

A LIFO QUEUE IN HEAVY TRAFFIC

BY VLADA LIMIC¹

Cornell University

This paper describes the heavy-traffic behavior of an M/G/1 *last-in–first-out preemptive resume* queue. An appropriate framework for the analysis is provided by measure-valued processes. In particular, the paper exploits the setting of recent works by Le Gall and Le Jan. Their finite-measure-valued *exploration* process corresponds to our *RES-measure* (residual services measure) process, that captures all the relevant information about the evolution of the queue, while their *height* process corresponds to the *queue-length* process. The heavy-traffic “diffusion” approximations for the RES-measure and the queue-length processes are derived under the usual second moment assumptions on the service distributions. The tightness of queue lengths argument uses estimates for the total size and height of large Galton–Watson trees.

1. Introduction. Imagine customers arriving to a queue according to a Poisson (rate λ) process, each customer requesting an amount of service time with distribution function F , independently of the arrival process and of the service times of other customers. Let F have finite mean m . The server devotes all of its service potential to the last customer to have arrived. Moreover, at the moment of each new arrival, the server switches instantaneously from serving the current customer c (if any) to the newest customer \bar{c} . Customer c stays waiting in queue and only after \bar{c} is served completely and exits the queue does the service of c resume. The above rule (or *service discipline*), that applies to serving any customer c , corresponds to a branching structure (cf. Section 1.2). This system is called the M/G/1 *last-in–first-out (LIFO) preemptive resume* queue. From now on we will omit the qualifier “preemptive resume” for brevity. Note that the server is busy whenever the queue is nonempty, which is usually referred to as a *nonidling* or *work-conserving* property. We assume the queue is empty at time 0.

For any two numbers x, y , let x^+ , $x \wedge y$ and $x \vee y$ denote the positive part of x , the minimum and the maximum of x and y , respectively. For any x , $\lfloor x \rfloor$ and $\lceil x \rceil$ denote the largest integer smaller than or equal to x , and the smallest integer larger than x , respectively. Also identify $G(t)$ with G_t whenever G is a stochastic process.

Suppose a customer arrives to the queue at time t and requests an amount v of service time. If we let $u(s)$, $s \geq t$, be its total amount of time in service by time s , the *residual service time* of this customer at time s is $v - u(s) = (v - u(s))^+$. Denote by $(A(t), t \geq 0)$ the Poisson (counting) process of arrivals, by $Z(t)$ the

Received May 1999; revised May 2000.

¹Research supported by NSF Mathematical Sciences Postdoctoral Research Fellowship.

AMS 2000 *subject classifications*. 60K25, 60J80.

Key words and phrases. LIFO queue, heavy traffic, measure-valued process, branching.

queue length at time t , that is, the number of individuals in queue at time t , and by $W(t)$ the (*immediate*) *workload* of the queue at time t , that is, the total amount of work still required by customers present in the system at time t (measured in units of server time). So the workload is equal to the total sum of all residual service times. The parameter $\rho = \lambda m$, called the *traffic intensity* of the queue, is the average amount of work arriving per unit time. It is a well-known (and easy) (cf. Section 2) fact that the workload process does not vary over work-conserving service disciplines. In particular, the workload process ($W(t), t \geq 0$) is the same for the *first-in-first-out* (FIFO) queue, where the customers are served in the order of their arrival. Therefore, W is a Markov process with respect to the filtration \mathcal{F}_t generated by arrivals and service times up to time t , and it is positive recurrent, null-recurrent and transient whenever $\rho < 1$, $\rho = 1$ and $\rho > 1$, respectively. From a practical point of view, it is desirable to “keep the server busy” most of the time without getting it overwhelmed with work. This corresponds to the situation $\rho = 1 - \varepsilon$ for some small $\varepsilon > 0$, and as $\varepsilon \searrow 0$, the queue approaches *heavy traffic*.

The pioneering works in heavy-traffic approximations to queues (Kingman [27]) and queueing networks (Iglehart and Whitt [21, 22], Harrison [19], Reiman [31] and Whitt [35]) appeared a while ago. A detailed overview of the enormous literature is given in Williams [36]. Recent papers by Bramson [7] and Williams [37] provide powerful tools for the analysis of multiclass queueing networks with feedback in heavy traffic. However, their techniques are developed for *head-of-the-line* (HL) service disciplines. It is intuitively clear what head-of-the-line means (see, e.g., [37] for precise definitions). The FIFO discipline is the simplest HL discipline, while the LIFO discipline is perhaps the simplest non-HL discipline. Recall that our LIFO discipline is preemptive resume, where the server switches to serving the newest customer immediately upon arrival. The queue-length process of the *non-preemptive* LIFO queue, where the server serves each customer completely, and immediately afterwards begins to serve the last customer to have arrived (if any), is equal (in distribution) to that of the FIFO queue, so its heavy-traffic approximation is given in [21, 22]. Although the LIFO discipline might seem “unfair,” and therefore less natural than the FIFO discipline, it naturally arises in applications (e.g., LIFO stack in computer science; see also [4, 25]). In fact, here is a natural “optimization” problem. Suppose that, for a queue close to heavy traffic, we have a server with the ability of serving in both FIFO and LIFO orders (equivalently, both non-preemptive and preemptive resume LIFO orders). Due to limited space (say), it is important to minimize the queue length, and the question is: which service discipline to use? Similar questions were considered by Coffman and Mitrani [9]. We discuss the answer in Section 3.4.2, where we see that typically one of the two disciplines is optimal.

Some important aspects of LIFO preemptive resume queues have been investigated previously. Shanthikumar and Sumita [33] study properties of invariant measure in the more general setting of renewal arrivals. Relation to risk processes (Sigman [34]) is another connection to applications. Abate and

Whitt [1] establish heavy-traffic limits for the steady-state waiting time in the M/G/1 LIFO queue.

The goal of this paper is to describe the heavy-traffic behavior of an M/G/1 LIFO queue (Theorems 1 and 5). The corresponding description of its FIFO counterpart is given in [21, 22]. An appropriate framework for the heavy-traffic analysis of the LIFO queue is provided by measure-valued processes. A random process is *measure-valued* if it takes values in a space of measures. Measure-valued processes have been actively studied in the past two decades (see, e.g., Dawson [10], Dynkin [12] for references). In particular, this paper exploits the setting of the recent papers by Le Gall and Le Jan [14, 15], who construct and study the finite-measure-valued *exploration process* (an analogue of which we describe in Section 1.2 and call the *RES-measure* process) as a step in their pathwise construction of superprocesses with general branching mechanism. The *height* process of [14, 15] corresponds to our queue-length process. The theorems in Section 3.2 are stated and proved for an interesting case (from the queueing perspective) where the approximation $X(t)$ to the load process (cf. Section 1.1) is a Brownian motion. By examining the argument, it is easily checked that Theorem 1 continues to hold when X belongs to a more general class of Lévy processes, and where the approximation $Z(t)$ to the queue-length process has continuous paths. In particular, it holds in the case of *heavy-tailed* service times, where the approximation to the load is a Lévy stable- α process with $\alpha \in (1, 2)$. Stable processes are common in queueing models (e.g., [24, 20]).

Theorem 5 provides the heavy-traffic approximation for the queue-length processes. Our tightness (in the Skorokhod topology) of queue-length argument rests on asymptotics for the distribution of a super-near-critical Galton–Watson tree (cf. Lemma 8) where the offspring distribution has finite variance. The weak convergence of queue-length processes for heavy-tailed service times remains an open problem (we discuss this briefly in Section 3.4.4). The Brownian motion approximation (Theorem 5) to queue length is analogous to the Brownian excursion approximation to *depth-first search* walk of a large (conditioned on total size) Galton–Watson tree, Aldous [2] (see Section 3.4.3). For some other interesting relations between queues and trees, we refer the reader to Kersting and Geiger [16] and Shalmon [32].

The paper is organized as follows. Sections 1.1 and 1.2 introduce basic concepts and some important relations. Section 2 is a brief analysis of the workload process in heavy-traffic. Section 3 is devoted to the heavy-traffic limit theorems for the RES-measure (cf. Section 1.2) and the queue-length processes. We discuss some consequences and related work in Section 3.4, and give the directions for further research in Section 4.

1.1. *M/G/1 LIFO queue load as a Lévy process.* Let $(X_t, t \geq 0)$ be the Lévy process obtained by superimposing positive discrete jumps on the shift $-at$, where $a > 0$. More precisely, the jumps occur at the times of increase of a counting Poisson (rate λ) process $A(t)$, the sizes of the jumps v_i , $i \geq 1$, are i.i.d. random variables with distribution F concentrated on $(0, \infty)$, and in

between the jumps the process is linear with constant negative drift $-a$. The Lévy characterization of X is

$$(1) \quad E \exp(-xX_t) = \exp\{tax + t\lambda \int_{(0,\infty)} (e^{-xr} - 1) F(dr)\}, \quad x > 0,$$

and a Lévy measure of X is $\pi(dx) = \lambda F(dx)$. See Section 3.1 for further definitions and Bertoin [5] for background on Lévy processes. We can assume by scaling that a equals 1. Then (as noted in [15]) X is the *load process*

$$(2) \quad X_t = \sum_{i=1}^{A(t)} v_i - t, \quad t \geq 0,$$

of an M/G/1 LIFO queue with customers arriving at the times of the jumps of X , and requesting service equal to sizes of the jumps. It is also clear that the load process X of any M/G/1 LIFO queue is a Lévy process of the above form.

Figure 1 shows a possible path of X over a finite time interval. Suppose X had a jump at some (random) time s and write $X_{s-} = \lim_{u \uparrow s} X_u$. Let $\gamma_s = \inf\{u \geq s: X_u \leq X_{s-}\}$. We identify the actual set of times when this customer is in service with the set $\mathcal{A}_s = \{u \in [s, \gamma_s]: X_u \in [X_{s-}, X_s] \text{ and } \inf_{t \in [s, u]} X_t \geq X_u\}$, indicated in bold on the time axis in the figure. At time γ_s , this customer exits the queue; in the meantime, its service might be interrupted several times due to jumps of X , that is, arrivals of new customers. The “gaps” in \mathcal{A}_s correspond to services of these customers. The customer who arrived (jumped) at time s will still be in queue at time $t > s$ if and only if $\gamma_s > t$, that is,

$$(3) \quad X_{s-} < \inf_{u \in [s, t]} X_u$$

(as it happens for s and t in the figure). The difference $(\inf_{u \in [s, t]} X_u - X_{s-})^+$ is its residual service time at t . Therefore, the queue-length process $Z_t = Z(t)$ satisfies

$$(4) \quad Z_t = \#\left\{s \leq t: X_{s-} < \inf_{s \leq u \leq t} X_u\right\}.$$

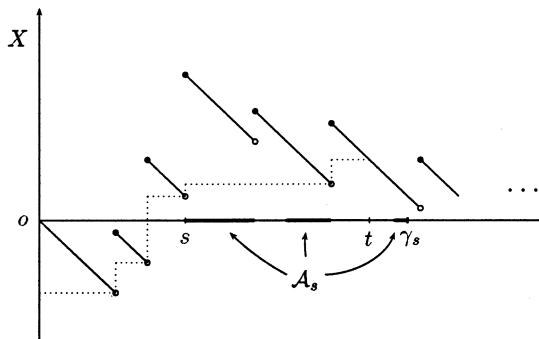


FIG. 1.

1.2. *Branching and the RES-measure.* The relation between queueing and branching goes back to Kendall [26]. Some highlights of the literature are given in [15]. Suppose we call a customer who arrives at time t a descendant of a customer that arrived at time s if the latter is still in queue at time t , that is, if (3) holds. Any customer either finds the queue empty upon arrival, in which case it becomes a *progenitor* (or *root*), or finds the queue nonempty, in which case it becomes a *child* (or *offspring*) of the customer being served previous to its arrival. This procedure yields a sequence of Galton–Watson trees, with each busy cycle corresponding to a different progenitor. The corresponding offspring distribution depends on the Lévy measure of X [15]; it is easy to see that the mean equals λm and the variance equals $\lambda^2 \beta + \lambda m - \lambda^2 m^2$, where β is the second moment of F . We return to this very useful characterization in the heavy-traffic analysis, Section 3.3.

Let $I_t = \inf_{s \in [0, t]} X_s$ be the past infimum process, and let $I_s^t = \inf_{u \in [s, t]} X_u$ be the *future infimum* (dotted line in Figure 1). Note that $-I_t$ equals the (*cumulative*) *idle time*, the amount of time for which there has been no customer in queue up to t . This is true since $-I_t$ increases at constant rate 1 during the time intervals with no customers in queue. The jumps of I^t may occur only at the times $s < t$ at which customers arrive, and the jump sizes $(\inf_{u \in [s, t]} X_u - X_{s-})^+$ are the residual service times (at time t) of the corresponding customers. The workload process is then given by $W_t = X_t - I_t$. The excursions of X above its past infimum, or equivalently, the excursions of W above 0, correspond to the busy cycles of the queue. The relation between excursions of random walks and branching goes back to Harris [18].

Let $\tau_x = \inf\{s \geq 0: I_s \leq -x\}$, and let $M_x = \#\{s \in [0, \tau_x]: X_{s-} < X_s \text{ and } W_s = Z_s = 0\}$ be the number of customers arriving to an empty queue during the interval $[0, \tau_x]$. So M_x is the number of busy cycles (i.e., the number of trees) started in the interval $[0, \tau_x]$. Observe that

$$(5) \quad M_x \stackrel{d}{=} \text{Poisson}(\text{rate } \lambda x),$$

where λ is the arrival rate.

One can think of a LIFO queue as a continuous-time process with values in the state space of finite lists of arbitrary length. At each time t , the state of the queue is the list of residual service times for all queued customers ordered by their arrival times. Of course, one can obtain the above list from $(X_s, s \leq t)$, the path of the load process up to t , or from the path of the workload W up to t . The state space of (finite and infinite) lists appeared in [15], and it seems convenient for certain types of analysis of evolution of the queue, e.g., for obtaining the stationary distribution and the dual process (cf. [15], Section 3). However, it is not convenient for the heavy-traffic analysis since the lists “become uncountable” in the limit, which is related to the fact that the above discrete branching mechanism converges under the heavy-traffic assumptions to a continuum branching mechanism driven by Brownian excursions. A crucial ingredient for this paper is the existence of a measure-valued encoding q_t of the state space, analogous to the *exploration process* of [14], that has a

natural extension in the limit. The queue length is recoverable from q by

$$(6) \quad Z_t = \sup(\text{Supp}(q_t)),$$

where $\text{Supp}(\mu)$ denotes the closed support of measure μ . Moreover, the list of residual service times at time t equals $(q_t(1), q_t(2), \dots, q_t(Z_t))$, the list of masses of atoms of q_t . Of course, the workload is then given by $W_t = \sum_{i=1}^{Z_t} q_t(i) = \int_0^t dI_s^t = \langle q_t, 1 \rangle$. The process q_t is defined by

$$(7) \quad \langle q_t, \varphi \rangle := \int_0^t \varphi(Z_s) dI_s^t = \int_0^t \varphi(Z_s^t) dI_s^t.$$

Here $\langle \mu, \varphi \rangle$ stands for $\int_{[0, \infty)} \varphi d\mu$, and $\varphi: R^+ \rightarrow R$ is continuous, with bounded support, and

$$(8) \quad Z_s^t = \#\left\{u \leq s: X_{u-} < \inf_{u \leq z \leq t} X_z\right\}$$

is the number of individuals in queue at time s that will still be in queue at time t . Note that the integrals in (7) are in fact finite sums. The equality in (7) is due to a simple fact: $Z_s \equiv Z_s^t$, $I^t(ds)$ -a.e., for each $t \geq 0$, where the step function I_s^t defines an atomic measure $I^t(ds)$ in the usual way. Note that $I_u^s = I_u^t$ for all $u \leq s$, $I^t(ds)$ almost everywhere. It is easily seen (and shown in [15]) that the process q_t is strong Markov. The process Z_s^t is nondecreasing in s for each fixed t , and this monotonicity will be essential in the heavy-traffic analysis (Section 3).

The following observation will be important later on for deriving the queue-length heavy-traffic approximation. If we fix any time t and time-reverse the load X from t back to 0 (or equivalently, rotate Figure 1 about the origin by 180 degrees), the future infimum I^t “gets mapped” onto the (past) supremum process of the time-reversed load process $\tilde{X}_s^t = X_t - X_{(t-s)-}$. In particular, the queue length Z_t , which equals the number of jumps of the future infimum (4), also equals the number of jumps of the time-reversed supremum process occurring in $[0, t]$. The precise statements and their generalizations are deferred until Section 3.1.

We prefer integrals to sums in (7) since the heavy-traffic limit Theorem 1 involves the convergence of rescaled q 's to a limit of the same form. The queue length is usually denoted by Z_t . In [15], the corresponding process is denoted by H and is called the *height*. (The two processes are analogous, though H is defined in discrete time and the walk H , unlike Z , never visits the same vertex twice.) The height process (the queue length) visits the vertices (the customers) of the sequence of trees (busy cycles) in the *depth-first search* order (children before siblings), recording their distance from the root. The exploration process “explores” these trees in a similar way, carrying a lot of additional information, and its advantage over the height process is the Markov property and relation (6). The process Z_t is Markov only if the service time distribution F is exponential. From now on we identify any queue with the corresponding measure-valued q , which we call the *RES-measure* (derived from *REsidual Services measure*) process.

2. The workload in heavy traffic. In this section we describe the framework of heavy traffic. We state some of the usual assumptions (e.g., [7, 37]) and discuss asymptotics of the workload processes. The workload does not depend on the service discipline and is a relatively simple object for analysis. At the same time, the workload is the simplest (interesting) process related to the queue, so any “natural” convergence of queues should comprise the convergence of corresponding workloads.

Let $q^r = (q_t^r, t \geq 0)$ be a family of RES-measure processes of M/G/1 LIFO queues, indexed by r . Here r ranges over real numbers; it is easiest to think of a sequence increasing to ∞ . The r th M/G/1 queue has the arrival rate λ^r and the service time distribution function F^r with mean m^r . Assume

$$(9) \quad \lambda^r \rightarrow \lambda \in (0, \infty), \quad m^r \rightarrow m \in (0, \infty) \quad \text{as } r \rightarrow \infty.$$

Let $\rho^r = \lambda^r m^r$. A usual heavy-traffic assumption is

$$(10) \quad r^{1/2}(1 - \rho^r) = r^{1/2}(1 - \lambda^r m^r) \rightarrow c \in R \quad \text{as } r \rightarrow \infty.$$

Let $A^r(\cdot)$, $W^r(\cdot)$ and $Z^r(\cdot)$ denote the corresponding arrival, workload and queue-length processes. We assume that, for each r , the queue is empty at time 0, or equivalently, $W^r(0) = 0$, so that notation of Sections 1.1 and 1.2 directly applies. More general initial conditions can be treated with additional work (cf. Section 3.4.1).

Let v_i^r be the service time requested by the i th customer who arrives to the queue. So $\{v_i^r, i \geq 1\}$ is an i.i.d. sequence with distribution F^r , and denote by $V^r(n) = \sum_{i=1}^n v_i^r$, $n \geq 1$, the *cumulative service time* process. The *workload equation* is (the same for all work-conserving disciplines)

$$(11) \quad W^r(t) = V^r(A^r(t)) - t - I^r(t) = X^r(t) - I^r(t),$$

where $X^r(t)$ is the load and $-I^r(t) = -\inf_{s \leq t} X^r(s)$ is the idle time.

If we assume in addition to (9), (10) that, for each r , the service times have second moment $\beta^r < \infty$ and

$$(12) \quad \beta^r \rightarrow \beta < \infty \quad \text{as } r \rightarrow \infty,$$

$$(13) \quad \sup_r E[(v_1^r)^2 1_{\{v_1^r \geq K\}}] \rightarrow 0 \quad \text{as } K \rightarrow \infty$$

(so a Lindeberg–Feller type of condition is satisfied), the rescaled load processes $(r^{-1/2}X^r(rt), t \geq 0)$ converge weakly to a Brownian motion X with drift $-c$, defined in (10), and variance $\lambda\beta$. The value $\lambda\beta$ for the asymptotic variance can be verified using standard arguments; intuitively, it is due to the fact that the infinitesimal drift $E(\Delta(X_t^r)^2 | \mathcal{F}_t^r)$ equals $\lambda^r E(v_1^r)^2 dt = \lambda^r \beta^r dt$ (in the obvious notation), so in the limit $X_t^2 - \lambda\beta t$ should be (and is) a martingale. By the continuous mapping theorem (e.g., Billingsley [6], Lemma 6.1), under the same scaling, the workload processes W^r converge to $W = X - I$ obtained by reflecting X above the past minimum $I_t = \inf_{s \leq t} X_s$.

REMARK. Since Brownian motion X is a stable process with index $\alpha = 2$, (10) reads

$$(14) \quad r^{1-\gamma}(1 - \rho^r) \rightarrow c,$$

where $\gamma = 1/\alpha = 1/2$. It is standard that $(A^r(rt)/r, t \geq 0) \Rightarrow (\lambda t, t \geq 0)$, as $r \rightarrow \infty$. Now consider a more general setting, where $\alpha \in (0, 2]$ and $\gamma = 1/\alpha$. Assume that $\lim_{r \rightarrow \infty} b^r/r^\gamma$ exists, that (9), (14) hold, and consider the rescaled processes

$$\begin{aligned} X^r(rt)/b^r &= \frac{V^r(A^r(rt)) - rt}{b^r} \\ &= \frac{\sum_{i=1}^{A^r(rt)} (v_i^r - m^r)}{b^r} + \frac{A^r(rt)m^r - rt}{b^r} \\ &\sim \frac{\sum_{i=1}^{r\lambda t} (v_i^r - m^r)}{b^r} + \frac{r(\rho^r - 1)t}{b^r}. \end{aligned}$$

Assume moreover that, for all large r , the service times have *heavy tails*, that is, F^r is in the domain of attraction of the stable- α law (cf. Breiman [8], page 207). Assume that $m < \infty$ so it must be $\alpha \in (1, 2]$. Then (e.g., Jacod and Shiryaev [23]) $X^r(rt)/b^r$ converges to a stable- α process X , and again the workload converges to X reflected above the past minimum. Under more general conditions on b^r and F^r (see [23], Theorem VII.2.35 and [15], Proposition 5.1), $X^r(rt)/b^r$ will converge to a Lévy process X with no negative jumps [i.e., the Lévy measure π in (16) is concentrated on $(0, \infty)$].

3. Heavy-traffic limits. In this section we state and prove the heavy-traffic limit theorems for M/G/1 LIFO queues. In Section 3.1 we define the limit processes and mention some of their properties from Le Gall and Le Jan [14, 15]. Sections 3.2 and 3.3 are devoted to the convergence, and Section 3.4 comments on some consequences and extensions, and relates our result to the existing literature.

3.1. *Limit processes.* We briefly describe the setting of [14, 15]. Let $X = (X_t, t \geq 0)$ be a Lévy process with no negative jumps such that $\liminf_{t \rightarrow \infty} X_t = -\infty$. Then

$$(15) \quad E \exp(-x X_t) = \exp(t\psi(x)), \quad x > 0,$$

where the *Laplace exponent* $\psi(x)$ is of the form

$$(16) \quad \psi(x) = cx + \frac{\sigma^2 x^2}{2} + \int_{(0, \infty)} (e^{-xr} - 1 + xr)\pi(dr), \quad c \geq 0,$$

and the *Lévy measure* $\pi(dr)$ satisfies

$$(17) \quad \int_{(0, \infty)} (r \wedge r^2)\pi(dr) < \infty.$$

In Section 1.1 we identified the load processes (2) of M/G/1 LIFO queue with a class of analogous (though much simpler) Lévy processes characterized by (1). Moreover, in Section 2 we saw how some of the processes characterized by (15)–(17) arise naturally as limits of rescaled load processes of queues approaching heavy traffic. Such X can be viewed as a “generalized” queue load, and it is plausible that a generalized queue length Z_t and a generalized RES-measure q_t can be obtained from the load by mimicking (4) and (7). Indeed, this has been done in [14, 15] [for X with general Laplace exponent (16), (17)], and we recall their definitions here briefly. We will mainly deal with the special Brownian case $\psi(x) = cx + \sigma^2 x^2/2$, $c \in (-\infty, \infty)$, $\sigma > 0$. If the drift $-c$ is strictly positive, then $\lim_{t \rightarrow \infty} X_t = \infty$ (not $-\infty$), and such processes were not considered in [14, 15]. However, the definitions below and the rest of the analysis extend naturally.

Let X^\bullet be a (not generalized) M/G/1 LIFO queue load process as in Section 1.1. In this section only, all the LIFO queue-related processes from Section 1 have additional “ \bullet ” in the superscript. Recall the related infimum processes $I_t^\bullet, I_s^{\bullet,t}$ from Section 1.2. For each fixed $t > 0$, denote by $(\tilde{X}_s^{\bullet,t}, 0 \leq s \leq t)$ the time-reversed process X^\bullet from t , that is,

$$\tilde{X}_s^{\bullet,t} = X_t^\bullet - X_{(t-s)^-}^\bullet, \quad 0 \leq s < t \quad \text{and} \quad \tilde{X}_t^{\bullet,t} = X_t^\bullet,$$

and let $\tilde{S}_s^{\bullet,t} = \sup_{u \in [0, s]} \tilde{X}_u^{\bullet,t}$. Rewrite identities (4) and (8) as

$$(18) \quad \begin{aligned} Z_t^\bullet &= \#\{z: z \in [0, t], \tilde{S}_z^{\bullet,t} > \tilde{S}_{z^-}^\bullet\}, \\ Z_s^{\bullet,t} &= \#\{z: z \in [t-s, t], \tilde{S}_z^{\bullet,t} > \tilde{S}_{z^-}^{\bullet,t}\}, \end{aligned}$$

since the future infimum $I_s^{\bullet,t}$ corresponds (in reversed time) to the past supremum $\tilde{S}_{t-s}^{\bullet,t}$.

Now fix a generalized queue load process X , and let $(\tilde{X}_s^t, 0 \leq s \leq t)$ be the corresponding time-reversed process, so that

$$\begin{aligned} \tilde{X}_s^t &= X_t - X_{(t-s)^-}, \quad 0 \leq s < t, \\ \tilde{X}_t^t &= X_t \quad \text{and} \quad \tilde{S}_s^t = \sup_{u \in [0, s]} \tilde{X}_u^t. \end{aligned}$$

Now the set of increase points of \tilde{S}^t is measured using local time. Let $(\tilde{L}_s^t, 0 \leq s \leq t)$ be a local time of the process $(\tilde{S}_s^t - \tilde{X}_s^t, 0 \leq s \leq t)$ at level 0. Then \tilde{L}_s^t is a continuous nondecreasing *additive* process with support on the zero set of $\tilde{S}_s^t - \tilde{X}_s^t$ (see, e.g., [5], Chapter IV). By analogy to (18), define

$$(19) \quad Z_s^t = \tilde{L}_s^t - \tilde{L}_{t-s}^t, \quad 0 \leq s \leq t \quad \text{and} \quad Z_t^t = Z_t^t = \tilde{L}_t^t.$$

Local time is unique ([5], Proposition IV.2.5) up to a multiplicative constant. For X a Brownian motion with drift $-c$ and variance σ^2 , we can choose \tilde{L}^t as

$$(20) \quad \tilde{L}_s^t = \frac{2}{\sigma^2} \tilde{S}_s^t,$$

which translates back to

$$(21) \quad \begin{aligned} Z_s^t &= \frac{2}{\sigma^2}(I_s^t - I_t), \quad t \geq 0, \\ Z_t &= \frac{2}{\sigma^2}(X_t - I_t) = \frac{2}{\sigma^2}W_t, \end{aligned}$$

where $I_t = \inf_{s \leq t} X_s$ and $I_s^t = \inf_{u \in [s, t]} X_u$ as always. The choice of the normalizing factor in (20) is motivated by heavy-traffic limits (e.g., Lemma 2). Note that the processes Z_t, Z_s^t in (21) are continuous. By (21), Z_s^t is nondecreasing in s for every t and $\lim_{s \uparrow t} Z_s^t = Z_t$. Moreover, it is straightforward to see that again $Z_s \equiv Z_s^t, I^t(ds)$ -a.e., $t \geq 0$. So definition (7) extends to the generalized setting. The finite-measure-valued process q_t given by

$$\langle q_t, \varphi \rangle = \int_0^t \varphi(Z_s) I^t(ds) = \int_0^t \varphi(Z^t(s)) I^t(ds)$$

is a strong Markov process, the identity (6) carries over, and moreover,

$$\text{Supp}(q_t) = [0, Z_t], \quad t \geq 0.$$

In fact, q_t is a constant ($\sigma^2/2$) multiple of the Lebesgue measure on $[0, Z_t]$. The above statements can be easily checked in our (Brownian) setting, and some have analogues in the more general setting (16), (17) of [14].

3.2. Convergence. We are now ready to state our main result. Denote by $M_f(\mathbb{R}_+)$ the complete, separable metric space of finite measures on $[0, \infty)$ (cf. Billingsley [6], or Dawson [10], Section 3), by $D_R[0, \infty), D_R[0, t], D_{M_f(\mathbb{R}_+)}[0, \infty)$ the usual Skorokhod spaces, and by \Rightarrow the corresponding weak convergence of processes.

Let $q^r = (q_t^r, t \geq 0), r \geq 1$, be a family of RES-measures with corresponding queue load processes $X^r = (X_t^r, t \geq 0)$, as in Sections 1 and 2. Similarly, denote by $Z^r, Z^{t,r}$ the queue-length processes in (6) and (8), and let $I^{t,r}$ be the future infimum processes of X^r .

Assume that (9), (10) and (12), (13) hold. Then we know (Section 2) that $\widehat{X}^r = (r^{-1/2}X^r(rt), t \geq 0) \Rightarrow X$, where X is a Brownian motion with variance $\lambda\beta$ and drift $-c$, and by the Skorokhod representation theorem, we may assume that

$$(22) \quad \widehat{X}^r \rightarrow X \quad \text{a.s. in } D_R[0, \infty) \quad \text{as } r \rightarrow \infty.$$

For each $r \geq 1$, let $\widehat{I}_s^{t,r} \equiv \widehat{I}^{t,r}(s) := r^{-1/2}I^{t,r}(rs)$, and

$$(23) \quad \widehat{Z}_t^r = r^{-1/2}Z_{rt}^r, \quad \widehat{Z}^{t,r}(s) = r^{-1/2}Z^{rt,r}(rs),$$

$$(24) \quad \langle \widehat{q}^r(t), \varphi \rangle = \int_0^t \varphi(\widehat{Z}_s^r) \widehat{I}^{t,r}(ds) = \int_0^t \varphi(\widehat{Z}^{t,r}(s)) \widehat{I}^{t,r}(ds).$$

Convergence in (22) implies that, for t fixed,

$$(25) \quad -\widehat{I}^{t,r}(\cdot) \rightarrow -I^t(\cdot) \quad \text{a.s. in } D_{R_+}[0, t] \quad \text{as } r \rightarrow \infty,$$

where $I^t(s) = \inf_{u \in [s, t]} X_u$.

Let $Z_t, Z^t(s) \equiv Z_s^t$ be as in (21). Note that $\widehat{X}^r, X, \widehat{Z}^r, Z \in D_{R_+}[0, \infty)$, $\widehat{Z}^{t,r}(\cdot), Z^t(\cdot) \in D_{R_+}[0, t]$ and $\widehat{q}^r, q \in D_{M_f(R_+)}[0, \infty)$.

THEOREM 1. *Under assumptions (9), (10), (12), (13), we have $\widehat{q}^r \Rightarrow q$, as $r \rightarrow \infty$.*

The proof is based on (25) and the following lemma:

LEMMA 2. *Let $\widehat{Z}^{t,r}(s)$ be as in (23). Then, for each fixed $t \geq 0$,*

$$P\left(\sup_{s \in [0, t]} |\widehat{Z}^{t,r}(s) - Z^t(s)| > \varepsilon\right) \rightarrow 0 \quad \text{as } r \rightarrow \infty.$$

PROOF. Assume (22), fix $t > 0$ and a finite subdivision $0 \leq s_1 < s_2 < \dots < s_k = t$ on $[0, t]$. We show

$$(26) \quad \begin{aligned} & (\widehat{Z}^{t,r}(s_1), \widehat{Z}^{t,r}(s_2), \dots, \widehat{Z}^{t,r}(s_k)) \\ & \xrightarrow{p} (Z^t(s_1), Z^t(s_2), \dots, Z^t(s_k)), \quad r \rightarrow \infty, \end{aligned}$$

where \xrightarrow{p} denotes convergence in probability. Recall $\widehat{Z}^{t,r}(s) = \#\{u: t - s \leq u \leq t, \widetilde{S}^{\widehat{t},r}(u) > \widetilde{S}^{\widehat{t},r}(u-)\} \cdot r^{-1/2}$, where $\widetilde{S}^{\widehat{t},r}(u) = \sup_{x \in [0, u]} \widetilde{X}^{\widehat{t},r}(x)$ is the supremum process of $\widetilde{X}^{\widehat{t},r}(s) = r^{-1/2}(X^r(rt) - X^r((rt - rs)-))$, the rescaled and time-reversed X^r .

Consider the time-reversed process $\widetilde{X}^{rt,r}(s) = X^r(rt) - X^r((rt - s)-)$, $0 \leq s \leq rt$. For $z \leq t$, denote by $M_t^r(z)$ the number of jumps of $\widetilde{X}^{rt,r}$ above its past maximum in the interval $[0, rz]$. Note that $\widehat{Z}^{t,r}(s) = r^{-1/2}(M_t^r(t) - M_t^r((t - s)-)) = r^{-1/2}(M_t^r(t) - M_t^r(t - s))$ since a.s. there is no jump at time $t - s$. Due to (19), (20), it suffices to show that, for each fixed $z \in [0, t]$,

$$(27) \quad r^{-1/2} M_t^r(z) \xrightarrow{p} \frac{2}{\lambda\beta} \widetilde{S}^t(z), \quad r \rightarrow \infty,$$

where $\widetilde{S}^t(z)$ is defined above (19). The lemma will then follow from (26), since $\widehat{Z}^{t,r}(\cdot)$ is nondecreasing for each r and t , and $Z^t(\cdot)$ is continuous and nondecreasing for each t .

In order to show (27), it is convenient to consider “extension” process $(\widetilde{X}^{rt,t}(s), s \geq 0)$ of $(\widetilde{X}^{rt,t}(s), 0 \leq s \leq rt)$, defined in the following way. Independently of the filtration generated by X^r , take a sequence $\{u_{-i}^r, i \geq 1\}$ of i.i.d. exponential (rate λ^r) random variables, and a sequence $\{v_{-i}^r, i \geq 1\}$ of

i.i.d. random variables with distribution F^r . Define $\tilde{X}^{rt,r}(rt+z) = \tilde{X}^{rt,r}(rt) + \sum_{i=1}^{A_-^r(z)} v_{-i}^r - z$, $z \geq 0$, where $A_-^r(z) = \sup\{j: \sum_{i=1}^j u_{-i}^r \leq z\}$. Then the extended process $\tilde{X}^{rt,r}$ has the distribution of the load process X^r .

Let $\tilde{S}^{rt,r}$ be the supremum process of the extended $\tilde{X}^{rt,t}$. Denote by $T_1^{r,t} < T_2^{r,t} < \dots$ the successive increase (jump) times of $\tilde{S}^{rt,r}$, and let $J_i^{r,t} = \tilde{X}^{rt,r}(T_i^{r,t}) - \tilde{S}^{rt,r}(T_i^{r,t}-)$, $i \geq 1$, be the corresponding *overshoots*. If $\lambda^r m^r \geq 1$, there are infinitely many overshoots almost surely, and they form an i.i.d. sequence of random variables. The distribution of J_1^r is known,

$$(28) \quad \frac{P(J_1^r = dz)}{dz} = \lambda^r \int_{[z, \infty)} \exp\{-\Phi^r(0)(y-z)\} F^r(dy), \quad z \geq 0,$$

where $\Phi^r(0) \geq 0$ (see [5], page 188 for interpretation) is such that the right-hand side of (28) defines a (proper) probability distribution. If $\lambda^r m^r < 1$, there are $N_{r,t}$ many overshoots, where $N_{r,t}$ is a geometric $(1 - \lambda^r m^r)$ random variable, and conditionally on $N_{r,t}$, the overshoots $(J_1^{r,t}, \dots, J_{N_{r,t}}^{r,t})$ are independent and identically distributed, with known distribution

$$(29) \quad \frac{P(J_1^r = dz | \text{overshoot occurs})}{dz} = \frac{1}{m^r} F^r([z, \infty)), \quad z \geq 0.$$

Therefore, $E(J_1^r | \text{overshoot occurs}) = \beta^r / (2m^r)$ in this case. Both (28) and (29) are special cases of [5], Theorem VII.17(ii).

We are mainly interested in the overshoots that occurred by (reversed) time $rz \leq rt$. Note that

$$(30) \quad \widehat{S}^{t,r}(z) = \sum_{i=1}^{M_i^r(z)} \frac{1}{r^{1/2}} J_i^r,$$

and the convergence in (22) implies

$$(31) \quad \widehat{S}^{t,r}(z) \rightarrow \tilde{S}^t(z) \quad \text{a.s.} \quad \text{as } r \rightarrow \infty.$$

Assertion (27) is now a consequence of (30), (31) and Lemma 3 below. \square

LEMMA 3.

$$\frac{1}{M_i^r(z)} \sum_{i=1}^{M_i^r(z)} J_i^r \xrightarrow{p} \frac{\lambda\beta}{2} \quad \text{as } r \rightarrow \infty.$$

PROOF. An adaptation of the law of large numbers. It suffices to consider subsequences r_k of r for which either $\lambda^{r_k} m^{r_k} \geq 1$ for all k , or $\lambda^{r_k} m^{r_k} < 1$ for all k . Denote such subsequences again by r .

Assume $\lambda^r m^r \geq 1$. Note that

$$(32) \quad \Phi^r(0) \text{ in (28) tends to 0} \quad \text{as } r \rightarrow \infty.$$

This is due to (9), (10), (12), (13). Suppose $\Phi^r(0) > \delta$ along a subsequence, for some $\delta > 0$. By (9) and the Markov inequality, the sequence of service

time distributions F^r is tight. So we may assume that F^r converges weakly to distribution F along the same subsequence. Distribution F has mean m and second moment β , due to (12), (13). Let $v_1^r =^d F^r$ and $v_1 =^d F$. Due to

$$1 = P(J_1^r \geq 0) = \lambda^r \int_{[0, \infty)} \exp\{-\Phi^r(0)x\} F^r([x, \infty)) dx,$$

$m^r \lambda^r \rightarrow 1$, and $1 \geq \exp\{-\delta x\} > \exp\{-\Phi^r(0)x\}$, an application of the sandwich theorem gives

$$\begin{aligned} \lim_r \lambda^r \int_{[0, \infty)} e^{-\delta x} F^r([x, \infty)) dx &= \frac{\lambda}{\delta} \lim_r E(1 - \exp(-\delta v_1^r)) \\ &= \frac{\lambda}{\delta} E(1 - \exp(-\delta v_1)) = 1, \end{aligned}$$

where the limit is taken along the subsequence above. The last equality is impossible, since $E(1 - \exp(-\delta v_1)) \leq 1 - \exp(-\delta E v_1) = 1 - e^{-\delta m}$ and $\lambda(1 - e^{-\delta m})/\delta$ is strictly smaller than $\lambda m = 1$.

Due to (13), random variables J_1^r in (28) are uniformly integrable. Namely, an application of Fubini's theorem gives

$$\begin{aligned} \lim_{K \rightarrow \infty} E J_1^r \mathbf{1}_{\{J_1^r \geq K\}} &\leq \lim_{K \rightarrow \infty} \lambda^r \int_{[K, \infty)} z F^r([z, \infty)) dz \\ &\leq \lim_{K \rightarrow \infty} \frac{1}{2} \int_{[K, \infty)} z^2 F^r(dz) = 0. \end{aligned}$$

Since

$$\begin{aligned} E J_1^r &= \int_0^\infty P(J_1^r \geq y) dy \\ &= \lambda^r \int_0^\infty F^r(dz) \int_0^\infty dy \int_0^\infty dx \exp(-\Phi^r(0)x) \mathbf{1}_{\{x \geq 0, y \geq 0, x+y \leq z\}}, \end{aligned}$$

by (32) and uniform integrability we have $E J_1^r \rightarrow (\lambda\beta)/2$ as $r \rightarrow \infty$. So, it suffices to show

$$\frac{1}{M_t^r(z)} \sum_{i=1}^{M_t^r(z)} (J_i^r - E J_1^r) \xrightarrow{p} 0 \quad \text{as } r \rightarrow \infty.$$

It is not hard to see that $M_t^r(z) \rightarrow \infty$ in probability, due to (30) and (31). Take $\varepsilon' > 0$. Since

$$\begin{aligned} P\left(\left|\frac{1}{M_t^r(z)} \sum_{i=1}^{M_t^r(z)} (J_i^r - E J_1^r)\right| > 4\varepsilon'\right) \\ \leq P(M_t^r(z) < l) + P\left(\sup_{k \geq l} \left|\frac{1}{k} \sum_{i=1}^k (J_i^r - E J_1^r)\right| > 4\varepsilon'\right), \end{aligned}$$

it suffices to show that

$$(33) \quad P\left(\sup_{k \geq l} \left|\frac{1}{k} \sum_{i=1}^k (J_i^r - E J_1^r)\right| > 4\varepsilon'\right) \leq o(r, l),$$

where $\lim_{l \rightarrow \infty} \sup_r o(r, l) = 0$. Note that (33) is just an extension of the usual strong law of large numbers. The proof is a variation on the classical proof. Fix large K and split the sum in (33) into

$$\sum_{i=1}^k \left(J_i^r \mathbf{1}_{\{J_i^r \leq K\}} - E J_i^r \mathbf{1}_{\{J_i^r \leq K\}} \right) + \sum_{i=1}^k \left(J_i^r \mathbf{1}_{\{J_i^r > K\}} - E J_i^r \mathbf{1}_{\{J_i^r > K\}} \right).$$

Random variables $J_i^r \mathbf{1}_{\{J_i^r \leq K\}}$ are bounded; in particular, they have uniformly (in r) bounded fourth moments and the simplest proof of the strong law of large numbers yields

$$P \left(\sup_{k \geq l} \left| \frac{1}{k} \sum_{i=1}^k \left(J_i^r \mathbf{1}_{\{J_i^r \leq K\}} - E J_i^r \mathbf{1}_{\{J_i^r \leq K\}} \right) \right| > 2\varepsilon' \right) \leq o(r, l, K)$$

where $\lim_{l \rightarrow \infty} \sup_r o(r, l, K) = 0$. By uniform integrability, $E J_i^r \mathbf{1}_{\{J_i^r > K\}} = \frac{1}{k} \sum_{i=1}^k E J_i^r \mathbf{1}_{\{J_i^r > K\}}$ can be made uniformly small in r for large enough K . So it suffices to show

$$(34) \quad \limsup_{l \rightarrow \infty} \sup_r P \left(\sup_{k \geq l} \frac{1}{k} \sum_{i=1}^k J_i^r \mathbf{1}_{\{J_i^r > K\}} > \varepsilon' \right) \leq o(K),$$

where $\lim_K o(K) = 0$. Perhaps the easiest way to show (34) is by imitating Etemadi's proof of SLLN [11], Theorem 1.8.4. Define truncated $H_i^{r, K} = J_i^r \mathbf{1}_{\{J_i^r > K\}} \mathbf{1}_{\{J_i^r \mathbf{1}_{\{J_i^r > K\}} \leq i\}}$. Then

$$\sup_r P \left(H_i^{r, K} \neq J_i^r \mathbf{1}_{\{J_i^r > K\}} \text{ for some } i \geq 1 \right) \leq \sup_r E J_i^r \mathbf{1}_{\{J_i^r > K\}} \leq o(K),$$

where $\lim_K o(K) = 0$, so it suffices to show

$$\limsup_{l \rightarrow \infty} \sup_r P \left(\sup_{k \geq l} \frac{1}{k} \sum_{i=1}^k H_i^{r, K} > \varepsilon' \right) \leq o(K).$$

Fix $\varepsilon'' > 0$ and $a > 1$. Let $T_n^{r, K} = \sum_{i=1}^n H_i^{r, K}$, and let $k(n) = \lfloor a^n \rfloor$. By the same calculation as in Etemadi's proof, we get

$$\begin{aligned} \sup_r \sum_{n=1}^{\infty} P \left(\left| T_{k(n)}^{r, K} - E T_{k(n)}^{r, K} \right| > \varepsilon'' k(n) \right) \\ \leq \sup_r \frac{16}{(1-a^2)\varepsilon''^2} E J_i^r \mathbf{1}_{\{J_i^r > K\}} \leq o(K). \end{aligned}$$

Since

$$\sup_K \sup_r \left| E H_i^{r, K} - E J_i^r \mathbf{1}_{\{J_i^r > K\}} \right| \rightarrow 0 \quad \text{as } i \rightarrow \infty,$$

then $\sup_r |ET_{k(n)}^{r,K}/k(n) - EJ_1^r 1_{\{J_1^r > K\}}| \rightarrow 0$, as $n \rightarrow \infty$, and therefore,

$$\lim_{l \rightarrow \infty} \sup_r P\left(\sup_{k(n) \geq l} \frac{1}{k(n)} T_{k(n)}^{r,K} > 2\varepsilon''\right) \leq o(K).$$

The rest is the same as in Etemadi's proof. Details are left to the reader.

Assume $\lambda^r m^r < 1$. By the same reasoning as above, again $M_t^r(z) \xrightarrow{P} \infty$, as $r \rightarrow \infty$. Similarly, the distributions in (29) are again uniformly integrable. Recall $E(J_1^r | \text{overshoot occurs}) = \beta^r / (2m^r) \rightarrow \lambda\beta/2$. Write

$$\begin{aligned} &P\left(\left|\frac{1}{M_t^r(z)} \sum_{i=1}^{M_t^r(z)} \left(J_i^r - \frac{\lambda^r \beta^r}{2}\right)\right| > 4\varepsilon'\right) \\ &\leq P(M_t^r(z) < l) + P\left(\sup_{l \leq k \leq N_{r,t}} \left|\frac{1}{k} \sum_{i=1}^k \left(J_i^r - \frac{\lambda^r \beta^r}{2}\right)\right| > 4\varepsilon', N_{r,t} \geq l\right) \\ &\leq P(M_t^r(z) < l) \\ &\quad + \sum_{j=l}^{\infty} P(N_{r,t} = j) P\left(\sup_{l \leq k \leq j} \left|\frac{1}{k} \sum_{i=1}^k \left(J_i^r - \frac{\lambda^r \beta^r}{2}\right)\right| > 4\varepsilon' \mid N_{r,t} = j\right). \end{aligned}$$

By the observation made above (29), the last term in the sum above is dominated by

$$P\left(\sup_{k \geq l} \left|\frac{1}{k} \sum_i (\bar{J}_i^r - E\bar{J}_1^r)\right| > 4\varepsilon'\right),$$

where $\bar{J}_i^r, i \geq 1$, are i.i.d. with distribution (29), and the rest of the proof is the same as in the supercritical case. \square

COROLLARY 4. For any fixed t and $0 \leq t_1 < t_2 \dots < t_k \leq t$,

$$\left(\widehat{Z}_{t_1}^r, \widehat{Z}_{t_2}^r, \dots, \widehat{Z}_{t_k}^r\right) \xrightarrow{P} \left(Z_{t_1}, Z_{t_2}, \dots, Z_{t_k}\right).$$

PROOF OF THEOREM 1. From Lemma 2 and (22)–(25), we easily obtain the convergence of finite-dimensional distributions

$$\left(\hat{q}_{t_1}^r, \dots, \hat{q}_{t_k}^r\right) \Rightarrow \left(q_{t_1}, \dots, q_{t_k}\right), \quad r \rightarrow \infty.$$

So it suffices to show the tightness of the family (\hat{q}^r) in $D_{M_f(R_+)}[0, \infty)$. Let $q = q^r$, fix a bounded and continuous function $\varphi: [0, \infty) \rightarrow R$ and calculate

$$\begin{aligned} &|\langle q(s+h), \varphi \rangle - \langle q(s), \varphi \rangle| \\ &= \left| \left(\int_0^s + \int_s^{s+h} \right) \varphi(Z_z^{s+h}) dI_z^{s+h} - \int_0^s \varphi(Z_z^s) dI_z^s \right| \\ &\leq \left| \int_0^s \varphi(Z_z^{s+h}) dI_z^s - \int_0^s \varphi(Z_z^{s+h}) dI_z^{s+h} \right| \\ &\quad + \left| \int_s^{s+h} \varphi(Z_z^{s+h}) dI_z^{s+h} \right| + \left| \int_0^s \varphi(Z_z^s) - \varphi(Z_z^{s+h}) dI_z^s \right|, \end{aligned}$$

where the inequality is an application of the triangle inequality, after convenient rearrangement of the terms. Use the fact that $I_z^s - I_z^{s+h}$ is nondecreasing in $z \in [0, s]$ in order to bound the first term from above by

$$\begin{aligned} & \|\varphi\|_\infty \left| \int_0^s dI_z^s - dI_z^{s+h} \right| \\ &= \|\varphi\|_\infty (I_s^s - I_0^s - (I_s^{s+h} - I_0^{s+h})) \\ &\leq \|\varphi\|_\infty (X_s - I_s^{s+h}). \end{aligned}$$

For the second term, use a similar identity $\int_s^{s+h} dI_z^{s+h} = I_{s+h}^{s+h} - I_s^{s+h} = X_{s+h} - I_s^{s+h}$. To bound the third term, note that Z_z^s and Z_z^{s+h} differ (clearly $Z_z^s \geq Z_z^{s+h}$ for all $z \in [0, s]$) only for $z \in [0, s]$ such that $I_z^s > I_z^{s+h}$, that is $I_z^s > I_s^{s+h}$. Now $\int_0^s (\varphi(Z_z^s) - \varphi(Z_z^{s+h})) dI_z^s = \int_{\{z \in [0, s]: I_z^s > I_s^{s+h}\}} (\varphi(Z_z^s) - \varphi(Z_z^{s+h})) dI_z^s$, so that

$$\left| \int_0^s (\varphi(Z_z^s) - \varphi(Z_z^{s+h})) dI_z^s \right| \leq 2 \|\varphi\|_\infty (I_s^s - I_s^{s+h}).$$

The above calculations imply

$$\begin{aligned} |\langle q(s+h), \varphi \rangle - \langle q(s), \varphi \rangle| &\leq \|\varphi\|_\infty (3(X_s - I_s^{s+h}) + (X_{s+h} - I_s^{s+h})) \\ &\leq 4\|\varphi\|_\infty \sup_{\theta \in [0, h]} |X_{s+\theta} - X_s|, \end{aligned}$$

and after scaling,

$$|\langle \hat{q}^r(s+h), \varphi \rangle - \langle \hat{q}^r(s), \varphi \rangle| \leq 4\|\varphi\|_\infty \sup_{\theta \in [0, h]} |\hat{X}_{s+\theta}^r - \hat{X}_s^r|.$$

The tightness of \hat{q}^r now follows from the tightness of \hat{X}^r , by combining Jakubowski and Aldous criteria, [10], Theorems 3.6.4 and 3.6.5 (see also [13], Theorem III.8.6).

REMARK. Since the mapping $\mu \mapsto \text{sup}(\text{Supp}(\mu))$ is clearly not continuous in the topology on $M_f(R_+)$, Theorem 5 below does not immediately follow from Theorem 1 and (6). If X in (22) is a stable- α ($\alpha \in (1, 2)$) and the processes Z_s^t, Z_t are appropriately defined (cf. Section 3.4.4), then assertion (26) becomes a consequence of the analysis in [15], Section 5. Therefore, an analogue of Theorem 1 exists in the stable- α setting, $\alpha \in (1, 2)$.

THEOREM 5. Under assumptions (9), (10), (12), (13), we have $\hat{Z}^r \Rightarrow Z$ as $r \rightarrow \infty$.

The difficulties in analyzing \hat{Z}^r and Z are related to their lack of Markov property. Fix some $\varepsilon > 0$ and $\eta > 0$. Fix time $T > 0$, and let $t_i = i(T/n)$, $0 \leq i \leq n$, be the subdivision of $[0, T]$ with mesh size T/n . For n large enough, we have

$$(35) \quad P\left(\sup_{1 \leq i \leq n} \sup_{u \in [t_{i-1}, t_i]} |Z_{t_i} - Z_u| > \varepsilon\right) \leq \eta,$$

by continuity of Z (cf. [13], page 122). Recall that, for each t , the process $(Z^t(s), 0 \leq s \leq t)$ given by (21) is continuous. The processes $Z^{t_i}(t_i) - Z^{t_i}(t_i - \theta) = \tilde{L}_\theta^{t_i}$, $\theta \in [0, T/n]$, $1 \leq i \leq n$, as defined in (19), are all identically distributed. The processes $\tilde{L}_\theta^{t_i}$ are also independent, due to independent increments of Brownian motion, but this property is not used in the argument. We claim that, for all large n ,

$$(36) \quad P\left(\sup_{1 \leq i \leq n} \sup_{\theta \in [0, T/n]} |Z^{t_i}(t_i) - Z^{t_i}(t_i - \theta)| > \varepsilon\right) \leq \eta.$$

Local time $\tilde{L}_\theta^{t_i}$ is increasing in the variable θ , so it suffices to show that, for any $\varepsilon > 0$, $nP(L_{T/n} > \varepsilon) \rightarrow 0$ as $n \rightarrow \infty$. This follows from (20) and the following easy lemma. The proof is left to the reader.

LEMMA 6. For $S_t = \sup_{u \in [0, t]} B_u$, the supremum process of Brownian motion B with variance 1 and drift $c \in R$, we have

$$(37) \quad nP(S_{1/n} > \varepsilon) \rightarrow 0 \quad \text{as } n \rightarrow \infty \quad \text{for all } \varepsilon > 0.$$

PROOF OF THEOREM 5. The finite-dimensional distributions of \hat{Z}^r are converging to those of Z due to Corollary 4. So it suffices to show the tightness of \hat{Z}^r , $r \geq 1$, with respect to the Skorokhod topology on $D_R[0, \infty)$. It suffices to show that any subsequence r_n has a further subsequence r_{n_k} so that $\hat{Z}^{r_{n_k}}$ is tight. For a given subsequence r_n , find a weakly converging further subsequence $F^{r_{n_k}} \rightarrow F$. This is possible again by tightness. Recall that F has mean m and second moment β . To simplify the notation, denote the subsequence r_{n_k} again by r .

Formally, the idea is to use $\hat{Z}_{t_i - \theta}^r \approx \hat{Z}^{t_i, r}(t_i - \theta) \approx \hat{Z}^{t_i, r}(t_i) = \hat{Z}_{t_i}^r \approx Z_{t_i}$ for small θ , and exploit the monotonicity of $\hat{Z}_s^{t_i, r}$ and Z_s^t in s . Let \mathcal{F}_t^r be the filtration generated by \hat{X}^r . Let $t_i = iT/n$ as above, and $t \in [t_{i-1}, t_i]$. Observe that, for each r ,

$$(38) \quad \hat{Z}_{t_{i-1}}^{t_i, r} \leq \hat{Z}_t^{t_i, r} \leq \hat{Z}_t^r \leq \hat{Z}_{t_{i-1}}^r + \hat{Z}_{t-t_{i-1}}^{r, i},$$

where $(\hat{Z}_u^{r, i}, u \in [0, T/n])$ has the same law as $(\hat{Z}_u^r, u \in [0, T/n])$, and is independent of $\mathcal{F}_{t_{i-1}}^r$. The first inequality in (38) is the monotonicity of $Z_s^{t_i, r}$ in s ; the second inequality trivially follows from the interpretation of $\hat{Z}_t^{t_i, r}$ as the (rescaled) number of individuals in queue at time rt whose service will not have been completed by time rt_i . The last inequality in (38) is a special case of [15], Lemma 4.5, though it can be argued using again queueing interpretation: $\hat{Z}_{t-t_{i-1}}^{r, i}$ is the (rescaled) number of customers who arrived to the queue in the time interval $[rt_{i-1}, rt]$ and who did not exit by time rt . In particular,

$$(39) \quad \hat{Z}_{t_i}^r = \hat{Z}_{t_{i-1}}^{t_i, r} + \hat{Z}_{t_i - t_{i-1}}^{r, i}, \quad 1 \leq i \leq n, \quad \text{almost surely.}$$

PROPOSITION 7. For any fixed $\varepsilon, \eta > 0$, there exist an integer $n \geq 1$, and $r_1 \geq 1$ so that

$$(40) \quad \sup_{r \geq r_1} P \left(\sup_{1 \leq i \leq n} \left| \widehat{Z}_{t_i}^r - \widehat{Z}_{t_{i-1}}^r \right| > 2\varepsilon \right) \leq 2\eta$$

and

$$(41) \quad \sup_{r \geq r_1} P \left(\sup_{1 \leq i \leq n} \sup_{s \in [t_{i-1}, t_i]} \left| \widehat{Z}_s^{t_i, r} - \widehat{Z}_{t_i}^{t_i, r} \right| > 2\varepsilon \right) \leq 2\eta,$$

$$(42) \quad \sup_{r \geq r_1} P \left(\sup_{1 \leq i \leq n} \sup_{u \in [0, T/n]} \widehat{Z}_u^{r, i} > 5\varepsilon \right) \leq 18\eta.$$

By Corollary 4, Lemma 2 and (35), (36), we can find $r_1 \geq 1$ so that both (40) and (41) hold. The hard estimate (42) will be shown in the next section. Estimates (40), (41) imply

$$\sup_{r \geq r_1} P \left(\sup_{1 \leq i \leq n} \left| \widehat{Z}_{t_{i-1}}^r - \widehat{Z}_{t_{i-1}}^{t_i, r} \right| > 4\varepsilon \right) \leq 4\eta,$$

so the left-most and the right-most side in (38) “typically differ” by at most $4\varepsilon + \sup_{u \in [0, T/n]} \widehat{Z}_u^{r, i}$. Combined with (42), this implies that, for any $0 < h < T/n$, we have

$$\sup_{r \geq r_1} P \left(\sup_{|s-t| < h} \left| \widehat{Z}_s^r - \widehat{Z}_t^r \right| > 20\varepsilon \right) \leq 22\eta.$$

The last estimate gives relative compactness of the sequence Z^r , for example, by [13], Corollary III.7.4, completing the proof of Theorem 5.

3.3. *Tree estimates.* This section proves assertion (42) in Proposition 7. The proof uses estimates for the joint total size and height distribution of a sequence of supercritical (near-critical) Galton–Watson trees.

Let \mathcal{T} be a Galton–Watson random tree with offspring distribution Ξ , where Ξ is concentrated on nonnegative integers. By this we mean a tree-valued random variable constructed from a sequence of i.i.d.- Ξ random variables. The root of the tree is the zero generation. In the first step, the root gives birth to ξ_0 children, where $\xi_0 \sim_d \Xi$. If $\xi_0 = 0$, then \mathcal{T} consists of the root only. If $\xi_0 \geq 1$, then the root has children s_j^1 , $1 \leq j \leq \xi_0$, that form the first generation. Each vertex s_j^1 is connected to the root by an edge. The tree is formed recursively. In the n th step, each vertex in the $(n-1)$ st generation gives birth according to Ξ , independently of others and the previous generations. Again an edge connects each child to its parent. The children of vertices in the $(n-1)$ st generation are the n th generation. Continue until *extinction* (no births from any vertex in the same generation) occurs, or forever, if no extinction occurs. For any vertex $\varsigma \in \mathcal{T}$, let $\text{gen}(\varsigma)$ denote its generation number. Let \mathcal{T}_ς denote the tree spanned by ς and all of its descendents (children, children of children, etc.). Then $\mathcal{T}_\varsigma =^d \mathcal{T}$ is an elementary consequence of the construction.

The Galton–Watson tree is called *strictly supercritical* if $E\xi_0 > 1$, *critical* if $E\xi_0 = 1$ and *strictly subcritical* if $E\xi_0 < 1$. Similarly, \mathcal{T} is *super (sub)-critical* if $E\xi_0 \geq 1 (\leq 1)$. Let $|\mathcal{T}|$ denote the *total size* (number of vertices) of \mathcal{T} , and let $\text{ht}(\mathcal{T})$ denote the *height* (the maximal generation number) of \mathcal{T} , respectively. The case of $P(\xi_0 = 1) = 1$ does not appear in the setting of this paper, so we exclude it from consideration. Then, it is well known (e.g., [3]) that subcritical trees have finite size (therefore height) with probability 1, whereas strictly supercritical trees have infinite size (height) with nonzero probability. One readily checks by induction that the number of vertices in the n th generation of \mathcal{T} has expectation $(E\xi_0)^n$.

Recall the branching interpretation for the queue length from Section 1.2. Each busy cycle of the queue corresponds to an excursion of the load (workload) process, and yields a Galton–Watson tree of customers who entered (and exited) the queue during this busy cycle. There is one-to-one and onto correspondence between the vertices of the tree and the customers of the busy cycle. A new customer that arrives at time s creates a new vertex s in the corresponding tree. If the queue was empty immediately before the arrival ($Z(s-) = 0$), then s becomes the root of the tree. Otherwise, s becomes a child of the customer whose service was interrupted, and $\text{gen}(s) = Z(s-) = Z(s) - 1$. The queue-length process Z is the height process (or depth-first search walk) that visits trees in chronological order (of busy cycles), and within each tree visits vertices in the depth-first search (LIFO) order. At each time u , $Z(u)$ records the generation number (plus 1) of the vertex corresponding to the customer currently in service (if any).

It is not hard to calculate the exact offspring distribution Ξ^r corresponding to the queue of index r . Consider a typical customer. At the moment of arrival, this customer requests $v^r \stackrel{d}{=} F^r$ amount of service time. Its service may be interrupted several times due to new arrivals, each such arrival producing a single offspring. Now it is easy to see that, after conditioning on v^r , the total number of offspring is Poisson $(\lambda^r v^r)$, so that

$$(43) \quad \Xi^r(i) = P(\xi^r = i) = E \left[\frac{\exp(-\lambda^r v^r)(\lambda^r v^r)^i}{i!} \right], \quad i = 0, 1, 2, \dots$$

Without loss of generality, we can assume supercriticality, $E\xi^r = \lambda^r E v^r = \lambda^r m^r \geq 1$, for all $r \geq 1$. If $E\xi^r < 1$, define $c^r = 1/E\xi^r = 1/\lambda^r m^r > 1$ and $\bar{v}^r = c^r v^r$ with distribution \bar{F}^r . If $E\xi^r \geq 1$, set $\bar{F}^r = F^r$. Then clearly we can couple the queue-length processes $Z_u^{r,i}$ and $\bar{Z}_u^{r,i}$ having service distributions F^r and \bar{F}^r , so that, for each r ,

$$Z_u^{r,i} \leq \bar{Z}_u^{r,i} \quad \text{for all } u, \quad 1 \leq i \leq n, \quad \text{a.s.}$$

The convergence $F^r \rightarrow F$ implies $\bar{F}^r \rightarrow F$, convergence relations (9), (10), (12) continue to hold and (13) implies

$$\sup_r E \left[(\bar{v}^r)^2 1_{\{\bar{v}^r \geq K\}} \right] \rightarrow 0 \quad \text{as } K \rightarrow \infty.$$

Henceforth, we assume $E\xi^r \geq 1$ so that $\bar{F}^r = F^r, r \geq 1$.

PROOF OF (42). Recall the setting of Proposition 7. For each i , we have $\widehat{Z}_u^{r,i} = (1/\sqrt{r})Z_{ru}^{r,i}$, $u \in [0, T/n]$, where $Z^{r,i}$ is a copy of the queue length Z^r , and relations (38), (39) are satisfied.

Let $\mathcal{T}_{i,j}^r$, $1 \leq j \leq M_i^r$, be the Galton–Watson trees of the busy cycles (corresponding to $Z^{r,i}$) started after time 0 and completed before time rT/n . Let \mathcal{T}_i^r be the tree of the busy cycle containing the customer present in service at time rT/n . Denote by $\mathcal{T}_i^{r,\dagger}$ the initial portion of the last tree \mathcal{T}_i^r , traversed by the queue length $Z^{r,i}$ up to time rT/n . If the queue is empty at time rT/n , set $\mathcal{T}_i^r = \mathcal{T}_i^{r,\dagger} = \emptyset$ to be the empty tree with $\text{ht}(\emptyset) = 0$. Due to the above reasoning, the maximal queue length $\sup_{u \in [0, rT/n]} Z_{ru}^{r,i}$ is dominated by $\max_{1 \leq j \leq M_i^r} \text{ht}(\mathcal{T}_{i,j}^r) \vee \text{ht}(\mathcal{T}_i^{r,\dagger}) + 1$, the maximal height of all trees (of busy cycles) started in $[0, rT/n]$. So assertion (42) will follow from

$$(44) \quad \sup_{r \geq r_1} P \left(\max_{1 \leq i \leq n} \left(\max_{1 \leq j \leq M_i^r} \text{ht}(\mathcal{T}_{i,j}^r) \right) \vee \text{ht}(\mathcal{T}_i^{r,\dagger}) > 5\varepsilon\sqrt{r} \right) \leq 18\varepsilon\eta.$$

It is well known (cf. Athreya and Ney [3], Theorem I.9.1) that

$$(45) \quad P(\text{ht}(\mathcal{T}) > r) \sim \frac{2}{\sigma^2 r} \quad \text{as } r \rightarrow \infty,$$

where \mathcal{T} is a critical Galton–Watson tree with offspring distribution ξ , $E\xi = 1$ and $0 < \text{var}(\xi) = \sigma^2 < \infty$. This is also the content of Kolchin [28], Theorem 2.1.2. Aldous ([2], Proposition 24) gives the following estimate for the joint height and total size distribution of the same tree:

$$(46) \quad r^{1/2} P(\text{ht}(\mathcal{T}) > \varepsilon r^{1/2}, |\mathcal{T}| < \delta r) \rightarrow \sigma^{-1} \delta^{-1/2} G(\varepsilon \delta^{-1/2} \sigma) \text{ as } r \rightarrow \infty,$$

where $G(x) \leq \kappa_1 \exp(-x/\kappa_2)$, $0 < x < \infty$, for some $0 < \kappa_1, \kappa_2 < \infty$. Since we allow the offspring distribution Ξ^r to vary with r , we will need the following analogous lemma.

LEMMA 8. *Let \mathcal{T}^r be a sequence of supercritical Galton–Watson trees with offspring distribution ξ^r such that:*

- (i) $\xi^r \rightarrow_d \xi \sim_d \Xi$, where $E\xi = 1$,
- (ii) $r^{1/2}(1 - E\xi^r) \rightarrow c \leq 0$,
- (iii) $0 < \text{var}(\xi^r) = (\sigma^r)^2 \rightarrow \sigma^2 \in (0, \infty)$,
- (iv) $\sup_r E((\xi^r)^2 \mathbf{1}_{\{\xi^r > K\}}) \rightarrow 0$, as $K \rightarrow \infty$.

Then

$$(47) \quad \limsup_r r^{1/2} P(\text{ht}(\mathcal{T}^r) > \varepsilon r^{1/2}, |\mathcal{T}^r| < \delta r) \leq \sigma^{-1} \delta^{-1/2} G(\varepsilon \delta^{-1/2} \sigma),$$

where $G(x) \leq \kappa_1 \exp(-x/\kappa_2)$, $0 < x < \infty$, for some $0 < \kappa_1, \kappa_2 < \infty$.

Note that $F^r \rightarrow F$, together with (43), (13) and (9), (10), (12), imply the conditions (i)–(iv) of the lemma. The uniform integrability condition (iv) is natural (cf. [30, 17]).

Assume Lemma 8 for now. Let $\varepsilon, \eta > 0$ and T be as in Proposition 7, let K_1 be a large number such that $e^{-4c\varepsilon}(1 + \eta)/K_1 < \eta$ and recall that λ is the asymptotic arrival rate. Choose n_1 large enough so that

$$(48) \quad K_1 \frac{\kappa_1(\lambda + \varepsilon)}{\sigma\sqrt{(\lambda + \varepsilon)T}} n^{3/2} \exp\left(-\frac{\varepsilon\sigma n^{1/2}}{\kappa_2\sqrt{(\lambda + \varepsilon)T}}\right) \leq \eta \quad \text{for all } n \geq n_1,$$

for κ_1, κ_2 in the lemma, and also large enough so that (35), (36) are satisfied for all $n \geq n_1$. Fix some $n \geq n_1$. Assume $r_1 \geq 1$ to be large enough so that both (40) and (41) hold. We use estimate (47) to bound the probabilities of events

$$\left\{ \max_{1 \leq i \leq n} \max_{1 \leq j \leq M_i^r} \text{ht}(\mathcal{S}_{i,j}^r) > \varepsilon\sqrt{r} \right\} \quad \text{and} \quad \left\{ \max_{1 \leq i \leq n} \text{ht}(\mathcal{S}_i^{r,\dagger}) > \varepsilon\sqrt{r} \right\}.$$

Recall M_i^r is the number of trees corresponding to the completed busy cycles of $Z^{r,i}$, and let $N_i^r = |\mathcal{S}_{i,1}^r| + |\mathcal{S}_{i,2}^r| + \dots + |\mathcal{S}_{i,M_i^r}^r| + |\mathcal{S}_i^{r,\dagger}|$ be the total number of vertices visited before time rT/n .

LEMMA 9. *For any fixed $n \geq n_1$, there exists $r_3 \geq 1$ such that*

$$(49) \quad \sup_{r \geq r_3} P\left(\max_{1 \leq i \leq n} N_i^r \geq (\lambda + \varepsilon)\frac{T}{n}r\right) \leq \eta,$$

$$(50) \quad \sup_{r \geq r_3} P\left(\max_{1 \leq i \leq n} M_i^r \geq \sqrt{r}(\lambda + \varepsilon)\right) \leq 3\eta.$$

PROOF. The first assertion is easy since $N_i^r =^d$ Poisson (rate $r\lambda^r T/n$), $1 \leq i \leq n$. For the second one, consider processes $X^{r,i} = (X_{t_{i-1}+s}^r - X_{t_{i-1}}^r, s \in [0, rT/n])$, and let $I_s^{r,i} = \inf_{0 \leq u \leq s} X_u^{r,i}$, and $\tau_x^{r,i} := \inf\{s \geq 0: I_s^{r,i} \leq -x\}$. Note that $I_s^{r,i} = -|\{u \in [0, s]: Z_u^{r,i} = 0\}|$. Recall how (37) implied (41). Since the asymptotic load X is a Brownian motion, the same assertion (37) implies that, for any $n \geq n_1$, we can find r_3 large enough so that

$$\sup_{r \geq r_3} P\left(\min_{1 \leq i \leq n} \tau_{\sqrt{r}}^{r,i} < rT/n\right) \leq \sup_{r \geq r_3} nP\left(-I_{rT/n}^{r,1} > \sqrt{r}\right) \leq 2\eta.$$

On the complement of $\{\min_{1 \leq i \leq n} \tau_{\sqrt{r}}^{r,i} < rT/n\}$, we have $M_i^r \leq M_{\sqrt{r}}^{r,i}$, $1 \leq i \leq n$, where $M_{\sqrt{r}}^{r,i}$ equals the number of busy cycles started, and completed, during $[0, \tau_{\sqrt{r}}^{r,i}]$. By (5), $M_{\sqrt{r}}^{r,i} =^d M_{\sqrt{r}} =^d$ Poisson (rate $\sqrt{r}\lambda$), and the second assertion of the lemma follows just like the first one. \square

Now for any fixed $n \geq n_1$ and any $r \geq \max\{r_1, r_3\}$, we get

$$\begin{aligned} & P\left(\max_{1 \leq i \leq n} \max_{1 \leq j \leq M_i^r} \text{ht}(\mathcal{T}_{i,j}^r) > \varepsilon\sqrt{r}\right) \\ & \leq 4\eta + \sum_{i=1}^n P\left(\max_{1 \leq j \leq M_i^r \leq \sqrt{r}(\lambda + \varepsilon)} \text{ht}(\mathcal{T}_{i,j}^r) > \varepsilon\sqrt{r}, \max_j |\mathcal{T}_{i,j}^r| \right. \\ & \qquad \left. < (\lambda + \varepsilon)\frac{T}{n}r, \text{ and } M_i^r \leq \sqrt{r}(\lambda + \varepsilon)\right) \\ & \leq 4\eta + n\sqrt{r}(\lambda + \varepsilon)P\left(\text{ht}(\mathcal{T}^r) > \varepsilon\sqrt{r}, |\mathcal{T}^r| < (\lambda + \varepsilon)\frac{T}{n}r\right), \end{aligned}$$

and by (47) there exists r_4 (possibly larger than r_3) so that, for each $r \geq r_4$,

$$\begin{aligned} & P\left(\max_{1 \leq i \leq n} \max_{1 \leq j \leq M_i^r} \text{ht}(\mathcal{T}_{i,j}^r) > \varepsilon\sqrt{r}\right) \\ & \leq 5\eta + \frac{\kappa_1(\lambda + \varepsilon)}{\sigma\sqrt{(\lambda + \varepsilon)T}}n^{3/2} \exp\left(-\frac{\varepsilon\sigma n^{1/2}}{\kappa_2\sqrt{(\lambda + \varepsilon)T}}\right) \\ & \leq 6\eta \quad (\text{by (48)}). \end{aligned}$$

LEMMA 10. For $n \geq n_1$ fixed as above, there exists some $r_5 \geq r_4$ such that

$$\sup_{r \geq r_5} P\left(\max_{1 \leq i \leq n} \text{ht}(\mathcal{T}_i^{r,\dagger}) > 5\varepsilon\sqrt{r}\right) \leq 12\eta.$$

PROOF. The following is an extension of the idea in the argument for (50). Let $X^{r,i}$, $I^{r,i}$ and $\tau^{r,i}$ be as in Lemma 9. Let $(Y_s^{r,i}, s \geq 0)$, $1 \leq i \leq n$, be mutually independent, distributed as the load process $(X_s^r, s \geq 0)$ and independent of X^r . One can construct a new copy $X^{*,r} =^d X^r$ from $X^{r,i}$ and $Y^{r,i}$ as described below. The point of the construction is that (typically) each tree $\mathcal{T}_{i,j}^r$, $1 \leq j \leq M_i^r$, $1 \leq i \leq n$, corresponding to a busy cycle of $Z^{r,i}$ (that is, $X^{r,i}$) reappears as a tree corresponding to a busy cycle of $X^{*,r}$. More importantly, each tree $\mathcal{T}_i^{r,\dagger}$, $1 \leq i \leq n$, reappears as the initial portion of a tree of some busy cycle of $X^{*,r}$. The idea is simple, but the notation could get messy, so sometimes we omit “ r ” in the superscript. Define stopping times $\tau_i = \tau_{\sqrt{r}}^{r,i} \wedge (rT/n)$, $1 \leq i \leq n$. Then $\bar{X}^{*,r,i}$, defined by

$$\bar{X}_s^{*,r,i} = X_s^{r,i}1_{\{s \leq \tau_i\}} + (X_{\tau_i}^{r,i} + Y_{s-\tau_i}^{r,i})1_{\{s > \tau_i\}},$$

equals X^r in distribution, due to independence of $X^{r,i}$ and $Y^{r,i}$, and the strong Markov property of X^r . Note that, moreover, $\bar{X}^{*,r,i}$, $1 \leq i \leq n$, are mutually independent as processes.

Let $\bar{\tau}_{\sqrt{r}}^i = \inf\{s \geq 0: \inf_{u \leq s} \bar{X}_u^{*,r,i} < -\sqrt{r}\}$. Then the processes

$$(51) \quad (\bar{X}_s^{*,r,i}, s \in [0, \bar{\tau}_{\sqrt{r}}^i]), \quad 1 \leq i \leq n,$$

are independent and identically distributed, where $(\bar{X}_s^{*,r,1}, s \in [0, \bar{\tau}_{\sqrt{r}}^1]) =^d (X_s^r, s \in [0, \tau_{\sqrt{r}}^r])$, and $\tau_{\sqrt{r}}^r = \inf\{u \leq s: X_u^r < -\sqrt{r}\}$. Now define $\tau_0^* = 0$, $\tau_{i\sqrt{r}}^* := \sum_{j \leq i} \bar{\tau}_{\sqrt{r}}^j$, and let

$$X_s^{*,r} = \sum_{i=1}^{n-1} \bar{X}_{(s-\tau_{(i-1)\sqrt{r}}^*) \wedge \bar{\tau}_{\sqrt{r}}^i}^{*,r,i} \mathbf{1}_{\{\tau_{(i-1)\sqrt{r}}^* \leq s\}} + \bar{X}_{s-\tau_{(n-1)\sqrt{r}}^*}^{*,r,n} \mathbf{1}_{\{s \geq \tau_{(n-1)\sqrt{r}}^*\}}.$$

So the path of $X^{*,r}$ is the concatenation of paths (51) for $1 \leq i \leq n - 1$ and the whole path $(\bar{X}_s^{*,r,n}, s \geq 0)$.

Again by Markov property, $X^{*,r}$ equals X^r in distribution. Note that $\tau_{i\sqrt{r}}^* = \inf\{s \geq 0: \inf_{u \leq s} X_u^{*,r} < -i\sqrt{r}\}$, which agrees with the usual notation. Moreover, on the event $\{\min_{1 \leq i \leq n} \tau_{\sqrt{r}}^{r,i} \geq rT/n\}$, we have $\tau_i = rT/n$ and $\bar{\tau}_{\sqrt{r}}^i \geq rT/n$. So on the same event, for each i , the path $(X_s^{r,i}, s \in [0, rT/n])$ is the initial part of the path $(\bar{X}_s^{*,r,i}, s \in [0, \bar{\tau}_{\sqrt{r}}^i])$, and therefore, $(X_s^{r,i}, s \in [0, rT/n]) = (X_{\tau_{(i-1)\sqrt{r}}^*+s}^{*,r} - X_{\tau_{(i-1)\sqrt{r}}^*}^{*,r}, s \in [0, rT/n])$ almost surely. Hence, on the event $\{\min_{1 \leq i \leq n} \tau_{\sqrt{r}}^{r,i} \geq rT/n\}$, the trees $\mathcal{T}_{i,j}^r, 1 \leq j \leq M_i^r$ (resp. $\mathcal{T}_i^{r,\dagger}$), $1 \leq i \leq n$, all reappear as trees (resp. initial parts of trees) corresponding to busy cycles of $(X_s^{*,r}, s \in [0, \tau_{n\sqrt{r}}^*])$.

Identity (39) together with bound (41) implies

$$\sup_{r \geq r_1} P\left(\max_{1 \leq i \leq n} \widehat{Z}_{t_i-t_{i-1}}^{r,i} > 4\varepsilon\right) \leq 4\eta.$$

On the event $\{\max_{1 \leq i \leq n} \widehat{Z}_{t_i-t_{i-1}}^{r,i} \leq 4\varepsilon\}$, the generation number of the last vertex visited by the queue-length process $Z^{r,i}$ on the interval $[0, rT/n]$ is smaller than or equal to $4\varepsilon\sqrt{r}$, for all $i \leq n$ simultaneously.

Now consider the intersection A^r of “good” events

$$\begin{aligned} A^r &= \left\{ \min_{1 \leq i \leq n} \tau_{\sqrt{r}}^{r,i} > rT/n \right\} \cap \left\{ \max_{1 \leq i \leq n} \widehat{Z}_{t_i-t_{i-1}}^{r,i} \leq 4\varepsilon \right\} \\ &\cap \left\{ \max_{1 \leq i \leq n} N_i^r \leq (\lambda + \varepsilon)rT/n \right\} \\ &\cap \left\{ \max_{1 \leq i \leq n} M_i^r \leq (\lambda + \varepsilon)\sqrt{r} \right\}. \end{aligned}$$

By previous considerations and Lemma 9, the probability of the complement of A^r is bounded from above by 8η for all r larger than $\max\{r_1, r_3\}$. The condition $\{\max_{1 \leq i \leq n} M_i^r \leq (\lambda + \varepsilon)\sqrt{r}\}$ will be used in later calculations, cf. (52).

Let $\mathcal{T}_1^{*,r}, \dots, \mathcal{T}_{M^*,r}^{*,r}$ be the sequence of Galton–Watson trees generated by $(X_s^{*,r}, s \in [0, \tau_{n\sqrt{r}}^*])$. Recall that $\mathcal{T}_s^{*,r}$ is the subtree spanned by vertex s and all of its descendants. Due to the above construction of $X^{*,r}$,

$$P\left(\left\{ \max_{1 \leq i \leq n} \text{ht}(\mathcal{T}_i^{r,\dagger}) > 5\varepsilon\sqrt{r} \right\} \cap A^r\right) \leq P(A_0^{*,r}),$$

where $A_0^{*,r} = \{\text{ht}(\mathcal{T}_s^{*,r}) > \varepsilon\sqrt{r} - 1 \text{ and } |\mathcal{T}_s^{*,r}| < r(\lambda + \varepsilon)T/n \text{ for some vertex } s^{*,r} \in \mathcal{T}_1^{*,r} \cup \dots \cup \mathcal{T}_{M^{*,r}}^{*,r}, \text{ gen}(s^{*,r}) = \lfloor 4\varepsilon\sqrt{r} \rfloor + 1\}$. The last statement is true by the “triangle inequality.” On the event $\{\max_{1 \leq i \leq n} \text{ht}(\mathcal{T}_i^{r,\dagger}) > 5\varepsilon\sqrt{r}\} \cap A^r$, the last vertex visited in each of the trees $\mathcal{T}_i^{r,\dagger}$ belongs to one of the first $\lfloor 4\varepsilon\sqrt{r} \rfloor$ generations, and the total size of $\mathcal{T}_i^{r,\dagger}$ is smaller than or equal to $(\lambda + \varepsilon)rT/n$, $1 \leq i \leq n$. However, there is at least one vertex in $\mathcal{T}_j^{r,\dagger}$, for some $1 \leq j \leq n$, with generation number larger than $\lceil 5\varepsilon\sqrt{r} \rceil$. The ancestor s^r of this vertex in generation $\lfloor 4\varepsilon\sqrt{r} \rfloor + 1$ must belong to the same tree $\mathcal{T}_j^{r,\dagger}$, due to the depth-first search order. Similarly, due to the depth-first search order, the whole tree $\mathcal{T}_{s^r}^{*,r}$ is contained in $\mathcal{T}_j^{r,\dagger}$ for this j . In the construction of $X^{*,r}$, vertex s^r and its tree of descendants $\mathcal{T}_{s^r}^{*,r}$ become $s^{*,r}$ and $\mathcal{T}_{s^{*,r}}^{*,r}$, where $\text{gen}(s^{*,r}) = \lfloor 4\varepsilon\sqrt{r} \rfloor + 1$, so that $A_0^{*,r}$ occurs.

Since $M^{*,r} \stackrel{d}{=} \text{Poisson}(\text{rate } \lambda^r n \sqrt{r})$ by (5), we have $\lim_r P(M^{*,r} \geq n(\lambda + \varepsilon)\sqrt{r}) = 0$, so one may assume r to be large enough so that $P(M^{*,r} \geq n(\lambda + \varepsilon)\sqrt{r}) \leq \eta$. The trees $\mathcal{T}_1^{*,r}, \dots, \mathcal{T}_{M^{*,r}}^{*,r}$ are, conditionally on $M^{*,r}$, independent and identically distributed as \mathcal{T}^r . So, given $M^{*,r} < n(\lambda + \varepsilon)\sqrt{r}$, the total expected number of vertices in generation $\lfloor 4\varepsilon\sqrt{r} \rfloor + 1$ is bounded from above by

$$n(\lambda + \varepsilon)\sqrt{r}(E\xi^r)^{4\varepsilon\sqrt{r}+1} \leq 4(1 + \eta)n(\lambda + \varepsilon)\sqrt{r}e^{-4c\varepsilon}, \quad r \rightarrow \infty,$$

due to assumption (10). Recall the large number K_1 from (48). If we denoted by $M^{*,r}(l) = \#\{s \in \mathcal{T}_1^{*,r} \cup \dots \cup \mathcal{T}_{M^{*,r}}^{*,r} : \text{gen}(s) = l\}$ the total size of generation l , Markov inequality implies

$$(52) \quad \begin{aligned} P\left(M^{*,r}(\lfloor 4\varepsilon\sqrt{r} \rfloor + 1) > K_1 n(\lambda + \varepsilon)\sqrt{r} \mid M^{*,r} \leq n(\lambda + \varepsilon)\sqrt{r}\right) \\ \leq \frac{(1 + \eta)e^{-4c\varepsilon}}{K_1} \leq \eta, \end{aligned}$$

for all large r . Each vertex s in generation $\lfloor 4\varepsilon\sqrt{r} \rfloor + 1$ has equal probability $P(\text{ht}(\mathcal{T}_s^{*,r}) > \varepsilon\sqrt{r}, |\mathcal{T}_s^{*,r}| < r(\lambda + \varepsilon)T/n) = P(\text{ht}(\mathcal{T}^r) > \varepsilon\sqrt{r}, |\mathcal{T}^r| < r(\lambda + \varepsilon)T/n)$ of contributing to event $A_0^{*,r}$. The above estimates put together with (47), (48) imply the existence of some large $r_5 \geq r_4$ such that

$$\begin{aligned} \sup_{r \geq r_5} P(A_0^{*,r}) &\leq 2\eta + \sup_{r \geq r_5} K_1 n(\lambda + \varepsilon)\sqrt{r} P(\text{ht}(\mathcal{T}^r) > \varepsilon\sqrt{r}, \\ &\quad |\mathcal{T}^r| < r(\lambda + \varepsilon)T/n) \\ &\leq 3\eta + K_1 \frac{\kappa_1(\lambda + \varepsilon)}{\sigma\sqrt{(\lambda + \varepsilon)T}} n^{3/2} \exp\left(-\frac{\varepsilon\sigma n^{1/2}}{\kappa_2\sqrt{(\lambda + \varepsilon)T}}\right) \leq 4\eta. \end{aligned}$$

Therefore,

$$\sup_{r \geq r_5} P\left(\max_{1 \leq i \leq n} \text{ht}(\mathcal{T}_i^{r,\dagger}) > 5\varepsilon\sqrt{r}\right) \leq \sup_{r \geq r_5} P(\text{not } A^r) + P(A_0^{*,r}) \leq 12\eta.$$

Now take r_5 to be r_1 in the statement of Proposition 7. \square

It remains to prove Lemma 8. The proof of (47) consists of adapting the corresponding arguments in [28], and then applying the reasoning of [2]. We sketch the proof, recalling the arguments of Kolchin along the way. Let $\xi_i, i \geq 1$, be i.i.d. integer-valued random variables with span 1 (that is, $P(\xi_1 = 1) > 0$), and with mean $E\xi_1 = a$ and $\text{var}(\xi_1) = \sigma^2$. Then the (standard) local central limit theorem (e.g., [11], Theorem 2.5.2 or [28], Theorem 1.4.2) gives

$$(53) \quad \sigma\sqrt{N}P(\xi_1 + \dots + \xi_N = m) - \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{(m - aN)^2}{2\sigma^2 N}\right\} \rightarrow 0,$$

uniformly in m , as $N \rightarrow \infty$. The identity in [28], Lemma 2.1.3, implies

$$(54) \quad P(|\mathcal{S}| = N) = \frac{1}{N}P(\xi_1 + \dots + \xi_N = N - 1), \quad N \geq 1,$$

and evaluating (53), (54) with $a = 1$ gives ([28], Lemma 2.1.4)

$$(55) \quad P(|\mathcal{S}| = N) \sim \frac{1}{\sqrt{2\pi}\sigma} N^{-3/2}.$$

It is important to note that the criticality assumption ($E\xi = 1$) gets used here only when applying (53), while the identity (54) holds for noncritical ξ 's as well.

Theorem 2.4.3 of [28] shows the convergence of heights conditioned on total size, which Aldous [2] recognizes in terms of the maximum W^* of the standard (unit length) Brownian excursion as

$$(56) \quad P\left(\text{ht}(\mathcal{S}) > \frac{x}{\sigma} N^{1/2} \mid |\mathcal{S}| = N\right) = P(2W^* > x)(1 + o_N(1)),$$

$$\lim_N o_N(1) = 0.$$

It suffices to show uniform (in r) analogues of (53), (56):

LEMMA 11. (i) As $N \rightarrow \infty$, uniformly in m ,

$$(57) \quad \sup_{r \geq 1} \left(\sigma^r \sqrt{N} P(\xi_1^r + \dots + \xi_N^r = m) - \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{(m - E\xi^r N)^2}{2N(\sigma^r)^2}\right\} \right) \rightarrow 0,$$

(ii)

$$(58) \quad P\left(\text{ht}(\mathcal{S}^r) > \frac{x}{\sigma^r} N^{1/2} \mid |\mathcal{S}^r| = N\right) = P(2W^* > x)(1 + o(r, N)),$$

where $\lim_{N \rightarrow \infty} \sup_{r \geq 1} |o(r, N)| = 0$.

Due to (57) and (54), we have

$$\begin{aligned}
 (59) \quad P(|\mathcal{T}^r| = \lfloor ur \rfloor) &= \frac{1}{\sqrt{2\pi\sigma}} (ur)^{-3/2} \exp(-c^2 u / (2\sigma^2)) (1 + o_r(u)) \\
 &\leq \frac{1}{\sqrt{2\pi\sigma}} (ur)^{-3/2} (1 + o_r(u)),
 \end{aligned}$$

where $\sup_{u \in [\varepsilon/\sqrt{r}, \delta]} |o_r(u)| \rightarrow 0$ as $r \rightarrow \infty$. The rest of the argument for (47) is identical to the one in [2], Proposition 24 for (46), using (59) and (58) in place of (55) and (56).

SKETCH OF THE PROOF OF LEMMA 11. Assertion (57) is proved in the same way as (53). Let $\phi^r(t) = E \exp(it\xi^r)$ be the characteristic function of ξ^r . One uses the inversion formula, splits the real-line into the same four regions and estimates the integrands, this time, uniformly over r . It is important that, by assumption (i) of Lemma 8, $\phi^r \rightarrow \phi$ uniformly (on R since all ϕ^r are 2π periodic), where ϕ is the characteristic function of ξ . We omit the details; the reader can easily check that the argument carries over, step by step.

To prove (58), it suffices to consider only critical trees. Suppose the moment generating function $G(z) = \sum_{i=0}^\infty P(\xi = i)z^i$, $z \in [0, 1]$, of the offspring distribution of \mathcal{T} satisfies

$$(60) \quad G(a) = aG'(a) < \infty,$$

for some real $a > 0$. It is easy to check (originally due to Kolchin) that the critical Galton–Watson tree \mathcal{T}_a with offspring distribution $P(\xi_a = i) = (a^i/G(a))P(\xi = i)$ and the original tree \mathcal{T} have the same distribution, when conditioned on their total size. In particular,

$$P(\text{ht}(\mathcal{T}) > xN^{1/2} \mid |\mathcal{T}| = N) = P(\text{ht}(\mathcal{T}_a) > xN^{1/2} \mid |\mathcal{T}_a| = N).$$

If the tree \mathcal{T} is critical, then $a = 1$. If the tree \mathcal{T} is supercritical, that is, $G'(1) > 1$, then the smallest fixed point $z_0 \in [0, 1]$ of G is strictly smaller than 1. By convexity of G , it must be $G'(z_0) < 1$, so

$$G'(1) > G(1) \quad \text{and} \quad z_0 G'(z_0) < G(z_0) = z_0.$$

By continuity and convexity, there exists a unique $a \in (z_0, 1)$ so that (60) holds. Using assumptions (ii) and (iv) of the lemma, it is easy to check that the sequence of smallest fixed points z_0^r of G^r satisfies $\lim_r z_0^r \rightarrow 1$; therefore, $\lim_r a^r = 1$. The sequence $\xi_{a^r}^r$ of integer random variables inherits all the properties in assumptions of the lemma. Thus, without loss of generality, we assume that all trees are critical, $E \xi^r = 1$, $r \geq 1$.

Due to our assumptions (i)–(iv), again the relevant estimates in [28] can be made uniform in r . For example, [28], Theorem 2.1.2 implies the analogue of (45),

$$(61) \quad P\left(\text{ht}(\mathcal{T}^r) > \frac{N}{\sigma^r}\right) = \frac{1}{\sqrt{2\pi N}} (1 + o_{r,N}(1)),$$

where $\lim_{N \rightarrow \infty} \sup_r |o_{r,N}(1)| = 0$, which is the first step in showing (58). The arguments in [28], Lemmas 2.4.3–2.4.5, Corollaries 2.4.1 and 2.4.2 and Theorems 2.4.1–2.4.3 extend in a similar way; we omit the details.

3.4. Discussion.

3.4.1. Initial condition. We first comment on the heavy-traffic limits under more general initial conditions. It is clear that convergence of initial load on diffusion scale $\widehat{X}^r(0) = X^r(0)/\sqrt{r} \rightarrow_d X_0$, where X_0 is a.s. a finite nonnegative random variable, would imply the convergence in distribution of the load processes \widehat{X}^r to a shifted Brownian motion started at X_0 . As before, this implies the convergence of the workload processes

$$(62) \quad \widehat{W}_t^r \Rightarrow W,$$

where $W_t = (X_t - \inf_{s \leq t} (X_s \wedge 0))$, $t \geq 0$. It is intuitively clear that, provided we have convergence of rescaled initial pure-atomic measures $q^r(0)$ (with atoms $\{1, 2, \dots, Z^r(0)\}$ and intensities equal to residual service times) to a fixed finite measure q_0 with finite support $[0, Z(0)] = [0, \lim_r \widehat{Z}^r(0)]$ and such that $\langle 1, q_0 \rangle = X_0$, then Theorems 1 and 5 should extend accordingly. We need to introduce some notation in order to identify these limits. As in [15], for any scalar $a \geq 0$ and any measure μ such that $\text{Supp}(\mu) \subset [0, \infty]$, let the *truncation* of μ at level a be the measure $\mu|_a$ defined by $\mu|_a[0, x] = \mu[0, x] \wedge a$. So, if $a \leq 0$, then $\mu|_a$ is the zero measure. Also if μ, ν are measures such that $\text{Supp}(\mu) \cup \text{Supp}(\nu) \subset [0, \infty]$, and $\sup(\text{Supp}(\mu)) = b < \infty$, define μ concatenated with ν as $(\mu \oplus \nu)([0, x]) = \mu([0, x \wedge b]) + \nu([0, (x - b)^+])$. Then it is easy to see that, at each level r , the RES-measure process $q^r(t) = q^r(0)|_{X^r(0)+I_t^{*,r}} \oplus q^{*,r}(t)$ encodes all the information, where $q^{*,r}(t)$ is a copy of the RES-measure process from Section 1.2 (started at zero measure) and $I_t^{*,r} = \inf_{u \in [0, t]} X_u^{r,*}$ is the corresponding infimum process. As $r \rightarrow \infty$, the rescaled $\hat{q}^r(t)$ should converge in the Skorokhod topology to $q_t = q_0|_{X_0+I_t^*} \oplus q_t^*$, where q^* is a copy of the generalized RES-measure process from Theorem 1.

The heavy-traffic limit for the queue length, on the other hand, depends on the finer properties of the asymptotic initial measure q_0 . By Theorems 1 and 5 and convergence (25), under certain regularity assumptions, one should get $\widehat{Z}^r(t) = \sup(\text{Supp}(\hat{q}_t^r)) \Rightarrow \sup(\text{Supp}(q_0|_{X_0+I_t^*})) + \sup(\text{Supp}(q_t^*))$. If $q_0(dx) = \frac{\lambda\beta}{2} dx$, $x \leq Z(0)$, this means

$$\widehat{Z}_t^r \Rightarrow Z_t = \frac{2}{\lambda\beta} (X(0) + I_t^*)^+ + \frac{2}{\lambda\beta} W_t^* = \frac{2}{\lambda\beta} W_t,$$

where W is the limit in (62).

3.4.2. LIFO vs. FIFO. Recall the optimization question from the Introduction. Assume a sequence of queues approaches heavy traffic (9), (10), (12), (13) and fix some large r . The two queues have the same workload process (11) W^r , which is approximated by W , a reflected Brownian motion (variance $\lambda\beta$ and drift $-c$). Denote by Z_{FI}^r and Z_{LI}^r the queue lengths under the FIFO and

the LIFO disciplines. For the FIFO queue, we have

$$(63) \quad W^r(t) = \sum_{i=1}^{Z_{\text{FI}}^r(t)} v_i^r + \varepsilon^r(t),$$

where $(v_i^r, i \geq 1)$ are i.i.d. with distribution F^r and $\varepsilon^r(t)$ is the residual service time of the customer currently in service. Due to the law of large numbers, $\widehat{Z}_{\text{FI}}^r \Rightarrow Z_{\text{FI}} = \lambda W$ in the limit; therefore,

$$Z_{\text{FI}}^r(\cdot) \approx \lambda W^r(\cdot),$$

while for the LIFO queue, (21) and Theorem 5 give $Z_{\text{LI}}^r(\cdot) \approx \frac{2}{\lambda\beta} W^r(\cdot)$, so in order to minimize the queue length in heavy traffic, the server should use the LIFO discipline iff $\lambda^2\beta > 2$ (equivalently, $\beta > 2m^2$ or $\sigma^2 > 2m$) and the FIFO discipline (alternatively, LIFO non-preemptive) otherwise. Note that in the special case, where both the arrival and the service times are exponential (rate λ^r), we can make the two queue lengths Z_{FI}^r and Z_{LI}^r coincide (as processes); therefore, their limits coincide, confirming $\lambda = \frac{2}{\lambda\beta}$. If all customers have constant service time m^r , $P(v = m^r) = 1$, then $\beta = m^2 < 2m^2$ and, of course, the FIFO discipline is optimal.

3.4.3. Random tree analogy. The argument in Lemma 2 uses an analogue of (63) for the LIFO case. It is not surprising that the mean residual service time depends on both the first and the second moment of the service time distribution F . Moreover, its exact value is in agreement with the analogous result in Aldous [2] about a “diffusion approximation” to a large Galton–Watson tree. To simplify the comparison, we assume that, for all large r , $\lambda^r = \lambda = 1$ and the service times have distribution $F^r = F$ with mean $m = 1$ and variance $\beta - m^2$. Then the busy cycles of the queue correspond to trees with critical offspring distribution $\xi^r = \xi$, $E\xi = 1$, $\text{var}(\xi) = \beta$. Denote by X^r the discrete-time depth-first search walk (from [2]) of a Galton–Watson tree \mathcal{T} with offspring distribution ξ , conditioned on $|\mathcal{T}| = r$. Theorem 23 in [2] states

$$(r^{-1/2} X^r([2rt]), t \in [0, 1]) \Rightarrow \left(\frac{2}{\sqrt{\beta}} W_t^*, t \in [0, 1] \right),$$

where W^* is standard Brownian excursion. The LIFO queue-length process Z^r is the depth-first search walk (continuous-time analogue) of an infinite sequence of critical Galton–Watson trees generated from queueing. Theorem 5 and (21) state

$$(r^{-1/2} Z^r(rt), t \geq 0) \Rightarrow \left(\frac{2}{\sqrt{\beta}} W^*(t), t \geq 0 \right),$$

where $W^*(t)$ is standard (mean 0, variance 1) Brownian motion.

3.4.4. *The heavy tails.* Recall the heavy tails setting from the Remark in Section 2. For X a stable- α process, $\alpha \in (1, 2)$, we can choose ([15], Proposition 4.3) the analogue of (21) as

$$(64) \quad \begin{aligned} Z_s^t &= \lim_{\varepsilon \rightarrow 0} \varepsilon^{\alpha-1} \#\{u \in (0, s], X_{u-} < I_u^t, \Delta X_u \geq \varepsilon\}, \\ s &\in [0, t] \quad \text{and} \quad Z_t = Z_t^t, \end{aligned}$$

where $\Delta X_u = X_u - X_{u-}$. Again Z_s^t and Z_t are continuous processes (cf. [15], Theorem 4.7). At each level r , define $\widehat{Z}_s^{t,r} = r^{1/\alpha-1} Z^{rt,r}(rs)$ and $\widehat{Z}_t^r = \widehat{Z}_t^{t,r}$. Then (26) is satisfied (cf. [15], Proposition 5.2) with Z_s^t in (64), so Theorem 1 extends in this case.

The finite-dimensional convergence of queue lengths (or heights) is a consequence of a more general result ([15], Proposition 5.2). The “tightness from below” for the queue length is again a consequence of (38) and (36). For the “tightness from above,” an analogue of Proposition 7 might be obtained using tree estimates analogous to those in Lemma 8. As remarked in Section 2, one can construct a triangular array of loads X^r converging (after scaling) to a general Lévy process X with Laplace exponent (16), (17). Duquesne and Le Gall (personal communication) consider this setting, where $\sigma = 0$ in (16), and obtain an analogue of Theorem 5 under suitable assumptions.

4. Directions for further research. Taking the FIFO queueing discipline as a paradigm, we list several natural ways to generalize the result of this paper. The full name of our queue, *feed-forward, single class, single server M/G/1 LIFO preemptive resume queue*, gives a list of assumptions that might be relaxed. Allowing renewal (non-Markovian) arrivals would be valuable extensions for applications. We consider the above setting in the forthcoming paper [29]. It turns out that LIFO preemptive resume service discipline induces an unconventional heavy-traffic behavior, in that the limit for the queue length depends on the type of arrivals (and services) in an intricate way.

Introducing feedback, or more customer classes, to the system (where the classes differ by their interarrival and/or service time distributions) or considering networks of LIFO queues, complicates the global arrival process and might result in additional “surprises” in heavy traffic.

Acknowledgment. I am grateful to Ruth Williams for suggesting this project, and for her numerous comments that improved the flow of the paper. Many thanks to the anonymous referee for detailed comments and for finding a flaw in the previous proof of Proposition 7.

REFERENCES

- [1] ABATE, J. and WHITT, W. (1997). Limits and approximations for the $M/G/1$ LIFO waiting-time distribution. *Oper. Res. Lett.* **20** 199–206.
- [2] ALDOUS, D. J. (1993). The continuum random tree III. *Ann. Probab.* **21** 248–289.

- [3] ATHREYA, K. B. and NEY, P. E. (1972). *Branching Processes*. Springer, New York.
- [4] BASKETT, F., CHANDY, K. M., MUNTZ, R. R. and PALACIOS, F. G. (1975). Open, closed and mixed networks of queues with different classes of customers. *J. ACM* **22** 248–260.
- [5] BERTOIN, J. (1996). *Lévy Processes*. Cambridge Univ. Press.
- [6] BILLINGSLEY, P. (1968). *Convergence of Probability Measures*. Wiley, New York.
- [7] BRAMSON, M. (1998). State space collapse with application to heavy traffic limits for multi-class queueing networks. *Queueing Systems: Theory Appl.* **30** 89–148.
- [8] BREIMAN, L. (1968). *Probability*. Addison-Wesley, Reading, MA.
- [9] COFFMAN, E. G. and MITRANI, I. (1988). Storage of the single-server queue. In *Queueing theory and Its Applications*. 193–205. North-Holland, Amsterdam.
- [10] DAWSON, D. (1993). *Measure-Valued Markov Processes. École d'Été de Probabilités de Saint-Flour XXI. Lecture Notes in Math.* **1541**. Springer, Berlin.
- [11] DURRETT, R. (1991). *Probability. Theory and Examples*. Wadsworth and Brooks/Cole, Pacific Grove, CA.
- [12] DYNKIN, E. B. (1994). *An Introduction to Branching Measure-Valued Processes*. Amer. Math. Soc., Providence, RI.
- [13] ETHIER, S. N. and KURTZ, T. G. (1986). *Markov Processes: Characterization and Convergence*. Wiley, New York.
- [14] LE GALL, J. F. and LE JAN, Y. (1998). Branching processes in Lévy processes: Laplace functionals of snakes and superprocesses. *Ann. Probab.* **26** 1407–1432.
- [15] LE GALL, J. F. and LE JAN, Y. (1998). Branching processes in Lévy processes: the exploration process. *Ann. Probab.* **26** 213–252.
- [16] GEIGER, J. and KERSTING, G. (1997). Depth-first search of random trees, and Poisson point processes. *IMA Vol. Math. Appl.* **84** 111–126.
- [17] GRIMVALL, A. (1974). On the convergence of sequences of branching processes. *Ann. Probab.* **2** 1027–1045.
- [18] HARRIS, T. E. (1974). First passage and recurrence distributions. *Trans. Amer. Math. Soc.* **73** 471–486.
- [19] HARRISON, J. M. (1978). The diffusion approximation for tandem queues in heavy traffic. *Adv. in Appl. Probab.* **10** 886–905.
- [20] HEATH, D., RESNICK, S. and SAMORODNITSKY, G. (1998). Heavy tails and long range dependence in on/off processes and associated fluid models. *Math. Oper. Res.* **23** 145–165.
- [21] IGLEHART, D. L. and WHITT, W. (1970). Multiple channel queues in heavy traffic I. *Adv. in Appl. Probab.* **2** 150–177.
- [22] IGLEHART, D. L. and WHITT, W. (1970). Multiple channel queues in heavy traffic II. *Adv. in Appl. Probab.* **2** 355–364.
- [23] JACOD, J. and SHIRYAEV, A. N. (1987). *Limit Theorems for Stochastic Processes*. Springer, Berlin.
- [24] JAGERMAN, D. L., MELAMED, B. and WILLINGER, W. (1997). Stochastic modeling of traffic processes. In *Frontiers in Queueing*. 271–320. CRC Press, Boca Raton, FL.
- [25] KELLY, F. P. (1975). Networks of queues with customers of different types. *J. Appl. Probab.* **12** 542–554.
- [26] KENDALL, D. G. (1951). Some problems in the theory of queues. *J. Roy. Statist. Soc. B* **13** 151–185.
- [27] KINGMAN, J. F. C. (1961). The single server queue in heavy traffic. *Proc. Cambridge Philos. Soc.* **57** 902–904.
- [28] KOLCHIN, V. F. (1986). *Random Mappings*. Optimization Software, New York. [Trans. of Russian original.]
- [29] LIMIC, V. (2000). On the behavior of LIFO preemptive resume queues in heavy traffic. *Electron. Comm. Probab.* **5** 13–27.
- [30] LINDVALL, T. (1974). Limit theorems for some functionals of certain Galton–Watson branching processes. *Adv. in Appl. Probab.* **6** 309–321.
- [31] REIMAN, M. I. (1984). Open queueing networks in heavy traffic. *Math. Oper. Res.* **9** 441–458.
- [32] SHALMON, M. (1988). Analysis of the GI/GI/1 queue and its variations via the LCFS preemptive resume discipline and its random walk interpretation. *Probab. Engng. Inform. Sci.* **2** 215–230.

- [33] SHANTHIKUMAR, J. G. and SUMITA, U. (1986). On $G/G/1$ queues with LIFO-P service discipline. *J. Oper. Res. Soc. Japan* **29** 220–231.
- [34] SIGMAN, K. (1996). Queues under preemptive LIFO and ladder height distributions for risk processes: a duality. *Comm. Statist. Stochastic Models* **12** 725–735.
- [35] WHITT, W. (1971). Weak convergence theorems for priority queues: preemptive-resume discipline. *J. Appl. Probab.* **8** 74–94.
- [36] WILLIAMS, R. J. (1996). On the approximation of queueing networks in heavy traffic. In *Stochastic Networks: Theory and Applications* (S. Zachary, F. P. Kelly and I. Ziedins, eds.) 35–56. Clarendon Press, Oxford.
- [37] WILLIAMS, R. J. (1998). Diffusion approximations for open multiclass queueing networks: sufficient conditions involving state space collapse. *Queueing Systems Theory Appl.* **30** 27–88.

DEPARTMENT OF MATHEMATICS
CORNELL UNIVERSITY
ITHACA, NEW YORK 14853–4201
E-MAIL: limic@math.cornell.edu