{(n)}(t))G(l)| \xrightarrow{P} 0.$$

But

$$(y \vee \widehat{F}^{(n)}(t)) \Rightarrow (y \vee -W^*(t)).$$

Thus, we have proved that  $\widehat{\mathcal{N}}^{(n)}(t)$  converges weakly to the product measure with the uniform distribution on  $[-W^*(t), 0]$  in one coordinate and the cumulative distribution function  $G$  in the other one. Let us recall that the first coordinate indicates the *negative* of the time in queue and the second one the initial lead times of customers.

The lead time of a customer is equal to the sum of its initial lead time and the negative of the time it has already spent in queue. The mapping from  $(-\infty, 0] \times \mathbb{R}$  into  $\mathbb{R}$  given by  $(y, l) \rightarrow l + y$  induces a continuous function from the space of measures on  $(-\infty, 0] \times \mathbb{R}$  to measures on the real line. Thus, we get the first part of (6.1). The proof of the second part is analogous, with  $v_j^{(n)}$  and  $\mu^{(n)}$  replaced by 1.  $\square$

As in the case of the EDF service discipline, the above result can be easily generalized to the case of  $K$  arrival processes. The corresponding result is

**THEOREM 6.2.** *Let  $\widehat{\mathcal{W}}_k^*$  be the measure-valued process defined by*

$$(6.8) \quad \widehat{\mathcal{W}}_k^*(t)(B) \triangleq \rho_k \int_B \left( G_k \left( \eta + \frac{W_k^*(t)}{\rho_k} \right) - G_k(\eta) \right) d\eta,$$

*that is, the process whose distribution is the convolution of  $dG_k$  and the uniform distribution on  $[-W_k^*(t)/\rho_k, 0]$  with density  $\rho_k$ . Also, let  $\widehat{\mathcal{Q}}_k^* \triangleq \mu_k \widehat{\mathcal{W}}_k^*$ ,  $\widehat{\mathcal{W}}^* \triangleq (\widehat{\mathcal{W}}_1^*, \dots, \widehat{\mathcal{W}}_K^*)$ ,  $\widehat{\mathcal{Q}}^* \triangleq (\widehat{\mathcal{Q}}_1^*, \dots, \widehat{\mathcal{Q}}_K^*)$ . We have, in the case of the first-in, first-out queue discipline,*

$$(\widehat{\mathcal{W}}_1^{(n)}, \dots, \widehat{\mathcal{W}}_K^{(n)}) \Rightarrow \widehat{\mathcal{W}}^*, \quad (\widehat{\mathcal{Q}}_1^{(n)}, \dots, \widehat{\mathcal{Q}}_K^{(n)}) \Rightarrow \widehat{\mathcal{Q}}^*.$$

The proof of Theorem 6.2 is an easy generalization of the argument proving Theorem 6.1, where the necessary results from [9] are replaced by their  $K$ -dimensional counterparts discussed in Section 5. In particular, Proposition 5.1 should be used in place of Proposition 3.10 from [9].

**COROLLARY 6.3.** *Under the first-in, first-out queue discipline, the  $K$ -dimensional scaled queue length process  $(\widehat{\mathcal{Q}}_k^{(n)})_{k=1, \dots, K}$  converges weakly to  $(\mu_k \widehat{\mathcal{W}}_k^*)_{k=1, \dots, K}$ .*

Corollary 6.3 is an example of “state space collapse,” a notion introduced by Reiman [18] to capture the idea that one limiting process (scaled queue length) is a deterministic function of another limiting process (scaled workload). Corollary 6.3 is close to the results of Bramson [3] and Williams [19], who obtain the relationship between individual queue lengths and total workload under head-of-the-line proportional processor sharing, but not the weighted processor sharing of this paper.

**7. Simulations.** 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 corresponding 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 Sections 2 and 3. We denoted the queue length for the  $k$ th input stream in the  $n$ th queueing system by  $Q_k^{(n)}$ , and the scaled queue length by  $\widehat{Q}_k^{(n)} \triangleq (1/\sqrt{n})Q_k^{(n)}$ . For large values of  $n$ , under both the earliest-deadline-first (EDF) and the first-in, first-out (FIFO) disciplines,  $\widehat{Q}_k^{(n)}$  is approximately equal to  $Q_k^* = \mu_k W_k^*$  (Corollaries 5.3 and 6.3), where  $W_k^*$  is the limit of  $\widehat{W}_k^{(n)}$  as in Assumption 3.4. Similarly, the measure-valued processes  $\widehat{Q}_k^{(n)}$  converges to  $\mu_k \widehat{W}_k^*$ , where (recalling equations from earlier sections)

$$(5.3) \quad \widehat{W}_k^*(t)(B) \triangleq \rho_k \int_{B \cap [F_k^*(t), \infty)} (1 - G_k(y)) dy,$$

under EDF (Theorem 5.2), and

$$(6.8) \quad \widehat{W}_k^*(t)(B) \triangleq \rho_k \int_B \left( G_k \left( \eta + \frac{W_k^*(t)}{\rho_k} \right) - G_k(\eta) \right) d\eta,$$

under FIFO (Theorem 6.2). In (5.3), the process  $F_k^*$  is given by

$$(5.2) \quad F_k^*(t) \triangleq H_k^{-1} \left( \frac{W_k^*(t)}{\rho_k} \right).$$

Recall that customers arrive with initial lead time (deadline) distribution

$$(2.4) \quad \mathbb{P}\{L_{k,j}^{(n)} \leq \sqrt{n}y\} = G_k(y).$$

We define  $G_k^{(n)}(x) = G_k(x/\sqrt{n})$ , so that  $\mathbb{P}\{L_{k,j}^{(n)} \leq x\} = G_k^{(n)}(x)$  is the cumulative distribution function of initial lead times in the  $k$ th input stream of the  $n$ th queueing system. The process  $F_k^*(t)$  of (5.2) is characterized in terms of the function

$$(5.1) \quad H_k(y) \triangleq \int_y^\infty (1 - G_k(\eta)) d\eta.$$

We introduce the function

$$(7.1) \quad H_k^{(n)}(x) \triangleq \sqrt{n} H_k \left( \frac{x}{\sqrt{n}} \right) = \int_x^\infty (1 - G_k^{(n)}(\xi)) d\xi.$$



Under EDF, for large values of  $n$ ,  $\widehat{Q}_k^{(n)}(y, \infty) \approx \lambda_k H_k(y \vee F_k^*)$ , and because  $H_k$  is nonincreasing,

$$\begin{aligned}
 \lambda_k H_k(y \vee F_k^*) &= \lambda_k H_k\left(y \vee H_k^{-1}\left(\frac{W_k^*}{\rho_k}\right)\right) \\
 &= \lambda_k H_k(y) \wedge \mu_k W_k^* \\
 (7.2) \qquad &= \left(\lambda_k \frac{1}{\sqrt{n}} H_k^{(n)}(\sqrt{n}y)\right) \wedge Q_k^* \\
 &\approx \left(\lambda_k \frac{1}{\sqrt{n}} H_k^{(n)}(\sqrt{n}y)\right) \wedge \frac{1}{\sqrt{n}} Q_k^{(n)}.
 \end{aligned}$$

Multiplying (7.2) by  $\sqrt{n}$ , we obtain

$$(7.3) \qquad Q_k^{(n)}(x, \infty) = \sqrt{n} \widehat{Q}_k^{(n)}\left(\frac{x}{\sqrt{n}}, \infty\right) \approx (\lambda_k H_k^{(n)}(x)) \wedge Q_k^{(n)}.$$

Relation (7.3) connects the unscaled queue length  $Q_k^{(n)}$  with the number of customers whose unscaled lead times exceed  $x$ , and the function  $H_k^{(n)}$  appearing in this relation can be computed from the cumulative distribution function  $G_k^{(n)}$  for the unscaled lead times. Relation (7.3) can be tested by simulation under the earliest-deadline-first discipline without knowledge of the parameter  $n$ .

Under FIFO, we have from (6.8) that

$$\begin{aligned}
 \widehat{Q}_k^{(n)}(y, \infty) &\approx \lambda_k \int_y^\infty \left(G_k\left(\eta + \frac{W_k^*}{\rho_k}\right) - G_k(\eta)\right) d\eta \\
 (7.4) \qquad &= \frac{\lambda_k}{\sqrt{n}} \int_{\sqrt{n}y}^\infty \left(G_k^{(n)}\left(\xi + \frac{\sqrt{n}W_k^*}{\rho_k}\right) - G_k^{(n)}(\xi)\right) d\xi.
 \end{aligned}$$

Replacing  $\sqrt{n}W_k^*$  in (7.4) by the approximation  $(1/\mu_k)Q_k^{(n)}$  and multiplying by  $\sqrt{n}$ , we obtain

$$\begin{aligned}
 Q_k^{(n)}(x, \infty) &= \sqrt{n} \widehat{Q}_k^{(n)}\left(\frac{x}{\sqrt{n}}, \infty\right) \\
 (7.5) \qquad &\approx \lambda_k \int_x^\infty \left(G_k^{(n)}\left(\xi + \frac{1}{\lambda_k} Q_k^{(n)}\right) - G_k^{(n)}(\xi)\right) d\xi.
 \end{aligned}$$

Once again, we have a relation that connects the unscaled queue length  $Q_k^{(n)}$  with the number of customers whose unscaled lead times exceed  $x$ , and the only function appearing in (7.5) is the unscaled lead time distribution  $G_k^{(n)}$ . Relation (7.5) can be tested by simulation under the first-in, first-out discipline without knowledge of the parameter  $n$ .

In each of Figures 1–8, there are eight input streams creating an overall traffic intensity of 0.96. In Figures 1, 3, 5 and 7, this traffic intensity is divided equally among the eight streams, whereas in Figures 2, 4, 6 and 8, the eighth stream is

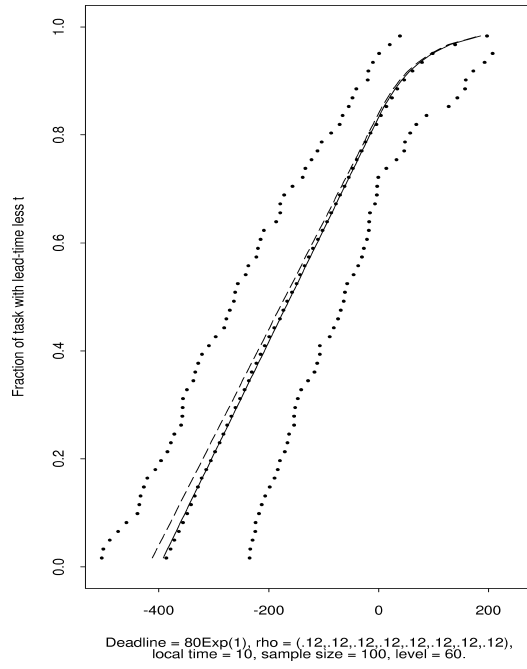


FIG. 1. *M/M/1, EDF, heavy-traffic queue.*

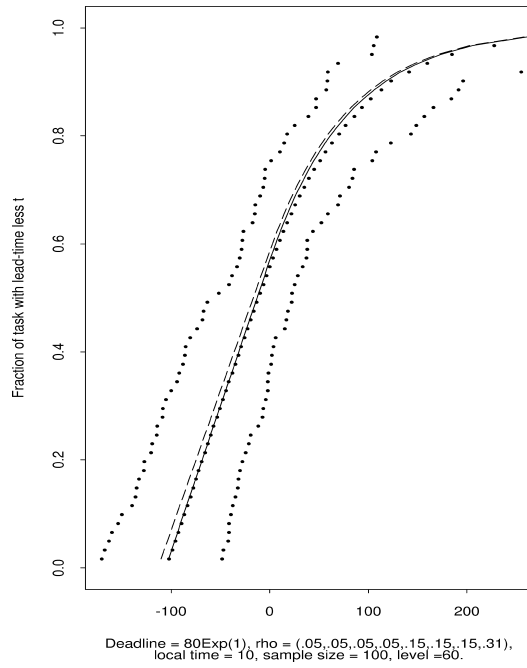
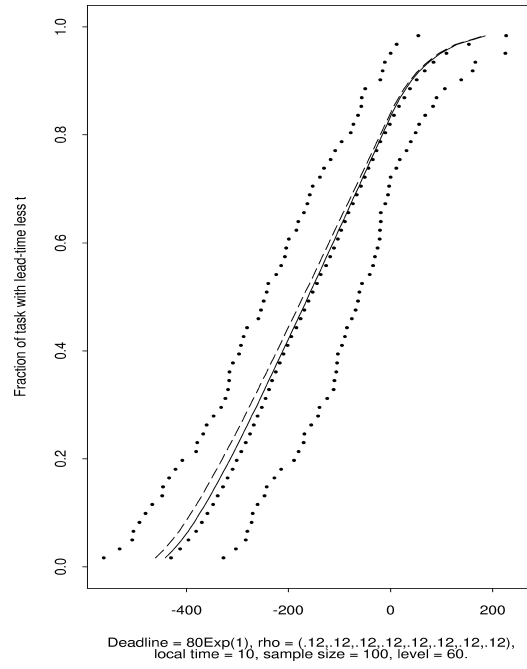
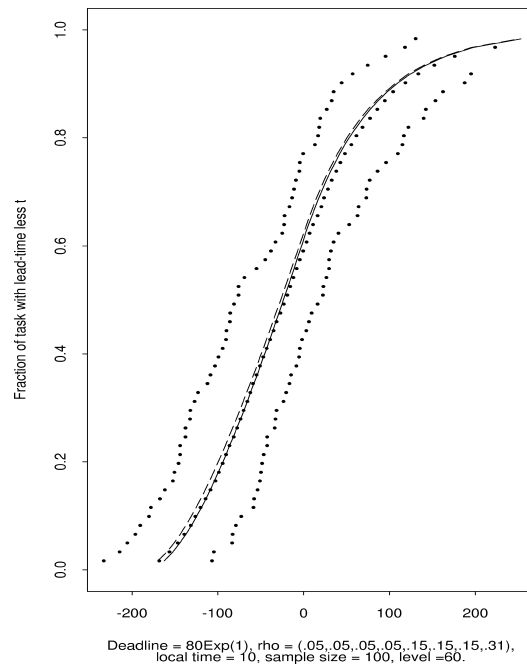


FIG. 2. *M/M/1, EDF, heavy-traffic queue.*

FIG. 3. *M/M/1, FIFO, heavy-traffic queue.*FIG. 4. *M/M/1, FIFO, heavy-traffic queue.*

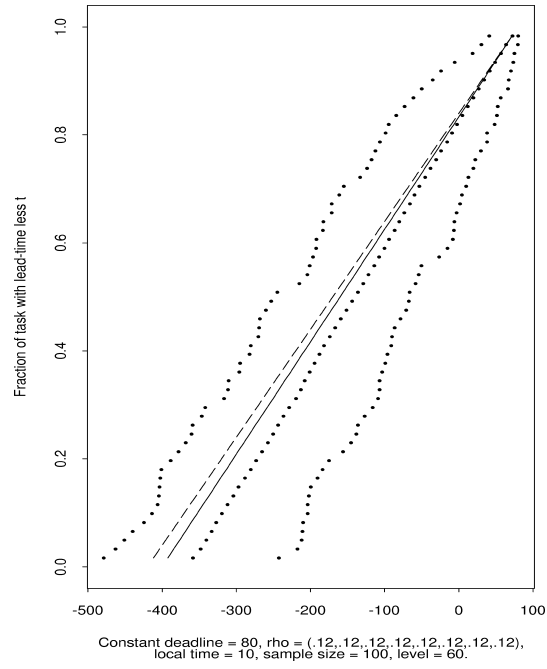


FIG. 5. *M/M/1, heavy-traffic queue.*

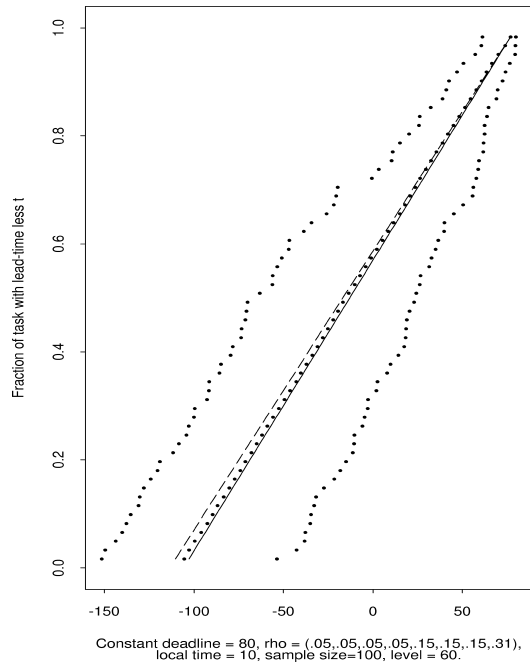
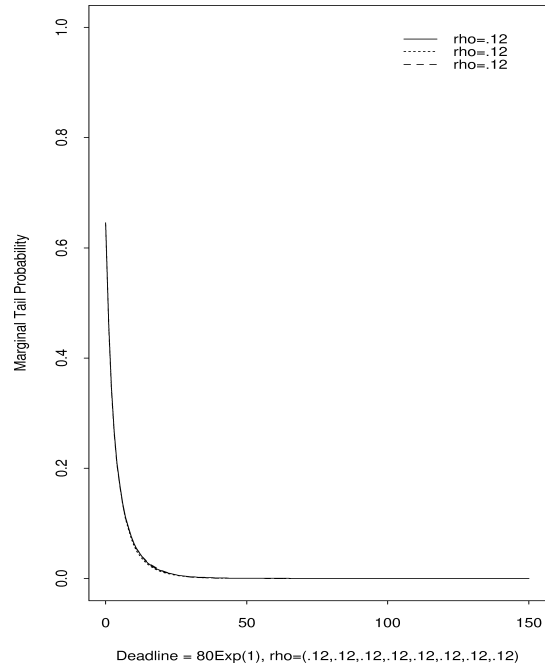
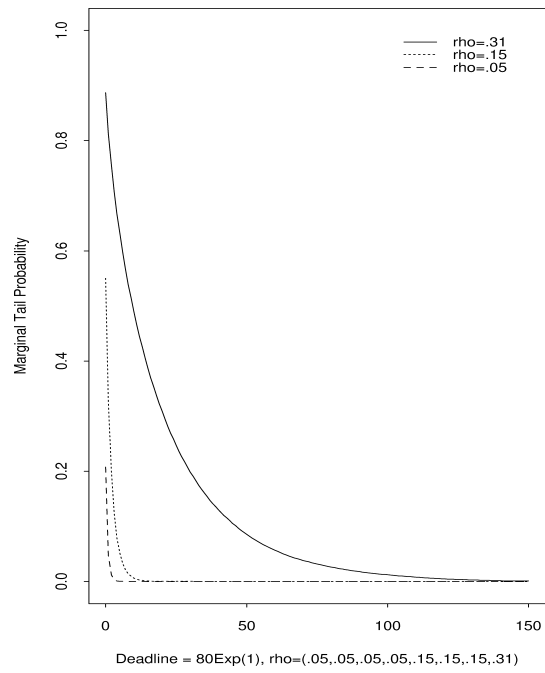


FIG. 6. *M/M/1, heavy-traffic queue.*

FIG. 7. *M/M/1, EDF.*FIG. 8. *M/M/1, EDF.*

dominant. The traffic intensities  $\rho_k^{(n)}$  are indicated in the captions; the service rates in all cases are  $\mu_k^{(n)} = 1$ , so that  $\lambda_k^{(n)} = \rho_k^{(n)}$ . In Figures 1–6, each simulation run is initiated with an empty queue and a local time is accumulated when the queue length of the eighth queue is exactly equal to 60. To avoid excessive dependence between successive samples, at the instant the local time reaches the value 10, a sample is taken and the local time counter is reset to 0. The simulation continues until 100 samples are recorded. In particular, in the formulas (7.3) and (7.5), we used the value  $Q_8^{(n)} = 60$ . In Figures 1–4,  $G_8(x) = (1 - e^{-x/80})\mathbb{I}_{\{x \geq 0\}}$ . In Figures 5 and 6,  $G_8(x) = \mathbb{I}_{\{x \geq 80\}}$ ; in this case, EDF and FIFO coincide, as do the right-hand sides of (7.3) and (7.5).

The leftmost dots indicate the pointwise minimum empirical cumulative distribution function of the lead time profile for these 100 samples, the rightmost dots indicate the pointwise maximum and the central dots are the average. As a function of  $x$ , the right-hand side of (7.3) is plotted as a dashed curve in Figures 1, 2, 5 and 6. In Figures 3, 4, 5 and 6, the right-hand side of (7.5) is plotted as a dashed curve.

We obtained the solid curves in Figures 1–6 by replacing  $\lambda_8^{(n)}$  in the right-hand sides of (7.3) and (7.5) by  $\lambda_8^{(n)}/0.96$ . This normalization by the total traffic intensity causes the theory to have better predictive value. Indeed, with this normalization, the theoretical cumulative distribution functions and the pointwise average empirical cumulative distribution functions are in almost perfect agreement, except in Figure 5. The cause of the disagreement in this figure, and indeed the reasonableness of multiplying the approximations (7.2) and (7.4) by  $\sqrt{n}$  in order to obtain the relations (7.3) and (7.5), is a subject for future research.

Finally, in Appendix A, it is shown that when there is a dominant traffic intensity, the scaled workloads and queue lengths in the nondominant queues converge to 0. Figures 7 and 8 graph the empirical marginal tail probabilities of the queue length processes. In Figure 7, one of the identical eight streams with common traffic intensity 0.12 is shown. In Figure 8, there are eight streams of three types, having traffic intensities 0.05, 0.15 or 0.31. The empirical marginal tail probabilities are shown for one stream of each type. The behavior in these two cases is quite different. Figure 8 shows that the nondominant traffic flows become a negligible part of the workload. Figure 7 shows that in the equal traffic intensity case, the workload will be evenly spread across the different streams, none of which becomes negligible.

## APPENDIX A

This appendix provides a sufficient condition for Assumption 3.4 to hold, the condition described immediately following that assumption. For this appendix, we adopt the notation and assumptions of Sections 2 and 3, except, of course, Assumption 3.4 and its consequence Corollary 3.5. In its place we impose

Assumption A.2. A technical “crushing lemma” needed for this appendix is contained in Appendix B.

Recall that, for  $k = 1, \dots, K$ ,  $R_k^{(n)}(t)$  defined by (2.17) is the work performed by the server on queue  $k$  up to time  $t$ . Recall also the notation  $e(t) = t$  for all  $t \geq 0$ .

LEMMA A.1. *For every  $k$ , we have*

$$\frac{1}{n} R_k^{(n)} \circ ne \Rightarrow \rho_k e.$$

PROOF. From Corollary 3.3, we see that  $\frac{1}{n} \sum_{k=1}^K W_k^{(n)} \circ ne \Rightarrow 0$ , and since each term in the sum is nonnegative, we have  $\frac{1}{n} W_k^{(n)} \circ ne \Rightarrow 0$  for each  $k$ . Equation (2.18) implies

$$\frac{1}{n} R_k^{(n)}(nt) = -\frac{1}{n} W_k^{(n)}(nt) + \frac{1}{n} [V_k^{(n)}(A_k^{(n)}(nt)) - \rho_k^{(n)} nt] + \rho_k^{(n)} t,$$

and the right-hand side converges to  $\rho_k t$  because of Corollary 3.2.  $\square$

To proceed further, we make the following assumption.

ASSUMPTION A.2. *There exist positive constants  $w_k$ ,  $k = 1, \dots, K$ , satisfying (3.18) and such that, for every  $T > 0$ , the sequence of random variables*

$$\sqrt{n} \sup_{0 \leq t \leq nT} |w_k^{(n)}(t) - w_k|,$$

$k = 1, \dots, K$ ,  $n = 1, 2, \dots$ , *is tight.*

This, of course, implies that  $w_k^{(n)} \Rightarrow w_k$  as  $n \rightarrow \infty$  for  $k = 1, \dots, K$ .

PROPOSITION A.3. *Under Assumption A.2, we have  $\widehat{W}_k^{(n)} \Rightarrow 0$  for all  $k < K$ .*

COROLLARY A.4.  *$\widehat{W}_K^{(n)}$  converges weakly to a reflected Brownian motion with drift.*

Corollary A.4 follows immediately from Proposition A.3 and Corollary 3.3. We devote the remainder of the appendix to the proof of the proposition.

PROOF OF PROPOSITION A.3. For simplicity, we shall prove convergence of the workloads  $\widehat{W}_k^{(n)}$ ,  $k < K$ , on  $[0, 1]$ . Equations (2.14) and (2.19) imply

$$(A.1) \quad \dot{R}_k^{(n)}(t) = \sum_{m=0}^{K-1} \sum_{\alpha \in A_k^c(m)} \frac{w_k^{(n)}(t)}{1 - \sum_{\ell \in \alpha} w_\ell^{(n)}(t)} j_\alpha^{(n)}(t).$$

Choose  $0 < c_0 < \min_{j=1, \dots, K} w_j$  and let

$$(A.2) \quad E_n = \{w_k^{(n)}(t) > c_0, t \in [0, n], k = 1, \dots, K\}.$$

By Assumption A.2,

$$(A.3) \quad \mathbb{P}(E_n) \rightarrow 1.$$

On  $E_n$ , we have, for  $k < K$ ,

$$\begin{aligned} \frac{\dot{R}_K^{(n)}(t)}{w_K^{(n)}(t)} - \frac{\dot{R}_k^{(n)}(t)}{w_k^{(n)}(t)} &= \sum_{m=1}^{K-1} \sum_{\alpha \in A_K^c(m) \cap A_k(m)} \frac{1}{1 - \sum_{l \in \alpha} w_l^{(n)}(t)} j_\alpha^{(n)}(t) \\ &\quad - \sum_{m=1}^{K-1} \sum_{\alpha \in A_K(m) \cap A_k^c(m)} \frac{1}{1 - \sum_{l \in \alpha} w_l^{(n)}(t)} j_\alpha^{(n)}(t) \\ &\leq \sum_{m=1}^{K-1} \sum_{\alpha \in A_k(m)} \frac{1}{1 - \sum_{l \in \alpha} w_l^{(n)}(t)} j_\alpha^{(n)}(t) \\ &\leq \frac{1}{c_0} \dot{i}_k^{(n)}(t). \end{aligned}$$

But then

$$\begin{aligned} \dot{U}_k^{(n)}(t) &= \sum_{m=1}^K \sum_{\alpha \in A_k(m)} \frac{w_k^{(n)}(t)}{w_k^{(n)}(t) + 1 - \sum_{l \in \alpha} w_l^{(n)}(t)} j_\alpha^{(n)}(t) \\ &\geq w_k^{(n)}(t) \sum_{m=1}^K \sum_{\alpha \in A_k(m)} j_\alpha^{(n)}(t) \\ &= w_k^{(n)}(t) \dot{i}_k^{(n)}(t) \\ &\geq c_0 w_k^{(n)}(t) \left( \frac{\dot{R}_K^{(n)}(t)}{w_K^{(n)}(t)} - \frac{\dot{R}_k^{(n)}(t)}{w_k^{(n)}(t)} \right) \end{aligned}$$

and

$$\begin{aligned} \dot{T}_k^{(n)}(t) &= \dot{R}_k^{(n)}(t) + \dot{U}_k^{(n)}(t) \\ &\geq \dot{R}_k^{(n)}(t) + c_0 w_k^{(n)}(t) \left( \frac{\dot{R}_K^{(n)}(t)}{w_K^{(n)}(t)} - \frac{\dot{R}_k^{(n)}(t)}{w_k^{(n)}(t)} \right). \end{aligned}$$

We define

$$(A.4) \quad \tilde{T}_k^{(n)}(t) \triangleq \begin{cases} R_k^{(n)}(t) + c_0 \int_0^t w_k^{(n)}(s) \left( \frac{\dot{R}_K^{(n)}(s)}{w_K^{(n)}(s)} - \frac{\dot{R}_k^{(n)}(s)}{w_k^{(n)}(s)} \right) ds, & \text{on } E_n, \\ T_k^{(n)}(t), & \text{on } E_n^c. \end{cases}$$



Then

$$(A.5) \quad \dot{T}_k^{(n)}(t) \geq \frac{d}{dt} \tilde{T}_k^{(n)}(t) \geq 0$$

for all  $t \geq 0$ , and, after integration, we get  $T_k^{(n)} \geq \tilde{T}_k^{(n)}$ . By Assumption A.2, (A.3) and Lemma A.1,

$$\frac{1}{n} \tilde{T}_k^{(n)} \circ ne \Rightarrow \rho_k e + c_0 w_k \left( \frac{\rho_K}{w_K} - \frac{\rho_k}{w_k} \right) e > \rho_k e.$$

We define

$$\tilde{N}_k^{(n)}(t) = V_k^{(n)}(A_k^{(n)}(t)) - \tilde{T}_k^{(n)}(t),$$

so that  $N_k^{(n)} \leq \tilde{N}_k^{(n)}$ . We set  $\tilde{U}_k^{(n)}(t) = -\min_{0 \leq s \leq t} \tilde{N}_k^{(n)}(s)$  and  $\tilde{W}_k^{(n)}(t) = \tilde{N}_k^{(n)}(t) + \tilde{U}_k^{(n)}(t)$ . We show that

$$(A.6) \quad W_k^{(n)} \leq \tilde{W}_k^{(n)}.$$

Indeed, let us define

$$\tau = \sup\{t \geq 0; W_k^{(n)}(s) \leq \tilde{W}_k^{(n)}(s) \text{ for all } s \in [0, t]\}.$$

We wish to prove  $\tau = \infty$ , so assume, on the contrary, that  $\tau < \infty$ . Note that, for every  $t$ , we have

$$\begin{aligned} N_k^{(n)}(t) - N_k^{(n)}(t-) &= \tilde{N}_k^{(n)}(t) - \tilde{N}_k^{(n)}(t-) \\ &= V_k^{(n)}(A_k^{(n)}(t)) - V_k^{(n)}(A_k^{(n)}(t-)) \geq 0, \end{aligned}$$

and, consequently,

$$W_k^{(n)}(t) - W_k^{(n)}(t-) = \tilde{W}_k^{(n)}(t) - \tilde{W}_k^{(n)}(t-) \geq 0.$$

By the definition of  $\tau$ , we have  $W_k^{(n)}(\tau-) \leq \tilde{W}_k^{(n)}(\tau-)$  and hence  $W_k^{(n)}(\tau) \leq \tilde{W}_k^{(n)}(\tau)$ . If  $W_k^{(n)}(\tau) > 0$ , then, for some  $\varepsilon > 0$ , we have  $W_k^{(n)}(t) > 0$  for  $\tau \leq t \leq \tau + \varepsilon$ . It follows that  $U_k^{(n)}(t) = U_k^{(n)}(\tau)$  for  $\tau \leq t \leq \tau + \varepsilon$  and so, for this range of  $t$ , we have

$$\begin{aligned} W_k^{(n)}(t) &= W_k^{(n)}(t) - W_k^{(n)}(\tau) + W_k^{(n)}(\tau) \\ &\leq N_k^{(n)}(t) - N_k^{(n)}(\tau) + \tilde{W}_k^{(n)}(\tau) \\ &\leq \tilde{N}_k^{(n)}(t) - \tilde{N}_k^{(n)}(\tau) + \tilde{W}_k^{(n)}(\tau) \leq \tilde{W}_k^{(n)}(t) \end{aligned}$$

[the second inequality follows from (A.5)]. This contradicts the definition of  $\tau$ . It follows that we must have  $W_k^{(n)}(\tau) = 0$ , and, in particular,  $W_k^{(n)}(\tau) - W_k^{(n)}(\tau-) = 0$ ; that is, there is no arrival at time  $\tau$ . Thus, there is an  $\varepsilon > 0$  with the property that there is no arrival in the interval  $[\tau, \tau + \varepsilon]$ , and, in particular,  $W_k^{(n)}(t) = 0$  for

$\tau \leq t \leq \tau + \varepsilon$ . Again, the definition of  $\tau$  is contradicted. We conclude that  $\tau = \infty$ , which is equivalent to (A.6).

Now we want to show that, for  $k = 1, \dots, K - 1$ ,

$$(A.7) \quad \widehat{W}_k^{(n)}(t) \Rightarrow 0,$$

that is, that the workloads corresponding to lower traffic intensities get crushed under the heavy-traffic scaling. We prove this by induction on  $k$ . Let  $k_0 \leq K - 1$  be given and suppose that (A.7) is true for all  $k = 1, \dots, k_0 - 1$  (in particular, no assumption is necessary to consider  $k_0 = 1$ ). The crushing lemma (see Lemma B.4) applied to

$$\begin{aligned} \overline{Z}^{(n)} &\triangleq \widehat{Z}_{k_0}^{(n)}, \\ Z^* &\triangleq V_{k_0}^* \circ \lambda_{k_0} e + \frac{1}{\mu_{k_0}} A_{k_0}^*, \\ \overline{T}^{(n)}(t) &\triangleq \frac{1}{n} \widetilde{T}_{k_0}^{(n)}(nt) - \rho_{k_0}^{(n)} t \end{aligned}$$

[which is Lipschitz-continuous with constant  $L := 1$  because  $0 \leq (d/dt) \widetilde{T}_{k_0}^{(n)}(t) \leq 1$  for all  $n, t$ ] and

$$T^*(t) \triangleq c_0 w_{k_0} \left( \frac{\rho_K}{w_K} - \frac{\rho_{k_0}}{w_{k_0}} \right) t$$

yields  $\widetilde{\tau}_{k_0}^{(n)}(t) \Rightarrow t$ , where  $\widetilde{\tau}_{k_0}^{(n)}(t) \triangleq \sup\{s \in [0, t] : (1/\sqrt{n}) \widetilde{W}_{k_0}^{(n)}(ns) = 0\}$ . Thus, because  $0 \leq W_k^{(n)} \leq \widetilde{W}_k^{(n)}$ ,

$$(A.8) \quad \tau_{k_0}^{(n)}(t) \Rightarrow t,$$

where  $\tau_{k_0}^{(n)}(t) \triangleq \sup\{s \in [0, t] : (1/\sqrt{n}) W_{k_0}^{(n)}(ns) = 0\}$ . By the definition of  $\tau_{k_0}^{(n)}$ , we also have

$$(A.9) \quad \frac{1}{\sqrt{n}} W_{k_0}^{(n)}(n\tau_{k_0}^{(n)}(t)-) = 0.$$

For all  $k = 1, \dots, K$ , we show that

$$(A.10) \quad \widehat{Z}_k^{(n)}(t) - \widehat{Z}_k^{(n)}(\tau_{k_0}^{(n)}(t)-) \Rightarrow 0.$$

Indeed,

$$\begin{aligned} &\widehat{Z}_k^{(n)}(t) - \widehat{Z}_k^{(n)}(\tau_{k_0}^{(n)}(t)-) \\ &= \left[ \widehat{Z}_k^{(n)}(t) - \widehat{Z}_k^{(n)}(\tau_{k_0}^{(n)}(t)) \right] + \left[ \widehat{Z}_k^{(n)}(\tau_{k_0}^{(n)}(t)) - \widehat{Z}_k^{(n)}(\tau_{k_0}^{(n)}(t)-) \right]. \end{aligned}$$

The first process on the right-hand side converges weakly to 0 by the differencing theorem, and the second one converges weakly to 0 because, by Corollary 3.2,

$\widehat{Z}_k^{(n)}$  has a continuous limit and thus the maximal jump of  $\widehat{Z}_k^{(n)}$  on every finite time horizon  $[0, T]$  converges weakly to 0. This proves (A.10).

Thus, for  $k \neq k_0$ , the process  $c_k^{(n)}$  defined by

$$(A.11) \quad c_k^{(n)}(t) \triangleq \left( \widehat{Z}_k^{(n)}(t) - \widehat{Z}_k^{(n)}(\tau_{k_0}^{(n)}(t)-) \right) - \frac{w_k}{w_{k_0}} \left( \widehat{Z}_{k_0}^{(n)}(t) - \widehat{Z}_{k_0}^{(n)}(\tau_{k_0}^{(n)}(t)-) \right)$$

converges weakly to 0. We have

$$\begin{aligned} c_k^{(n)}(t) &= \frac{1}{\sqrt{n}} \left( V_k^{(n)}(A_k^{(n)}(nt)) - V_k^{(n)}(A_k^{(n)}(n\tau_{k_0}^{(n)}(t)-)) \right) \\ &\quad - \frac{w_k}{\sqrt{n}w_{k_0}} \left( V_{k_0}^{(n)}(A_{k_0}^{(n)}(nt)) - V_{k_0}^{(n)}(A_{k_0}^{(n)}(n\tau_{k_0}^{(n)}(t)-)) \right) \\ &\quad - \sqrt{n} \left( \rho_k^{(n)} - \frac{w_k}{w_{k_0}} \rho_{k_0}^{(n)} \right) (t - \tau_{k_0}^{(n)}(t)). \end{aligned}$$

For every  $s \in [n\tau_{k_0}^{(n)}(t), nt]$ , queue  $k_0$  is nonempty, and hence every nonzero term  $1/(1 - \sum_{l \in \alpha} w_l^{(n)}(t)) j_\alpha^{(n)}(s)$  in the sum

$$\frac{\dot{R}_k^{(n)}(s)}{w_k^{(n)}(s)} = \sum_{m=0}^{K-1} \sum_{\alpha \in A_k^c(m)} \frac{1}{1 - \sum_{l \in \alpha} w_l^{(n)}(t)} j_\alpha^{(n)}(s)$$

also appears in the sum

$$\frac{\dot{R}_{k_0}^{(n)}(s)}{w_{k_0}^{(n)}(s)} = \sum_{m=0}^{K-1} \sum_{\alpha \in A_{k_0}^c(m)} \frac{1}{1 - \sum_{l \in \alpha} w_l^{(n)}(t)} j_\alpha^{(n)}(s).$$

In other words,

$$(A.12) \quad \frac{\dot{R}_k^{(n)}(s)}{w_k^{(n)}(s)} \leq \frac{\dot{R}_{k_0}^{(n)}(s)}{w_{k_0}^{(n)}(s)}$$

for all  $k = 1, \dots, K$  and all  $s \in [n\tau_{k_0}^{(n)}(t), nt]$  for which  $w_k^{(n)}(s) \neq 0$  and  $w_{k_0}^{(n)}(s) \neq 0$ .

As a final preparatory step, we let  $\varepsilon_0 > 0$  be given and choose  $\overline{C} > 0$  so large that  $\mathbb{P}(G_n) \geq 1 - \varepsilon_0$  for all  $n$ , where

$$(A.13) \quad G_n = \left\{ \sqrt{n} \sup_{0 \leq t \leq n} \max_{j=1, \dots, K} |w_k^{(n)}(t) - w_k| \leq \overline{C} \right\}.$$

(A choice is possible by Assumption A.2.) Let  $n_1$  be so large that  $\min_{j=1, \dots, K} w_j > \overline{C}/\sqrt{n_1}$ . In particular, on  $G_n$  we have  $w_k^{(n)}(t) > 0$  for all  $n \geq n_1$ ,  $k = 1, \dots, K$

and  $0 \leq t \leq n$ . Then there exists a constant  $C$  such that, on  $G_n$  for all  $n \geq n_1$ ,  $k = 1, \dots, K$  and  $0 \leq t \leq n$ , we have

$$(A.14) \quad \sqrt{n} \left| \frac{w_k^{(n)}(t)}{w_{k_0}^{(n)}(t)} - \frac{w_k}{w_{k_0}} \right| \leq C.$$

Consider an arbitrary  $k > k_0$ . By (A.12) and (A.14), we have, on  $G_n$ ,

$$(A.15) \quad \begin{aligned} & \widehat{W}_k^{(n)}(t) - \widehat{W}_k^{(n)}(\tau_{k_0}^{(n)}(t)-) \\ &= \frac{1}{\sqrt{n}} \left( V_k^{(n)}(A_k^{(n)}(nt)) - V_k^{(n)}(A_k^{(n)}(n\tau_{k_0}^{(n)}(t)-)) \right) \\ & \quad - \frac{1}{\sqrt{n}} \left( R_k^{(n)}(nt) - R_k^{(n)}(n\tau_{k_0}^{(n)}(t)) \right) \\ & \geq \frac{1}{\sqrt{n}} \left( V_k^{(n)}(A_k^{(n)}(nt)) - V_k^{(n)}(A_k^{(n)}(n\tau_{k_0}^{(n)}(t)-)) \right) \\ & \quad - \frac{1}{\sqrt{n}} \int_{n\tau_{k_0}^{(n)}(t)}^{nt} \frac{w_k^{(n)}(s)}{w_{k_0}^{(n)}(s)} \dot{R}_{k_0}^{(n)}(s) ds \\ & \geq \frac{1}{\sqrt{n}} \left( V_k^{(n)}(A_k^{(n)}(nt)) - V_k^{(n)}(A_k^{(n)}(n\tau_{k_0}^{(n)}(t)-)) \right) \\ & \quad - \frac{w_k}{\sqrt{n}w_{k_0}} \left( R_{k_0}^{(n)}(nt) - R_{k_0}^{(n)}(n\tau_{k_0}^{(n)}(t)) \right) - C(t - \tau_{k_0}^{(n)}(t)) \\ &= \frac{w_k}{\sqrt{n}w_{k_0}} \left( V_{k_0}^{(n)}(A_{k_0}^{(n)}(nt)) - V_{k_0}^{(n)}(A_{k_0}^{(n)}(n\tau_{k_0}^{(n)}(t)-)) \right) \\ & \quad + \sqrt{n} \left( \rho_k^{(n)} - \frac{w_k}{w_{k_0}} \rho_{k_0}^{(n)} \right) (t - \tau_{k_0}^{(n)}(t)) + c_k^{(n)}(t) \\ & \quad - \frac{w_k}{\sqrt{n}w_{k_0}} \left( R_{k_0}^{(n)}(nt) - R_{k_0}^{(n)}(n\tau_{k_0}^{(n)}(t)) \right) - C(t - \tau_{k_0}^{(n)}(t)) \\ &= \frac{w_k}{w_{k_0}} \left( \widehat{W}_{k_0}^{(n)}(t) - \widehat{W}_{k_0}^{(n)}(\tau_{k_0}^{(n)}(t)-) \right) \\ & \quad + \sqrt{n} \left( \rho_k^{(n)} - \frac{w_k}{w_{k_0}} \rho_{k_0}^{(n)} \right) (t - \tau_{k_0}^{(n)}(t)) + c_k^{(n)}(t) - C(t - \tau_{k_0}^{(n)}(t)) \\ & \geq \frac{w_k}{w_{k_0}} \left( \widehat{W}_{k_0}^{(n)}(t) - \widehat{W}_{k_0}^{(n)}(\tau_{k_0}^{(n)}(t)-) \right) + d_k^{(n)}(t), \end{aligned}$$

where

$$d_k^{(n)}(t) = c_k^{(n)}(t) - \left( \left( 1 + \frac{w_k}{w_{k_0}} \right) c + C \right) (t - \tau_{k_0}^{(n)}(t)) \Rightarrow 0,$$

with constants  $c$  and  $C$  entering the definition of  $d_k^{(n)}$  being the same as in (3.4) and (A.14), respectively. The last inequality in (A.15) follows from (3.18) and (3.4) and uses the fact that  $k_0 \leq K - 1$ , to wit,

$$\begin{aligned} & \sqrt{n} \left( \rho_k^{(n)} - \frac{w_k}{w_{k_0}} \rho_{k_0}^{(n)} \right) \\ &= \sqrt{n} w_k \left( \frac{\rho_k^{(n)}}{w_k} - \frac{\rho_{k_0}^{(n)}}{w_{k_0}} \right) \\ &= \sqrt{n} w_k \left( \frac{\rho_k^{(n)} - \rho_k}{w_k} - \frac{\rho_{k_0}^{(n)} - \rho_{k_0}}{w_{k_0}} \right) + \sqrt{n} w_k \left( \frac{\rho_k}{w_k} - \frac{\rho_{k_0}}{w_{k_0}} \right) \\ &\geq - \left( 1 + \frac{w_k}{w_{k_0}} \right) c. \end{aligned}$$

By Corollary 3.3 and the differencing theorem once again,

$$\sum_{k=1}^K (\widehat{W}_k^{(n)}(t) - \widehat{W}_k^{(n)}(\tau_{k_0}^{(n)}(t)-)) \Rightarrow 0.$$

Using our induction assumption, we actually have

$$p_k^{(n)}(t) \triangleq \sum_{k=k_0}^K (\widehat{W}_k^{(n)}(t) - \widehat{W}_k^{(n)}(\tau_{k_0}^{(n)}(t)-)) \Rightarrow 0.$$

However, by our previous considerations, on  $G_n$ ,

$$\begin{aligned} p_k^{(n)}(t) &\geq \left( \sum_{k=k_0}^K \frac{w_k}{w_{k_0}} \right) (\widehat{W}_{k_0}^{(n)}(t) - \widehat{W}_{k_0}^{(n)}(\tau_{k_0}^{(n)}(t)-)) + \sum_{k=k_0+1}^K d_k^{(n)}(t) \\ &= \left( \sum_{k=k_0}^K \frac{w_k}{w_{k_0}} \right) \widehat{W}_{k_0}^{(n)}(t) + \sum_{k=k_0+1}^K d_k^{(n)}(t) \end{aligned}$$

owing to (A.9). Thus, on the set  $G_n$  with  $\mathbb{P}(G_n) \geq 1 - \varepsilon_0$ , we have the following bound on  $\widehat{W}_{k_0}^{(n)}$ :

$$0 \leq \widehat{W}_{k_0}^{(n)}(t) \leq \frac{1}{\sum_{k=k_0}^K (w_k/w_{k_0})} \left( p_k^{(n)}(t) - \sum_{k=k_0+1}^K d_k^{(n)}(t) \right) \Rightarrow 0.$$

This proves (A.7) for  $k = k_0$ , because  $\varepsilon_0 > 0$  is arbitrary.  $\square$

**COROLLARY A.5.** *If, instead of the dominance assumption (3.18), we have*

$$0 < \frac{\rho_1}{w_1} \leq \frac{\rho_2}{w_2} \leq \dots \leq \frac{\rho_{k_0}}{w_{k_0}} < \frac{\rho_{k_0+1}}{w_{k_0+1}} = \dots = \frac{\rho_K}{w_K},$$

then  $\widehat{W}_k^{(n)}(t) \Rightarrow 0$  for all  $k \leq k_0$ .

Indeed, the proof given above applies, without any modifications, also to this more general situation.

## APPENDIX B

The purpose of this appendix is to prove the crushing lemma (Lemma B.4) used in the proof of Proposition A.3.

**B.1. Convergence of inverse processes.** Let  $\phi : [0, 1) \rightarrow [0, \infty)$  be a strictly increasing, continuous function with  $\phi(0) = 0$ ,  $\phi(1) \triangleq \lim_{x \uparrow 1} \phi(x) = \infty$ . For  $x, y \in [0, \infty]$ , define

$$\rho(x, y) = |\phi^{-1}(x) - \phi^{-1}(y)|.$$

Then  $\rho$  is a metric on  $[0, \infty]$  and  $([0, \infty], \rho)$  is separable and compact.

Let  $D_{[0, \infty]}[0, \infty)$  be the space of RCLL functions on  $[0, \infty)$ , taking values in  $([0, \infty], \rho)$ . Let  $D_{\mathbb{R}}[0, \infty)$  be the space of RCLL functions on  $[0, \infty)$ , taking values in  $\mathbb{R}$ . On both these spaces, we impose the Skorohod topology discussed in Chapter 3, Section 5, of [10] and briefly reviewed in Appendix A of [9].

For  $f \in D_{\mathbb{R}}[0, \infty)$  and  $t \geq 0$ , define

$$\Lambda f(t) = \inf\{s \geq 0; f(s) > t\},$$

where  $\inf \emptyset = \infty$ . It is clear from this definition that  $\Lambda f(t)$  is a nondecreasing function of  $t$ .

**LEMMA B.1.**  $\Lambda$  maps  $D_{\mathbb{R}}[0, \infty)$  into  $D_{[0, \infty]}[0, \infty)$ .

**PROOF.** Let  $f \in D_{\mathbb{R}}[0, \infty)$  be given. Since  $\Lambda f$  is nondecreasing, it has left limits. We need to show that  $\Lambda f$  is right-continuous.

Let  $t \geq 0$  be given and define  $s = \Lambda f(t)$ . If  $s = \infty$ , then  $\Lambda f$  is right-continuous at  $t$ .

Assume  $s < \infty$ . By the definition of  $\Lambda f(t)$ , there is a sequence  $\{s_n\}$  with  $s < s_n$  for every  $n$ ,  $s_n \downarrow s$ , and  $f(s_n) > t$  for all  $n$ . Choose a sequence  $\{t_n\}$  with  $t_n \downarrow t$  and  $t < t_n < f(s_n)$  for every  $n$ . Then  $\Lambda f(t_n) \leq s_n$  and, consequently,

$$\liminf_{t' \downarrow t} \Lambda f(t') \leq \liminf_{n \rightarrow \infty} \Lambda f(t_n) \leq s = \Lambda f(t).$$

Because  $\Lambda f$  is nondecreasing,

$$\liminf_{t' \downarrow t} \Lambda f(t') = \lim_{t' \downarrow t} \Lambda f(t') \geq \Lambda f(t),$$

and we have established the right-continuity of  $\Lambda f$ .  $\square$

**LEMMA B.2.** Let  $\{f_n\}$  be a sequence of functions in  $D_{\mathbb{R}}[0, \infty)$  converging to a continuous, strictly increasing function  $f$ . Then  $\Lambda f$  is continuous (taking values in  $[0, \infty]$ ) and  $\Lambda f_n$  converges to  $\Lambda f$  in  $D_{[0, \infty]}[0, \infty)$ .

PROOF. We first show that  $\Lambda f$  is continuous. Since  $\Lambda f$  is nondecreasing and right-continuous, it suffices to show

$$\limsup_{t' \uparrow t} \Lambda f(t') \geq \Lambda f(t)$$

for every  $t > 0$ . Let  $t > 0$  be given and set  $s = \Lambda f(t)$ .

*Case I.*  $s = \infty$ . Define  $M = \sup_{u \geq 0} f(u)$ . We have  $M \leq t < \infty$ . Since  $f$  is strictly increasing,  $f$  cannot attain the value  $M$ . We may choose a sequence  $0 < s_1 < s_2 < \dots$  such that  $s_n \rightarrow \infty$ ,  $t_n = f(s_n) < M$  and  $t_n \rightarrow M$ . Then  $s_n = \Lambda f(t_n)$  and

$$\limsup_{t' \uparrow t} \Lambda f(t') \geq \lim_{n \rightarrow \infty} \Lambda f(t_n) = \lim_{n \rightarrow \infty} s_n = \infty = \Lambda f(t).$$

*Case II.*  $0 < s < \infty$ . Since  $f$  is strictly increasing,  $f(s) = t$ . Let  $0 < s_1 < s_2 < \dots$  be such that  $s_n \uparrow s$  and set  $t_n = f(s_n)$ . Then  $t_n \uparrow t$ ,  $\Lambda f(t_n) = s_n$  and

$$\limsup_{t' \uparrow t} \Lambda f(t') \geq \lim_{n \rightarrow \infty} \Lambda f(t_n) = \lim_{n \rightarrow \infty} s_n = s = \Lambda f(t).$$

*Case III.*  $s = 0$ . Then  $\Lambda f \equiv 0$  on  $[0, t]$ .

We next show that  $\Lambda f_n$  converges to  $\Lambda f$ . It suffices to show that for every  $T$  the restriction of  $\Lambda f_n$  to  $[0, T]$  converges in the Skorohod metric to the restriction of  $\Lambda f$  to  $[0, T]$ .

Fix  $T > 0$ . If  $\Lambda f(T) < \infty$ , set  $S = \Lambda f(T)$ , so that  $T = f(S)$ . If  $\Lambda f(T) = \infty$ , let  $\varepsilon > 0$  be given and choose  $S < \infty$ , so that  $\rho(S - 2, \infty) < \frac{\varepsilon}{2}$ .

Next, choose  $\delta_1 \in (0, 1)$ , so that  $|x - y| \leq 2\delta_1$  implies  $\rho(x, y) \leq \frac{\varepsilon}{2}$  for every  $x, y \in [0, S + 2]$ .

The continuous positive function  $f(s + \delta_1) - f(s)$  has a positive minimum  $\delta_2$  over  $[0, S]$ ; that is, for all  $s \in [0, S]$ , we have

$$f(s + \delta_1) - f(s) \geq \delta_2.$$

Because  $f_n \rightarrow f$  in  $D_{\mathbb{R}}[0, \infty)$ , we may choose  $N$  so that for every  $n \geq N$  there is a strictly increasing mapping  $\lambda_n$  mapping  $[0, S + 2\delta_1]$  onto itself with

$$\sup_{0 \leq s \leq S + 2\delta_1} |f_n(s) - f(\lambda_n(s))| + \sup_{0 \leq s \leq S + 2\delta_1} |s - \lambda_n(s)| < \delta_1 \wedge \delta_2.$$

Let  $t \in [0, f(S)]$  be given and set  $s = \Lambda f(t)$ , so  $0 \leq s \leq S$  and  $t = f(s)$ . We have

$$f_n(s + 2\delta_1) > f(\lambda_n(s + 2\delta_1)) - \delta_2 \geq f(s + \delta_1) - \delta_2 \geq f(s) = t,$$

which shows that

$$\Lambda f_n(t) \leq s + 2\delta_1 = \Lambda f(t) + 2\delta_1.$$

On the other hand, for  $0 \leq u \leq s - 2\delta_1$ ,

$$f_n(u) < f(\lambda_n(u)) + \delta_2 \leq f(u + \delta_1) + \delta_2 \leq f(s - \delta_1) + \delta_2 \leq f(s) = t,$$

which shows that

$$\Lambda f_n(t) \geq s - 2\delta_1 = \Lambda f(t) - 2\delta_1.$$

We conclude that, for every  $t \in [0, f(S)]$ ,

$$|\Lambda f_n(t) - \Lambda f(t)| \leq 2\delta_1.$$

By the choice of  $\delta_1$  and the fact that  $\Lambda f(t) \in [0, S]$  for  $t \in [0, f(S)]$ , we have

$$\rho(\Lambda f_n(t), \Lambda f(t)) \leq \frac{\varepsilon}{2}$$

for all  $t \in [0, f(S)]$ .

If  $\Lambda f(T) < \infty$ , so  $f(S) = T$ , we are done. In the event that  $\Lambda f(T) = \infty$ , we must consider  $\rho(\Lambda f_n(t), \Lambda f(t))$  for  $t \in (f(S), T]$ . But, for  $t > f(S)$ ,

$$\Lambda f(t) \geq \Lambda f(f(S)) = S$$

and

$$\Lambda f_n(t) \geq \Lambda f_n(f(S)) \geq \Lambda f(f(S)) - 2\delta_1 = S - 2\delta_1 > S - 2.$$

Therefore, for  $t > f(S)$ ,

$$\begin{aligned} \rho(\Lambda f_n(t), \Lambda f(t)) &= |\phi^{-1}(\Lambda f_n(t)) - \phi^{-1}(\Lambda f(t))| \\ &\leq |\phi^{-1}(\Lambda f_n(t)) - 1| + |1 - \phi^{-1}(\Lambda f(t))| \\ &\leq |1 - \phi^{-1}(S - 2)| + |1 - \phi^{-1}(S)| \\ &= \rho(S - 2, \infty) + \rho(S, \infty) < \varepsilon. \end{aligned} \quad \square$$

**REMARK B.3.** Assume the hypotheses of Lemma B.2 and let  $s \geq 0$  be given. Fix  $n$ . If  $f_n(s)$  is positive, we may choose a sequence  $t_m \uparrow f_n(s)$  with  $0 \leq t_m < f_n(s)$  for all  $m$ . Then  $\Lambda f_n(t_m) \leq s$ . Letting  $m \rightarrow \infty$ , we obtain

$$\Lambda f_n(f_n(s)-) \leq s$$

for all  $s \geq 0$  such that  $f_n(s) > 0$ . Fix  $S > 0$ . Set  $T = f(S)$  and define

$$M_n(T) = \sup_{0 \leq t \leq T+1} |\Lambda f_n(t) - \Lambda f_n(t-)|.$$

Because  $\Lambda f_n \rightarrow \Lambda f$ , which is finite and continuous, and  $M_n(T)$  is an upper semicontinuous function of  $\Lambda f_n$ , we have

$$(B.1) \quad 0 \leq \limsup_{n \rightarrow \infty} M_n(T) \leq \sup_{0 \leq t \leq T+1} |\Lambda f(t) - \Lambda f(t-)| = 0.$$



Let  $0 \leq s \leq S$  be given and suppose that  $f_n(s) > 0$ . Then assuming, as we can, that  $n$  is large enough to guarantee  $\sup_{0 \leq s \leq S} f_n(s) \leq T + 1$ , we have

$$s \geq \Delta f_n(f_n(s)-) \geq \Delta f_n(f_n(s)) - M_n(T) = \Delta f_n(\max\{f_n(s), 0\}) - M_n(T).$$

In fact, under an additional assumption that  $\Delta f_n(0) = 0$ , for every  $0 \leq s \leq S$  we have

$$(B.2) \quad s \geq \Delta f_n(\max\{f_n(s), 0\}) - M_n(T)$$

also in the event that  $f_n(s) \leq 0$ .

### B.2. The crushing lemma.

LEMMA B.4. *Let  $\{\bar{Z}^{(n)}\}_{n=1}^\infty$  be a sequence of RCLL processes converging weakly to a continuous limit  $Z^*$ . Assume  $\bar{Z}^{(n)}(0) = 0$  for every  $n$ . Let  $\{\bar{T}^{(n)}\}_{n=1}^\infty$  be a sequence of RCLL processes satisfying*

$$(B.3) \quad \bar{T}^{(n)}(t) \geq \bar{T}^{(n)}(t-)$$

for all  $t \geq 0$ ,  $\bar{T}^{(n)}(0) = 0$  for every  $n$ . Assume also that there exists a finite (possibly negative) constant  $L$  and a sequence of strictly positive numbers  $\{\varepsilon_n\}_{n=1}^\infty$  such that

$$(B.4) \quad \bar{T}^{(n)}(t) \geq Lt$$

for all  $t \in [0, \varepsilon_n]$ .

Assume  $\bar{T}^{(n)} \Rightarrow T^*$ , where  $T^*$  is nonrandom, continuous and strictly increasing. Define

$$\bar{N}^{(n)}(t) = \bar{Z}^{(n)}(t) - \sqrt{n} \bar{T}^{(n)}(t),$$

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

$$\bar{W}^{(n)} = \bar{N}^{(n)} + \bar{I}^{(n)}$$

and let, for  $t \geq 0$ ,

$$(B.5) \quad \tau^{(n)}(t) = \sup\{s \in [0, t]; \bar{W}^{(n)}(s) = 0\}.$$

Then  $\tau^{(n)}(t) \Rightarrow t$ .

PROOF. First, note that  $\bar{W}^{(n)}(0) = 0$ , so the supremum in (B.5) is not over the empty set. We shall prove convergence of  $\tau^{(n)}(t)$  for  $0 \leq t \leq 1$ .

We can always assume that the constant  $L$  entering (B.4) is actually 1. Indeed, if  $L < 1$  (this is the only case we need to analyze), consider, for each  $n$ , an auxiliary

function  $g_n(t)$ , which is equal to  $t$  on  $[0, \frac{1}{n}]$ ,  $\frac{2}{n} - t$  on  $[\frac{1}{n}, \frac{2}{n}]$  and is identically 0 for  $t > \frac{2}{n}$ . Then define

$$\begin{aligned}\bar{S}^{(n)}(t) &= \bar{T}^{(n)}(t) + (1 - L)g_n(t), \\ \delta_n &= \frac{1}{n} \wedge \varepsilon_n.\end{aligned}$$

It is clear that  $\bar{S}^{(n)}$  satisfies all the assumptions about  $\bar{T}^{(n)}$  with  $L$  and  $\varepsilon_n$  changed to 1 and  $\delta_n$ , respectively. Moreover,  $\bar{S}^{(n)} \geq \bar{T}^{(n)}$  and  $\bar{S}^{(n)}(t) = \bar{T}^{(n)}(t)$  for  $t \geq \frac{2}{n}$ . Thus, if we put

$$\begin{aligned}\bar{M}^{(n)}(t) &= \bar{Z}^{(n)}(t) - \sqrt{n}\bar{S}^{(n)}(t), \\ \bar{J}^{(n)}(t) &= -\min_{0 \leq s \leq t} \bar{M}^{(n)}(s),\end{aligned}$$

then  $\bar{M}^{(n)} \leq \bar{N}^{(n)}$ , which implies  $\bar{J}^{(n)} \geq \bar{I}^{(n)}$ . This yields, for  $t \geq \frac{2}{n}$ ,

$$0 \leq \bar{W}^{(n)}(t) \leq \bar{M}^{(n)}(t) + \bar{J}^{(n)}(t).$$

For  $0 \leq t \leq \frac{2}{n}$ , we have

$$0 \leq \tau^{(n)}(t) \leq \tau^{(n)}\left(\frac{2}{n}\right) \leq \frac{2}{n},$$

so

$$|t - \tau^{(n)}(t)| \leq \frac{2}{n}.$$

Thus, the validity of our lemma for  $\bar{S}^{(n)}$  clearly implies its validity for  $\bar{T}^{(n)}$ . Therefore, without loss of generality, we assume  $L = 1$ .

On  $(\tau^{(n)}(t), t]$ , the process  $\bar{W}^{(n)}$  is strictly positive and so  $\bar{I}^{(n)}$  is constant. Therefore,

$$\begin{aligned}\bar{W}^{(n)}(\tau^{(n)}(t)) &= \bar{Z}^{(n)}(\tau^{(n)}(t)) - \sqrt{n}\bar{T}^{(n)}(\tau^{(n)}(t)) + \bar{I}^{(n)}(\tau^{(n)}(t)) \\ &= \bar{Z}^{(n)}(\tau^{(n)}(t)) - \sqrt{n}\bar{T}^{(n)}(\tau^{(n)}(t)) + \bar{I}^{(n)}(t) \\ \text{(B.6)} \quad &\geq \bar{Z}^{(n)}(\tau^{(n)}(t)) - \sqrt{n}\bar{T}^{(n)}(\tau^{(n)}(t)) - \bar{N}^{(n)}(t) \\ \text{(B.7)} \quad &= \bar{Z}^{(n)}(\tau^{(n)}(t)) - \bar{Z}^{(n)}(t) + \sqrt{n}(\bar{T}^{(n)}(t) - \bar{T}^{(n)}(\tau^{(n)}(t))).\end{aligned}$$

We may rewrite this as

$$\bar{T}^{(n)}(t) - \bar{T}^{(n)}(\tau^{(n)}(t)) \leq \frac{1}{\sqrt{n}}\bar{W}(\tau^{(n)}(t)) + \frac{1}{\sqrt{n}}[\bar{Z}^{(n)}(t) - \bar{Z}^{(n)}(\tau^{(n)}(t))].$$

Now, by (B.3),

$$\frac{1}{\sqrt{n}} \bar{W}(\tau^{(n)}(t)) \leq \frac{1}{\sqrt{n}} \max_{0 \leq s \leq 1} [\bar{Z}^{(n)}(s) - \bar{Z}^{(n)}(s-)],$$

and the right-hand side has limit 0. Furthermore,

$$\frac{1}{\sqrt{n}} [\bar{Z}^{(n)}(t) - \bar{Z}^{(n)}(\tau^{(n)}(t))] \leq \frac{1}{\sqrt{n}} \left[ \bar{Z}^{(n)}(t) - \min_{0 \leq s \leq 1} \bar{Z}^{(n)}(s) \right],$$

and again the right-hand side has limit 0. On the other hand, by the continuous mapping theorem,

$$\begin{aligned} \bar{T}^{(n)}(t) - \bar{T}^{(n)}(\tau^{(n)}(t)) &\geq \bar{T}^{(n)}(t) - \sup_{0 \leq s \leq t} \bar{T}^{(n)}(s) \\ &\Rightarrow T^*(t) - \sup_{0 \leq s \leq t} T^*(s) = T^*(t) - T^*(t) = 0, \end{aligned}$$

because  $T^*$  is increasing. We conclude that

$$(B.8) \quad \bar{T}^{(n)} - \bar{T}^{(n)} \circ \tau^{(n)} \Rightarrow 0,$$

and hence

$$(B.9) \quad \bar{T}^{(n)} \circ \tau^{(n)} = \bar{T}^{(n)} + (\bar{T}^{(n)} \circ \tau^{(n)} - \bar{T}^{(n)}) \Rightarrow T^*.$$

We want to show that

$$(B.10) \quad \tau^{(n)}(t) \Rightarrow t.$$

Lemma B.2 and the continuous mapping theorem imply

$$(B.11) \quad \Lambda \bar{T}^{(n)} \Rightarrow \Lambda T^* = (T^*)^{-1}.$$

Set

$$M_n = \sup_{0 \leq t \leq T^*(1)+1} |\Lambda \bar{T}^{(n)}(t) - \Lambda \bar{T}^{(n)}(t-)|.$$

Define, for a given  $n$ ,

$$F_n = \left\{ \sup_{0 \leq t \leq 1} \bar{T}^{(n)}(t) \leq T^*(1) + 1 \right\}.$$

Our assumptions imply

$$(B.12) \quad \lim_{n \rightarrow \infty} P(F_n) = 1.$$

By (B.4) with  $L = 1$ , we see that  $\Lambda \bar{T}^{(n)}(0) = 0$  for every  $n$ . Thus, according to Remark B.3, on  $F_n$  we have, for  $0 \leq t \leq 1$ ,

$$(B.13) \quad t \geq \tau^{(n)}(t) \geq \Lambda \bar{T}^{(n)}(\max\{\bar{T}^{(n)}(\tau^{(n)}(t)), 0\}) - M_n$$

[compare (B.2)]. By the same remark [see especially (B.1)] and the continuous mapping theorem, we have  $M_n \Rightarrow 0$ . Thus, if we show

$$(B.14) \quad \Lambda \bar{T}^{(n)}(\max\{\bar{T}^{(n)}(\tau^{(n)}(t)), 0\}) \Rightarrow t,$$

then (B.12) and (B.13) imply (B.10) by the squeeze theorem.

Consider arbitrary constants  $\eta > 0$ ,  $\varepsilon > 0$ . Let  $c \in (0, 1]$  be chosen in such a way that the value  $A \triangleq T^*(1) + c$  is attained by the function  $T^*$ ; that is,  $A = T^*(S_1)$  for some  $S_1 > 1$ . Choose  $\delta > 0$  such that

$$(B.15) \quad w_{(T^*)^{-1}}(\delta, [0, A]) < \frac{\varepsilon}{4},$$

where the last symbol represents the modulus of continuity of  $(T^*)^{-1}$  on  $[0, A]$ :  $w_{(T^*)^{-1}}(\delta, [0, A]) \triangleq \max_{0 \leq s, t \leq A, |s-t| \leq \delta} |(T^*)^{-1}(s) - (T^*)^{-1}(t)|$ . Let, for any real function  $f$  on  $[0, A]$ ,

$$\|f\|_{[0, A]} \triangleq \sup_{s \in [0, A]} |f(s)|$$

and let

$$(B.16) \quad A_n = \left\{ \|\Lambda \bar{T}^{(n)} - (T^*)^{-1}\|_{[0, A]} \leq \frac{\varepsilon}{4} \right\}.$$

By (B.11) and the continuity of  $(T^*)^{-1}$ , we can choose  $n_0$  such that, for all  $n \geq n_0$ ,

$$(B.17) \quad P(A_n) \geq 1 - \frac{\eta}{3}.$$

Indeed, convergence in the Skorohod metric to a continuous limit implies uniform convergence (see [2], page 124). Let

$$(B.18) \quad B_n = \left\{ w_{\Lambda \bar{T}^{(n)}}(\delta, [0, A]) \leq \frac{3\varepsilon}{4} \right\}.$$

By (B.15),  $A_n \subseteq B_n$  for each  $n$ . Now let

$$(B.19) \quad C_n = \{\|\bar{T}^{(n)}\|_{[0, 1]} \leq A\}$$

and suppose that  $n_1$  is such that

$$(B.20) \quad P(C_n) \geq 1 - \frac{\eta}{3}$$

for  $n \geq n_1$  [see the comment after (B.17)]. Define also

$$(B.21) \quad D_n = \{\|\max\{\bar{T}^{(n)}(\tau^{(n)}), 0\} - T^*\|_{[0, 1]} \leq \delta\}.$$

Choose  $n_2$  such that, for all  $n \geq n_2$ , we have

$$(B.22) \quad P(D_n) \geq 1 - \frac{\eta}{3}.$$

Such a choice is possible because, by (B.9) and the continuous mapping theorem,

$$\max\{\bar{T}^{(n)}(\tau^{(n)}(t)), 0\} \implies \max\{T^*(t), 0\} = T^*(t)$$

and  $T^*$  [hence  $\max\{T^*(t), 0\}$ ] is continuous. Take

$$E_n = A_n \cap C_n \cap D_n, \quad n_3 = n_0 \vee n_1 \vee n_2.$$

By (B.17), (B.20) and (B.22), for  $n \geq n_3$ , we have  $P(E_n) \geq 1 - \eta$ . Moreover, on  $E_n$ , we have, for all  $t \in [0, 1]$ ,

$$\begin{aligned} & |\Lambda \bar{T}^{(n)}(\max\{\bar{T}^{(n)}(\tau^{(n)}(t)), 0\}) - t| \\ &= |\Lambda \bar{T}^{(n)}(\max\{\bar{T}^{(n)}(\tau^{(n)}(t)), 0\}) - (T^*)^{-1}(T^*(t))| \\ &\leq |\Lambda \bar{T}^{(n)}(\max\{\bar{T}^{(n)}(\tau^{(n)}(t)), 0\}) - \Lambda \bar{T}^{(n)}(T^*(t))| \\ &\quad + |\Lambda \bar{T}^{(n)}(T^*(t)) - (T^*)^{-1}(T^*(t))| \\ &\leq \frac{3\varepsilon}{4} + \frac{\varepsilon}{4} = \varepsilon \end{aligned}$$

by (B.16), (B.18), (B.19) and (B.21). This proves (B.14), and (B.10) follows.  $\square$

## REFERENCES

- [1] BILLINGSLEY, P. (1986). *Probability and Measure*, 2nd ed. Wiley, New York.
- [2] BILLINGSLEY, P. (1999). *Convergence of Probability Measures*, 2nd ed. Wiley, New York.
- [3] BRAMSON, M. (1998). State space collapse with applications to heavy traffic limits for multiclass queueing networks. *Queueing Systems Theory Appl.* **30** 89–148.
- [4] CHANG, C.-S. (1994). Stability, queue length, and delay of deterministic and stochastic queueing networks. *IEEE Trans. Automat. Control* **39** 913–931.
- [5] 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.
- [6] CRUZ, R. L. (1991). A calculus for network delay. I. Network elements in isolation. *IEEE Trans. Inform. Theory* **37** 114–131.
- [7] CRUZ, R. L. (1991). A calculus for network delay. II. Network analysis. *IEEE Trans. Inform. Theory* **37** 132–141.
- [8] DEMERS, A., KESHAV, S. and SHENKER, S. (1989). Design and analysis of a fair queueing system. In *Proceedings of ACM SIGCOMM'89*. ACM Press, New York.
- [9] DOYTCHINOV, B., LEHOCZKY, J. P. and SHREVE, S. E. (2001). Real-time queues in heavy traffic with earliest-deadline-first queue discipline. *Ann. Appl. Probab.* **11** 332–378.
- [10] ETHIER, S. N. and KURTZ, T. G. (1985). *Markov Processes: Characterization and Convergence*. Wiley, New York.
- [11] IGLEHART, D. and WHITT, W. (1970). Multiple channel queues in heavy traffic I. *Adv. in Appl. Probab.* **2** 150–177.
- [12] KESHAV, S. (1997). *An Engineering Approach to Computer Networking*. Addison–Wesley, Reading, MA.

- [13] KRUK, L., LEHOCZKY, J., SHREVE, S. and YEUNG, S. N. (2002). Earliest-deadline-first service in heavy traffic acyclic networks. Working paper, Dept. Mathematical Sciences, Carnegie Mellon Univ.
- [14] PAREKH, A. K. and GALLAGER, R. G. (1993). A generalized processor sharing approach to flow control in integrated services networks: The single node case. *IEEE/ACM Trans. Networking* **1** 344–357.
- [15] PAREKH, A. K. and GALLAGER, R. G. (1994). A generalized processor sharing approach to flow control in integrated services networks: The multiple node case. *IEEE/ACM Trans. Networking* **2** 137–150.
- [16] PROKHOROV, YU. (1956). Convergence of random processes and limit theorems in probability theory. *Theory Probab. Appl.* **1** 157–214.
- [17] RAMANAN, K. and REIMAN, H. (2002). Fluid and heavy traffic diffusion limits for a generalized processor sharing model. Working paper, Bell Labs, Murray Hill, NJ.
- [18] REIMAN, M. (1984). Some diffusion approximations with state space collapse. *Modelling and Performance Evaluation Methodology. Lecture Notes in Control and Information Systems* **60** 209–240. Springer, New York.
- [19] WILLIAMS, R. J. (1998). Diffusion approximations for open multiclass queueing networks: sufficient conditions involving state space collapse. *Queueing Systems Theory Appl.* **30** 27–88.
- [20] YEUNG, S. N. and LEHOCZKY, J. P. (2002). Real-time queueing networks in heavy traffic with EDF and FIFO queue discipline. Working paper, Dept. Statist., Carnegie Mellon Univ.
- [21] ZHANG, H. (1995). Service disciplines for guaranteed performance service in packet-switching networks. *Proc. IEEE* **83** 1374–1396.

L. KRUK  
INSTITUTE OF MATHEMATICS  
MARIA CURIE-SKŁODOWSKA UNIVERSITY  
MARIJ CURIE-SKŁODOWSKIEJ 1  
20-031 LUBLIN  
POLAND  
E-MAIL: lkruk@golem.umcs.lublin.pl

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

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

S.-N. YEUNG  
AT&T LABS  
180 PARK AVENUE  
FLORHAM PARK, NEW JERSEY 07932  
E-MAIL: syeung@homer.att.com