

## DIFFUSION LIMITS FOR SHORTEST REMAINING PROCESSING TIME QUEUES UNDER NONSTANDARD SPATIAL SCALING

BY AMBER L. PUHA<sup>1</sup>

*California State University San Marcos*

We develop a heavy traffic diffusion limit theorem under nonstandard spatial scaling for the queue length process in a single server queue employing shortest remaining processing time (SRPT). For processing time distributions with unbounded support, it has been shown that standard diffusion scaling yields an identically zero limit. We specify an alternative spatial scaling that produces a nonzero limit. Our model allows for renewal arrivals and i.i.d. processing times satisfying a rapid variation condition. We add a corrective spatial scale factor to standard diffusion scaling, and specify conditions under which the sequence of unconventionally scaled queue length processes converges in distribution to the same nonzero reflected Brownian motion to which the sequence of conventionally scaled workload processes converges. Consequently, this corrective spatial scale factor characterizes the order of magnitude difference between the queue length and workload processes of SRPT queues in heavy traffic. It is determined by the processing time distribution such that the rate at which it tends to infinity depends on the rate at which the tail of the processing time distribution tends to zero. For Weibull processing time distributions, we restate this result in a manner that makes the resulting state space collapse more apparent.

**1. Introduction.** We study the heavy traffic behavior of the queue length process in a shortest remaining processing time (SRPT) queue. We consider a single server queue with renewal arrivals and independent and identically distributed processing times that are also independent of the arrival process. Jobs are served in a nonidling fashion such that at each instant the job with the shortest remaining processing time is served at rate one. This is done with preemption so that if a job arrives for which the total processing time is smaller than that remaining of the job in service, the job in service is placed on hold and the arriving job enters service. Therefore, in order to adequately track the state of the system, it is necessary to keep track of all remaining processing times of all jobs in the system. We do this using a measure valued process that at each time has a unit atom at the remaining processing time of each job in the system. This is introduced formally in Section 2.

---

Received January 2014; revised July 2014.

<sup>1</sup>Supported in part by ViaSat Inc. and an ROA supplement to NSF Grant DMS-12-06772.

*MSC2010 subject classifications.* Primary 60K25, 60F17; secondary 60G57, 68M20, 90B22.

*Key words and phrases.* Heavy traffic, queueing, shortest remaining processing time, diffusion limit, nonstandard scaling, rapidly varying processing times.

Optimality of shortest remaining processing time, in the sense that it is the queue length minimizer over all nonidling service disciplines, has been known since the 1960s [21, 24]. One anticipates that this results at the expense of lengthy delays for jobs with large total processing times. Hence, sojourn times are naturally of interest for SRPT. For Markovian arrivals, the early work Schrage and Miller [22] develops a formula for the mean response time in steady state with extended results available in Schassberger [20] and Perera [18]; see Schreiber [23] for a survey. Also see results in Pavlov [16] and Pechinkin [17] on steady state queue length distributions. Recently, Lin, Wierman and Zwart [13] followed up on the work in [22] by characterizing the asymptotic behavior of the steady state mean sojourn time as the traffic intensity approaches one for a large class of processing time distributions. Interestingly, the rate at which the mean sojourn time tends to infinity depends on the tail behavior of the processing time distribution. Results in this spirit were also an outcome of [5, 6], where a fluid model (formal functional law of large numbers limit) was proposed and an associate weak convergence result (functional law of large numbers result) was stated and proved. There the rate at which a fluid analog of the sojourn time as a function of the initial processing time tends to infinity depends on the tail behavior of the processing distribution. Other somewhat recent studies of SRPT have focused on fairness (e.g., [1, 26]) or tail behavior [14, 15].

Here we focus on further developing existing diffusion limit results (functional central limit theorems) for SRPT. The paper [9] contains a diffusion limit theorem for the sequence of measure valued state descriptors under standard heavy traffic conditions and standard diffusion scaling. For limiting processing time distributions with bounded support, the main result in [9] indicates that the limiting measure valued process is a single atom supported at the supremum of the support of the limiting processing time distribution for all time. The height of that atom varies randomly in time as determined by the limiting workload process. More specifically, under standard heavy traffic conditions, the sequence of conventionally diffusion scaled workload processes converges in distribution to a semimartingale reflected Brownian motion [11]. The height of the atom for the measure valued diffusion limit is then given by the limiting workload process divided by the supremum of the support of the limiting processing time distribution. This is analogous to early results for strict priority queues where in the heavy traffic diffusion limit work piles up in the lowest priority class [25]. The main result in [9] goes on to state that for limiting processing time distributions with unbounded support, the limiting measure valued process is identically equal to the zero measure. In particular, the limiting queue length process is identically equal to zero. Such behavior had not been observed prior to this for other nonidling service disciplines. The limiting queue length process is typically recovered from the limiting workload process via multiplication by a positive constant, a phenomenon known as state space collapse; see [8, 11, 12, 25], for instance. The fact that the sequence of rescaled queue length processes is of lower order magnitude than the sequence of rescaled workload processes quantifies the extreme queue length minimizing nature of SRPT.

A natural follow-up question to the work in [9] is whether or not there is an alternative scaling that can be employed to yield a nontrivial limit for this unconventionally rescaled queue length process. If such a limit exists, it would be of interest to describe how that limit is related to the limiting workload process that arises under standard diffusion scaling. Here we identify such a nonstandard scaling for continuous processing time distributions with unbounded support for which the tails satisfy a rapid variation condition. The main theorem in this paper, Theorem 3.1, specifies that the spatial scaling must be modified by multiplying by a certain inverse function related to the tails of the first moment of the processing time distribution. In particular, there is a multiplicative correction factor that must be applied to standard diffusion scaling to obtain a nontrivial limit. The order of magnitude of that correction factor depends on the rate at which the tails of the first moment tend to zero. With this corrective scaling, the limiting process is identically equal to the limiting workload process that arises under standard diffusion scaling. Hence, with this corrective scaling factor, a generalized version of state space collapse holds.

The corrective scaling identified here was inspired by fluid limit results in [6]. The order of magnitude agrees with that of the left edge of the support of fluid model solutions as time approaches infinity. This seems to be the first result in the queuing theory literature where the nature of the scaling depends on the tail behavior of the processing time distribution. In fact, we only know of one previous result [10] that employs nonstandard scaling. The scaling in [10] is a mixture of conventional fluid (functional law of large numbers) and conventional diffusion (functional central limit theorem) scaling.

It is interesting to note that the result in Theorem 3.1 is consistent with the rapid variation case of [13], Theorem 3, as follows. Theorem 3 in [13] specifies an asymptotic formula for the mean sojourn time in steady state as the traffic intensity increases to one. By using the rate at which the traffic intensity approaches one for standard heavy traffic conditions [see (3.2)], one can informally translate their asymptotic formula into one indexed by the sequence of systems here. This results in an asymptotic formula that has the same order of magnitude as the spatial scaling specified by Theorem 3.1. Of course the former is for the steady-state mean response time, and the latter is for the unconventionally rescaled queue length process. But, due to Little's law, the queue length and response time should be of the same order of magnitude.

In general, the inverse function that produces the corrective scaling is not available in closed form. Hence, the multiplicative constant contained within it is not immediately available. However, for Weibull processing time distributions, explicit calculations can be done to separate the order of magnitude and multiplicative constant. The order of magnitude is determined by the shape parameter and the multiplicative constant is given by the scale parameter. This is stated precisely in Corollary 3.2, which provides an interesting illustration of the resulting generalized state space collapse.

This raises the next natural question of what happens when the processing time distributions satisfy a regular variation condition. The work here does not address that case. The works [6] and [13] suggest that the same function might provide an appropriate corrective scaling. However, the proof of Theorem 3.1 does not generalize to that case. The slowly varying nature of the inverse function that specifies the corrective scaling factor plays an important role in the proof of Theorem 3.1; see (2.2) and (2.3). Determining the behavior in the case of regular variation is work in progress.

In the next section, we precisely define the model and associated measure valued state descriptor. Then we specify the sequence of systems and associated asymptotic conditions that they must satisfy. This allows us to state the main result of the paper, Theorem 3.1, and its corollary for Weibull processing time distributions, Corollary 3.2. The remainder of the paper contains the proof of the main result.

1.1. *Notation.* Throughout  $\mathbb{R}$  denotes the real numbers, and  $\mathbb{R}_+$  denotes the nonnegative real numbers. Similarly,  $\mathbb{Z}$  denotes the integers, and  $\mathbb{Z}_+$  denotes the nonnegative integers. Then  $\mathbb{N}$  denotes the positive integers. For  $a, b \in \mathbb{R}$ ,  $a \wedge b$  and  $a \vee b$ , respectively, denote the minimum and maximum of  $a$  and  $b$ . Also, for  $a \in \mathbb{R}$ ,  $|a| = (-a) \vee a$  denotes the absolute value of  $a$ .

We define  $\mathbf{C}(\mathbb{R}_+)$  to be the set of continuous real valued functions with domain  $\mathbb{R}_+$ . Then  $\mathbf{C}_b(\mathbb{R}_+)$  denotes those elements of  $\mathbf{C}(\mathbb{R}_+)$  that are bounded. We use the notation  $\mathbf{1}(\cdot)$  for the function in  $\mathbf{C}_b(\mathbb{R}_+)$  that is identically equal to one and  $\chi(\cdot)$  for the identity function in  $\mathbf{C}(\mathbb{R}_+)$ .

For a Polish space  $\mathcal{S}$ , we let  $\mathbf{D}([0, \infty), \mathcal{S})$  denote the set of functions of time taking values in  $\mathcal{S}$  that are right continuous with finite left limits. We endow this space with the Skorohod  $J_1$ -topology. Then  $\mathbf{D}([0, \infty), \mathcal{S})$  is also a Polish space [7]. We denote the function in  $\mathbf{D}([0, \infty), \mathbb{R})$  that is identically equal to zero by  $\mathbf{0}(\cdot)$ .

We use the notation  $\mathcal{M}$  for the set of finite, nonnegative Borel measures on  $\mathbb{R}_+$ . The zero measure in  $\mathcal{M}$  is denoted by  $\mathbf{0}$ . Traditionally, for  $x \in \mathbb{R}_+$ ,  $\delta_x \in \mathcal{M}$  is the unit atom at  $x$ . For  $x \in \mathbb{R}_+$ , we also define  $\delta_x^+$  to be the measure in  $\mathcal{M}$  that is  $\delta_x$  if  $x > 0$  and  $\mathbf{0}$  otherwise. Given a Borel measurable function  $f: \mathbb{R}_+ \rightarrow \mathbb{R}$  and  $\zeta \in \mathcal{M}$ , we let  $\langle f, \zeta \rangle = \int_{\mathbb{R}_+} f(x)\zeta(dx)$ , when the integral exists. Then  $\langle \mathbf{1}, \zeta \rangle$  is the total mass of  $\zeta$ . We refer to  $\langle \chi, \zeta \rangle$  as the first moment of  $\zeta$ . The set  $\mathcal{M}$  is endowed with the topology of weak convergence. In particular, for  $\{\zeta_n\}_{n \in \mathbb{N}} \subset \mathcal{M}$  and  $\zeta \in \mathcal{M}$ ,  $\zeta_n \xrightarrow{w} \zeta$  as  $n \rightarrow \infty$  if and only if  $\lim_{n \rightarrow \infty} \langle g, \zeta_n \rangle = \langle g, \zeta \rangle$  for all  $g \in \mathbf{C}_b(\mathbb{R}_+)$ . With this topology,  $\mathcal{M}$  is a Polish space. We denote the function in  $\mathbf{D}([0, \infty), \mathcal{M})$  that is identically equal to the zero measure by  $\mathbf{0}(\cdot)$ .

We use “ $\Rightarrow$ ” to denote convergence in distribution of random elements of a metric space. Following Billingsley [2], we use  $\mathbb{P}$  and  $\mathbb{E}$ , respectively, to denote the probability measure and expectation operator associated with whatever space the relevant random element is defined on. Unless otherwise specified, all stochastic processes used in this paper are assumed to have paths that are right continuous with finite left limits (r.c.l.l.).

Finally, following [3], we say that a measurable function  $g : (0, \infty) \rightarrow (0, \infty)$  is rapidly varying of index  $\infty$  if for all  $\varepsilon > 0$ ,

$$(1.1) \quad \lim_{x \rightarrow \infty} \frac{g((1 + \varepsilon)x)}{g(x)} = \infty,$$

and is rapidly varying of index  $-\infty$  if for all  $\varepsilon > 0$ ,

$$(1.2) \quad \lim_{x \rightarrow \infty} \frac{g((1 + \varepsilon)x)}{g(x)} = 0.$$

Together the functions in these two classes are called *rapidly varying*. In addition,  $g : (0, \infty) \rightarrow (0, \infty)$  is slowly varying if for all  $c > 0$ ,

$$(1.3) \quad \lim_{x \rightarrow \infty} \frac{g(cx)}{g(x)} = 1.$$

**2. The stochastic model and state descriptor.** We consider a  $GI/GI/1$  SRPT queue such that the processing time distribution is continuous, has unbounded support and the tails satisfy a rapid variation condition. In particular, jobs arrive according to a delayed renewal process  $E(\cdot)$  with rate  $\lambda \in (0, \infty)$  such that the interarrival times have finite standard deviation  $\sigma_a$  and  $E(0) = 0$ . Then, for  $t \in [0, \infty)$ ,  $E(t)$  denotes the number of jobs that have arrived to the system exogenously by time  $t$ . Processing times for these jobs are independent and identically distributed positive random variables with common continuous cumulative distribution function  $F(\cdot)$ , finite mean and finite standard deviation  $\sigma_s \in (0, \infty)$ . The sequence  $\{v_i\}_{i \in \mathbb{N}}$  of processing times is also assumed to be independent of the arrival process. For  $i \in \mathbb{N}$ , the  $i$ th job to arrive to the system has total processing time  $v_i$ . For simplicity, we refer to the  $i$ th job to arrive to the system as job  $i$ , or the  $i$ th job. We use the notation  $v$  to denote a random variable that is equal in distribution to a generic processing time. Specifically,

$$\bar{F}(x) = 1 - F(x) = \mathbb{P}(v > x), \quad x \in \mathbb{R}_+.$$

Our assumptions include that  $\bar{F}(\cdot)$  is rapidly varying with index minus infinity; see (1.2). We restrict attention to a subset of such processing time distributions that includes, for example, Weibull distributions. For this, given  $x \in \mathbb{R}_+$ , let

$$S(x) = \frac{1}{\mathbb{E}[v 1_{\{v > x\}}]}.$$

In [6],  $s(\cdot)$  is the fluid analog of the sojourn time of initial jobs as a function of the remaining processing time at time zero. The notation  $S(\cdot)$  is chosen here to highlight its similarity with  $s(\cdot)$  in [6] (they differ by factor that tends to a positive constant as  $x$  tends to infinity). Note that  $S(0) = 1/\mathbb{E}[v]$ . Further,  $S(\cdot)$  is positive, nondecreasing, continuous, unbounded and rapidly varying with index plus infinity. Set

$$(2.1) \quad S^{-1}(y) = \inf\{x \in \mathbb{R}_+ : S(x) > y\}, \quad y \in \mathbb{R}_+.$$

Hence,  $S^{-1}(\cdot)$  is positive, strictly increasing, right continuous, unbounded and slowly varying. Further,  $S(S^{-1}(y)) = y$  for all  $y \in \mathbb{R}_+$ . We assume that for some  $c > 1$ ,

$$(2.2) \quad \lim_{y \rightarrow \infty} \left( \frac{S^{-1}(cy)}{S^{-1}(y)} - 1 \right) \ln(S^{-1}(y)) = 0.$$

This is an assumption about the rate at which the ratio associated with the slowly varying function  $S^{-1}(\cdot)$  converges to one. Rate of convergence conditions such as this and their implications are discussed more fully in [3], Section 2.3.1. Here we note that (2.2) is not satisfied by all slowly varying functions. For example, as noted in [3], page 78, (2.2) does not hold for slowly varying functions of the form  $\exp((\ln(\cdot))^\delta)$ , where  $1/2 \leq \delta < 1$ . However, it does hold for many processing distributions, including Weibull processing time distributions. That Weibull processing time distributions satisfy (2.2) is demonstrated in Section 3.

The reason for assuming (2.2) is that by [3], Theorem 2.3.3 (originally stated in [4]), it follows that for all  $\delta \in \mathbb{R}$ ,

$$(2.3) \quad \lim_{y \rightarrow \infty} \frac{S^{-1}((S^{-1}(y))^\delta y)}{S^{-1}(y)} = 1.$$

Recall that  $S^{-1}(\cdot)$  is slowly varying so that  $\lim_{y \rightarrow \infty} S^{-1}(cy)/S^{-1}(y) = 1$  for all  $c > 0$ . Then (2.3) says that one can replace the constant  $c > 0$  with  $(S^{-1}(y))^\delta$ ,  $y \in \mathbb{R}_+$ . As  $y$  tends to infinity, this tends to infinity if  $\delta > 0$  and to zero if  $\delta < 0$ . In Section 3, (2.3) is used to obtain (3.8), which is in turn used to prove the main theorem of the paper, Theorem 3.1.

As far as the initial state of the system is concerned, there are  $Q(0)$  jobs in the system at time zero. Here  $Q(0)$  is assumed to be a random variable taking values in  $\mathbb{Z}_+$ . The time zero remaining processing times for such jobs are the first  $Q(0)$  elements of the sequence  $\{\tilde{v}_i\}_{i \in \mathbb{N}} \subset \mathbb{R}_+$ . Each member of the sequence  $\{\tilde{v}_i\}_{i \in \mathbb{N}}$  is assumed to be a positive random variable. For  $1 \leq i \leq Q(0)$ , we refer to the job in the system at time zero with remaining processing time  $\tilde{v}_i$  at time zero as initial job  $i$ , or the  $i$ th initial job. Let  $W(0) = \sum_{i=1}^{Q(0)} \tilde{v}_i$ , which is a random variable taking values in  $\mathbb{R}_+$ . Then  $W(0)$  corresponds to the total work (measured in units of processing time) in the system at time zero. Finally, let  $\mathcal{Z}(0) \in \mathcal{M}$  be given by

$$\mathcal{Z}(0) = \sum_{i=1}^{Q(0)} \delta_{\tilde{v}_i}^+.$$

Note that

$$Q(0) = \langle 1, \mathcal{Z}(0) \rangle \quad \text{and} \quad W(0) = \langle \chi, \mathcal{Z}(0) \rangle.$$

Jobs are served in a nonidling fashion. In particular, the server does not idle if there are jobs in the system. At any given instance at which the system is nonempty,

the job with the shortest remaining processing time is served at rate one. This is done with preemption so that when a job arrives to the system that requires less processing time than that remaining for the job currently in service, the job in service is placed on hold and the arriving job enters service immediately. For  $1 \leq i \leq Q(0)$  and  $t \in [0, \infty)$ ,  $\tilde{v}_i(t)$  denotes the remaining processing time of initial job  $i$  at time  $t$ . For  $1 \leq i \leq E(t)$  and  $t \in [0, \infty)$ ,  $v_i(t)$  denotes the remaining processing time of job  $i$  at time  $t$ . So then, for  $t \in [0, \infty)$ , let

$$\mathcal{Z}(t) = \sum_{i=1}^{Q(0)} \delta_{\tilde{v}_i(t)}^+ + \sum_{i=1}^{E(t)} \delta_{v_i(t)}^+.$$

In particular,  $\mathcal{Z}(\cdot) \in \mathbf{D}([0, \infty), \mathcal{M})$  is the associated measure valued state descriptor. For  $t \in [0, \infty)$ , let

$$Q(t) = \langle 1, \mathcal{Z}(t) \rangle \quad \text{and} \quad W(t) = \langle \chi, \mathcal{Z}(t) \rangle.$$

Then  $Q(\cdot)$  and  $W(\cdot)$ , respectively, denote the queue length and workload processes.

**3. Statement of the main result.** Let  $\mathcal{R}$  be a sequence taking values in  $(1, \infty)$  tending to infinity. Fix a sequence of  $GI/GI/1$  SRPT queues indexed by  $\mathcal{R}$  for which the initial conditions and stochastic primitive inputs satisfy the conditions specified in Section 2. We further require that the processing time distributions do not depend on  $r$  and have common cumulative distribution function  $F(\cdot)$ . We place a superscript  $r$  on all parameters and processes associated with the  $r$ th system. So then for each  $r \in \mathcal{R}$ , we have  $\lambda^r, \sigma_a^r, E^r(\cdot), \mathcal{Z}^r(\cdot), Q^r(\cdot)$  and  $W^r(\cdot)$ , which may depend on  $r$ , but  $F^r(\cdot) = F(\cdot)$  for all  $r \in \mathcal{R}$ . Also, for  $r \in \mathcal{R}$ , set

$$\rho^r = \lambda^r \mathbb{E}[v].$$

For convenience later on, for  $r \in \mathcal{R}$  and  $x \in \mathbb{R}_+$ , we also define

$$\rho_x^r = \lambda^r \mathbb{E}[v 1_{\{v \leq x\}}].$$

Then, for  $r \in \mathcal{R}$  and  $x \in \mathbb{R}_+$ ,

$$(3.1) \quad \rho^r - \rho_x^r = \frac{\lambda^r}{S(x)}.$$

We assume that the stochastic primitive inputs satisfy the following asymptotic heavy traffic conditions. For some  $\kappa \in \mathbb{R}$ , as  $r \rightarrow \infty$ ,

$$(3.2) \quad \sigma_a^r \rightarrow \sigma_a \quad \text{and} \quad r(\rho^r - 1) \rightarrow \kappa.$$

Then it follows that  $\lambda^r \rightarrow \lambda$  as  $r \rightarrow \infty$ , where  $\lambda = 1/\mathbb{E}[v]$ . For  $r \in \mathcal{R}$  and  $t \in [0, \infty)$ , let

$$\bar{E}^r(t) = \frac{E^r(r^2 t)}{r^2} \quad \text{and} \quad \hat{E}^r(t) = \frac{E^r(r^2 t) - \lambda^r r^2 t}{r}.$$

Also assume that as  $r \rightarrow \infty$ ,

$$(3.3) \quad \widehat{E}^r(\cdot) \Rightarrow E^*(\cdot),$$

where  $E^*(\cdot)$  is a Brownian motion starting from zero with drift zero and variance  $(\lambda)^3(\sigma_a)^2$ . This implies a functional weak law of large numbers for the exogenous arrival process. Specifically, set  $\lambda(t) = \lambda t$  for  $t \in [0, \infty)$ . Then, as  $r \rightarrow \infty$ ,

$$(3.4) \quad \overline{E}^r(\cdot) \Rightarrow \lambda(\cdot).$$

Given  $r \in \mathcal{R}$ , let

$$(3.5) \quad c^r = S^{-1}(r).$$

For  $r \in \mathcal{R}$  and  $t \in [0, \infty)$ , set

$$\tilde{Q}^r(t) = \frac{c^r Q^r(r^2 t)}{r} \quad \text{and} \quad \tilde{Z}^r(t) = \frac{c^r Z^r(r^2 t)}{r}.$$

Also, for  $r \in \mathcal{R}$  and  $t \in [0, \infty)$ , set

$$\widehat{W}^r(t) = \frac{W^r(r^2 t)}{r} \quad \text{and} \quad \widehat{Z}^r(t) = \frac{Z^r(r^2 t)}{r}.$$

Then, the “hat” notation corresponds to processes under standard diffusion scaling and the “tilde” notation corresponds to processes under the nonstandard scaling consisting of standard diffusion scaling multiplied by the spatial correction factor  $c^r$ ,  $r \in \mathcal{R}$ . Note that  $\lim_{r \rightarrow \infty} c^r = \infty$ . Assume that for some random variable  $W_0$  that is finite almost surely, as  $r \rightarrow \infty$ ,

$$(3.6) \quad (\widehat{W}^r(0), \tilde{Q}^r(0)) \Rightarrow (W_0, W_0).$$

For  $r \in \mathcal{R}$  and  $\varepsilon > 0$ , let

$$(3.7) \quad l_\varepsilon^r = S^{-1}(r(c^r)^{-2-\varepsilon}) \quad \text{and} \quad u_\varepsilon^r = S^{-1}(r(c^r)^{2+\varepsilon}).$$

Then, for  $\varepsilon > 0$  and  $r \in \mathcal{R}$ , we have that  $0 < l_\varepsilon^r < c^r < u_\varepsilon^r < \infty$ . Also, for all  $\varepsilon > 0$ ,  $\lim_{r \rightarrow \infty} l_\varepsilon^r = \lim_{r \rightarrow \infty} u_\varepsilon^r = \infty$ . Further, by (2.3), (3.5) and (3.7), for each  $\varepsilon > 0$ ,

$$(3.8) \quad \lim_{r \rightarrow \infty} \frac{c^r}{l_\varepsilon^r} = 1 \quad \text{and} \quad \lim_{r \rightarrow \infty} \frac{c^r}{u_\varepsilon^r} = 1.$$

The proof of the main result (Theorem 3.1) will proceed by demonstrating for any given  $\varepsilon > 0$ , the contribution to the total mass under the unconventional scaling and to the work under the conventional scaling asymptotically concentrates in  $(l_\varepsilon^r, u_\varepsilon^r]$  as  $r \rightarrow \infty$ . Therefore, we further assume that for all  $\varepsilon > 0$ , as  $r \rightarrow \infty$ ,

$$(3.9) \quad \langle (1 \vee \chi) 1_{[0, l_\varepsilon^r]}, \tilde{Z}^r(0) \rangle \Rightarrow 0 \quad \text{and} \quad \langle \chi 1_{(u_\varepsilon^r, \infty)}, \widehat{Z}^r(0) \rangle \Rightarrow 0.$$



THEOREM 3.1. *Assume that (3.2), (3.3), (3.6) and (3.9) hold. As  $r \rightarrow \infty$ ,*

$$(\tilde{Q}^r(\cdot), \widehat{W}^r(\cdot)) \Rightarrow (W^*(\cdot), W^*(\cdot)),$$

where  $W^*(\cdot)$  is a reflected Brownian motion with drift  $\kappa$  and variance  $\lambda((\sigma_a)^2 + (\sigma_s)^2)$  such that  $W^*(0)$  is equal in distribution to  $W_0$ .

For the class of processing time distributions that satisfy the rapid variation condition (2.2), Theorem 3.1 implies that the asymptotic order of magnitude difference between the  $\mathcal{R}$  indexed queue length and workload processes in heavy traffic is given by  $c^r = S^{-1}(r)$ ,  $r \in \mathcal{R}$ . Through (2.1), the order of magnitude of the correction factor  $c^r$ ,  $r \in \mathcal{R}$ , is determined by the rate at which the tail of the first moment of the processing time distribution tends to zero.

One can view Theorem 3.1 as a generalized state space collapse result with a multiplicative lifting factor of one; that is, the heavy traffic limit of the unconventionally rescaled queue length process is one times the heavy traffic limit of the conventionally rescaled workload process. The proof of Theorem 3.1 given in Section 4 provides insight into how this phenomenon manifests itself. We give an informal overview there as well. Another way to view this result is that the sequence of spatial correction factors  $\{c^r\}_{r \in \mathcal{R}}$ , has embedded in it both the order of magnitude difference between the  $\mathcal{R}$  indexed queue length and workload processes in heavy traffic and the reciprocal of the multiplicative lifting map. For many processing time distributions that are of interest in practice, one can compute these explicitly. We illustrate this in the following corollary.

In the following corollary, we consider Weibull processing time distributions with positive shape parameter  $\alpha$  and positive rate parameter  $\beta$ . For these processing time distributions, the corollary precisely identifies the order of magnitude of the corrective spatial scaling factor as  $\sqrt[\alpha]{\ln r}$ . It also identifies what can be viewed as a state space collapse lifting map that obtains the limit of the sequence of diffusion scaled queue length processes with the  $r$ th member multiplied by  $\sqrt[\alpha]{\ln r}$  from the limit of the sequence of diffusion scaled workload process via multiplication by the rate parameter  $\beta$ . In this regard, it is interesting to note that multiplication of the limiting workload process by  $\beta$  is not the same as division by the mean processing time, except in the exponential case  $\alpha = 1$ . Indeed, the mean processing time is given by  $\Gamma(1 + \alpha)/\beta$ , where  $\Gamma(t) = \int_{\mathbb{R}_+} x^{t-1} \exp(-x) dx$ ,  $t \in (0, \infty)$ , denotes the gamma function. Note that  $\Gamma(1 + \alpha) < 1$  for  $0 < \alpha < 1$  and  $\Gamma(1 + \alpha) > 1$  for  $\alpha > 1$ . Then, under this nonstandard spatial scaling, the limiting residual processing time per job in the system  $1/\beta$  exceeds the mean processing time for  $0 < \alpha < 1$ . The opposite is true for  $\alpha > 1$ .

COROLLARY 3.2. *Let  $\alpha, \beta > 0$ . Assume that (3.2), (3.3), (3.6) and (3.9) hold and that  $\bar{F}(x) = \exp(-(\beta x)^{-\alpha})$ ,  $x \in \mathbb{R}_+$  (so that the processing time distribution*

is Weibull distributed with rate parameter  $\beta > 0$  and shape parameter  $\alpha > 0$ ). Then, as  $r \rightarrow \infty$ ,

$$\frac{\sqrt[\alpha]{\ln(r)} Q^r(r^2 \cdot)}{r} \Rightarrow \beta W^*(\cdot),$$

where  $W^*(\cdot)$  is a reflected Brownian motion with drift  $\kappa$  and variance  $\lambda((\sigma_a)^2 + (\sigma_s)^2)$  such that  $W^*(0)$  is equal in distribution to  $W_0$ .

PROOF. Fix  $\alpha, \beta > 0$ . We begin by more precisely determining the asymptotic behavior of  $S^{-1}(\cdot)$ ; see (3.10) below. Then we use this asymptotic behavior to verify (2.2) so that we may apply Theorem 3.1. The continuous mapping theorem together with (3.10), then allows us to replace  $c^r = S^{-1}(r)$  with  $\sqrt[\alpha]{\ln r}/\beta$  and then to multiply by the constant  $\beta$  to obtain the desired conclusion.

For  $x \in \mathbb{R}_+$ ,

$$\frac{1}{S(x)} = \mathbb{E}[v1_{\{v>x\}}] = x\bar{F}(x) + \int_x^\infty \bar{F}(y) dy \geq \frac{x}{\exp((\beta x)^\alpha)}.$$

Using L'Hopital's rule, one can verify that

$$\lim_{x \rightarrow \infty} \frac{\exp((\beta x)^\alpha)}{xS(x)} = \lim_{x \rightarrow \infty} \frac{\mathbb{E}[v1_{\{v>x\}}]}{x\bar{F}(x)} = 1.$$

Fix  $\delta \in (0, 1)$ . Then there exists  $X \in \mathbb{R}_+$  such that for all  $x > X$ ,

$$\begin{aligned} (1 - \delta) \exp(((1 - \delta)\beta x)^\alpha) &\leq \frac{(1 - \delta) \exp((\beta x)^\alpha)}{x} \\ &\leq S(x) \\ &\leq \frac{\exp((\beta x)^\alpha)}{x} \\ &\leq \exp((\beta x)^\alpha). \end{aligned}$$

So then it follows that there exists  $Y \in \mathbb{R}_+$  such that for  $y > Y$ ,

$$\frac{\sqrt[\alpha]{\ln(y)}}{\beta} \leq S^{-1}(y) \leq \frac{\sqrt[\alpha]{\ln(y/(1 - \delta))}}{(1 - \delta)\beta}.$$

Hence

$$(3.10) \quad \lim_{y \rightarrow \infty} \frac{\beta S^{-1}(y)}{\sqrt[\alpha]{\ln y}} = 1.$$

Fix  $c > 1$ . For  $y > 1$ , we have

$$\begin{aligned} &\left( \frac{\sqrt[\alpha]{\ln(cy)}}{\sqrt[\alpha]{\ln(y)}} - 1 \right) \ln\left( \frac{\sqrt[\alpha]{\ln(y)}}{\beta} \right) \\ &= \left( \sqrt[\alpha]{1 + \frac{\ln(c)}{\ln(y)}} - \sqrt[\alpha]{1} \right) \left( \frac{\ln(\ln(y))}{\alpha} - \ln \beta \right). \end{aligned}$$

Set  $h(z) = \sqrt[\alpha]{1+z}$ ,  $z \in (-1, \infty)$ . Using Taylor’s remainder theorem and the fact that  $h'(\cdot)$  is continuous in a neighborhood of the origin, there exists  $B, \delta > 0$  such that for all  $|z| < \delta$ ,

$$1 - B|z| \leq h(z) \leq 1 + B|z|.$$

So then for all  $y$  sufficiently larger than 1,

$$0 \leq \left( \sqrt[\alpha]{1 + \frac{\ln(c)}{\ln(y)}} - \sqrt[\alpha]{1} \right) \ln \left( \frac{\sqrt[\alpha]{\ln(y)}}{\beta} \right) \leq \frac{B \ln(c)}{\ln(y)} \left( \frac{\ln(\ln(y))}{\alpha} - \ln \beta \right).$$

Hence

$$(3.11) \quad \lim_{y \rightarrow \infty} \left( \frac{\sqrt[\alpha]{\ln(cy)}}{\sqrt[\alpha]{\ln(y)}} - 1 \right) \ln \left( \frac{\sqrt[\alpha]{\ln(y)}}{\beta} \right) = 0.$$

Combining (3.10) and (3.11) implies (2.2) for  $S^{-1}(\cdot)$ . Hence the result follows from Theorem 3.1, (3.10) and the continuous mapping theorem.  $\square$

**4. Proof of Theorem 3.1.** Here we state the main facts that will be proved in subsequent sections in order to verify Theorem 3.1. Then we prove Theorem 3.1 using these facts.

Henceforth, we assume that we have a sequence of  $\mathcal{R}$  indexed  $GI/GI/1$  SRPT queues satisfying the conditions in Section 3 and that  $W^*(\cdot)$  denotes a semi-martingale reflected Brownian motion with drift  $\kappa$  and variance  $\lambda((\sigma_a)^2 + (\sigma_s)^2)$  such that  $W^*(0)$  is equal in distribution to  $W_0$ . Then, by [11], as  $r \rightarrow \infty$ ,

$$(4.1) \quad \widehat{W}^r(\cdot) \Rightarrow W^*(\cdot).$$

In Section 5.1, we state and prove Lemma 5.1. This together with the fact that  $c^r < u_\varepsilon^r$  for all  $\varepsilon > 0$  and  $r \in \mathcal{R}$  implies that for all  $\varepsilon > 0$ , as  $r \rightarrow \infty$ ,

$$(4.2) \quad \langle 1_{(u_\varepsilon^r, \infty)}, \widetilde{Z}^r(\cdot) \rangle \Rightarrow 0(\cdot) \quad \text{and} \quad \langle \chi 1_{(u_\varepsilon^r, \infty)}, \widehat{Z}^r(\cdot) \rangle \Rightarrow 0(\cdot).$$

In Section 5.2.2, we state and prove Lemma 5.2. This implies that for all  $\varepsilon > 0$ , as  $r \rightarrow \infty$ ,

$$(4.3) \quad \langle 1_{[0, l_\varepsilon^r]}, \widetilde{Z}^r(\cdot) \rangle \Rightarrow 0(\cdot) \quad \text{and} \quad \langle \chi 1_{[0, l_\varepsilon^r]}, \widehat{Z}^r(\cdot) \rangle \Rightarrow 0(\cdot).$$

The asymptotic behavior summarized in (4.2) and (4.3) is used below in the proof of Theorem 3.1.

Before proceeding to prove Theorem 3.1, we provide an overview, which provides some insight into how the state space collapse that it implies arises. For this, let  $\varepsilon > 0$ . Then (4.2) and (4.3) imply that in heavy traffic the jobs that contribute to the unconventionally rescaled queue length process or to the conventionally rescaled workload process have residual processing times that asymptotically concentrate in  $(l_\varepsilon^r, u_\varepsilon^r]$  as  $r \rightarrow \infty$ . For each  $r \in \mathcal{R}$ , this interval contains the scale factor  $c^r$ . The interval itself is shifting out to infinity as  $r \rightarrow \infty$ . However, since

the workload process converges to a nondegenerate limit under diffusion scaling, the number of jobs with residual service time in this interval must tend to zero on diffusion scale. That the diffusion scaled queue length has a zero limit was shown rigorously in [9], which implies that the diffusion scaled measure valued state descriptor has a zero limit as well. However, due to (3.8), all members of this interval are of order  $c^r$ . In particular, each job with residual processing time in this interval contributes order  $c^r$  to the diffusion scaled workload process. So then, since jobs with residual service time outside of  $(l_\varepsilon^r, u_\varepsilon^r]$  do not asymptotically contribute to the unconventionally rescaled queue length process, it should follow that as  $r \rightarrow \infty$

$$c^r \widehat{Q}^r(\cdot) \approx \widehat{W}^r(\cdot).$$

The proof of Theorem 3.1 given next demonstrates this in precise terms, and thereby validates this line of reasoning.

**PROOF OF THEOREM 3.1.** We have that for all  $\varepsilon > 0, r \in \mathcal{R}$  and  $t \in [0, \infty)$ ,

$$l_\varepsilon^r \langle 1_{(l_\varepsilon^r, u_\varepsilon^r]}, \widehat{Z}^r(t) \rangle \leq \langle \chi 1_{(l_\varepsilon^r, u_\varepsilon^r]}, \widehat{Z}^r(t) \rangle \leq u_\varepsilon^r \langle 1_{(l_\varepsilon^r, u_\varepsilon^r]}, \widehat{Z}^r(t) \rangle.$$

Then, for all  $\varepsilon > 0, r \in \mathcal{R}$  and  $t \in [0, \infty)$ ,

$$(4.4) \quad \frac{c^r}{u_\varepsilon^r} \langle \chi 1_{(l_\varepsilon^r, u_\varepsilon^r]}, \widehat{Z}^r(t) \rangle \leq \langle 1_{(l_\varepsilon^r, u_\varepsilon^r]}, \widehat{Z}^r(t) \rangle \leq \frac{c^r}{l_\varepsilon^r} \langle \chi 1_{(l_\varepsilon^r, u_\varepsilon^r]}, \widehat{Z}^r(t) \rangle.$$

Fix  $T, \varepsilon, \eta, \delta > 0$ . Given  $r \in \mathcal{R}$ , let

$$\begin{aligned} \Omega_1^r &= \left\{ \sup_{t \in [0, T]} \langle 1_{(u_\varepsilon^r, \infty)}, \widehat{Z}^r(t) \rangle < \delta/3 \right\} \cap \left\{ \sup_{t \in [0, T]} \langle \chi 1_{(u_\varepsilon^r, \infty)}, \widehat{Z}^r(t) \rangle < \delta/3 \right\}, \\ \Omega_2^r &= \left\{ \sup_{t \in [0, T]} \langle 1_{[0, l_\varepsilon^r]}, \widehat{Z}^r(t) \rangle < \delta/3 \right\} \cap \left\{ \sup_{t \in [0, T]} \langle \chi 1_{[0, l_\varepsilon^r]}, \widehat{Z}^r(t) \rangle < \delta/3 \right\}. \end{aligned}$$

By (4.2) and (4.3),

$$(4.5) \quad \lim_{r \rightarrow \infty} \mathbb{P}(\Omega_1^r \cap \Omega_2^r) = 1.$$

By (4.4), for each  $r \in \mathcal{R}$ , on  $\Omega_1^r \cap \Omega_2^r$ , for all  $t \in [0, T]$ ,

$$\frac{c^r}{u_\varepsilon^r} \widehat{W}^r(t) - \frac{2\delta}{3} \leq \widehat{Q}^r(t) \leq \frac{c^r}{l_\varepsilon^r} \widehat{W}^r(t) + \frac{2\delta}{3}.$$

Then, for each  $r \in \mathcal{R}$ , on  $\Omega_1^r \cap \Omega_2^r$ , for all  $t \in [0, T]$ ,

$$(4.6) \quad \left( \frac{c^r}{u_\varepsilon^r} - 1 \right) \widehat{W}^r(t) - \frac{2\delta}{3} \leq \widehat{Q}^r(t) - \widehat{W}^r(t) \leq \left( \frac{c^r}{l_\varepsilon^r} - 1 \right) \widehat{W}^r(t) + \frac{2\delta}{3}.$$

Given  $r \in \mathcal{R}$  and  $M \in \mathbb{N}$ , let

$$\Omega^r(M) = \left\{ \sup_{t \in [0, T]} \widehat{W}^r(t) < M \right\} \quad \text{and} \quad \Omega(M) = \left\{ \sup_{t \in [0, T]} W^*(t) < M \right\}.$$

Since  $W^*(\cdot)$  is continuous almost surely,

$$\mathbb{P}\left(\bigcup_{M \in \mathbb{N}} \Omega(M)\right) = 1.$$

Hence, there exists  $M_\eta \in \mathbb{N}$  such that

$$\mathbb{P}(\Omega(M_\eta)) \geq 1 - \eta.$$

Then, by (4.1) and the Portmanteau theorem,

$$\liminf_{r \rightarrow \infty} \mathbb{P}(\Omega^r(M_\eta)) \geq 1 - \eta.$$

This together with (4.5) implies that

$$(4.7) \quad \liminf_{r \rightarrow \infty} \mathbb{P}(\Omega_1^r \cap \Omega_2^r \cap \Omega^r(M_\eta)) \geq 1 - \eta.$$

Further, by (4.6), for each  $r \in \mathcal{R}$ , on  $\Omega_1^r \cap \Omega_2^r \cap \Omega^r(M_\eta)$ ,

$$\sup_{t \in [0, T]} |\tilde{Q}^r(t) - \widehat{W}^r(t)| \leq \max\left(\frac{c^r}{l_\varepsilon^r} - 1, 1 - \frac{c^r}{u_\varepsilon^r}\right) M_\eta + \frac{2\delta}{3}.$$

By (3.8), there exists  $R \in \mathcal{R}$  such that for all  $r > R$ ,

$$\max\left(\frac{c^r}{l_\varepsilon^r} - 1, 1 - \frac{c^r}{u_\varepsilon^r}\right) \leq \frac{\delta}{3M_\eta}.$$

Then, for each  $r > R$ , on  $\Omega_1^r \cap \Omega_2^r \cap \Omega^r(M_\eta)$ ,

$$\sup_{t \in [0, T]} |\tilde{Q}^r(t) - \widehat{W}^r(t)| \leq \delta.$$

Hence, by (4.7),

$$\liminf_{r \rightarrow \infty} \mathbb{P}\left(\sup_{t \in [0, T]} |\tilde{Q}^r(t) - \widehat{W}^r(t)| \leq \delta\right) \geq 1 - \eta.$$

Since  $T, \eta, \delta > 0$  were arbitrary,

$$\tilde{Q}^r(\cdot) - \widehat{W}^r(\cdot) \Rightarrow 0(\cdot).$$

This together with (4.1) and the converging together lemma completes the proof. □

**5. Verification of (4.2) and (4.3).** Theorem 3.1 was proved in Section 4 as a consequence of (4.2) and (4.3) and other facts already established in the paper. The remainder of the paper is devoted to stating and proving the two lemmas that imply (4.2) and (4.3), namely Lemmas 5.1 and 5.2.

5.1. *Workload process tail behavior.* In this section we prove Lemma 5.1, which implies (4.2). The tail behavior asserted here is relatively easy to verify since it is simply a manifestation of the scaling. This is evident in the proof given below.

LEMMA 5.1. *For all  $\varepsilon > 0$ , as  $r \rightarrow \infty$ ,*

$$(5.1) \quad \langle \chi 1_{(u_\varepsilon^r, \infty)}, \widehat{\mathcal{Z}}^r(\cdot) \rangle \Rightarrow 0(\cdot).$$

PROOF. Fix  $\varepsilon > 0$ . For  $r \in \mathcal{R}$  and  $t \in [0, \infty)$ ,  $\tilde{v}_i^r(t) \leq \tilde{v}_i^r$  for all  $1 \leq i \leq Q^r(0)$  and  $v_i^r(t) \leq v_i$  for all  $1 \leq i \leq E^r(r^2t)$ . Hence, for  $r \in \mathcal{R}$ ,

$$(5.2) \quad \langle \chi 1_{(u_\varepsilon^r, \infty)}, \widehat{\mathcal{Z}}^r(\cdot) \rangle \leq \langle \chi 1_{(u_\varepsilon^r, \infty)}, \widehat{\mathcal{Z}}^r(0) \rangle + \frac{1}{r} \sum_{i=1}^{r^2 \overline{E}^r(\cdot)} v_i 1_{\{v_i > u_\varepsilon^r\}}.$$

Further, for  $r \in \mathcal{R}$ ,

$$\frac{1}{r} \sum_{i=1}^{r^2 \overline{E}^r(\cdot)} v_i 1_{\{v_i > u_\varepsilon^r\}} = \frac{1}{r} \left( \sum_{i=1}^{r^2 \overline{E}^r(\cdot)} v_i 1_{\{v_i > u_\varepsilon^r\}} - r^2 \lambda^r(\cdot) \mathbb{E}[v 1_{\{v > u_\varepsilon^r\}}] \right) + \frac{r \lambda^r(\cdot)}{S(u_\varepsilon^r)}.$$

By (3.7), (3.2) and  $\lim_{r \rightarrow \infty} c^r = \infty$ ,

$$(5.3) \quad \lim_{r \rightarrow \infty} \frac{r \lambda^r(\cdot)}{S(u_\varepsilon^r)} = \lim_{r \rightarrow \infty} \frac{\lambda^r(\cdot)}{(c^r)^{2+\varepsilon}} = 0(\cdot).$$

Further, as  $r$  tends to infinity,  $\mathbb{E}[v 1_{\{v > u_\varepsilon^r\}}]$  and  $\mathbb{E}[v^2 1_{\{v > u_\varepsilon^r\}}]$  converge to zero since  $\lim_{r \rightarrow \infty} u_\varepsilon^r = \infty$ . Hence, by Proposition A.1, as  $r \rightarrow \infty$ ,

$$\frac{1}{r} \left( \sum_{i=1}^{r^2 \overline{E}^r(\cdot)} v_i 1_{\{v_i > u_\varepsilon^r\}} - r^2 \lambda^r(\cdot) \mathbb{E}[v 1_{\{v > u_\varepsilon^r\}}] \right) \Rightarrow 0(\cdot).$$

Therefore, as  $r \rightarrow \infty$ ,

$$\frac{1}{r} \sum_{i=1}^{r^2 \overline{E}^r(\cdot)} v_i 1_{\{v_i > u_\varepsilon^r\}} \Rightarrow 0(\cdot).$$

Combining this with (3.9) and (5.2) implies (5.1).  $\square$

5.2. *Behavior in large neighborhoods of the origin.* In this section, we prove the following lemma, which implies (4.3).

LEMMA 5.2. *For all  $\varepsilon > 0$ , as  $r \rightarrow \infty$ ,*

$$(5.4) \quad \langle (1 \vee \chi) 1_{[0, u_\varepsilon^r]}, \widetilde{\mathcal{Z}}^r(\cdot) \rangle \Rightarrow 0(\cdot).$$

The behavior asserted in Lemma 5.2 is more subtle than that asserted in Lemma 5.1 since it relies on the SRPT processing dynamics. Key elements used in verifying this result are asymptotics obtained for the duration of busy periods for large neighborhoods of the origin; see Lemmas 5.3 and 5.4. Such results are refinements of [9], (4.9), where the neighborhood of the origin does not grow with  $r \in \mathcal{R}$ , and a slower rate of convergence to zero is verified for fixed width neighborhoods of the origin. Equations (5.6) and (5.7) developed below play a central role in proving these rate of convergence results. They exploit the nonidling nature of SRPT as well as the order in which jobs are processed.

Once Lemmas 5.3 and 5.4 are established, we verify that the total mass in a fixed width neighborhood of the origin converges to zero; see Lemma 5.5. The proof of Lemma 5.5 utilizes an inequality similar in spirit to (5.6), but for total mass rather than the total amount of work; see (5.15). This inequality is less precise than (5.6) since knowing how many time units the server has spent processing work does not exactly prescribe the number of jobs that exit the system during that timeframe. However, by fixing the width of the neighborhood of the origin, one can utilize this dynamic inequality together with the result in Lemma 5.3 to obtain the desired conclusion.

The final step is to verify that the total amount of work in a growing neighborhood of the origin tends to zero; see Lemma 5.6. For this, we return to (5.7) multiplied by the corrective spatial scaling factor  $c^r$  and with  $x$  taken to be  $l_\varepsilon^r$ ,  $\varepsilon > 0$  and  $r \in \mathcal{R}$ . This yields an upper bound on the desired quantity. Then we need to verify that all terms on the right-hand side tend to zero. In particular, we must verify that the net change over certain busy periods of what could be referred to as centered truncated load processes tends to zero sufficiently fast. This is addressed by Lemma 5.7. Since these centered, truncated load processes converge to Brownian motion (as noted in the Appendix), the proof strategy is to use Hölder continuity of Brownian motion to bound such differences by quantities involving the duration of the busy period. This allows one to utilize the asymptotics obtained in Lemma 5.4 to prove Lemma 5.7. The result in Lemma 5.7 is combined with other facts in order to prove Lemma 5.6 at the end of Section 5.2.3.

For completeness, we write out the proof of Lemma 5.2 as a consequence of Lemmas 5.5 and 5.6 here.

PROOF OF LEMMA 5.2. Fix  $\varepsilon > 0$ . Then, given  $r \in \mathcal{R}$ ,

$$\langle (1 \vee \chi)1_{[0, l_\varepsilon^r]}, \tilde{Z}^r(\cdot) \rangle \leq \langle 1_{[0, 1]}, \tilde{Z}^r(\cdot) \rangle + \langle \chi 1_{[0, l_\varepsilon^r]}, \tilde{Z}^r(\cdot) \rangle.$$

This together with Lemmas 5.5 and 5.6 immediately implies (5.4).  $\square$

The remainder of this section contains the statements and proofs of Lemmas 5.5 and 5.6.

5.2.1. *Asymptotics for busy period durations.* For  $x \in \mathbb{R}_+$ ,  $r \in \mathcal{R}$  and  $t \in [0, \infty)$ , let

$$\tau^r(t, x) = \sup\{s \in [0, t] : \langle 1_{[0,x]}, \tilde{\mathcal{Z}}^r(s) \rangle = 0\} \quad \text{and} \quad \theta^r(t, x) = t - \tau^r(t, x).$$

Given  $x \in \mathbb{R}_+$ ,  $r \in \mathcal{R}$  and  $t \in [0, \infty)$ ,  $\theta^r(t, x)$  represents the amount of time that has elapsed since the  $r$ th system had no jobs with residual processing time in  $[0, x]$ . In particular, given  $x \in \mathbb{R}_+$ ,  $r \in \mathcal{R}$  and  $t \in [0, \infty)$ ,  $\langle \chi 1_{[0,x]}, \mathcal{Z}^r(r^2s) \rangle > 0$  for all  $s \in (\tau^r(t, x), t]$ . Hence, during the time interval  $(r^2\tau^r(t, x), r^2t]$  the server in the  $r$ th system is busy and devoted to serving jobs with remaining processing time in  $[0, x]$ . Hence, for each  $x \in \mathbb{R}_+$ ,  $r \in \mathcal{R}$  and  $t \in [0, \infty)$ ,

$$\begin{aligned} \langle \chi 1_{[0,x]}, \mathcal{Z}^r(r^2t) \rangle &= \langle \chi 1_{[0,x]}, \mathcal{Z}^r(r^2\tau^r(t, x)) \rangle \\ &+ \sum_{i=E^r(r^2\tau^r(t, x))+1}^{E^r(r^2t)} v_i 1_{\{v_i \leq x\}} - r^2\theta^r(t, x). \end{aligned}$$

For  $x \in \mathbb{R}_+$ ,  $r \in \mathcal{R}$  and  $t \in [0, \infty)$ , set

$$\begin{aligned} V_x^r(t) &= \sum_{i=1}^{E^r(t)} v_i 1_{\{v_i \leq x\}}, \\ \bar{V}_x^r(t) &= \frac{V_x^r(r^2t)}{r^2}, \\ \hat{V}_x^r(t) &= \frac{1}{r}(V_x^r(r^2t) - \rho_x^r r^2t). \end{aligned}$$

Here, given  $r \in \mathcal{R}$  and  $x \in \mathbb{R}_+$ ,  $V_x^r(\cdot)$  is referred to as a truncated load process. Then, for  $x \in \mathbb{R}_+$ ,  $r \in \mathcal{R}$  and  $t \in [0, \infty)$ ,

$$\begin{aligned} \langle \chi 1_{[0,x]}, \hat{\mathcal{Z}}^r(t) \rangle &= \langle \chi 1_{[0,x]}, \hat{\mathcal{Z}}^r(\tau^r(t, x)) \rangle + r(\bar{V}_x^r(t) - \bar{V}_x^r(\tau^r(t, x))) \\ &\quad - r\theta^r(t, x), \\ \langle \chi 1_{[0,x]}, \hat{\mathcal{Z}}^r(t) \rangle &= \langle \chi 1_{[0,x]}, \hat{\mathcal{Z}}^r(\tau^r(t, x)) \rangle + \hat{V}_x^r(t) - \hat{V}_x^r(\tau^r(t, x)) \\ &\quad + (\rho_x^r - 1)r\theta^r(t, x). \end{aligned}$$

Given  $x \in \mathbb{R}_+$ ,  $r \in \mathcal{R}$  and  $t \in [0, \infty)$ , either  $\tau^r(t, x) = 0$  or  $\tau^r(t, x) > 0$ . If the latter, then at time  $\tau^r(t, x)$ , either a job with total processing time in  $[0, x]$  arrives exogenously or a job with total processing time greater than  $x$  was in service immediately before time  $\tau^r(t, x)$ , and its remaining processing time at time  $\tau^r(t, x)$  is  $x$ . Hence, for  $x \in \mathbb{R}_+$ ,  $r \in \mathcal{R}$  and  $t \in [0, \infty)$ ,

$$\begin{aligned} \langle \chi 1_{[0,x]}, \hat{\mathcal{Z}}^r(\tau^r(t, x)) \rangle &\leq \langle \chi 1_{[0,x]}, \hat{\mathcal{Z}}^r(0) \rangle \\ (5.5) \quad &+ \frac{1}{r}(V_x^r(r^2\tau^r(t, x)) - V_x^r(r^2\tau^r(t, x)-)) + \frac{x}{r} \end{aligned}$$



$$\begin{aligned} &= \langle \chi 1_{[0,x]}, \widehat{\mathcal{Z}}^r(0) \rangle \\ &\quad + r(\overline{V}_x^r(\tau^r(t,x)) - \overline{V}_x^r(\tau^r(t,x)-)) + \frac{x}{r} \\ &= \langle \chi 1_{[0,x]}, \widehat{\mathcal{Z}}^r(0) \rangle \\ &\quad + \widehat{V}_x^r(\tau^r(t,x)) - \widehat{V}_x^r(\tau^r(t,x)-) + \frac{x}{r}. \end{aligned}$$

Therefore, for  $x \in \mathbb{R}_+$ ,  $r \in \mathcal{R}$  and  $t \in [0, \infty)$ ,

$$(5.6) \quad \begin{aligned} \langle \chi 1_{[0,x]}, \widehat{\mathcal{Z}}^r(t) \rangle &\leq \langle \chi 1_{[0,x]}, \widehat{\mathcal{Z}}^r(0) \rangle + r(\overline{V}_x^r(t) - \overline{V}_x^r(\tau^r(t,x)-)) \\ &\quad - r\theta^r(t,x) + \frac{x}{r}, \end{aligned}$$

$$(5.7) \quad \begin{aligned} \langle \chi 1_{[0,x]}, \widehat{\mathcal{Z}}^r(t) \rangle &\leq \langle \chi 1_{[0,x]}, \widehat{\mathcal{Z}}^r(0) \rangle + \widehat{V}_x^r(t) - \widehat{V}_x^r(\tau^r(t,x)-) \\ &\quad + (\rho_x^r - 1)r\theta^r(t,x) + \frac{x}{r}. \end{aligned}$$

We use (5.6) to prove the next lemma, which specifies the asymptotic behavior of  $\theta^r(\cdot, x)$  as  $r \rightarrow \infty$ . We use (5.7) to prove the subsequent lemma, which specifies the asymptotic behavior of  $\theta^r(\cdot, l'_\varepsilon)$  as  $r \rightarrow \infty$ .

LEMMA 5.3. *For each  $x \in \mathbb{R}_+$ , as  $r \rightarrow \infty$ ,*

$$(5.8) \quad c^r r\theta^r(\cdot, x) \Rightarrow 0(\cdot).$$

PROOF. Given  $x \in \mathbb{R}_+$ , let  $\rho_x = \lambda \mathbb{E}[v 1_{\{v \leq x\}}]$  and  $\rho_x(t) = \rho_x t$  for all  $[0, \infty)$ . Then, (A.1) implies that, for each  $x \in \mathbb{R}_+$ , as  $r \rightarrow \infty$ ,

$$(5.9) \quad \overline{V}_x^r(\cdot) \Rightarrow \rho_x(\cdot).$$

Fix  $x \in \mathbb{R}_+$ ,  $T > 0$  and  $\gamma > 0$ . Note that  $\rho_x < 1$ . Let  $\delta > 0$  be such that  $(1 + \delta)\rho_x < 1$ . For  $r \in \mathcal{R}$ , let

$$\begin{aligned} \Omega_0^r &= \left\{ \langle \chi 1_{[0,x]}, \widehat{\mathcal{Z}}^r(0) \rangle \leq \frac{\gamma}{2} \right\}, \\ \Omega_1^r &= \left\{ \sup_{0 \leq s \leq t \leq T} \overline{V}_x^r(t) - \overline{V}_x^r(s-) < (1 + \delta)\rho_x(t - s) \right\}, \\ \Omega^r &= \Omega_0^r \cap \Omega_1^r. \end{aligned}$$

By (3.9) and (5.9),

$$\lim_{r \rightarrow \infty} \mathbb{P}(\Omega^r) = 1.$$

By (5.6), for each  $r \in \mathcal{R}$ , on  $\Omega^r$ , for each  $t \in [0, T]$ ,

$$\langle \chi 1_{[0,x]}, \widehat{\mathcal{Z}}^r(t) \rangle \leq \langle \chi 1_{[0,x]}, \widehat{\mathcal{Z}}^r(0) \rangle + ((1 + \delta)\rho_x - 1)r\theta^r(t,x) + \frac{x}{r}.$$

But for each  $r \in \mathcal{R}$ ,  $\langle \chi 1_{[0,x]}, \widehat{\mathcal{Z}}^r(t) \rangle \geq 0$  for each  $t \in [0, T]$ . Hence, for each  $r \in \mathcal{R}$ , on  $\Omega^r$ , for each  $t \in [0, T]$ ,

$$(1 - (1 + \delta)\rho_x)c^r r\theta^r(t, x) \leq \langle \chi 1_{[0,x]}, \widehat{\mathcal{Z}}^r(0) \rangle + \frac{c^r x}{r}.$$

Recall that  $S^{-1}(\cdot)$  is slowly varying so that  $\lim_{y \rightarrow \infty} S^{-1}(y)/y = 0$ . Hence  $\lim_{r \rightarrow \infty} c^r/r = \lim_{r \rightarrow \infty} S^{-1}(r)/r = 0$ . Then for  $r \in \mathcal{R}$  sufficiently large, on  $\Omega^r$ , for each  $t \in [0, T]$ ,

$$c^r r\theta^r(t, x) \leq \frac{\gamma}{(1 - (1 + \delta)\rho_x)}.$$

Since  $\lim_{r \rightarrow \infty} \mathbb{P}(\Omega^r) = 1$ , (5.8) holds.  $\square$

One feature of the SRPT discipline that is utilized in the above proof is that by restricting to jobs with remaining processing time in  $[0, x]$  for a fixed  $x$ , the workload process truncated to jobs with remaining processing time in  $[0, x]$  effectively behaves as a subcritical queue. We wish to obtain a version of Lemma 5.3 on  $[0, l'_\varepsilon]$  for fixed  $\varepsilon > 0$  with  $r \rightarrow \infty$ . Note that for  $\varepsilon > 0$ ,  $\lim_{r \rightarrow \infty} l'_\varepsilon = \infty$ . Therefore, on such time intervals, the truncated workload process approaches that of a critical queue. This makes the verification of Lemma 5.4 a bit more delicate, and the rate of convergence result obtained is not as rapid. For this, for  $\varepsilon > 0$ ,  $r \in \mathcal{R}$  and  $t \in [0, \infty)$ , we adopt the shorthand notation

$$\tau_\varepsilon^r(t) = \tau^r(t, l'_\varepsilon) \quad \text{and} \quad \theta_\varepsilon^r(t) = \theta^r(t, l'_\varepsilon).$$

LEMMA 5.4. For  $\varepsilon > 0$ , as  $r \rightarrow \infty$ ,

$$(c^r)^{2+\varepsilon} \theta_\varepsilon^r(\cdot) \Rightarrow 0(\cdot).$$

PROOF. Fix  $\varepsilon > 0$  and  $t \in [0, \infty)$ . Given  $r \in \mathcal{R}$ , we take  $x = l'_\varepsilon$  in (5.7), and then we subtract and add  $\rho^r r\theta_\varepsilon^r(t)$ , and use (3.1) and the fact that  $S(l'_\varepsilon) = r(c^r)^{-2-\varepsilon}$  to obtain that for  $r \in \mathcal{R}$ ,

$$\begin{aligned} \langle \chi 1_{[0,l'_\varepsilon]}, \widehat{\mathcal{Z}}^r(t) \rangle &\leq \langle \chi 1_{[0,l'_\varepsilon]}, \widehat{\mathcal{Z}}^r(0) \rangle + \widehat{V}_{l'_\varepsilon}^r(t) - \widehat{V}_{l'_\varepsilon}^r(\tau_\varepsilon^r(t)-) \\ &\quad - \lambda^r (c^r)^{2+\varepsilon} \theta_\varepsilon^r(t) + (\rho^r - 1)r\theta_\varepsilon^r(t) + \frac{l'_\varepsilon}{r}. \end{aligned} \tag{5.10}$$

We have that  $\langle \chi 1_{[0,l'_\varepsilon]}, \widehat{\mathcal{Z}}^r(t) \rangle \geq 0$  and  $\theta_\varepsilon^r(t) \geq 0$  for all  $r \in \mathcal{R}$ . This together with the fact that  $l'_\varepsilon < c^r$  implies that, for all  $r \in \mathcal{R}$ ,

$$\begin{aligned} 0 &\leq \lambda^r (c^r)^{2+\varepsilon} \theta_\varepsilon^r(t) \\ &\leq \langle \chi 1_{[0,l'_\varepsilon]}, \widehat{\mathcal{Z}}^r(0) \rangle + \widehat{V}_{l'_\varepsilon}^r(t) - \widehat{V}_{l'_\varepsilon}^r(\tau_\varepsilon^r(t)-) + (\rho^r - 1)r\theta_\varepsilon^r(t) + \frac{c^r}{r}. \end{aligned} \tag{5.11}$$

Upon dividing by  $(c^r)^{2+\varepsilon}$  and using  $\lim_{r \rightarrow \infty} c^r = \infty$ , (3.2), (3.9) and (A.2), we see that, as  $r \rightarrow \infty$ ,

$$(5.12) \quad \theta_\varepsilon^r(\cdot) \Rightarrow 0(\cdot).$$

Hence, by (A.2) and the fact that  $V^*(\cdot)$  is continuous, as  $r \rightarrow \infty$ ,

$$(5.13) \quad \widehat{V}_{l_\varepsilon}^r(\cdot) - \widehat{V}_{l_\varepsilon}^r(\tau_\varepsilon^r(\cdot)-) \Rightarrow 0(\cdot).$$

Then letting  $r \rightarrow \infty$  in (5.11) and using (3.2), (3.9), (5.12), (5.13) and the fact that  $c^r = S^{-1}(r)$  and  $S^{-1}(\cdot)$  is slowly varying completes the proof.  $\square$

5.2.2. *Truncated queue length process asymptotics.* We are prepared to use Lemma 5.3 to verify that the total mass in a fixed width neighborhood of the origin vanishes as  $r$  tends to infinity.

LEMMA 5.5. *For all  $x \in \mathbb{R}_+$ , as  $r \rightarrow \infty$ ,*

$$(5.14) \quad \langle 1_{[0,x]}, \tilde{Z}^r(\cdot) \rangle \Rightarrow 0(\cdot).$$

PROOF. Fix  $x \in \mathbb{R}_+$  and  $T > 0$ . By ignoring any processing that occurs in  $(r^2\tau^r(t,x), r^2t]$ , for  $r \in \mathcal{R}$  and  $t \in [0, \infty)$ , we have that

$$\langle 1_{[0,x]}, \tilde{Z}^r(t) \rangle \leq \langle 1_{[0,x]}, \tilde{Z}^r(\tau^r(t,x)) \rangle + c^r r (\overline{E}^r(t) - \overline{E}^r(\tau^r(t,x))).$$

Further, by using arguments similar to those that yielded (5.5), for  $r \in \mathcal{R}$  and  $t \in [0, \infty)$ ,

$$\langle 1_{[0,x]}, \tilde{Z}^r(\tau^r(t,x)) \rangle \leq \langle 1_{[0,x]}, \tilde{Z}^r(0) \rangle + c^r r (\overline{E}^r(\tau^r(t,x)) - \overline{E}^r(\tau^r(t,x)-)) + \frac{c^r}{r}.$$

Then, for  $r \in \mathcal{R}$  and  $t \in [0, \infty)$ , we have that

$$(5.15) \quad \langle 1_{[0,x]}, \tilde{Z}^r(t) \rangle \leq \langle 1_{[0,x]}, \tilde{Z}^r(0) \rangle + c^r r (\overline{E}^r(t) - \overline{E}^r(\tau^r(t,x)-)) + \frac{c^r}{r}.$$

Fix  $\gamma > 0$ . For  $r \in \mathcal{R}$ , let

$$\begin{aligned} \Omega_0^r &= \left\{ \langle 1_{[0,x]}, \tilde{Z}^r(0) \rangle \leq \frac{\gamma}{3} \right\}, \\ \Omega_1^r &= \left\{ \sup_{0 \leq s \leq t \leq T} (\overline{E}^r(t) - \overline{E}^r(s-)) < 2\lambda(t-s) \right\}, \\ \Omega_2^r &= \left\{ \sup_{t \in [0,T]} \theta^r(t,x) < \frac{\gamma}{6\lambda c^r r} \right\}, \\ \Omega^r &= \Omega_0^r \cap \Omega_1^r \cap \Omega_2^r. \end{aligned}$$

By (3.4), (3.9) and (5.8),  $\lim_{r \rightarrow \infty} \mathbb{P}(\Omega^r) = 1$ . Then since  $c^r = S^{-1}(r)$  and  $S^{-1}(\cdot)$  is slowly varying, it follows that, on  $\Omega^r$ , for  $r$  sufficiently large,

$$\sup_{t \in [0,T]} \langle 1_{[0,x]}, \tilde{Z}^r(t) \rangle \leq \gamma.$$

Since  $\gamma > 0$  was arbitrary, the proof is complete.  $\square$

5.2.3. *Truncated workload process asymptotics.* We are prepared to use Lemma 5.4 to verify that the total work in a growing neighborhood of the origin vanishes as  $r$  tends to infinity.

LEMMA 5.6. *For all  $\varepsilon > 0$ , as  $r \rightarrow \infty$ ,*

$$\langle \chi 1_{[0, l_\varepsilon^r]}, \tilde{Z}^r(\cdot) \rangle \Rightarrow 0(\cdot).$$

Before proving Lemma 5.6, we begin with an observation. By (5.10), for each  $\varepsilon > 0$ ,  $r \in \mathcal{R}$  and  $t \in [0, \infty)$ ,

$$(5.16) \quad \langle \chi 1_{[0, l_\varepsilon^r]}, \tilde{Z}^r(t) \rangle \leq \langle \chi 1_{[0, l_\varepsilon^r]}, \tilde{Z}^r(0) \rangle + c^r (\widehat{V}_{l_\varepsilon^r}^r(t) - \widehat{V}_{l_\varepsilon^r}^r(\tau_\varepsilon^r(t)-)) \\ - \lambda^r (c^r)^{3+\varepsilon} \theta_\varepsilon^r(t) + c^r (\rho^r - 1)r\theta_\varepsilon^r(t) + \frac{c^r l_\varepsilon^r}{r}.$$

We argue that each term on the right-hand side converges in distribution to the zero process. We begin by proving the following lemma.

LEMMA 5.7. *For each  $\varepsilon > 0$ , as  $r \rightarrow \infty$ ,*

$$c^r (\widehat{V}_{l_\varepsilon^r}^r(\cdot) - \widehat{V}_{l_\varepsilon^r}^r(\tau_\varepsilon^r(\cdot)-)) \Rightarrow 0(\cdot).$$

PROOF. Fix  $T, \varepsilon > 0$ . Recall that Brownian motion is Hölder continuous with exponent  $\gamma$  for any  $0 < \gamma < 1/2$ . Fix  $0 < \gamma < 1/2$  such that  $\gamma(2 + \varepsilon) > 1$ . For  $M \in \mathbb{N}$ , let

$$\Omega(M) = \{|V^*(t) - V^*(s-)| < M(t - s)^\gamma \text{ for all } 0 \leq s \leq t \leq T\}.$$

We have that  $\Omega(M) \subset \Omega(M + 1)$  for all  $M \in \mathbb{N}$  and

$$\mathbb{P}\left(\bigcup_{M \in \mathbb{N}} \Omega(M)\right) = 1.$$

Hence given  $\eta > 0$ , there exists  $M_\eta \in \mathbb{N}$  such that

$$\mathbb{P}(\Omega(M_\eta)) \geq 1 - \eta.$$

Given  $r \in \mathcal{R}$  and  $M \in \mathbb{N}$ , let

$$\Omega^r(M) = \{|\widehat{V}_{l_\varepsilon^r}^r(t) - \widehat{V}_{l_\varepsilon^r}^r(s-)| < M(t - s)^\gamma \text{ for all } 0 \leq s \leq t \leq T\}.$$

For each  $M \in \mathbb{N}$ , the set  $A(M)$ , given by

$$A(M) = \{f \in \mathbf{D}([0, T], \mathbb{R}) : |f(t) - f(s-)| < M(t - s)^\gamma \text{ for all } 0 \leq s \leq t \leq T\},$$

is open in the uniform topology. Hence, (A.2) and the Portmanteau theorem imply that

$$\liminf_{r \rightarrow \infty} \mathbb{P}(\Omega^r(M_\eta)) \geq \mathbb{P}(\Omega(M_\eta)) \geq 1 - \eta.$$

For  $r \in \mathcal{R}$ , let

$$\Omega_1^r = \left\{ \sup_{t \in [0, T]} \theta_\varepsilon^r(t) < \frac{1}{\sqrt[\gamma]{M_\eta(c^r)^{2+\varepsilon}}} \right\}.$$

By Lemma 5.4,

$$\lim_{r \rightarrow \infty} \mathbb{P}(\Omega_1^r) = 1.$$

Then

$$\liminf_{r \rightarrow \infty} \mathbb{P}(\Omega^r(M_\eta) \cap \Omega_1^r) \geq 1 - \eta.$$

Given  $r \in \mathcal{R}$ , set

$$\Omega_2^r = \left\{ \sup_{t \in [0, T]} c^r |\widehat{V}^r(t) - \widehat{V}^r(\tau_\varepsilon^r(t)-)| < \frac{c^r}{(c^r)^{(2+\varepsilon)\gamma}} = \frac{1}{(c^r)^{(2+\varepsilon)\gamma-1}} \right\}.$$

Then, for  $r \in \mathcal{R}$ ,

$$\Omega^r(M_\eta) \cap \Omega_1^r \subset \Omega_2^r.$$

Hence

$$\liminf_{r \rightarrow \infty} \mathbb{P}(\Omega_2^r) \geq 1 - \eta.$$

But, for  $r \in \mathcal{R}$ ,  $\Omega_2^r$  does not depend on  $\eta$ . Therefore, we may let  $\eta$  decrease to zero so that

$$\liminf_{r \rightarrow \infty} \mathbb{P}(\Omega_2^r) = 1.$$

Fix  $\delta > 0$ . Given  $r \in \mathcal{R}$ , let

$$\Omega_3^r = \left\{ \sup_{t \in [0, T]} c^r |\widehat{V}^r(t) - \widehat{V}^r(\tau_\varepsilon^r(t)-)| < \delta \right\}.$$

Since  $\lim_{r \rightarrow \infty} c^r = \infty$  and  $(2 + \varepsilon)\gamma - 1 > 0$ , it follows that for  $r$  sufficiently large  $\Omega_2^r \subset \Omega_3^r$ . Therefore,  $\liminf_{r \rightarrow \infty} \mathbb{P}(\Omega_3^r) = 1$ . Since  $T, \varepsilon, \delta > 0$  were arbitrary, Lemma 5.7 holds.  $\square$

**COROLLARY 5.8.** *For each  $\varepsilon > 0$ , as  $r \rightarrow \infty$ ,*

$$\lambda^r (c^r)^{3+\varepsilon} \theta_\varepsilon^r(\cdot) \Rightarrow 0(\cdot).$$

**PROOF.** Fix  $\varepsilon > 0$ . By (5.16), the fact that  $\langle \chi 1_{[0, l_\varepsilon^r]}, \widetilde{Z}^r(t) \rangle \geq 0$  and  $\theta_\varepsilon^r(t) \geq 0$  for all  $r \in \mathcal{R}$  and  $t \in [0, \infty)$  and  $l_\varepsilon^r < c^r$  for all  $r \in \mathcal{R}$ , we have that, for all  $r \in \mathcal{R}$  and  $t \in [0, \infty)$ ,

$$\begin{aligned} 0 &\leq \lambda^r (c^r)^{3+\varepsilon} \theta_\varepsilon^r(t) \\ (5.17) \quad &\leq \langle \chi 1_{[0, l_\varepsilon^r]}, \widetilde{Z}^r(0) \rangle + c^r (\widehat{V}_{l_\varepsilon^r}^r(t) - \widehat{V}_{l_\varepsilon^r}^r(\tau_\varepsilon^r(t)-)) \\ &\quad + c^r (\rho^r - 1) r \theta_\varepsilon^r(t) + \frac{(c^r)^2}{r}. \end{aligned}$$

The result follows from this, (3.9), Lemma 5.7, (3.2), Lemma 5.4 and the fact that  $c^r = S^{-1}(r)$  and  $S^{-1}(\cdot)$  is slowly varying.  $\square$

PROOF OF LEMMA 5.6. Fix  $\varepsilon > 0$ . The result follows by combining (5.16), (3.9), Lemma 5.7, Corollary 5.8, (3.2), Lemma 5.4,  $l'_\varepsilon < c^r$  for  $r \in \mathcal{R}$ ,  $c^r = S^{-1}(r)$  for  $r \in \mathcal{R}$  and  $S^{-1}(\cdot)$  is slowly varying.  $\square$

APPENDIX: BEHAVIOR OF TRUNCATED LOAD PROCESSES

The following result is well known and follows from [19], Theorem 3.1, used to extend [2], Section 17.3.

PROPOSITION A.1. For each  $r \in \mathcal{R}$ , let  $\{x_k^r\}_{k=1}^\infty$  be an independent and identically distributed sequence of nonnegative random variables with finite mean  $m^r$  and finite standard deviation  $\sigma^r$  that is independent of  $E^r(\cdot)$ . Suppose that for some finite nonnegative constants  $m$  and  $\sigma$ ,  $\lim_{r \rightarrow \infty} m^r = m$  and  $\lim_{r \rightarrow \infty} \sigma^r = \sigma$ . Further assume that for each  $\delta > 0$ ,

$$\lim_{r \rightarrow \infty} \mathbb{E}[(x_1^r - m^r)^2 | x_1^r - m^r | > r\delta] = 0.$$

For  $r \in \mathcal{R}$ ,  $n \in \mathbb{N}$  and  $t \in [0, \infty)$ , let

$$X^r(n) = \sum_{k=1}^n x_k^r \quad \text{and} \quad \widehat{X}^r(t) = \frac{X^r(\lfloor r^2 t \rfloor) - \lfloor r^2 t \rfloor m^r}{r}.$$

Then, as  $r \rightarrow \infty$ ,  $(\widehat{E}^r(\cdot), \widehat{X}^r(\cdot)) \Rightarrow (E^*(\cdot), X^*(\cdot))$ , where  $E^*(\cdot)$  is given by (3.3), and  $X^*(\cdot)$  is a Brownian motion starting from zero with zero drift and variance  $\sigma^2$  per unit time, that is independent of  $E^*(\cdot)$ . Furthermore, as  $r \rightarrow \infty$ ,

$$\frac{X^r(r^2 \overline{E}^r(\cdot)) - r^2 \lambda^r(\cdot) m^r}{r} \Rightarrow X^*(\lambda(\cdot)) + m E^*(\cdot),$$

where for each  $r \in \mathcal{R}$  and  $t \in [0, \infty)$ ,  $\lambda^r(t) = \lambda^r t$  and  $\lambda(t) = \lambda t$ .

Recall that, for  $r \in \mathcal{R}$  and  $x \in \mathbb{R}_+$ ,

$$\widehat{V}_x^r(\cdot) = \frac{\sum_{i=1}^{r^2 \overline{E}^r(\cdot)} v_i 1_{\{v_i \leq x\}} - r^2 \lambda^r(\cdot) \mathbb{E}[v 1_{\{v \leq x\}}]}{r}.$$

Proposition A.1 implies that for each  $x \in \mathbb{R}_+$ , as  $r \rightarrow \infty$ ,

$$(A.1) \quad \widehat{V}_x^r(\cdot) \Rightarrow V_x^*(\cdot),$$

where  $V_x^*(\cdot)$  is a Brownian motion starting from zero with drift zero and finite variance per unit time. Similarly, Proposition A.1 together with  $0 \leq \mathbb{E}[v 1_{\{v \leq l'_\varepsilon\}}] \leq$

$\mathbb{E}[v]$  and  $0 \leq \mathbb{E}[v^2 1_{\{v \leq r'_\varepsilon\}}] \leq \mathbb{E}[v^2]$  for all  $\varepsilon > 0$  and  $r \in \mathcal{R}$  and the monotone convergence theorem implies that for each  $\varepsilon > 0$ , as  $r \rightarrow \infty$ ,

$$(A.2) \quad \widehat{V}_{r'_\varepsilon}^r(\cdot) \Rightarrow V^*(\cdot),$$

where  $V^*(\cdot)$  is a Brownian motion starting from zero with drift zero and finite variance per unit time.

**Acknowledgment.** The author would like to thank ViaSat Inc. for generously funding undergraduate research assistants Richard Hunperger and Sean Malter who developed code and performed simulations that helped the author formulate the statement of Theorem 3.1.

## REFERENCES

- [1] BANSAL, N. and HARCHOL-BALTER, M. (2001). Analysis of SRPT scheduling: Investigating unfairness. *ACM SIGMETRICS Performance Evaluation Review* **29** 279–290.
- [2] BILLINGSLEY, P. (1968). *Convergence of Probability Measures*. Wiley, New York. [MR0233396](#)
- [3] BINGHAM, N. H., GOLDIE, C. M. and TEUGELS, J. L. (1987). *Regular Variation. Encyclopedia of Mathematics and Its Applications* **27**. Cambridge Univ. Press, Cambridge. [MR0898871](#)
- [4] BOJANIĆ, R. and SENETA, E. (1971). Slowly varying functions and asymptotic relations. *J. Math. Anal. Appl.* **34** 302–315. [MR0274676](#)
- [5] DOWN, D. G., GROMOLL, H. C. and PUHA, A. L. (2009). State-dependent response times via fluid limits for shortest remaining processing time queues. *ACM SIGMETRICS Performance Evaluation Review* **37** 75–76.
- [6] DOWN, D. G., GROMOLL, H. C. and PUHA, A. L. (2009). Fluid limits for shortest remaining processing time queues. *Math. Oper. Res.* **34** 880–911. [MR2573501](#)
- [7] ETHIER, S. N. and KURTZ, T. G. (1986). *Markov Processes: Characterization and Convergence*. Wiley, New York. [MR0838085](#)
- [8] GROMOLL, H. C. (2004). Diffusion approximation for a processor sharing queue in heavy traffic. *Ann. Appl. Probab.* **14** 555–611. [MR2052895](#)
- [9] GROMOLL, H. C., KRUK, Ł. and PUHA, A. L. (2011). Diffusion limits for shortest remaining processing time queues. *Stoch. Syst.* **1** 1–16. [MR2948916](#)
- [10] HARRISON, J. M. and WILLIAMS, R. J. (1996). A multiclass closed queueing network with unconventional heavy traffic behavior. *Ann. Appl. Probab.* **6** 1–47. [MR1389830](#)
- [11] IGLEHART, D. L. and WHITT, W. (1970). Multiple channel queues in heavy traffic. I. *Adv. in Appl. Probab.* **2** 150–177. [MR0266331](#)
- [12] LIMIC, V. (2001). A LIFO queue in heavy traffic. *Ann. Appl. Probab.* **11** 301–331. [MR1843048](#)
- [13] LIN, M., WIERMAN, A. and ZWART, B. (2011). The heavy-traffic analysis of mean response time under shortest remaining processing time. *Performance Evaluation* **68** 955–966.
- [14] NÚÑEZ-QUEIJA, R. (2002). Queues with equally heavy sojourn time and service requirement distributions. *Ann. Oper. Res.* **113** 101–117. [MR1960684](#)
- [15] NUYENS, M. and ZWART, B. (2006). A large-deviations analysis of the  $GI/GI/1$  SRPT queue. *Queueing Syst.* **54** 85–97. [MR2268054](#)
- [16] PAVLOV, A. V. (1984). A system with Schrage servicing discipline in the case of a high load. *Engng. Cybernetics* **21** 114–121; translated from *Izv. Akad. Nauk SSSR Tekhn. Kibernet.* **6** (1983) 59–66 (Russian).

- [17] PECHINKIN, A. V. (1986). Heavy traffic in a system with a discipline of priority servicing for the job of shortest remaining length with interruption. *Mat. Issled.* **89** 85–93. [MR0836668](#)
- [18] PERERA, R. (1993). The variance of delay time in queueing system  $M/G/1$  with optimal strategy SRPT. *Archiv für Elektronik und Uebertragungstechnik* **47** 110–114.
- [19] PROHOROV, YU. V. (1956). Convergence of random processes and limit theorems in probability theory. *Theory Probab. Appl.* **1** 157–214.
- [20] SCHASSBERGER, R. (1990). The steady-state appearance of the  $M/G/1$  queue under the discipline of shortest remaining processing time. *Adv. in Appl. Probab.* **22** 456–479. [MR1053240](#)
- [21] SCHRAGE, L. E. (1968). A proof of the optimality of the shortest remaining processing time discipline. *Oper. Res.* **16** 687–690.
- [22] SCHRAGE, L. E. and MILLER, L. W. (1966). The queue  $M/G/1$  with the shortest remaining processing time discipline. *Oper. Res.* **14** 670–684. [MR0195173](#)
- [23] SCHREIBER, F. (1993). Properties and applications of the optimal queueing strategy SRPT: A survey. *Archiv für Elektronik und Übertragungstechnik* **47** 372–378.
- [24] SMITH, D. R. (1978). A new proof of the optimality of the shortest remaining processing time discipline. *Oper. Res.* **26** 197–199. [MR0471112](#)
- [25] WHITT, W. (1971). Weak convergence theorems for priority queues: Preemptive-resume discipline. *J. Appl. Probab.* **8** 74–94. [MR0307389](#)
- [26] WIERMAN, A. and HARCHOL-BALTER, M. (2003). Classifying scheduling policies with respect to unfairness in an  $M/GI/1$ . In *Proceedings of the 2003 ACM SIGMETRICS International Conference on Measurement and Modeling of Computer Systems* 238–249. ACM, New York.

DEPARTMENT OF MATHEMATICS  
CALIFORNIA STATE UNIVERSITY SAN MARCOS  
333 S. TWIN OAKS VALLEY ROAD  
SAN MARCOS, CALIFORNIA 92096-0001  
USA  
E-MAIL: [apuha@csusm.edu](mailto:apuha@csusm.edu)