



A Poisson limit for the departure process from a queue with many busy servers



Ward Whitt

Department of Industrial Engineering and Operations Research, Columbia University, New York, NY 10027-6699, USA

ARTICLE INFO

Article history:

Received 20 October 2015

Received in revised form

15 May 2016

Accepted 18 June 2016

Available online 6 July 2016

Keywords:

Poisson approximations

Departure processes

Output processes

Nonhomogeneous Poisson processes

Queueing networks

Many-server heavy-traffic limits for queues

ABSTRACT

We establish a limit theorem supporting a Poisson approximation for the departure process from a multi-server queue that tends to have many busy servers. This limit can support approximating a flow out of such a queue in a complex queueing network by an independent Poisson source. The main ideas are: (i) to scale time so that previous many-server heavy-traffic limits can be applied and (ii) for time-varying arrival-rate functions, to scale (spread out) time by a large factor about each fixed time.

© 2016 Elsevier B.V. All rights reserved.

1. Introduction

Complex queueing systems are typically networks of queues, with arrival processes at individual queues being composed of departures and overflows from other queues, with the service-time cumulative distribution functions (cdf's) often being not nearly exponential. Thus an arrival process at an internal queue usually cannot be assumed to be exactly a Poisson process; e.g., see [3]. Nevertheless, a Poisson approximation may be reasonable.

Example 1.1 (*Final Checkout in Online Shopping*). Suppose that we want to develop a stochastic arrival process model for the final checkout in a complex online shopping system. Many separate people shop online until they are ready for final checkout. To illustrate, we model the checkout as the second queue in a two-queue $G_t/GI/\infty \rightarrow \cdot/GI/1$ network, in which the first queue is an infinite-server (IS) model with a general arrival process having a time-varying arrival-rate function $\lambda(t)$, which is independent of service times that are independent and identically distributed (i.i.d.) with a general cdf F having a continuous probability density function (pdf) f with $F(t) = \int_0^t f(s) ds$, $t \geq 0$. The output of the IS queue is the arrival process to a final single-server (SS) checkout queue, with general service cdf, unlimited waiting room and

service in order of arrival. The exact form of the departure-rate function from the IS queue is

$$\delta(t) = \int_0^\infty f(y)\lambda(t-y) ds, \quad (1)$$

as given in Theorem 1 of [4]; it is the same for G_t as for M_t ; see §5 of [9]. In this setting we provide support for approximating the final SS queue by an $M_t/GI/1$ queue, where the arrival process is a nonhomogeneous Poisson process (NHPP) with arrival-rate function $\delta(t)$ in (1). An efficient algorithm to calculate performance measures when $\lambda(t)$ is periodic is given in [16].

For a concrete simulation, consider the stationary $GI/GI/\infty \rightarrow \cdot/GI/1$ model in which all service times are i.i.d. and the external arrival process is a renewal process. To introduce extra variability, we assume that all three GI components have the hyperexponential cdf (H_2 , mixture of two exponentials) with squared coefficient of variation (scv, variance divided by the square of the mean) $c^2 = 4$ and balanced means as in p. 137 of [21]; that leaves only the mean or its reciprocal, the rate, to be specified. We let the arrival rate be $\lambda = 100$ and the service rates at the two queues be $\mu_1 = 1$ and $\mu_2 = 200$. By Little's law, these rates make the mean steady-state number of busy servers in the IS queue be 100, which we regard as moderately large scale. In actual online checkout, the mean number of busy shoppers is likely to be much larger, and the difference between the two service rates is likely to be even greater.

In this context, we suggest that the performance at the final SS queue can be approximated by the $M/H_2/1$ model, for which

E-mail address: ww2040@columbia.edu.

the mean steady-state waiting time before starting service has the Pollaczek–Khintchine (PK) formula $EW = \rho\mu_2^{-1}(1 + c^2)/2(1 - \rho) = 0.0125$ for $\rho = 0.50$, $\mu_2 = 200$ and $c^2 = 4$. The intuition is that, with many busy servers, the departure process from the IS queue is much like the superposition of i.i.d. renewal processes, one for each server, for which the limit is Poisson, as discussed in §9.8 of [23]. Of course, the servers do not remain busy all the time and the number of busy servers is random, varying over time, so that representation is only approximate. Thus, there remains something to prove for departure processes.

A simulation experiment was conducted for this example. It shows that the interarrival-time cdf at the second queue is approximately exponential with mean 0.01 and that the estimated mean wait EW is only 8% above the PK formula for M arrivals; see the appendix for more details.

We conclude this example by mentioning that part of the justification for the $M/H_2/1$ approximation with a Poisson arrival process for the SS queue is the relatively low traffic intensity at the SS queue, because the departure process from the $H_2/H_2/\infty$ IS queue with many busy servers is only approximately Poisson over a short time scale. For example, the central limit theorem for the departure process will not have the same variability parameter as for a Poisson process. As discussed in §9.8 of [23], there is different variability at different time scales. As $\rho \uparrow 1$, the ratio of the actual mean $EW(\rho)$ to the mean with Poisson arrivals increases. We found that the $M/H_2/1$ approximation for the mean EW was 27% low when the service rate at the second queue was decreased so that $\rho_2 = 0.90$. See [20] for a related superposition process example. ■

In [22] we previously established a limit theorem supporting the Poisson approximation for the departure process in the simulated example; our purpose here is to extend the result to a larger class of models. First, for infinite-server models, we extend the result established for the $GI/Ph/\infty$ model in [24] to the $G_t/GI/\infty$ model, having a general service-time distribution (the GI) instead of Ph and from a renewal arrival process (GI) to general (allowing non-renewal) arrival process with a time-varying rate (the G_t). The proof is similar, except now we apply the two-parameter MSHT FWLLN for the $G_t/GI/\infty$ model reviewed in [18] instead of the single-parameter FWLLN for the $GI/Ph/\infty$ model in [24].

We are also interested in establishing a result that applies to models with finitely many servers, perhaps including customer abandonment and feedback. A concrete example of a closed network of two $\cdot/GI/s$ queues which could be used in this way is contained in [12]. In that model there is one SS station with state-dependent service rate and one IS station. In the same spirit, our approach provides the basis for an alternate proof of a Poisson limit for a queue with delayed feedback (which can be regarded as a $\cdot/GI/\infty$ IS queue) in [19]; they established the Poisson limit using a coupling technique.

The Poisson limit in [22] was established using martingale methods. The “martingale method” means that we focus on the stochastic departure rate or intensity of the departure process and its integral, called the compensator, which depends on a specification of the history or filtration; see [2,17] for introductions and [5,8] for advanced accounts. We will establish the Poisson limit, independent of the history of the queueing system, by showing that the compensators approach a deterministic limit; e.g., see Theorem VIII.4.10 in [8] and Problem 1 on p. 360 of [5].

We have special interest in many-server queues with time-varying arrival-rate functions. To obtain useful Poisson limits for those models, we will introduce a new scaling method, spreading out time about a fixed reference time. The Poisson limit then provides support for approximating the departure process by an

NHPP. For the required MSHT FWLLN's in $G_t/GI/\infty$ and $G_t/GI/s_t + GI$ models with general nonstationary arrival processes, we can apply [11,18,10,15], respectively. These limits exploit a random-measure or two-parameter framework. We present our results with minimum technicalities; we refer to those papers for the details.

In Section 2 we review the MSHT FWLLN in a $G_t/GI/\infty$ model and establish the required FWLLN for the departure rate process in Theorem 2.1. In Section 3 we establish the main result, Theorem 3.1, which provides general conditions for the desired Poisson limit in terms of associated MSHT limits. We present additional supplementary material on the simulation for Example 1.1 and a direct NHPP approximation for the departure process in an appendix, which is available from the author's website (<http://www.columbia.edu/~ww2040/allpapers.html>).

2. Review of the MSHT FWLLN for $G_t/GI/\infty$ queues

We start by reviewing the MSHT FWLLN in Theorem 3.1 in [18], because we will use established properties as conditions in our new theorem for other models.

Let \Rightarrow denote convergence in distribution and let $D \equiv D(I, \mathbb{R})$ be the usual Skorohod space of right-continuous real-valued functions with left limits on a subinterval I of the entire real line \mathbb{R} , possibly \mathbb{R} itself [5,8,23]. In our setting with a continuous limits, convergence in the Skorohod J_1 topology is equivalent to uniform convergence over bounded subintervals of I .

We consider a sequence of queueing models indexed by n . Let the arrival process have a well-defined arrival rate for each n ; i.e., let $A_n(t_1, t_2)$ be the number of arrivals in model n in the time interval $(t_1, t_2]$ and assume that

$$E[A_n(t_1, t_2)] = n\Lambda(t_1, t_2), \quad \text{where } \Lambda(t_1, t_2) \equiv \int_{t_1}^{t_2} \lambda(s) ds \quad (2)$$

for $-\infty < t_1 < t_2 < +\infty$, with \equiv denoting equality by definition. This can be achieved by scaling (accelerating) time in a fixed arrival process. Thus, the arrival rate in model n is

$$\lambda_n(t) = n\lambda(t), \quad -\infty < t < +\infty. \quad (3)$$

As a regularity condition, we also assume that $0 \leq \lambda(t) \leq \lambda_U < \infty$. We furthermore assume that the system starts empty at time $-t_0 \leq 0$. That avoids having to carefully treat the initial conditions, but for a way to do so, see [1]. Let $\bar{A}_n(t_1, t_2) \equiv n^{-1}A_n(t_1, t_2)$. We assume a FWLLN is valid for the arrival processes; i.e.,

$$\sup_{t_L \leq t_1 < t_2 \leq t_U} |\bar{A}_n(t_1, t_2) - \Lambda(t_1, t_2)| \Rightarrow 0 \quad \text{as } n \rightarrow \infty$$

for all t_L and t_U with $-\infty < -t_0 \leq t_L < t_U < \infty$ (weak convergence uniformly over bounded intervals).

Assumption 1 of [18] allows a general sequence of arrival processes, but they are required to satisfy a functional central limit theorem (FCLT) because the primary concern was establishing the MSHT FCLT. That FCLT condition can be weakened to having only a FWLLN, because Theorem 3.1 only requires the MSHT FWLLN conclusion. The proof of the FWLLN for the number of busy servers under the weaker FWLLN condition is not discussed in [18], but it is discussed in [17]; see Theorem 3.6 and §§3.4, 4.3, 5.2, 6.1 and 6.2.

Assumption 2 of [18] stipulates that the service times come from a single i.i.d. sequence, independent of n and the arrival processes, distributed as a random variable S having a general cdf F . In addition, we require that the cdf F have a continuous pdf f in terms of which we can write $F(t) = \int_0^t f(s) ds$, $t \geq 0$, for $F^c(t) \equiv 1 - F(t)$, and a failure-rate function $h(t) \equiv f(t)/F^c(t)$ that is bounded over finite intervals. In [18] the system starts empty at time 0. Without loss of generality, we assume that the system

starts empty at time $-t_0 < 0$. We then can let $t_0 \rightarrow \infty$ to obtain the simple approximation formula in (1).

Let $N_n^e(t, y)$ be the number of customers in service at time t in model n that have been so for at most time y . Let \bar{N}_n^e be the FWLLN-scaled version $\bar{N}_n^e(t, y) \equiv n^{-1}N_n^e(t, y)$. A variant of (3.5) and (3.7) of Theorem 3.1 of [18] then implies that

$$\sup_{t_L \leq t \leq t_U, y_L \leq y \leq y_U} |\bar{N}_n^e(t, y) - N^e(t, y)| \Rightarrow 0 \quad \text{as } n \rightarrow \infty \quad (4)$$

for all t_L and t_U with $-\infty < -t_0 \leq t_L < t_U < \infty$ and for all y_L and y_U with $-\infty < y_L < y_U < \infty$ (again weak convergence uniformly over bounded intervals), where

$$N^e(t, y) \equiv \int_0^y F^c(s)\lambda(t-s) ds. \quad (5)$$

Let $D_n(t) \equiv A_n(t) - N_n^e(t, t + t_0)$ be the associated departure counting process in model n and let $\bar{D}_n(t) \equiv n^{-1}D_n(t)$ be the fluid-scaled version. Along with (4), we also have the limit

$$\bar{D}_n \Rightarrow D \quad \text{in } D[-t_0, \infty) \text{ as } n \rightarrow \infty, \quad (6)$$

where

$$D(t) \equiv \Lambda(t) - N^e(t, t + t_0) = \int_0^{t+t_0} F(s)\lambda(t-s) ds, \quad t \geq -t_0. \quad (7)$$

For the new part, let $\Delta_n(t)$ be the stochastic departure rate at time t in model n . The departure rate can be expressed as a stochastic integral (which is just a random sum) via

$$\Delta_n(t) = \int_0^{t+t_0} h(y) dy N_n^e(t-, y) dy, \quad t \geq -t_0. \quad (8)$$

As in (2.1) of [22], we use the left limit t -in (8) to make $\Delta_n(t)$ be the predictable stochastic intensity with respect to the appropriate history that includes the ages of all the customers in service and the history of the arrival process at each time t ; see §1.3 of [2] and [17]. That can be understood and justified by a discretization argument, dividing the interval $[-t_0, t]$ into k subintervals, doing a discrete-time analysis and then letting $k \rightarrow \infty$. A detailed proof is given in §5.2 of [11]; see Lemma 5.4.

To elaborate, $\Delta_n(t)$ being a stochastic intensity means that the centered process $D_n(t) - C_n(t)$ is a martingale with compensator

$$C_n(t) = \int_{-t_0}^t \Delta_n(s) ds, \quad t \geq -t_0, \quad (9)$$

again with respect to the full system history at time t .

Let $\bar{\Delta}_n \equiv n^{-1}\Delta_n$ for in (8) be the FWLLN-scaled departure rate process. We first establish a bound on the expectations.

Lemma 2.1 (Expectation Bound). *Under the assumptions above for the sequence of $G_t/GI/\infty$ models,*

$$E[\bar{\Delta}_n(t)] \leq K \max\{1, t + t_0\} \sup_{0 \leq s \leq t+t_0} \{h(s)\} < \infty \quad (10)$$

for all n and t .

Proof. Since $N_n(t) \equiv N_n^e(t, \infty) \leq A_n(-t_0, t)$ we can apply (2). Since the failure rate function h is bounded over bounded intervals, we can replace it by a constant outside the integral. ■

Theorem 2.1 (MSHT Limit For the Departure Rate). *For the $G_t/GI/\infty$ model under the assumptions above,*

$$\bar{\Delta}_n \equiv n^{-1}\Delta_n \Rightarrow \delta \quad \text{in } D([-t_0, \infty), \mathbb{R}) \text{ as } n \rightarrow \infty, \quad (11)$$

where

$$\delta(t) \equiv \int_0^{t+t_0} h(y) dy N(t, y), \quad t \geq -t_0, \quad (12)$$

so that

$$\delta(t) = \int_0^{t+t_0} f(y)\lambda(t-y) dy \quad \text{and} \quad D(t) = \int_{-t_0}^t \delta(s) ds, \quad t \geq -t_0. \quad (13)$$

Proof. We first apply Lemma 2.1 to get bounded expectations. Then we apply the Skorohod representation theorem, Theorem 3.2.2 of [23], to reduce the argument to a deterministic one, but use the same notation. We establish the desired uniform convergence over bounded intervals by showing, for any t in a bounded interval and any sequence $\{t_n\}$ with $t_n \rightarrow t$ as $n \rightarrow \infty$, that $n^{-1}\Delta_n(t_n) \rightarrow \delta(t)$ as $n \rightarrow \infty$. To do that, we exploit the fact that the convergence in (4) corresponds to the weak convergence of finite measures, where we regard $\bar{N}_n^e(t, y)$ as a function of y as a cdf. Hence, we can show, for each $t \geq -t_0$ that we have the associated convergence of the integrals

$$n^{-1}\Delta_n(t_n) = \int_0^{t_n+t_0} h(y) dy \bar{N}_n^e(t_n, y) \rightarrow \int_0^{t+t_0} h(y) F^c(y)\lambda(t-y) dy \quad \text{as } n \rightarrow \infty.$$

We use the fact that h is continuous and bounded on the interval $[0, t + t_0]$. The limiting integral simplifies, yielding

$$\int_0^{t+t_0} h(y) F^c(y)\lambda(t-y) dy = \int_0^{t+t_0} f(y)\lambda(t-y) dy$$

by the simple relation $h(y)F^c(y) = f(y)$. That convergence implies that $\bar{\Delta}_n \rightarrow \delta$ in $D(\mathbb{R}, \mathbb{R})$ as $n \rightarrow \infty$, which implies the weak convergence for the original processes. ■

Remark 2.1 (Starting Empty in the Distant Past). In many papers on IS queues, the system is assumed to start empty in the distant past (at $-\infty$). That is tantamount to letting $t_0 \rightarrow \infty$. As $t_0 \rightarrow \infty$, $\delta(t)$ in Theorem 2.1 approaches (1), the departure rate $E[\lambda(t - S)]$ in the $M_t/GI/\infty$ model in equation (4) of Theorem 1 in [4] and in the associated $G_t/GI/\infty$ fluid model; see §4 of [14].

3. The supporting limit for a Poisson approximation

We now establish the Poisson limit for the departure process from a general $G_t/GI/\infty$ model. At the same time, we provide a framework for treating many other models. To do so, we assume some of the conclusions deduced for the $G_t/GI/\infty$ model is Section 2 rather than specify the detailed model. Thus, we now consider a more general multi-server queue. As before, we assume that the servers work independently in parallel having an individual remaining service-time failure rate function h . However, the queue may be in the middle of a complex network and there may be customer abandonment and feedback.

As in Section 2, we consider a sequence of models indexed by n in a MSHT framework. That typically means that the arrival rate is allowed to grow without bound as in (2) and if there are finitely many servers, then that number is allowed to grow as well. We directly assume that the processes $N_n^e(t, y)$, $D_n(t)$, $C_n(t)$ and $\Delta_n(t)$ are well defined with the same meaning as in Section 2, but we do not fully specify the system; e.g., we do not specify the arrival process. We directly assume that the stochastic departure rate can be defined by the stochastic integral in (8) and that $D_n(t) - C_n(t)$

is a martingale with respect to the system history up to time t , where $C_n(t)$ is the compensator and is the integral of $\Delta_n(t)$ as in (9). We also assume that the limits in (4) and (8) hold, but without assuming the explicit form of the limits $N^e(t, y)$ and $D(t)$ in (5) and (7). Finally, we assume that the bound in (10) holds. Under these assumptions, we also have the conclusions of Theorem 2.1 with the limit in (12), but without the explicit limit in (13), because the same proof applies. For example, these assumptions apply to the $G_t/GI/s + GI$ model with finitely many servers and customer abandonment, for which a FWLLN was established in [14,13].

Paralleling [22], we will do an additional slow-time scaling in order to establish the supporting Poisson limit. However, in order to capture the time-varying arrival rate appropriately, instead of simply undoing the MSHT scaling in (2), we do the time scaling about an arbitrary time t , which we regard as fixed.

For this purpose, we introduce two-parameter processes

$$\begin{aligned} D_n(t, u_2) - D_n(t, u_1) &\equiv D_n(t + u_2/n) - D_n(t + u_1/n), \\ C_n(t, u_2) - C_n(t, u_1) &\equiv C_n(t + u_2/n) - C_n(t + u_1/n), \\ \Delta_n(t, u) &\equiv \Delta_n(t + u/n)/n, \\ -\infty < u_1 < u_2 < +\infty. \end{aligned} \tag{14}$$

Note that the definitions for $C_n(t, u)$ and $\Delta_n(t, u)$ follow from the definition for $D_n(t, u)$. With these definitions and the assumptions above,

$$\begin{aligned} C_n(t, u_2) - C_n(t, u_1) &= \int_{u_1}^{u_2} \Delta_n(t, v) dv, \\ -\infty < u_1 < u_2 < +\infty, \end{aligned} \tag{15}$$

$\{D_n(t, s) - C_n(t, s) : s \geq u_1\}$ is a martingale and $\Delta_n(t, u)$ is a predictable stochastic intensity with respect to the system history.

With this preparation, we are able to establish our desired result. In our setting, weak convergence of the processes with nondecreasing sample paths to a Poisson process in $D(I, \mathbb{R})$ is equivalent to convergence of all finite-dimensional distributions; see VI.3.37 of [8].

Theorem 3.1 (Poisson Limit). *Under the assumptions in this section above,*

$$D_n(t, \cdot) \Rightarrow \Pi_{\delta(t)}(\cdot) \text{ in } D(\mathbb{R}, \mathbb{R}) \text{ as } n \rightarrow \infty, \tag{16}$$

where Π_c is a homogeneous Poisson process with constant rate c and $\delta(t)$ is the limit in (12); i.e., for any integer k , any k -tuple of disjoint subintervals $((u_{i,1}, u_{i,2}] : 1 \leq i \leq k)$ and any k -tuple of nonnegative integers $(j_i : 1 \leq i \leq k)$,

$$P(D_n(t, u_{i,2}) - D_n(t, u_{i,1}) = j_i : 1 \leq i \leq k)$$

$$\rightarrow \prod_{i=1}^k \frac{e^{-\mu_i(t)} \mu_i(t)^{j_i}}{j_i!}$$

as $n \rightarrow \infty$, where $\mu_i(t) \equiv \delta(t)(u_{i,2} - u_{i,1})$.

Proof. The proof is similar to the proof of Theorem 2 in [22]. The limit in (11) implies that

$$\sup_{u_L < u < u_U} |n^{-1} \Delta_n(t + (u/n)) - \delta(t)| \Rightarrow 0 \text{ as } n \rightarrow \infty$$

for all u_L and u_U with $-\infty < u_L < u_U < +\infty$. Then, paralleling the proof of Theorem 2 in [22], we write

$$\begin{aligned} C_n(t + (u_2/n)) - C_n(t + (u_1/n)) &= \int_{u_1/n}^{u_2/n} \Delta_n(t + v) dv \\ &= \int_{u_1}^{u_2} n^{-1} \Delta_n(t + v/n) dv \\ &\Rightarrow \int_{u_1}^{u_2} \delta(t) dv = \delta(t)(u_2 - u_1) \text{ as } n \rightarrow \infty. \end{aligned} \tag{17}$$

Combining (17) with (14), we have the analog of Corollary 2 of [22], i.e.,

$$C_n(t, u_2) - C_n(t, u_1) \Rightarrow \delta(t)(u_2 - u_1) \text{ as } n \rightarrow \infty.$$

That implies that the limit (16) holds, as claimed, by Theorem VIII.4.10 of [8]. ■

Remark 3.1 (Supporting an NHPP Approximation). The statement of Theorem 3.1 may seem a bit paradoxical, because it states that the departure process is asymptotically a homogeneous Poisson process but with the time-varying rate $\delta(t)$ in (12). That dichotomy arises because of our scaling about the fixed time t . For applications, we interpret the limit as supporting an NHPP approximation with time-varying rate $\delta(t)$.

Remark 3.2 (The Stationary Case). For a stable stationary model without abandonment, the rate out equals the rate in, so that the departure rate must equal the constant arrival rate. Consistent with that basic property, we see that $\delta(t) = \lambda$ for all t if the arrival process has a constant arrival rate λ .

Remark 3.3 (Models With Finitely Many Servers). For the stationary $GI/M/s$ and the $M/M/s + M$ models, the papers [7,6] can be applied to establish analogs of Theorem 2.1. For the quality-and-efficiency-driven (QED) and efficiency-driven (ED) MSHT regimes, $\delta(t) = \mu s$ for all t . The FWLLN follows immediately from the MSHT FCLTs established in those papers. These result can be extended to general arrival processes using §7.3 of [17]. Extensions to the $G/G/s$ and $G/GI/s + GI$ follow from [10,11].

We can also apply [15] to obtain the analog of Theorem 2.1 for the $G_t/M/s_t + GI$ Model with customer abandonment, which alternates between overloaded intervals and underloaded intervals. With exponential service times, it suffices to look at $N(t)$, the number of customers in service at each time, instead of the more complicated two-parameter process $N^e(t, y)$. The departure rate at time t is simply $\mu \min\{X(t), s(t)\}$, where μ is the fixed service rate, $X(t)$ is the number of customers in the system and $s(t)$ is the number of servers at time t . The FWLLN is given for overloaded intervals in (4.2) of Theorem 4.1 and §3 of [15]; then $\delta(t) = s(t)\mu$. The FWLLN is given for underloaded intervals in (5.1) and (5.2) of Theorem 5.1 of [15]; except for the initial conditions, $\delta(t)$ is the same as in an IS system. Extensions to GI service follow from [13].

Acknowledgments

The author thanks Vahid Sarhangian for suggesting that it would be good to extend [22], Guodong Pang and an anonymous referee for helpful comments, Jingtong Zhao for conducting the supporting simulation, and NSF for research support (CMMI 1265070).

Appendix A. Supplementary data

Supplementary material related to this article can be found online at <http://dx.doi.org/10.1016/j.orl.2016.06.006>.

References

- [1] K. Aras, Y. Liu, W. Whitt, Heavy-traffic limit for the initial content process, Columbia University, 2014. <http://www.columbia.edu/~ww2040/allpapers.html>.
- [2] P. Bremaud, Point Processes and Queues: Martingale Dynamics, Springer, New York, 1981.
- [3] R.L. Disney, D. Konig, Queueing networks: a survey of their random processes, SIAM Rev. 27 (3) (1985) 335–403.
- [4] S.G. Eick, W.A. Massey, W. Whitt, The physics of the $M_t/G/\infty$ queue, Oper. Res. 41 (1993) 731–742.

- [5] S.N. Ethier, T.G. Kurtz, *Markov Processes: Characterization and Convergence*, Wiley, New York, 1986.
- [6] O. Garnett, A. Mandelbaum, M.I. Reiman, Designing a call center with impatient customers, *Manuf. Serv. Oper. Manage.* 4 (3) (2002) 208–227.
- [7] S. Halfin, W. Whitt, Heavy-traffic limits for queues with many exponential servers, *Oper. Res.* 29 (3) (1981) 567–588.
- [8] J. Jacod, A.N. Shiryaev, *Limit Theorems for Stochastic Processes*, Springer, New York, 1987.
- [9] O.B. Jennings, A. Mandelbaum, W.A. Massey, W. Whitt, Server staffing to meet time-varying demand, *Manage. Sci.* 42 (1996) 1383–1394.
- [10] W. Kang, K. Ramanan, Law of large number limit for many-server queues, *Ann. Appl. Probab.* 21 (1) (2010) 33–114.
- [11] H. Kaspi, K. Ramanan, Law of large number limit for many-server queues, *Ann. Appl. Probab.* 20 (6) (2011) 2204–2260.
- [12] E.V. Krichagina, A.A. Puhalskii, A heavy-traffic analysis of a closed queueing system with a GI/∞ service center, *Queueing Syst.* 25 (1997) 235–280.
- [13] Y. Liu, W. Whitt, A many-server fluid limit for the $G_t/GI/s_t + GI$ queueing model experiencing periods of overloading, *Oper. Res. Lett.* 40 (2012) 307–312.
- [14] Y. Liu, W. Whitt, The $G_t/GI/s_t + GI$ many-server fluid queue, *Queueing Syst.* 71 (2012) 405–444.
- [15] Y. Liu, W. Whitt, Many-server heavy-traffic limits for queues with time-varying parameters, *Ann. Appl. Probab.* 24 (1) (2014) 378–421.
- [16] N. Ma, W. Whitt, A performance algorithm for periodic queues, Columbia University, Working paper, 2016.
- [17] G. Pang, R. Talreja, W. Whitt, Martingale proofs of many-server heavy-traffic limits for Markovian queues, *Probab. Surv.* 4 (2007) 193–267.
- [18] G. Pang, W. Whitt, Two-parameter heavy-traffic limits for infinite-server queues, *Queueing Syst.* 65 (2010) 325–364.
- [19] E.A. Pekoz, N. Joglekar, Poisson traffic flow in a general feedback queue, *J. Appl. Probab.* 39 (2002) 630–636.
- [20] K. Sriram, W. Whitt, Characterizing superposition arrival processes in packet multiplexers for voice and data, *IEEE J. Sel. Areas Commun. SAC* 4 (6) (1986) 833–846.
- [21] W. Whitt, Approximating a point process by a renewal process: two basic methods, *Oper. Res.* 30 (1982) 125–147.
- [22] W. Whitt, Departures from a queue with many busy servers, *Math. Oper. Res.* 9 (4) (1984) 534–544.
- [23] W. Whitt, *Stochastic-Process Limits*, Springer, New York, 2002.
- [24] W. Whitt, On the heavy-traffic limit theorem for $GI/G/\infty$ queue, *Adv. Appl. Probab.* 14 (1) (1982) 171–190.