

QUEUES WITH MANY SERVERS AND IMPATIENT CUSTOMERS

AVISHAI MANDELBAUM AND PETAR MOMČILOVIĆ

ABSTRACT. The many-server queue with abandonments, G/GI/N+GI, is considered in the Quality- and Efficiency-Driven (QED) regime. Here the number of servers N and the offered-load R are related via the square-root rule $N = R + \beta\sqrt{R} + o(\sqrt{R})$, for some constant β , as the number of servers increases indefinitely. QED performance entails short waiting times and scarce abandonments (high Quality) *jointly* with high servers' utilization (high Efficiency), which is feasible when many servers cater to a single queue. For the G/GI/N+GI queue, we derive diffusion approximations for both its queue-length and virtual-waiting-time processes. Special cases, for which closed-form analysis is provided, are the G/M/N+GI and G/D/N+GI queues, thus expanding and generalizing existing results.

1. INTRODUCTION

The Quality-and-Efficiency Driven (QED) regime achieves, jointly, high levels of system's *efficiency*, as manifested by servers' high utilization, and service *quality*, namely customer's short waiting times and hence scarce abandonments. QED performance is achievable in carefully-balanced queueing systems with many servers – indeed, with few servers, efficiency and quality must be traded off against each other. Within the G/GI/N framework, or more precisely G/GI/N+GI, the QED regime with abandonments is characterized by the relation $N = R + \beta\sqrt{R} + o(\sqrt{R})$, for some scalar β ; here N is number of servers and R is the offered load, namely the arrival rate multiplied by average service time. (This *square-root staffing* relation also characterizes the QED regime without abandonments, but β must be then taken positive to ensure stability.)

Relevance. Recent interest in multi-server queues with impatient customers is due to their applicability in modeling medium-to-large-scale customer call/contact centers. In such service operations, abandonments arise naturally and, in fact, *must* be accounted for in models (see Section 2 in [12] for an elaboration). Additionally, well-run call centers are QED [11] or some relatives of it (e.g. ED+QED, as in [24]). But QED queues also arise beyond call centers. To wit, waiting time in QED call centers is naturally measured in seconds while service times in minutes. This one-order time-reduction (from minutes to seconds in the case of call centers) is a QED characteristic; indeed, it also arises in transportation (searching for parking takes minutes while parking time is hours) and in healthcare (sojourn times in emergency departments take hours while hospitalization is days). Significantly, the abandonment phenomenon is relevant in all these examples, which is perhaps surprising for the latter: yet, a non-negligible fraction of patients leave emergency departments without being seen by a doctor [13].

Related research. Although the QED regime (without abandonments) can be traced back to Erlang [9] and Jagerman [16], the regime was first formalized by Halfin and Whitt in [14]; for recent result on the QED regime see [21, 27, 19, 26] and references therein, with Reed's

2000 *Mathematics Subject Classification.* Primary: 60K25; Secondary: 90B22.

Key words and phrases. Multi-server queue, abandonment, QED regime, Halfin-Whitt regime, diffusion approximation.

A. Mandelbaum's research was supported in part by BSF (Binational Science Foundation) Grants 2005175/2008480, ISF (Israeli Science Foundation) Grant 1357/08, and by the Technion funds for the promotion of research and sponsored research.

framework [27] for the G/GI/N queue being especially relevant. The M/M/N+M (Erlang-A) system in the QED regime (with abandonments) was considered in [12]. Extensions to the model with generally distributed abandonments can be found in [34, 33, 24, 29]. The M/M/N+G system in the Efficiency-Driven regime was analyzed in [32]; for a summary of performance measures of this system see [23]. Recently, fluid limits of many-server queues with abandonment were considered in [18]. Independently of our work, many-server queues with customer abandonment were investigated in [7], where the focus parallels our Lemma 9 (the authors establish a relation between the abandonment-count and queue-length processes). The literature on queues with abandonments is extensive and includes models with various features; we refer the reader to the discussions in [12, 33].

Contributions. We consider a G/GI/N+GI system in the QED regime. The limiting scaled number-of-customers-in-system process is described in terms of a non-linear operator (Corollary 1); a corresponding result for the limiting scaled waiting-time processes (virtual and offered) is obtained as well (Corollary 2). In the case when there is no abandonment, our operator coincides with the one earlier obtained in [27].

The proofs of our main results are based on two relations: (i) between corresponding systems with and without abandonment (Proposition 1), and (ii) between the queue-length and offered-waiting-time processes (Lemma 9). The first relation, in conjunction with results from [27], is used to obtain upper bounds for the queue-length and waiting-time processes; the bounds are tight enough to yield the proper orders of magnitude for QED convergence. This enables us to approximate the abandonment process, properly centered and scaled (see (14) and Lemma 5). The approximation in (14), in turn, reveals the QED limit of the abandonment process, which paves the way to other QED limits. The second relation is central in our paper as it allows one to circumvent the complex relation between the queue-length and the abandonment processes, which is necessary for obtaining limits of queue lengths. Indeed, the relation is complex since abandonments are determined by (offered) waiting times which, in turn, depend on queue length via a “first-passage-time” operator, as in Puhalskii [25] and Talreja and Whitt [30]. Still, one can not directly use [25, 30] to deduce limits of waiting times from those of queue lengths since we could not analyze the latter in isolation. This is in contrast to [27], which first derives limits of queue-length, then uses [25] to deduce directly limits of waiting times.

To summarize, Lemma 9 shows that queue-length and waiting times are asymptotically “close”. Therefore, via Lemma 5 and (14), we express the abandonment process in terms of queue length. This results in representing the queue length in terms of itself, which is resolved through a sample-path mapping, as introduced in (25).

The above provides a justification for the need to develop techniques and tools that are not required in Reed [27], who analyzed the QED G/GI/N queue. More specifically, for the QED G/GI/N+GI queue, queue length and waiting times must be analyzed jointly, as already discussed. Then, the mapping (25) that characterizes the limiting queue-length processes generalizes the one in [27]. (A mapping corresponding to waiting times is introduced in (34).) Our general result can be made explicit (see Section 5) in two special cases that have not yet been analyzed: Example 2 provides functional limit theorems that correspond and further generalize [33, 34] (which considered M/M/N+GI in steady state); Example 3 adds abandonments to [21] (which considered process limits for G/D/N). Finally, our results demonstrate that the role of the patience distribution in the QED regime is captured merely by the value of its density at the origin. This is practically important since patience data is censored (only lower bounds for patience are available for served customers), and possibly highly censored (e.g., 3% abandoning). We thus suggest an estimator for the patience density at the origin, based on a transient analogue of the steady-state relation between the probability of abandonment and the expected waiting time [22].

Contents. The paper is organized as follows. In the next section we describe the model and introduce the QED regime. Section 3 contains preliminaries; in particular, we discuss a relationship between the systems with and without abandonment, the infinite service-process associated with the arrival and service processes, processes associated with initial conditions and abandonment, the queue-length process, and a relationship between the queue length and offered waiting time. The main results of the paper are presented in Section 4. Examples are discussed in Section 5, future research is outlined in Section 6, and some proofs are provided in Section 7.

Notation. Denote by $D[0, \infty)$ the space of all real-valued functions on $[0, \infty)$ that are right-continuous with left limits (r.c.l.l.), endowed with the standard Skorohod J_1 topology. The J_1 metric is denoted by $d_{J_1}(\cdot, \cdot)$ and the uniform metric u is defined by the uniform norm:

$$\|x\|_T = \sup_{0 \leq t \leq T} |x(t)|,$$

for $x \in D[0, \infty)$ and $T \geq 0$; similarly, the L^1 metric is defined by the L^1 norm:

$$d_{L^1}^T(x, y) = \int_0^T |x(t) - y(t)| dt, \quad (1)$$

for $x, y \in D[0, \infty)$ and $T \geq 0$. Product metric spaces $(D^k[0, \infty), d_{J_1}^k)$ and $(D^k[0, \infty), u^k)$ are defined by $(D[0, \infty) \times \cdots \times D[0, \infty), d_{J_1} \times \cdots \times d_{J_1})$ and $(D[0, \infty) \times \cdots \times D[0, \infty), u \times \cdots \times u)$, respectively; $d_{J_1} \times \cdots \times d_{J_1}$ and $u \times \cdots \times u$ refer to the corresponding maximum metrics. Let \Rightarrow denote convergence in distribution – for stochastic processes in $D[0, \infty)$, as well as for random variables in \mathbb{R} . Let $1_{\{\cdot\}}$ be the usual indicator function and $e = \{e(t) = t, t \geq 0\}$ be the identity map. The composition map is denoted by \circ , i.e., for $(x, y) \in D[0, \infty) \times D[0, \infty)$, $x \circ y$ is defined by $(x \circ y)(t) = x(y(t))$, $t \geq 0$. For $x, y \in \mathbb{R}$, x^+ denotes the positive part of x , and $x \wedge y = \min\{x, y\}$.

2. ASSUMPTIONS

2.1. The model. We consider a sequence of first-come first-served (FCFS) $G/GI/N+GI$ queues indexed by the number of servers N . Customers arriving after $t = 0$ are indexed by natural numbers in an increasing order of their arrival times. Customer i arrives to the system at time $t_i > 0$ and two quantities are associated with it: the service requirement s_i and patience p_i . The service requirements of customers, $\{s_i, i \geq 1\}$, are independent and identically distributed (i.i.d.), characterized by a distribution function F , which does not vary with N (set $\bar{F} = 1 - F$). The sequence $\{p_i, i \geq 1\}$ is i.i.d. with a distribution G^N for the N th system. For simplicity of notation, we shall not index arrival times, service requirements and customers' patience by N – this dependency will be implicit.

Define $A^N(t)$, $t \geq 0$, to be the number of arrivals in the N th system over the time interval $[0, t]$. The process $A^N = \{A^N(t), t \geq 0\}$ is r.c.l.l., nondecreasing, nonnegative, integer-valued, with jumps of size 1 such that $A^N(0) = 0$ and $A^N(t) < \infty$ for all $t \geq 0$, almost surely (a.s.). The arrival process is related to the customer arrival times $\{t_i, i \geq 1\}$ by $t_i = \inf\{t \geq 0 : A^N(t) \geq i\}$, $i \geq 1$. Define $\tau^N = \{\tau^N(t), t \geq 0\}$ by $\tau^N(t) = t_{A^N(t)}$ for $t \geq t_1$ and $\tau^N(t) = 0$ for $t \leq t_1$; $\tau^N(t)$ is the time of the last arrival prior to t .

At time $t = 0$, there are q_0^N initial customers in the system, labeled by $-q_0^N, -q_0^N + 1, \dots, -1$. Those with indices $-q_0^N, -q_0^N + 1, \dots, -(q_0^N - N)^+ - 1$ are in service with i.i.d. service requirements drawn from the distribution F_* , the residual distribution associated with F :

$$F_*(x) = \mu \int_0^x \bar{F}(u) du, \quad (2)$$

where $\mu^{-1} = \mathbb{E}s_1$ is the mean service, which we assume exists (also set $\bar{F}_* = 1 - F_*$). The remaining $(q_0^N - N)^+$ initial customers (indexed by $-(q_0^N - N)^+, -(q_0^N - N)^+ + 1, \dots, -1$, if exist) have independent service requirements distributed according to F . However, their patience is infinite, i.e., $p_{-i} \equiv \infty$ for $i = 1, 2, \dots, (q_0^N - N)^+$. This assumption is convenient for the analysis while being non-restrictive, as argued at the end of this section.

Let v_i denote the *offered* waiting time of the i th customer – the amount of time the customer awaits service if the customer were infinitely patient ($p_i = \infty$). The *virtual* waiting time $V^N(t)$ at time $t \geq 0$ is the amount of time (measured beyond t) until one of the servers becomes idle, provided no new arrivals occur; by definition, $V^N(t) = 0$ if there exists an idle server at time t . The random variable $V^N(t)$ captures the amount of work in the queue at time t . (Note that a service completion which is immediately followed by a new service initiation does not render a server idle.) We set $V^N = \{V^N(t), t \geq 0\}$. The *actual* waiting time of the i th customer is then given by $v_i \wedge p_i$. That is, if customer i eventually enters service then v_i is equal to its actual waiting time and $p_i > v_i$; on the other hand, if customer i abandons the system then $v_i = V^N(t_i^-)$ (note that only customers with positive indices can abandon) and $v_i \geq p_i$. We use $V_{\leftarrow}^N = \{V_{\leftarrow}^N(t), t \geq 0\}$ to denote the offered-waiting-time process, with $V_{\leftarrow}^N(t) = v_{A^N(t)}$, for $t \geq t_1$, and $V_{\leftarrow}^N(t) = v_{-q_0^N}$ for $0 \leq t < t_1$. The offered-waiting-time process is defined in such a way that if customer i arrives at time t_i then $v_i = V_{\leftarrow}^N(t_i)$ rather than $v_i = V^N(t_i^-)$. Both V^N and V_{\leftarrow}^N are r.c.l.l. processes.

Define $Q^N = \{Q^N(t), t \geq 0\}$ where $Q^N(t)$ is the total number of customers in the system at time $t \geq 0$; this number includes customers receiving service, customers awaiting service that eventually receive service and customers awaiting service that eventually abandon. For the purpose of analysis, it is convenient to consider an alternative model in which customers who abandon the system, do so upon arrival (based on p_i 's). Namely, customers upon arrival, “compare” their p_i with v_i and immediately abandon the system if $p_i \leq v_i$; in this model, all customers awaiting service receive service eventually. Such dynamics is easier to analyze, and it turns out asymptotically equivalent to the original system. In order to distinguish between the two models, we introduce $H^N = \{H^N, t \geq 0\}$, where $H^N(t)$ is the number of customers at time $t \geq 0$ in the system with abandonment upon arrival.

2.2. The QED regime. We assume that the sequence of processes $\{A^N\}$ satisfies: (i) a functional strong law of large numbers (FSLLN):

$$A^N/\lambda^N \rightarrow e \quad (3)$$

u.o.c. a.s., as $N \rightarrow \infty$, where λ^N is the arrival rate in the N th system, and (ii) a functional central limit theorem (FCLT):

$$\hat{A}^N := \frac{1}{\sqrt{N}}(A^N - \lambda^N e) \Rightarrow \hat{A}, \quad (4)$$

as $N \rightarrow \infty$, where \hat{A} is a stochastic process with a.s. continuous sample paths.

The offered load to the N th system is λ^N/μ and the traffic intensity is $\rho^N = \lambda^N/(\mu N)$. In the QED regime, the number of servers N and traffic intensity ρ^N are related, in the limit as $N \rightarrow \infty$, via

$$\sqrt{N}(1 - \rho^N) \rightarrow \beta, \quad (5)$$

for some $-\infty < \beta < \infty$. In this regime, it is expected that the (virtual) waiting time vanishes as $N \rightarrow \infty$, hence only the behavior of G^N around the origin is relevant in the limit. To this end, we assume $G^N(0) = 0$ and

$$\hat{G}^N \rightarrow \theta e, \quad (6)$$

u.o.c., as $N \rightarrow \infty$, for some $0 \leq \theta < \infty$, where $\hat{G}(t) := \sqrt{N}G^N(t/\sqrt{N})$. The condition (6) is satisfied, for example, when $G^N = G$ for all N , and $G(t)/t \rightarrow \theta$ as $t \downarrow 0$ (or, equivalently, θ is the right-hand derivative of G at the origin).

The scaled and centered versions of Q^N and H^N are defined by

$$\hat{Q}^N = \{\hat{Q}^N(t), t \geq 0\} = \frac{1}{\sqrt{N}}(Q^N - N)$$

and

$$\hat{H}^N = \{\hat{H}^N(t), t \geq 0\} = \frac{1}{\sqrt{N}}(H^N - N),$$

respectively. As will be shown (see Theorem 1 and Corollary 1 in Section 4), the difference between \hat{Q}^N and \hat{H}^N vanishes in the limit, as $N \rightarrow \infty$. The scaled versions of the waiting time processes are given by

$$\hat{V}^N = \{\hat{V}^N(t), t \geq 0\} = \mu\sqrt{N}V^N$$

and

$$\hat{V}_{\leftarrow}^N = \{\hat{V}_{\leftarrow}^N(t), t \geq 0\} = \mu\sqrt{N}V_{\leftarrow}^N;$$

note that we use μ in the scaling for waiting time processes, which amounts to measuring wait in units of average service time.

2.3. Initial conditions. The number of customers in the system, at time $t = 0$, is given by $Q^N(0) = H^N(0) = q_0^N$. It is assumed that a scaled and centered version of q_0^N converges in distribution:

$$\hat{q}_0^N = \frac{1}{\sqrt{N}}(q_0^N - N) \Rightarrow \hat{q}_0, \tag{7}$$

as $N \rightarrow \infty$. This condition (together with the assumption that the residual service times of customers in service at $t = 0$ are i.i.d. with distribution F_*) is identical to the assumptions made in [27]. Although our initial condition is appealing in its simplicity, it is not the unique initial condition that induces the QED regime; e.g., see [21].

Next, we discuss an alternative model for patience of the initial customers. Namely, suppose that initial customers (at $t = 0$) do not have infinite patience but rather the sequence $\{p_{-i}, 1 \leq i \leq (q_0^N - N)^+\}$ is i.i.d., drawn from G^N . We argue that this variation does not impact our asymptotic results. To this end, let r_0^N be the number of initial customers that abandon the system:

$$r_0^N = \sum_{i=1}^{(q_0^N - N)^+} 1_{\{p_{-i} \leq v_{-i}\}};$$

i.e., $(q_0^N - r_0^N)$ initial customers awaiting service end up receiving service. Then, the following lemma holds.

Lemma 1. $r_0^N/\sqrt{N} \Rightarrow 0$, as $N \rightarrow \infty$.

Proof. See Section 7.1. □

As a consequence, we have

$$(q_0^N - N - r_0^N)/\sqrt{N} \Rightarrow \hat{q}_0,$$

as $N \rightarrow \infty$. Thus, the two models are asymptotically equivalent since (7) is the only assumption on the initial number of customers in the system and our limiting results are impacted via \hat{q}_0 only (see Theorem 1, Corollaries 1 and 2 in Section 4).

3. PRELIMINARIES

3.1. No abandonment. Consider the sequence of queues indexed by N , as introduced in the previous section. We next describe a corresponding sequence of systems without customer abandonment; entities associated with the systems without abandonment are appended by the “dot” symbol. Namely, for the N th system without abandonment, we set the initial and input parameters equal to those of the N th system with abandonment, except that all customers have infinite patience in the new system: $\dot{A}^N = A^N$; $\dot{q}_0^N = q_0^N$; $\dot{s}_i = s_i$, $i \geq -\dot{q}_0^N$; and $\dot{p}_i = \infty$, $i \geq -\dot{q}_0^N$. To obtain upper bounds on the offered waiting times $\{v_i, i \geq 1\}$, the following proposition [2] is used in conjunction with the results for the system without abandonment [27] (see Proposition 2 below). The process $\dot{V}_{\leftarrow}^N = \{\dot{V}_{\leftarrow}^N(t), t \geq 0\}$ is now a waiting-time process (as opposed to V_{\leftarrow}^N , which is an offered wait): if \dot{v}_i is the waiting time of customer i , then $\dot{V}_{\leftarrow}^N(t) = \dot{v}_{\dot{A}^N(t)}$, for $t \geq t_1$, and $\dot{V}_{\leftarrow}^N(t) = \dot{v}_{-1}$, for $0 \leq t < t_1$.

Proposition 1 (Bhattacharya & Ephremides [2]). $V_{\leftarrow}^N(t) \leq \dot{V}_{\leftarrow}^N(t)$ and $H^N(t) \leq \dot{H}^N(t) = \dot{Q}^N(t)$, for $t \geq 0$.

Proof. For completeness, we provide a proof in Section 7.2, which is verified within the setup of the present paper. \square

The following result is a consequence of the preceding proposition and Proposition 5.3 in [27].

Proposition 2. $V_{\leftarrow}^N \Rightarrow 0$, as $N \rightarrow \infty$.

3.2. Infinite-server processes. For each N , we consider a corresponding infinite-server process $X^N = \{X^N(t), t \geq 0\}$, defined by the original arrival process A^N and the sequence of service times $\{s_i, i \geq 1\}$, as follows:

$$\begin{aligned} X^N(t) &= \sum_{i=1}^{A^N(t)} 1_{\{t_i + s_i > t\}} \\ &= \sum_{i=1}^{A^N(t)} (1_{\{s_i > t - t_i\}} - \bar{F}(t - t_i)) + \int_0^t \bar{F}(t - s) dA^N(s). \end{aligned}$$

In addition, we introduce $\hat{X}^N = \{\hat{X}^N(t), t \geq 0\}$ to be a scaled and centered version of X^N :

$$\hat{X}^N = \frac{1}{\sqrt{N}}(X^N - N\rho^N F_*), \quad (8)$$

namely, for $t \geq 0$,

$$\hat{X}^N(t) = \frac{1}{\sqrt{N}} \sum_{i=1}^{A^N(t)} (1_{\{s_i > t - t_i\}} - \bar{F}(t - t_i)) + \int_0^t \bar{F}(t - s) d\hat{A}^N(s). \quad (9)$$

The following lemma, due to Krichagina and Puhalskii (see Theorem 3 in [20]), characterizes the limiting infinite-server process. Earlier results on the infinite-server process were obtained by Borovkov [4] and Iglehart [15]; for a recent measure-valued approach see [8, 28]. Define $U = \{U(t, x), t \geq 0, x \in [0, 1]\}$ to be a Keifer process, that is, a two-parameter continuous centered Gaussian process on $\mathbb{R}_+ \times [0, 1]$, with covariance function $\mathbb{E}[U(s, x)U(t, y)] = (s \wedge t)(x \wedge y - xy)$.

Lemma 2 (Krichagina & Puhalskii [20]). *The sequence of infinite-server processes $\{\hat{X}^N\}$ converges in distribution in $D[0, \infty)$, as $N \rightarrow \infty$, to the process $\hat{X} = \{\hat{X}(t), t \geq 0\}$ defined by*

$$\hat{X}(t) = \int_0^t \bar{F}(t - s) d\hat{A}(s) + \int_0^t \int_0^t 1_{\{s+x \leq t\}} dU(\mu s, F(x)), \quad t \geq 0,$$

where U is a Keifer process, \hat{A} and U are independent, and the first integral is understood as the result of integration by parts.

Recall the definition of offered waiting times $\{v_i, i \geq 1\}$ from Section 2. It will turn out convenient to define a (scaled and centered) process $\hat{X}_\Delta^N = \{\hat{X}_\Delta^N(t), t \geq 0\}$ by

$$\hat{X}_\Delta^N(t) = \frac{1}{\sqrt{N}} \sum_{i=1}^{A^N(t)} (1_{\{t-t_i-v_i < s_i \leq t-t_i\}} - \bar{F}(t-t_i-v_i) + \bar{F}(t-t_i)), \quad t \geq 0, \quad (10)$$

since this process relates to the (scaled and centered) number of customer with positive indices (those with arrival times $t_i > 0$) in the system at time $t \geq 0$, via the following equality:

$$\hat{X}^N(t) + \hat{X}_\Delta^N(t) - \int_0^t \bar{F}(t-s) d\hat{A}^N(s) = \frac{1}{\sqrt{N}} \sum_{i=1}^{A^N(t)} (1_{\{s_i > t-t_i-v_i\}} - \bar{F}(t-t_i-v_i)).$$

3.3. Initial-customers processes. In this subsection we consider the infinite-server processes associated with the customers initially in the system (at time $t = 0$). The process $I^N = \{I^N(t), t \geq 0\}$ is defined by

$$I^N(t) = \sum_{i=(q_0^N - N)^+ + 1}^{q_0^N} 1_{\{s_{-i} > t\}} + \sum_{i=1}^{(q_0^N - N)^+} 1_{\{s_{-i} > t\}},$$

for $t \geq 0$; recall that the random variables s_{-i} in the two sums are distributed according to F_* and F , respectively. Hence, the scaled and centered version $\hat{I}^N = \{\hat{I}^N(t), t \geq 0\}$ is defined by

$$\hat{I}^N = \frac{1}{\sqrt{N}} [I^N - (q_0^N \wedge N) \bar{F}_* - (q_0^N - N)^+ \bar{F}],$$

namely, for $t \geq 0$,

$$\hat{I}^N(t) = \frac{1}{\sqrt{N}} \sum_{i=(q_0^N - N)^+ + 1}^{q_0^N} (1_{\{s_{-i} > t\}} - \bar{F}_*(t)) + \frac{1}{\sqrt{N}} \sum_{i=1}^{(q_0^N - N)^+} (1_{\{s_{-i} > t\}} - \bar{F}(t)); \quad (11)$$

recall that $\mathbb{E}1_{\{s_{-i} > t\}} = \bar{F}_*(t)$, $(q_0^N - N)^+ < i \leq q_0^N$, and $\mathbb{E}1_{\{s_{-i} > t\}} = \bar{F}(t)$, $1 \leq i \leq (q_0^N - N)^+$. The following lemma characterizes the limiting behavior (as $N \rightarrow \infty$) of \hat{I}^N .

Lemma 3. $\hat{I}^N \Rightarrow \hat{I} = W \circ F_*$, as $N \rightarrow \infty$, where $W = \{W(t), t \in [0, 1]\}$ is a (standard) Brownian bridge, that is, a centered Gaussian process with covariance function $\mathbb{E}[W(t)W(s)] = t \wedge s - ts$.

Proof. Define $\hat{I}_1^N = \{\hat{I}_1^N(t), t \geq 0\}$ and $\hat{I}_2^N = \{\hat{I}_2^N(t), t \geq 0\}$ such that $\hat{I}^N = \hat{I}_1^N + \hat{I}_2^N$, i.e., $\hat{I}_1^N(t)$ and $\hat{I}_2^N(t)$ correspond to the first and second summand in (11), respectively. From Lemma 3.1 in [20], the random time change theorem and (7), it follows that $\hat{I}_1^N \Rightarrow \hat{I}$, as $N \rightarrow \infty$. By the same argument $\hat{I}_2^N \Rightarrow 0$, as $N \rightarrow \infty$, since (7) implies $(q_0^N - N)^+/N \Rightarrow 0$, as $N \rightarrow \infty$. \square

Next we introduce $\hat{I}_\Delta^N = \{\hat{I}_\Delta^N(t), t \geq 0\}$, where

$$\hat{I}_\Delta^N(t) = \frac{1}{\sqrt{N}} \sum_{i=1}^{(q_0^N - N)^+} (1_{\{t-v_{-i} < s_{-i} \leq t\}} - \bar{F}(t-v_{-i}) + \bar{F}(t)), \quad t \geq 0. \quad (12)$$

The relationship between \hat{I}_Δ^N and \hat{I}^N is similar to the relationship between \hat{X}_Δ^N and \hat{X}^N :

$$\hat{I}^N(t) + \hat{I}_\Delta^N(t) = \frac{1}{\sqrt{N}} \sum_{i=(q_0^N - N)^{++1}}^{q_0^N} (1_{\{s_{-i} > t\}} - \bar{F}_*(t)) + \frac{1}{\sqrt{N}} \sum_{i=1}^{(q_0^N - N)^+} (1_{\{s_{-i} > t - v_i\}} - \bar{F}(t - v_i)), \quad (13)$$

for $t \geq 0$, is the (scaled and centered) number of customers with negative indices that are in the system at time $t \geq 0$. Note that the sum in (12) contains elements corresponding to customers awaiting service at time $t = 0$ only; this is due to the fact that $v_{-i} = 0$ for $(q_0^N - N)^+ < i \leq q_0^N$ by definition. The following lemma states that the process $\hat{X}_\Delta^N + \hat{I}_\Delta^N$ diminishes as $N \rightarrow \infty$.

Lemma 4. $\hat{X}_\Delta^N + \hat{I}_\Delta^N \Rightarrow 0$, as $N \rightarrow \infty$.

Proof. See Section 7.3. □

3.4. Abandonment. Throughout the present section, we consider processes that correspond to the system with *abandonment upon arrival* – see the discussion in Section 2. This system is easier to analyze than the one where customers abandon after waiting. However, as already noted, the two systems are equivalent in the QED regime.

The infinite-server process X^N was constructed from all arriving customers. Now, let $Z^N = \{Z^N(t), t \geq 0\}$ be the infinite-server process induced only by arrivals that do abandon, i.e.,

$$Z^N(t) = \sum_{i=1}^{A^N(t)} 1_{\{t_i + s_i > t\}} 1_{\{p_i \leq v_i\}}.$$

Consequently, the scaled and centered version $\hat{Z}^N = \{\hat{Z}^N(t), t \geq 0\}$ is defined by

$$\hat{Z}^N(t) = \frac{1}{\sqrt{N}} \left[Z^N(t) - \int_0^t \bar{F}(t-s) G^N(V_{\leftarrow}^N(s)) dA^N(s) \right], \quad t \geq 0; \quad (14)$$

the independence of service requirements, customer patience and the arrival process, together with the independence of (s_i, p_i) and v_i , yield

$$\begin{aligned} \mathbb{E} \hat{Z}^N(t) &= \frac{1}{\sqrt{N}} \mathbb{E} \sum_{i=1}^{A^N(t)} (1_{\{t_i + s_i > t\}} 1_{\{p_i \leq v_i\}} - \bar{F}(t - t_i) G^N(v_i)) \\ &= \frac{1}{\sqrt{N}} \mathbb{E} \sum_{i=1}^{A^N(t)} (\mathbb{E}[1_{\{s_i > t - t_i\}} | t_i] \mathbb{E}[1_{\{p_i \leq v_i\}} | v_i] - \bar{F}(t - t_i) G^N(v_i)) = 0, \end{aligned} \quad (15)$$

where the second equality is due to $\mathbb{E}[1_{\{s_i > t - t_i\}} | t_i] = \bar{F}(t - t_i)$ and $\mathbb{E}[1_{\{p_i \leq v_i\}} | v_i] = G^N(v_i)$. The next lemma states that the process \hat{Z}^N is negligible in the limit, as $N \rightarrow \infty$. The lemma is based on the assumptions $G^N(0) = 0$ and $\theta < \infty$. Namely, during time intervals when the offered waiting time is positive, the rate at which customers abandon is proportional to \sqrt{N} (for large N), which is negligible relative to the total arrival rate λ^N , the latter being linear in N .

Lemma 5. $\hat{Z}^N \Rightarrow 0$, as $N \rightarrow \infty$.

Proof. See Section 7.4. □

Similarly, the infinite-server process due to customers who do not abandon will be denoted by $Y^N = \{Y^N(t), t \geq 0\} = X^N - Z^N$, with

$$Y^N(t) = \sum_{i=1}^{A^N(t)} 1_{\{t_i + s_i > t\}} 1_{\{p_i > v_i\}}.$$

Since customers abandon the system at a rate proportional to \sqrt{N} , the scaling and centering for Y^N is the same as for the process X^N in (8). Thus, $\hat{Y}^N = \{\hat{Y}^N(t), t \geq 0\}$ with

$$\hat{Y}^N = \frac{1}{\sqrt{N}}(Y^N - N\rho^N F_*), \quad (16)$$

which yields

$$\hat{Y}^N(t) = \hat{X}^N(t) - \hat{Z}^N(t) - \int_0^t \bar{F}(t-s) \sqrt{N} G^N(V_{\leftarrow}^N(s)) d\check{A}^N(s), \quad t \geq 0, \quad (17)$$

where $\check{A}^N = \{\check{A}^N(t), t \geq 0\}$ is a linearly-scaled arrival process:

$$\check{A}^N = A^N/N.$$

In parallel with \hat{X}_{Δ}^N and \hat{I}_{Δ}^N , define $\hat{Y}_{\Delta}^N = \{\hat{Y}_{\Delta}^N(t), t \geq 0\}$ by $\hat{Y}_{\Delta}^N = \hat{X}_{\Delta}^N - \hat{Z}_{\Delta}^N$, where $\hat{Z}_{\Delta}^N = \{\hat{Z}_{\Delta}^N(t), t \geq 0\}$ is given by

$$\hat{Z}_{\Delta}^N(t) = \frac{1}{\sqrt{N}} \sum_{i=1}^{A^N(t)} (1_{\{t-t_i-v_i < s_i \leq t-t_i\}} - \bar{F}(t-t_i-v_i) + \bar{F}(t-t_i)) 1_{\{p_i \leq v_i\}}, \quad t \geq 0. \quad (18)$$

Lemma 6. $\hat{Y}_{\Delta}^N + \hat{I}_{\Delta}^N \Rightarrow 0$, as $N \rightarrow \infty$.

Proof. See Section 7.5. □

Finally, we introduce $A_{\triangleright}^N = \{A_{\triangleright}^N(t), t \geq 0\}$ – the arrival process of customers that do not abandon the system, i.e., the customers that are eventually served; the process at time $t \geq 0$ is given by

$$A_{\triangleright}^N(t) = \sum_{i=1}^{A^N(t)} 1_{\{p_i > v_i\}}.$$

A corresponding scaled and centered version $\hat{A}_{\triangleright}^N = \{\hat{A}_{\triangleright}^N(t), t \geq 0\}$ is defined by

$$\hat{A}_{\triangleright}^N = \frac{1}{\sqrt{N}}(A_{\triangleright}^N - \lambda^N e);$$

this process is also used in the proof of Lemma 4 (see Section 7.3). The last lemma in this subsection stems from the fact that \hat{A} has a.s. continuous sample paths and V^N vanishes in the limit, as $N \rightarrow \infty$. The process $(\tau^N + V_{\leftarrow}^N)$ arises when the relation between H^N and V_{\leftarrow}^N is considered. In particular, $V_{\leftarrow}^N(t) = V_{\leftarrow}^N(\tau^N(t)) = V_{\leftarrow}^N \circ \tau^N(t)$, for $t \geq 0$, is defined not by $H^N \circ \tau^N(t)$ only but rather by $H^N \circ \tau^N(t)$ and $H^N \circ (\tau^N + V_{\leftarrow}^N)(t)$ jointly – see the proof of Lemma 9.

Lemma 7. $\hat{A}_{\triangleright}^N \circ (\tau^N + V_{\leftarrow}^N) - \hat{A}_{\triangleright}^N \circ \tau^N \Rightarrow 0$, as $N \rightarrow \infty$.

Proof. The value of the process $\hat{A}_{\triangleright}^N$, at time $t \geq 0$, is given by

$$\hat{A}_{\triangleright}^N(t) = \hat{A}^N(t) - \frac{1}{\sqrt{N}} \sum_{i=1}^{A^N(t)} 1_{\{p_i \leq v_i\}},$$

and, thus, $\hat{A}_{\triangleright}^N \circ (\tau^N + V_{\leftarrow}^N) - \hat{A}_{\triangleright}^N \circ \tau^N = \hat{A}^N \circ (\tau^N + V_{\leftarrow}^N) - \hat{A}^N \circ \tau^N - \hat{A}_{\Delta}^N$, where $\hat{A}_{\Delta}^N = \{\hat{A}_{\Delta}^N(t), t \geq 0\}$ which, for $t \geq 0$, satisfies

$$\begin{aligned} \hat{A}_{\Delta}^N(t) &= \frac{1}{\sqrt{N}} \sum_{i=A^N(\tau^N(t))+1}^{A^N(\tau^N(t)+V_{\leftarrow}^N(t))} 1_{\{p_i \leq v_i\}} \\ &\leq \frac{1}{\sqrt{N}} \sum_{i=A^N(t)+1}^{A^N(t+V_{\leftarrow}^N(t))} 1_{\{p_i \leq v_i\}}; \end{aligned}$$

the inequality follows from the monotonicity of $A^N(\cdot)$, $\tau^N(t) \leq t$ and $A^N(\tau^N(t)) = A^N(t)$. Assumption (4) and Proposition 2 imply $\hat{A}^N \circ (e + V_{\leftarrow}^N) - \hat{A}^N \Rightarrow 0$, as $N \rightarrow \infty$. These same two propositions yield $\hat{A}_{\Delta}^N \Rightarrow 0$, as $N \rightarrow \infty$, and the statement follows. \square

3.5. Queue length. The number of customers in the system at time $t \geq 0$ can be expressed as the sum of indicator functions [4, 20, 27]:

$$H^N(t) = \sum_{i=1}^{A^N(t)} 1_{\{t_i+s_i+v_i>t\}} 1_{\{p_i>v_i\}} + \sum_{i=(q_0^N-N)^++1}^{q_0^N} 1_{\{s_{-i}>t\}} + \sum_{i=1}^{(q_0^N-N)^+} 1_{\{s_{-i}+v_{-i}>t\}}. \quad (19)$$

On the other hand, Proposition 2.1 in [27] renders

$$\int_0^t (H^N(t-s)-N)^+ dF(s) = \sum_{i=1}^{A^N(t)} (\bar{F}(t-t_i-v_i) - \bar{F}(t-t_i)) 1_{\{p_i>v_i\}} + \sum_{i=1}^{(q_0^N-N)^+} (\bar{F}(t-v_{-i}) - \bar{F}(t)).$$

Then, combining the preceding equality and (19) yields, for $t \geq 0$,

$$\begin{aligned} H^N(t) &= \sum_{i=1}^{A^N(t)} (1_{\{t_i+s_i+v_i>t\}} - \bar{F}(t-t_i-v_i) + \bar{F}(t-t_i)) 1_{\{p_i>v_i\}} \\ &\quad + \sum_{i=(q_0^N-N)^++1}^{q_0^N} (1_{\{s_{-i}>t\}} - \bar{F}_*(t)) + \sum_{i=1}^{(q_0^N-N)^+} (1_{\{s_{-i}+v_{-i}>t\}} - \bar{F}(t-v_{-i})) \\ &\quad + (q_0^N - N)^+ \bar{F}(t) + (q_0^N \wedge N) \bar{F}_*(t) + \int_0^t (H^N(t-s) - N)^+ dF(s), \end{aligned}$$

or, equivalently, in terms of scaled processes (see (7), (10), (13), (16), (18)), for $t \geq 0$:

$$\begin{aligned} \hat{H}^N(t) &= (\hat{q}_0^N)^+ (\bar{F}(t) - \bar{F}_*(t)) + \hat{q}_0^N \bar{F}_*(t) + \hat{I}^N(t) + \hat{I}_{\Delta}^N(t) \\ &\quad + \hat{Y}^N(t) + \hat{Y}_{\Delta}^N(t) + \int_0^t (\hat{H}^N(t-s))^+ dF(s) - \sqrt{N}(1 - \rho^N) F_*(t). \quad (20) \end{aligned}$$

The following operator, $\varphi : D[0, \infty) \rightarrow D[0, \infty)$, was introduced in [27] – it plays a fundamental role in the analysis of QED queues without abandonment.

Definition 1 (Reed [27]). *For each $x \in D[0, \infty)$, let $\varphi(x)$ be the unique solution y to*

$$y(t) = x(t) + \int_0^t y^+(t-s) dF(s), \quad t \geq 0.$$

Then, (20) can be rewritten in terms of the operator φ :

$$\hat{H}^N = \varphi \left((\hat{q}_0^N)^+ (\bar{F} - \bar{F}_*) + \hat{q}_0^N \bar{F}_* + \hat{I}^N + \hat{I}_{\Delta}^N + \hat{Y}^N + \hat{Y}_{\Delta}^N - \sqrt{N}(1 - \rho^N) F_* \right). \quad (21)$$

The next proposition establishes L^1 -continuity of φ . In [27], only continuity of φ in the topology of uniform convergence was considered. The additional mode of L^1 -continuity is needed to relate \hat{H}^N and \hat{V}_{\leftarrow}^N in Lemma 9 (via Lemma 8). In particular, due to (14) (see also (17)), rather than approximating \hat{V}_{\leftarrow}^N by \hat{H}^N directly one needs only to relate integrals of these processes over finite intervals.

Proposition 3. *The function $\varphi : D[0, \infty) \rightarrow D[0, \infty)$ is Lipschitz continuous in the L^1 topology over bounded intervals.*

Proof. See Section 7.6. □

We now proceed to show that the scaled number-in-system process \hat{H}^N does not change significantly (in the L^1 sense, as $N \rightarrow \infty$) over time intervals during which individual customers await service. Note that $t = \tau^N(s) + V_{\leftarrow}^N(s)$ is the time when the last arriving customer before $t = s$ were to enter service if it had infinite patience (recall that V_{\leftarrow}^N is the offered-waiting-time process).

Lemma 8. *We have, as $N \rightarrow \infty$,*

$$\left\{ \int_0^t |\hat{H}^N \circ (\tau^N + V_{\leftarrow}^N)(s) - \hat{H}^N(s)| ds, t \geq 0 \right\} \Rightarrow 0.$$

Proof. See Section 7.7. □

3.6. Offered waiting time. The following lemma relates the (limiting and scaled) queue-length and offered-waiting-time processes in the QED regime. Recall that waiting is measured in units of average service-time.

Lemma 9. *We have, as $N \rightarrow \infty$,*

$$\left\{ \int_0^t |(\hat{H}^N(s))^+ - \hat{V}_{\leftarrow}^N(s)| ds, t \geq 0 \right\} \Rightarrow 0.$$

Remark 1. The lemma relates the queue-length and offered-waiting-time processes without a priori requiring that either of the processes converges weakly.

Proof. For $t \geq 0$, let $D^N(t)$ be the number of service completions during the time interval $[0, t]$. First, by definition, $V_{\leftarrow}^N(t)$ satisfies, for $t \geq 0$,

$$(H^N(\tau) + 1_{\{p \leq V_{\leftarrow}^N(t)\}} - N)^+ = D^N(\tau + V_{\leftarrow}^N(t)) - D^N(\tau), \quad (22)$$

where $\tau \equiv \tau^N(t)$ is the time of the last arrival prior to time t and $p \equiv p^N(t) = p_{A^N(t)}$ is the patience of the corresponding customer (set $p_0 = \infty$). The presence of the indicator function in (22) is due to the fact that the customer arriving at time τ might abandon the system on arrival (if $p \leq V_{\leftarrow}^N(\tau)$). Recall that, by definition, $V_{\leftarrow}^N(t) = V_{\leftarrow}^N(\tau)$ is the offered waiting time of the customer with index $A^N(t)$, i.e., the waiting time the customer would experience if it were not to abandon. The sum $H^N(\tau) + 1_{\{p \leq V_{\leftarrow}^N(t)\}}$ represents the number of customer in the system at time τ if the patience of the arriving customer is infinite. Second, the number of the customers in the system at time $\tau + V_{\leftarrow}^N(t) = \tau + V_{\leftarrow}^N(\tau)$ can be expressed as a linear combination of arrivals and departures:

$$\begin{aligned} H^N(\tau + V_{\leftarrow}^N(t)) &= H^N(\tau) + A_{\triangleright}^N(\tau + V_{\leftarrow}^N(t)) - A_{\triangleright}^N(\tau) - D^N(\tau + V_{\leftarrow}^N(t)) + D^N(\tau) \\ &= H^N(\tau) + A_{\triangleright}^N(\tau + V_{\leftarrow}^N(t)) - A_{\triangleright}^N(\tau) - (H^N(\tau) - N + 1_{\{p \leq V_{\leftarrow}^N(t)\}})^+, \end{aligned}$$

where the second equality is due to (22). Considering whether $V^N(t) > 0$ or $V^N(t) = 0$ in the preceding equation results in

$$(H^N(\tau + V_{\leftarrow}^N(t)) - N)^+ = A_{\triangleright}^N(\tau + V_{\leftarrow}^N(t)) - A_{\triangleright}^N(\tau) - 1_{\{p \leq V_{\leftarrow}^N(t)\}} 1_{\{H^N(\tau + V_{\leftarrow}^N(t)) = N\}}. \quad (23)$$

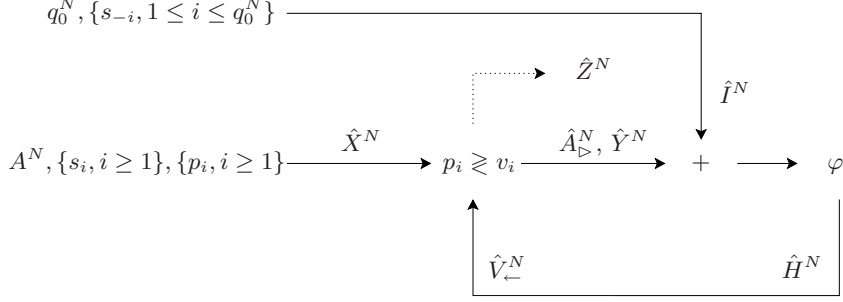


FIGURE 1. Relations between various processes/variables.

Third, rescaling the quantities in (23) gives rise to

$$\begin{aligned} (\hat{H}^N(t))^+ - \hat{V}_{\leftarrow}^N(t) &= (\hat{H}^N(t))^+ - (\hat{H}^N(\tau + V_{\leftarrow}^N(t)))^+ + \hat{A}_{\triangleright}^N(\tau + V_{\leftarrow}^N(t)) - \hat{A}_{\triangleright}^N(\tau) \\ &\quad - (1 - \rho^N) \hat{V}_{\leftarrow}^N(t) - 1_{\{p \leq V_{\leftarrow}^N(t)\}} 1_{\{H^N(\tau + V_{\leftarrow}^N(t)) = N\}} / \sqrt{N}. \end{aligned} \quad (24)$$

Next, note that (5) and Proposition 2 yield, as $N \rightarrow \infty$,

$$(1 - \rho^N) \hat{V}_{\leftarrow}^N \Rightarrow 0.$$

Finally, the statement follows from (24), the preceding limit and Lemmas 7, 8. \square

3.7. Summary of notation. We find it helpful to summarize, in Figure 1, various relations among the processes that have been introduced in this section. Process \hat{I}^N corresponds to customers that are initially in the system at time $t = 0$, while \hat{X}^N is the infinite-server process that corresponds to the customers that arrive after $t = 0$. Based on a comparison of customer patience and offered waiting times, process \hat{X}^N splits into the abandonment process \hat{I}^N and the infinite-server process \hat{Y}^N due to customers that receive service (do not abandon). Reed's operator φ provides a description of the queue-length process \hat{H}^N in terms of \hat{Y}^N and \hat{I}^N . Finally, the queue-length process is closely related to the (offered) waiting-time process \hat{V}_{\leftarrow}^N .

4. RESULTS

This section contains the main results of the paper. A central role is played by a mapping ϕ , applicable to the model with abandonment, which is a generalized version of the mapping φ in [27]. The two mappings coincide for $\theta = 0$ (no abandonment in the limit). Recall that the waiting time vanishes in the limit (Proposition 2) and, hence, the sequence of patience distributions $\{G^N\}$ manifests itself only through the parameter θ (cf. (6)).

Definition 2. *The mapping $\phi : D[0, \infty) \rightarrow D[0, \infty)$ is such that $\phi(x)$, for each $x \in D[0, \infty)$, is the unique solution y to*

$$y(t) = x(t) + \int_0^t y^+(t-s) dF(s) - \frac{\theta}{\mu} \int_0^t y^+(t-s) dF_*(s), \quad t \geq 0. \quad (25)$$

The next proposition guarantees that ϕ is well defined and summarizes some of its properties.

Proposition 4. *For each $x \in D[0, \infty)$ there exists a unique solution $\phi(x)$ to (25). The function $\phi : D[0, \infty) \rightarrow D[0, \infty)$ is Lipschitz continuous in the topology of uniform convergence over bounded intervals and measurable with respect to the Borel σ -field generated by the Skorohod J_1 topology.*

Proof. See Section 7.8. \square

4.1. Queue length. The following is the main result of our paper.

Theorem 1. *For the QED $G/GI/N+GI$ queue, with abandonments upon arrivals, we have, as $N \rightarrow \infty$,*

$$\hat{H}^N \Rightarrow \phi(\hat{q}_0^+(\bar{F} - \bar{F}_*) + \hat{q}_0\bar{F}_* + \hat{I} + \hat{X} - \beta F_*).$$

Remark 2. In the context of Theorem 1, the last term in (25) captures the effect of customers abandonment in the QED regime; note that the integration is with respect to the residual distribution F_* rather than the service distribution F . Namely, ϕ quantifies the negative feedback due to the abandonment effect (note that $\theta/\mu > 0$). The higher the number in the system, the higher the offered waiting time, the higher the abandonment rate, the lower the effective arrival rate of customers that eventually receive service, and the lower the number in the system; on the other hand, the lower the number in the system, the lower the offered waiting time, the lower the abandonment rate, the higher the arrival rate of customers that eventually receive service, and the higher the number in the system.

Proof. Using (17) and Definition 2, equality (21) can be rewritten as

$$\hat{H}^N = \phi(\hat{M}^N + \hat{I}^N + \hat{X}^N + \hat{\Delta}^N),$$

where

$$\hat{M}^N = (\hat{q}_0^N)^+(\bar{F} - \bar{F}_*) + \hat{q}_0^N\bar{F}_* - \sqrt{N}(1 - \rho^N)F_*, \quad (26)$$

and $\hat{\Delta}^N = \{\hat{\Delta}^N(t), t \geq 0\}$ is given by

$$\hat{\Delta}^N(t) = \hat{I}_{\Delta}^N(t) + \hat{Y}_{\Delta}^N(t) - \hat{Z}^N(t) - \int_0^t \bar{F}(t-s)\sqrt{N}G^N(V_{\leftarrow}^N(s))d\check{A}^N(s) + \frac{\theta}{\mu} \int_0^t (\hat{H}^N(t-s))^+ dF_*(s).$$

Combining Lemmas 5, 6 together with (2) and Lemma 9 yields, as $N \rightarrow \infty$,

$$\hat{\Delta}^N \Rightarrow 0. \quad (27)$$

From (5) and (7) it follows that, as $N \rightarrow \infty$,

$$\hat{M}^N \Rightarrow \hat{M} = \hat{q}_0^+(\bar{F} - \bar{F}_*) + \hat{q}_0\bar{F}_* - \beta F_*. \quad (28)$$

Now, we argue that, as $N \rightarrow \infty$, jointly

$$(\hat{M}^N, \hat{I}^N, \hat{X}^N, \hat{\Delta}^N) \Rightarrow (\hat{M}, \hat{I}, \hat{X}, 0); \quad (29)$$

note that the convergence of marginals is due to (28), Lemma 2, Lemma 3 and (27). To this end, introduce $\check{I}^N = \{\check{I}^N(t), t \geq 0\}$ with

$$\check{I}^N(t) = \frac{1}{\sqrt{N}} \sum_{i=q_0^N-N+1}^{q_0^N} (1_{\{\check{s}_{-i}>t\}} - \bar{F}_*(t)),$$

where $\check{s}_{-i} = s_{-i}$ for $(q_0^N - N)^+ < i \leq q_0^N$, and $\{\check{s}_{-i}, q_0^N - N < i \leq (q_0^N - N)^+\}$ is an i.i.d. sequence drawn from F_* and independent of all service requirements, arrival processes and q_0^N . Observe that the preceding sum contains exactly N elements (rather than a random number that depends on q_0^N), and the N -element sequence $\{\check{s}_{-i}, q_0^N - N < i \leq q_0^N\}$ is independent of q_0^N by construction (q_0^N is just an index in this case, and the elements of the sequence are independent of q_0^N); as a consequence, \check{I}^N and q_0^N are independent. Then the definitions of \check{I}^N and \hat{I}^N render, for $t \geq 0$,

$$\check{I}^N(t) - \hat{I}^N(t) = \frac{1}{\sqrt{N}} \sum_{i=q_0^N-N+1}^{(q_0^N-N)^+} (1_{\{\check{s}_{-i}>t\}} - \bar{F}_*(t)) - \frac{1}{\sqrt{N}} \sum_{i=1}^{(q_0^N-N)^+} (1_{\{s_{-i}>t\}} - \bar{F}(t)),$$

that, in turn, leads to (see the proof of Lemma 3), as $N \rightarrow \infty$,

$$(\check{I}^N, \hat{I}^N) \Rightarrow (\hat{I}, \hat{I}). \quad (30)$$

The limit $(\hat{M}^N, \check{I}^N, \hat{X}^N, 0) \Rightarrow (\hat{M}, \hat{I}, \hat{X}, 0)$, as $N \rightarrow \infty$, is due to the convergence of marginals and the independence of the pre-limit processes \hat{M}^N , \check{I}^N and \hat{X}^N [31, Theorem 11.4.4]; the independence is due to the fact that \hat{M}^N depends on q_0^N only (see (26)), \hat{X}^N depends only on the quantities associated with customers that are not initially in the system (see (9)), and \check{I} is independent of both q_0^N and A^N , $\{s_i, i \geq 1\}$. Furthermore, the following holds:

$$\begin{aligned} d_{J_1}((\hat{M}^N, \check{I}^N, \hat{X}^N, 0), (\hat{M}^N, \hat{I}^N, \hat{X}^N, \hat{\Delta}^N)) &\leq d_{J_1}(\check{I}^N, \hat{I}^N) + d_{J_1}(0, \hat{\Delta}^N) \\ &\Rightarrow 0, \end{aligned}$$

as $N \rightarrow \infty$, where the limit is due to (30), (27) and Theorem 11.4.8 in [31]. Finally, (29) follows from the preceding limit and Theorem 11.4.7 in [31].

The rest of the proof is almost identical to the corresponding part of the proof of Theorem 5.1 in [27]. Specifically, the space $D^4[0, \infty)$ is separable under the product topology (e.g., see Theorem 11.4.1 in [31]); therefore, due to (29) and the Skorohod representation theorem (e.g., see Theorem 3.2.2 in [31]), there exists an alternative probability space with $\{(\tilde{M}^N, \tilde{I}^N, \tilde{X}^N, \tilde{\Delta}^N)\}_{N \geq 1}$ and $(\tilde{M}, \tilde{I}, \tilde{X}, 0)$ defined on it with the following properties:

$$\begin{aligned} (\tilde{M}^N, \tilde{I}^N, \tilde{X}^N, \tilde{\Delta}^N) &\stackrel{d}{=} (\hat{M}^N, \hat{I}^N, \hat{X}^N, \hat{\Delta}^N), \\ (\tilde{M}, \tilde{I}, \tilde{X}, 0) &\stackrel{d}{=} (\hat{M}, \hat{I}, \hat{X}, 0), \\ (\tilde{M}^N, \tilde{I}^N, \tilde{X}^N, \tilde{\Delta}^N) &\rightarrow (\tilde{M}, \tilde{I}, \tilde{X}, 0) \text{ a.s.}, \end{aligned} \quad (31)$$

as $N \rightarrow \infty$. It should be noted that the last limit also holds under the uniform metric (not just J_1 metric) since both \hat{I} and \hat{X} have continuous sample paths and the set of discontinuity points of M^N is a subset of discontinuity points of F for all N . Hence, we have, as $N \rightarrow \infty$,

$$\tilde{M}^N + \tilde{I}^N + \tilde{X}^N + \tilde{\Delta}^N \rightarrow \tilde{M} + \tilde{I} + \tilde{X} \text{ a.s.} \quad (32)$$

under the uniform metric.

Define $\tilde{H}^N = \phi(\tilde{M}^N + \tilde{I}^N + \tilde{X}^N + \tilde{\Delta}^N)$ and note that, due to the measurability property of ϕ (Proposition 4) and (31), we have

$$\tilde{H}^N \stackrel{d}{=} \hat{H}^N. \quad (33)$$

Moreover, (32) and Proposition 4 (continuity part) yield, as $N \rightarrow \infty$,

$$\tilde{H}^N = \phi(\tilde{M}^N + \tilde{I}^N + \tilde{X}^N + \tilde{\Delta}^N) \rightarrow \phi(\tilde{M} + \tilde{I} + \tilde{X}) \text{ a.s.}$$

The fact that almost sure convergence implies convergence in distribution and convergence in the uniform metric implies convergence in the J_1 metric, together with (33), Proposition 4 (the measurability part) and the preceding limit yield

$$\hat{H}^N \Rightarrow \phi(\hat{M} + \hat{I} + \hat{X}),$$

as $N \rightarrow \infty$. The statement of the theorem now follows. \square

Recall that Q^N is the process of the total number of customers in the system when abandonments occur after waiting (as opposed to upon arrival). In view of Theorem 1, the following results indicates that, in the QED regime, the scaled number of customers awaiting service that eventually abandon becomes negligible (relative to the scaled total number of customers awaiting service) as the number of servers increases.

Corollary 1. *For the QED $G/GI/N+GI$ queue, with abandonments after waiting, we have, as $N \rightarrow \infty$,*

$$\hat{Q}^N \Rightarrow \phi(\hat{q}_0^+(\bar{F} - \bar{F}_*) + \hat{q}_0\bar{F}_* + \hat{I} + \hat{X} - \beta F_*),$$

where the limit coincides with that in Theorem 1.

Proof. The processes \hat{Q}^N and \hat{H}^N are related via $\hat{Q}^N = \hat{H}^N + \hat{R}^N$, where $\hat{R}^N = \{\hat{R}^N(t), t \geq 0\}$ is given by

$$\hat{R}^N(t) = \frac{1}{\sqrt{N}} \sum_{i=1}^{A^N(t)} 1_{\{t_i + p_i > t\}} 1_{\{p_i \leq v_i\}}.$$

Thus, in view of Theorem 1, it is sufficient to prove $\hat{R}^N \Rightarrow 0$, as $N \rightarrow \infty$. To this end, for any positive c and δ , the following inequality holds for all sufficiently large N :

$$\mathbb{P}[\|\hat{R}^N\|_T > \varepsilon] \leq \mathbb{P}[\|\hat{R}_{(c,\delta)}^N\|_T > \varepsilon] + \mathbb{P}[\|\hat{V}_{\leftarrow}^N\|_T > c],$$

where $\hat{R}_{(c,\delta)}^N = \{\hat{R}_{(c,\delta)}^N(t), t \geq 0\}$ is an infinite-server process with deterministic service times and

$$\hat{R}_{(c,\delta)}^N(t) = \frac{1}{\sqrt{N}} \sum_{i=1}^{A^N(t)} 1_{\{t_i + \delta > t\}} 1_{\{p_i \leq c/(\mu\sqrt{N})\}}.$$

Now, Theorem 3 in [20] implies $\hat{R}_{(c,\delta)}^N \Rightarrow \{c\theta(t \wedge \delta), t \geq 0\}$, as $N \rightarrow \infty$. On the other hand, Proposition 2 yields

$$\lim_{c \rightarrow \infty} \limsup_{N \rightarrow \infty} \mathbb{P}[\|\hat{V}_{\leftarrow}^N\|_T > c] = 0.$$

Therefore, given $\varepsilon > 0$, for any $\xi > 0$ it is possible to select c and δ such that $\mathbb{P}[\|\hat{R}^N\|_T > \varepsilon] < \xi$, for all N large enough. Consequently, $\hat{R}^N \Rightarrow 0$, as $N \rightarrow \infty$, and the corollary follows. \square

4.2. Waiting time. Now we introduce a mapping ψ that is an analogue of ϕ for the virtual-waiting-time process.

Definition 3. *The mapping $\psi : D[0, \infty) \rightarrow D[0, \infty)$ is such that $\psi(x)$, for each $x \in D[0, \infty)$, is the unique solution y to*

$$y(t) = \left[x(t) + \int_0^t y(t-s) dF(s) - \frac{\theta}{\mu} \int_0^t y(t-s) dF_*(s) \right]^+, \quad t \geq 0. \quad (34)$$

Remark 3. Note that, for $x \in D[0, \infty)$, if $y = \phi(x)$ then $y^+ = \psi(x)$, i.e., $\psi(x) = (\phi(x))^+$.

The next corollary characterizes the limiting waiting-time processes. Let $L^N = \{L^N(t), t \geq 0\}$ be the abandonment process in the N th system, that is, $L^N(t)$ is the number of customers that abandon by time t .

Corollary 2. *For the QED $G/GI/N+GI$ queue we have, as $N \rightarrow \infty$,*

$$\hat{V}^N \Rightarrow \hat{V} = \psi(\hat{q}_0^+(\bar{F} - \bar{F}_*) + \hat{q}_0\bar{F}_* + \hat{I} + \hat{X} - \beta F_*)$$

and

$$\hat{V}_{\leftarrow}^N \Rightarrow \hat{V}.$$

Remark 4. Note that, in view of Remark 3, the virtual waiting time \hat{V} and the queue length \hat{Q} are related via

$$\hat{V} = \hat{Q}^+.$$

(Recall that $\hat{V}^N = \mu\sqrt{N}V^N$ and $\hat{V}_{\leftarrow}^N = \mu\sqrt{N}V_{\leftarrow}^N$. In some of the literature, scaling that does not include the pre-factor μ is used, resulting in $\hat{V} = \hat{Q}^+/\mu$ rather than $\hat{V} = \hat{Q}^+$ as in the present paper.)

Proof. Recall that, from the definition of the process A_{\triangleright} , it follows that, for $t \geq 0$,

$$\hat{A}_{\triangleright}^N(t) = \hat{A}^N(t) - \frac{1}{\sqrt{N}} \int_0^t G^N(V_{\leftarrow}^N(s)) dA^N(t).$$

The preceding, (3), (4), (6), Lemma 9 and Theorem 1 yield, as $N \rightarrow \infty$,

$$\hat{A}_{\triangleright}^N \Rightarrow \left\{ \hat{A}(t) - \theta \int_0^t \hat{H}^+(s) ds, t \geq 0 \right\}.$$

Now, let $D^N = \{D^N(t), t \geq 0\}$ be the departure process in the N th system, i.e., $D^N(t)$ is the number of customers that receive service by time t . Then D^N and L^N can be expressed in terms of arrival and queue-length processes:

$$D^N(t) = A_{\triangleright}^N(t) - H^N(t) + q_0^N$$

and

$$L^N(t) = A^N(t) - A_{\triangleright}^N(t) + H^N(t) - Q^N(t),$$

where $t \geq 0$. These representations, (4.2), (4), Theorem 1 and Corollary 1 imply, as $N \rightarrow \infty$, $(D^N - \lambda^N e)/\sqrt{N} \Rightarrow \{\hat{D}(t), t \geq 0\}$, $L^N/\sqrt{N} \Rightarrow \{\hat{L}(t), t \geq 0\}$ and $L^N/N \rightarrow 0$ u.o.c. a.s., where

$$\hat{D}(t) = \hat{A}(t) - \theta \int_0^t \hat{H}^+(s) ds - \hat{H}(t) + \hat{q}_0, t \geq 0,$$

and

$$\hat{L}(t) = \theta \int_0^t \hat{H}^+(s) ds, t \geq 0. \quad (35)$$

Given the preceding limits, (3), (4), Corollary 1, the continuity of sample paths of \hat{A} and the Lipschitz continuity of ϕ (Proposition 4), the virtual-waiting-time process \hat{V}^N converges due to [30]: $\hat{V}^N \Rightarrow \hat{Q}^+ = \{Q^+(t), t \geq 0\}$, as $N \rightarrow \infty$, where \hat{Q} is such that $\hat{Q}^N \Rightarrow \hat{Q}$, as $N \rightarrow \infty$. However, from Corollary 1 it follows that $\hat{Q} = \phi(\hat{q}_0^+(\bar{F} - \bar{F}_*) + \hat{q}_0 \bar{F}_* + \hat{I} + \hat{X} - \beta F_*)$. The convergence of the offered-waiting-time processes \hat{V}_{\leftarrow}^N can be deduced from [25] since, in addition to convergence of the queue-length process, we have convergence of the arrival processes of customers that eventually receive service. \square

5. EXAMPLES

Example 1 (Estimating patience). In any application of models with abandonment there is the need to estimate the patience distribution [11]. Our results indicate that, in the QED regime, it suffices to merely estimate θ , the density of patience at the origin. The following corollary provides a theoretical justification for our proposed estimator.

Corollary 3. For the $G/GI/N+GI$ queue we have, as $N \rightarrow \infty$,

$$\left\{ \sqrt{N} \frac{L^N(t)}{A^N(t)}, t \geq 0 \right\} \Rightarrow \left\{ \frac{\theta}{\mu t} \int_0^t (\hat{Q}(s))^+ ds, t \geq 0 \right\}.$$

Proof. The statement follows from (4), Corollary 1 and (35). \square

The corollary suggests that an estimator for θ , $\hat{\theta}$, can be obtained in the following manner:

$$\hat{\theta} = \frac{L^N(t)/A^N(t)}{\frac{1}{\mu N t} \int_0^t (Q^N(s) - N)^+ ds}.$$

The numerator is simply the fraction of customers abandoning up to time t ; a practical approximation for the denominator can be the average waiting time up to time t . The accuracy of such estimators remains an open problem. \square

We next consider two specific examples. Both correspond to systems that have not yet been analyzed. The first example generalizes [33, 34] and the second expands on [21].

Example 2 (G/M/N+GI). Consider a system with exponential service times, noting that $F_* = F$. In addition, suppose that the sequences of random arrival times $\{t_i^N\}$ obey, for some $c > 0$,

$$\left\{ \frac{t_{\lfloor Nt \rfloor}^N - t/\mu}{\sqrt{N}c/\mu}, t \geq 0 \right\} \Rightarrow \hat{B}, \quad (36)$$

as $N \rightarrow \infty$, where \hat{B} is a standard Brownian motion. (Note that there exists a sequence of arrival times for each N , namely, the jump times of A^N .) Then, $\hat{H}^N \Rightarrow \hat{H}$ and $\hat{Q}^N \Rightarrow \hat{Q}$, as $N \rightarrow \infty$, due to Theorem 1 and Corollary 1, respectively. Here, $\hat{Q} = \hat{H}$ is the unique solution to

$$\hat{Q}(t) = \hat{I}(t) + \hat{X}(t) + (\hat{q}_0 + \beta) \exp\{-\mu t\} - \beta + (\mu - \theta) \int_0^t \hat{Q}^+(t-s) \exp\{-\mu s\} ds, \quad t \geq 0, \quad (37)$$

in which \hat{q}_0 is given in (7), \hat{I} in Lemma 3, and \hat{X} in Lemma 2. Similarly, due to Corollary 2, $\hat{V}^N \Rightarrow \hat{V}$, as $N \rightarrow \infty$, where \hat{V} is the unique solution to

$$\hat{V}(t) = \left[\hat{I}(t) + \hat{X}(t) + (\hat{q}_0 + \beta) \exp\{-\mu t\} - \beta + (\mu - \theta) \int_0^t \hat{V}(t-s) \exp\{-\mu s\} ds \right]^+, \quad t \geq 0.$$

The definitions of I^N and X^N give rise to

$$\frac{1}{\sqrt{N}}(I^N + X^N - \rho^N N) = \hat{I}^N + \hat{X}^N + (\hat{q}_0^N + \sqrt{N}(1 - \rho^N))\bar{F}, \quad (38)$$

where $\bar{F}(t) = \exp\{-\mu t\}$, $t \geq 0$. Since the service times are exponential, the process on the left-hand side of the preceding equality weakly converges to $\hat{S} = \{\hat{S}(t), t \geq 0\}$ that satisfies $\hat{S}(0) = \hat{q}_0 + \beta$ and

$$d\hat{S}(t) = -\mu\hat{S}(t) dt + \sqrt{\mu(1+c^2)} dB(t), \quad t \geq 0, \quad (39)$$

where $\{B(t), t \geq 0\}$ is a standard Brownian motion [20]. Now, (37) and (38) result in

$$d\hat{Q}(t) = d\hat{S}(t) - [\mu(\hat{Q}(t) - \hat{S}(t) + \beta) - (\mu - \theta)\hat{Q}^+(t)] dt,$$

which, combined with (39), yields

$$\begin{aligned} d\hat{Q}(t) &= \begin{cases} [\mu(\hat{S}(t) - \beta) - \theta\hat{Q}(t)] dt + d\hat{X}(t), & \hat{Q}(t) > 0, \\ [\mu(\hat{S}(t) - \beta) - \mu\hat{Q}(t)] dt + d\hat{X}(t), & \hat{Q}(t) \leq 0, \end{cases} \\ &= \begin{cases} -[\mu\beta + \theta\hat{Q}(t)] dt + \sqrt{\mu(1+c^2)} dB(t), & \hat{Q}(t) > 0, \\ -[\mu\beta + \mu\hat{Q}(t)] dt + \sqrt{\mu(1+c^2)} dB(t), & \hat{Q}(t) \leq 0; \end{cases} \end{aligned} \quad (40)$$

the initial condition for \hat{Q} is $\hat{Q}(0) = \hat{q}_0$.

Finally, in the special case $\mu = \theta$, the operator ϕ simplifies to $\phi(x) = x$ (see (25)), and, therefore, $\hat{Q} = \hat{H} = \hat{q}_0 + \hat{I} + \hat{X} - (\hat{q}_0 + \beta)F$, with $F(t) = 1 - \exp\{-\mu t\}$, $t \geq 0$. Note that $\hat{q}_0 + \hat{I} + \hat{X} - (\hat{q}_0 + \beta)F = \hat{S} - \beta$ is the limiting scaled and centered infinite-server process with the initial condition defined by \hat{q}_0 ; for the QED M/M/N+M system, this relation holds even in the pre-limit. \square

Example 3 (G/D/N+GI). The deterministic service distribution $F(s) = 1_{\{s \geq 1/\mu\}}$ implies a uniform residual distribution $F_*(s) = (\mu s)^+ \wedge 1$. Theorem 1 and Corollary 1 guarantee $\hat{H}^N \Rightarrow \hat{H}$ and $\hat{Q}^N \Rightarrow \hat{Q}$, as $N \rightarrow \infty$, where $\hat{Q} = \hat{H}$ satisfies, for $0 \leq t < 1/\mu$,

$$\hat{Q}(t) = \hat{q}_0^+ \mu t + \hat{q}_0(1 - \mu t) + \hat{I}(t) + \hat{X}(t) - \beta \mu t - \theta \int_0^t \hat{Q}^+(t-s) ds,$$

while, for $t \geq 1/\mu$,

$$\hat{Q}(t) = \hat{X}(t) - \beta + \hat{Q}^+(t - 1/\mu) - \theta \int_0^{1/\mu} \hat{Q}^+(t - s) ds;$$

as in the previous example, \hat{q}_0 is given in (7), \hat{I} in Lemma 3, and \hat{X} in Lemma 2. On the other hand, Corollary 2 implies $\hat{V}^N \Rightarrow \hat{V}$, as $N \rightarrow \infty$, where \hat{V} is the unique solution to

$$\hat{V}(t) = \left[\hat{q}_0^+ \mu t + \hat{q}_0(1 - \mu t) + \hat{I}(t) + \hat{X}(t) - \beta \mu t - \theta \int_0^t \hat{V}(t - s) ds \right]^+, \quad t \in [0, 1/\mu),$$

and

$$\hat{V}(t) = \left[\hat{X}(t) - \beta + \hat{V}(t - 1/\mu) - \theta \int_0^{1/\mu} \hat{V}(t - s) ds \right]^+, \quad t \geq 1/\mu.$$

When comparing the present example with the QED G/D/N queue (no abandonment) [17], one observes that having abandonments results in more complex dynamics. Specifically, while in [17] the distribution of $\hat{Q}(t)$ depends only on $\hat{Q}(t - 1/\mu)$ (as far as \hat{Q} is concerned), here $\hat{Q}(t)$ depends on all values of \hat{Q} during the time interval $[t - 1/\mu, t)$. This is due to the presence of the residual service distribution in the operators ϕ and ψ . \square

6. FUTURE RESEARCH: STATIONARY DISTRIBUTION

Our analysis addresses the transient behavior of a QED system with impatient customers. The *stationary* distributions of the queue length and the waiting time remain unknown, as is the case for the corresponding system without abandonment; note that the system with impatient customers remains stable (as $t \rightarrow \infty$) for all finite values of the capacity parameter β . (A large-deviation characterization of the stationary distributions for a QED queue without abandonments can be found in [10].)

We observe that Example 2 is consistent with the results in [12] on the *stationary* number-in-system process (for the M/M/N+M system in the QED regime). Namely, based on (37), it is tempting to conjecture that, for the G/M/N+GI system, the stationary versions of number-in-system processes converge weakly, in the QED regime, as $N \rightarrow \infty$, to the process $\tilde{Q} = \{\tilde{Q}(t), t \in \mathbb{R}\}$, where \tilde{Q} is the unique stationary process that solves

$$\tilde{Q}(t) = \tilde{X}(t) - \beta + (\mu - \theta) \int_{-\infty}^t \tilde{Q}^+(s) \exp\{-\mu(t - s)\} ds;$$

here $\tilde{X} = \{\tilde{X}(t), t \in \mathbb{R}\}$ is the stationary version of the infinite-server process \hat{X} (see also Lemma 3). Under assumption (36), \tilde{X} satisfies $d\tilde{X}(t) = -\mu\tilde{X}(t) dt + \sqrt{\mu(1 + c^2)} dB(t)$, where $\{B(t), t \in \mathbb{R}\}$ is a standard Brownian motion (since \hat{I} vanishes as $t \rightarrow \infty$ – see Lemma 3, and Example 2 in Section 5). An example where these assumptions (\tilde{X} stationary and (36)) prevail is when the arrival process is stationary renewal and \hat{q}_0 has the corresponding stationary distribution. A conjecture for the latter is provided in (41) below; in the case of Poisson arrivals the (diffusion) stationary distribution of \hat{q}_0 was calculated in [33]. Consequently, \tilde{Q} is a (piecewise) Ornstein-Uhlenbeck process (\tilde{Q} satisfies (40) where \tilde{Q} substitutes for \hat{Q}), as derived earlier in [12] for the case $c = 1$ (Poisson arrivals). Based on the preceding and [5], one can calculate the probability density function of $\tilde{Q}(t)$ (see also [12, 33]):

$$f_{\tilde{Q}(t)}(q) = \chi f_-(q) 1_{\{q \leq 0\}} + (1 - \chi) f_+(q) 1_{\{q > 0\}}, \quad (41)$$

where $f_-(q) = \tilde{c} \Phi'(\tilde{c}(q + \beta)) / \Phi(\tilde{c}\beta)$, $f_+(q) = \tilde{c} \sqrt{\theta/\mu} \Phi'(\tilde{c}(q\sqrt{\theta/\mu} + \beta\sqrt{\mu/\theta})) / \Phi(-\tilde{c}\beta\sqrt{\mu/\theta})$,

$$\tilde{c} = \sqrt{\frac{2}{1 + c^2}},$$

Φ and Φ' are the distribution and density functions of the standard normal random variable, respectively, and $\chi = f_+(0)/(f_+(0) + f_-(0))$. Furthermore, from the stochastic differential equation for \tilde{Q} one deduces directly that the stationary distribution of $\tilde{c}\tilde{Q}(t)$ is equal to the stationary distribution of the identically scaled limiting queue length $\tilde{Q}_A(t)$ in the Erlang-A model, but with the the load parameter $\tilde{c}\beta$. This makes results on the Erlang-A model, documented for example in [23], directly applicable to the QED G/M/N+GI queue. For example,

$$\mathbb{P}[\text{wait} > 0] = \mathbb{P}[\tilde{Q}(t) > 0] = \left[1 + \sqrt{\frac{\theta}{\mu} \frac{h(\tilde{c}\beta\sqrt{\mu/\theta})}{h(-\tilde{c}\beta)}} \right]^{-1}$$

and

$$\mathbb{E}\tilde{Q}^+(t) = \frac{\mu}{\tilde{c}\theta} [h(\tilde{c}\beta\sqrt{\mu/\theta}) - \tilde{c}\beta\sqrt{\mu/\theta}] \left[\sqrt{\frac{\mu}{\theta}} + \frac{h(\tilde{c}\beta\sqrt{\mu/\theta})}{h(-\tilde{c}\beta)} \right]^{-1},$$

where $h(q) = \Phi'(q)/(1 - \Phi(q))$ is the hazard rate. Corollary 2 and Remark 4 now provide a recipe for calculating also performance measures that involve waiting time. In particular, it is well known that $\mathbb{P}[\text{abandon}] = \theta\mathbb{E}[\text{wait}]$ when the patience distribution is exponential.

7. PROOFS

7.1. Proof of Lemma 1. Let $\{\hat{s}_{-i}, i \geq 1\}$ and $\{\hat{p}_{-i}, i \geq 1\}$ be two i.i.d. sequences defined by distributions F_* and G^N , respectively. The FCFS policy implies $v_{-i-1} \leq v_{-i}$, for $1 \leq i < (q_0^N - N)^+$ and, hence, for $\varepsilon > 0$, $v \geq 0$ and $c \geq 0$, we have

$$\begin{aligned} \mathbb{P}[r_0^N/\sqrt{N} > \varepsilon] &\leq \mathbb{P} \left[\sum_{i=1}^{(q_0^N - N)^+} 1_{\{p_{-i} \leq v\}} > \varepsilon\sqrt{N} \right] + \mathbb{P}[v_{-1} > v] \\ &\leq \mathbb{P} \left[\sum_{i=1}^{\lceil c\sqrt{N} \rceil} 1_{\{\hat{p}_{-i} \leq v\}} > \varepsilon\sqrt{N} \right] + \mathbb{P}[\hat{q}_0^N > c] + \mathbb{P} \left[\sum_{i=(q_0^N - N)^+ + 1}^{q_0^N} 1_{\{s_{-i} \leq v\}} < q_0^N - N \right] \\ &\leq \frac{\lceil c\sqrt{N} \rceil}{\varepsilon\sqrt{N}} G^N(v) + \mathbb{P} \left[\sum_{i=1}^N (F_*(v) - 1_{\{\hat{s}_{-i} \leq v\}}) > NF_*(v) - c\sqrt{N} \right] + 2\mathbb{P}[\hat{q}_0^N > c], \end{aligned}$$

where the second inequality is due to the fact that the event $\{v_{-1} > v\}$ implies that the number of service completions in the time interval $[0, v]$ is less than $(q_0^N - N)$; in addition, the number of service completions in $[0, t]$ is lower bounded by the sum in the last term in the second inequality; Markov inequality is used to obtain the third inequality. Setting $v = d/\sqrt{N}$, with $d = d(c)$ large enough such that $\sqrt{N}F_*(d/\sqrt{N}) - c > \varepsilon$ for all N large enough (which is feasible due to definition (2) of F_*) and applying Markov inequality result in

$$\mathbb{P}[r_0^N/\sqrt{N} > \varepsilon] \leq \frac{\lceil c\sqrt{N} \rceil}{\varepsilon\sqrt{N}} G^N(d/\sqrt{N}) + \frac{F_*(d/\sqrt{N})}{(\sqrt{N}F_*(d/\sqrt{N}) - c)^2} + 2\mathbb{P}[\hat{q}_0^N > c].$$

Finally, letting first $N \rightarrow \infty$, recalling (2), (6), and (7), and then letting $c \rightarrow \infty$ yields the statement of the lemma. \square

7.2. Proof of Proposition 1. It is sufficient to prove the statement for offered waiting times since it implies the result for queue lengths:

$$\begin{aligned} H^N(t) &= \sum_{i=1}^{A^N(t)} \mathbf{1}_{\{t_i + s_i + V_{\leftarrow}^N(t_i) > t\}} \mathbf{1}_{\{V_{\leftarrow}^N(t_i) \leq p_i\}} + \sum_{i=1}^{q_0^N \wedge N} \mathbf{1}_{\{s_{-i} > t\}} + \sum_{i=q_0^N \wedge N + 1}^{q_0^N} \mathbf{1}_{\{s_{-i} + v_{-i} > t\}} \\ &\leq \sum_{i=1}^{A^N(t)} \mathbf{1}_{\{t_i + s_i + \dot{V}_{\leftarrow}^N(t_i) > t\}} + \sum_{i=1}^{q_0^N \wedge N} \mathbf{1}_{\{s_{-i} > t\}} + \sum_{i=q_0^N \wedge N + 1}^{q_0^N} \mathbf{1}_{\{s_{-i} + \dot{v}_{-i} > t\}} = \dot{H}^N(t), \end{aligned}$$

for $t \geq 0$; note that $\dot{v}_i = v_i$ and $\dot{p}_i = p_i = \infty$, for $-q_0^N \leq i < -q_0^N \wedge N$, by construction. Furthermore, one must consider V_{\leftarrow}^N and \dot{V}_{\leftarrow}^N only at the moments of arrivals ($t = t_i$ for some $i \geq 0$) and $t = 0$, due to the fact that, between arrivals, both V_{\leftarrow}^N and \dot{V}_{\leftarrow}^N remain constant.

Now, consider the closely related shortest-workload-first routing policy (that can be conveniently described by the Kiefer-Wolfowitz recurrence, e.g., see [1, p. 91]), and let $W_n^N(t)$ and $\dot{W}_n^N(t)$, $1 \leq n \leq N$, be the n th smallest server workload in the system with and without abandonment, respectively. Then, it is well known that $V_{\leftarrow}^N(t_i) = W_1^N(t_i -)$ and $\dot{V}_{\leftarrow}^N(t_i) = \dot{W}_1^N(t_i -)$. Starting an induction, assume

$$W_n^N(t_i) \leq \dot{W}_n^N(t_i) \quad (42)$$

for some $i \geq 1$ and all $1 \leq n \leq N$; the base of the induction is due the assumption on the initial states (at $t = 0$). Let \mathcal{R} be the standard reorder operator. Then, since the vectors of W_n^N 's and \dot{W}_n^N 's satisfy the Kiefer-Wolfowitz recurrence, it follows that

$$\begin{aligned} &(W_1^N(t_{i+1}), W_2^N(t_{i+1}), \dots, W_N^N(t_{i+1})) \\ &= \mathcal{R}(W_1^N(t_i) + s_{i+1} \mathbf{1}_{\{p_i \leq W_1^N(t_{i+1}-)\}} - t_{i+1} + t_i, W_2^N(t_i) - t_{i+1} + t_i, \dots, W_N^N(t_i) - t_{i+1} + t_i)^+ \\ &\leq \mathcal{R}(\dot{W}_1^N(t_i) + s_{i+1} - t_{i+1} + t_i, \dot{W}_2^N(t_i) - t_{i+1} + t_i, \dots, \dot{W}_N^N(t_i) - t_{i+1} + t_i)^+ \\ &= (\dot{W}_1^N(t_{i+1}), \dot{W}_2^N(t_{i+1}), \dots, \dot{W}_N^N(t_{i+1})), \end{aligned}$$

where the inequality is due the inductive assumption (42); the operator $(\cdot)^+$ is applied elementwise. Therefore, (42) holds for all $i \geq 1$ and the proposition prevails. \square

7.3. Proof of Lemma 4. Let $A_{\triangleright}^N = \{A_{\triangleright}^N(t), t \geq 0\}$, where, for $t \geq 0$,

$$A_{\triangleright}^N(t) = \sum_{i=1}^{A^N(t)} \mathbf{1}_{\{p_i > v_i\}}$$

represents the number of customers with arrival times in $[0, t]$ that eventually receive service (do not abandon); the process A_{\triangleright}^N was also considered in Section 3.4 (see Lemma 7). Define a two-dimensional process $\{E^N(t, s), t \geq 0, s \geq 0\}$ by

$$E^N(t, s) = \sum_{i=A_{\triangleright}^N(t) - (H^N(t) - N)^+ + 1}^{A_{\triangleright}^N(t)} \mathbf{1}_{\{\tilde{s}_i \leq s\}},$$

where $\tilde{s}_i = s_{i-1}$, $-(q_0^N - N)^+ < i \leq 0$, and $\tilde{s}_i = s_{A^N(\tilde{t}_i)}$, $i \geq 1$ with $\tilde{t}_i = \inf\{t \geq 0 : A_{\triangleright}^N(t) = i\}$. The value of $E^N(t, s)$ is equal to the number of customers awaiting service at time t with service requirement at most s (recall that customers abandon upon arrival, if at all). Let $w_i = v_i \mathbf{1}_{\{p_i > v_i\}}$ for $i \geq 0$, $w_{-i} = v_{-i}$ for $1 \leq i \leq (q_0^N - N)^+$, and $w_{-i} = 0$ for $(q_0^N - N)^+ < i \leq q_0^N$;

note that $w_i = 0$ for all customers that abandon the system. Alternatively, $E^N(t, s)$ can be expressed as a sum over all customer indices:

$$E^N(t, s) = \sum_{i=-q_0^N}^{A^N(t)} \mathbf{1}_{\{t_i \leq t < t_i + w_i\}} \mathbf{1}_{\{s_i \leq s\}}, \quad (43)$$

where $t_{-i} = 0$ for $i = 1, \dots, q_0^N$, and the element of the sum corresponding to $i = 0$ does not exist. Furthermore, we define $\{F^N(t, s), t \geq 0, s \geq 0\}$ by

$$F^N(t, s) := \mathbf{1}_{\{H^N(t) > N\}} \frac{E^N(t, s)}{(H^N(t) - N)^+}, \quad (44)$$

and note that

$$E^N(t, s) = E^N(t, s) \mathbf{1}_{\{H^N(t) > N\}} = (H^N(t) - N)^+ F^N(t, s); \quad (45)$$

on the event $\{H^N(t) > N\}$, $F^N(t, \cdot)$ can be interpreted as the (empirical) distribution function of customers awaiting service at time t . Observe that, for $\delta > 0$, (43) renders

$$E^N(t - s, s + \delta) - E^N(t - s, s) = \sum_{i=-q_0^N}^{A^N(t)} \mathbf{1}_{\{t_i \leq t - s < t_i + w_i\}} \mathbf{1}_{\{s_i - \delta \leq s < s_i\}}.$$

In view of the preceding equality, the change in the order of summation results in

$$\begin{aligned} \int_0^t E^N(t - s, ds) &= \sum_{i=-q_0^N}^{A^N(t)} \int_0^t \mathbf{1}_{\{t_i \leq t - s < t_i + w_i\}} d\mathbf{1}_{\{s_i \leq s\}} \\ &= \sum_{i=-q_0^N}^{A^N(t)} \mathbf{1}_{\{t_i \leq t - s_i < t_i + w_i\}}, \end{aligned}$$

and, thus, due to (45), we have

$$\int_0^t (H^N(t - s) - N)^+ F^N(t - s, ds) = \sum_{i=-q_0^N}^{A^N(t)} \mathbf{1}_{\{t - t_i - w_i < s_i \leq t - t_i\}}. \quad (46)$$

On the other hand, for any $t \geq 0$, Proposition 2.1 in [27] yields

$$\int_0^t (H^N(t - s) - N)^+ dF(s) = \sum_{i=-q_0^N}^{A^N(t)} (\bar{F}(t - t_i - w_i) - \bar{F}(t - t_i)), \quad (47)$$

due to the fact that only customers that do not abandon potentially contribute to the sum on the right-hand side of (47). Therefore, (46) and (47) imply (see (10) and (12)), for $t \geq 0$,

$$(\hat{X}_\Delta^N + \hat{I}_\Delta^N)(t) = \int_0^t (\hat{H}^N(t - s))^+ (F^N(t - s, ds) - F(ds)). \quad (48)$$

Next, extend the i.i.d. sequence $\{\tilde{s}_i, i > -(q_0^N - N)^+\}$ to all integer indices (by letting $\{\tilde{s}_i, i \leq -(q_0^N - N)^+\}$ be an i.i.d. sequence, independent of $\{\tilde{s}_i, i > -(q_0^N - N)^+\}$, with its elements distributed according to F); observe that $\{\tilde{s}_i, i \in \mathbb{Z}\}$ is an i.i.d. sequence since both subsequences are i.i.d (defined by F) and independent of each other. Now, define a family of empirical distribution functions $F_{i,j} = \{F_{i,j}(s), s \geq 0\}$:

$$F_{i,j}(s) = \frac{1}{j} \sum_{k=i-j+1}^i \mathbf{1}_{\{\tilde{s}_k \leq s\}}, \quad (49)$$

where $i \geq 0$ and $j \geq 1$. In what follows, we estimate $\|F_{i,j} - F\|_\infty$ for a range of indices i and j . To this end, for any $\varepsilon > 0$ and $s \geq 0$, there exist constants $\theta(\varepsilon, s) > 0$ and $\gamma(\varepsilon, s) < \infty$ (e.g., see [3, p. 151]) such that, for all $j \geq 1$ (and all i),

$$\mathbb{P}[|F_{i,j}(s) - F(s)| > \varepsilon] \leq \gamma(\varepsilon, s) \exp\{-j\theta(\varepsilon, s)\}. \quad (50)$$

Moreover, by the same argument, replacing $1_{\{\bar{s}_k \leq s\}}$ with $1_{\{\bar{s}_k < s\}}$ in the definition of $F_{i,j}(s)$ yields

$$\mathbb{P}[|F_{i,j}(s-) - F(s-)| > \varepsilon] \leq \gamma(\varepsilon, s) \exp\{-j\theta(\varepsilon, s)\}, \quad (51)$$

where $F(s-) = \mathbb{E}1_{\{s_i < s\}}$, $i \geq 1$; the constants in (50) and (51) may differ in general. Given the distribution function F , for any $\varepsilon > 0$ there exists a finite sequence of nonnegative reals $\{a_l, 1 \leq l \leq L\}$ such that

$$\bigcap_{l=1}^L \{|F_{i,j}(a_l) - F(a_l)| \leq \varepsilon\} \cap \{|F_{i,j}(a_l-) - F(a_l-)| \leq \varepsilon\} \subseteq \{\|F_{i,j} - F\|_\infty \leq 2\varepsilon\}.$$

This relationship, (50) and (51) imply the existence of $\theta(\varepsilon) > 0$ and $\gamma(\varepsilon) < \infty$ such that

$$\mathbb{P}[\|F_{i,j} - F\|_\infty > \varepsilon] \leq \gamma(\varepsilon) \exp\{-j\theta(\varepsilon)\}, \quad (52)$$

for $i \geq 0$ and $j \geq 1$. Now, we introduce a nonnegative real that characterizes a distance between F and $F_{i,j}$ for multiple indices i and j :

$$f_{k,l,n} = \sup_{0 \leq i \leq k} \sup_{l \leq j \leq k+n} \|F_{i,j} - F\|_\infty, \quad (53)$$

where $l \geq 1$. Then, for $\varepsilon > 0$, the union bound and (52) yield $\mathbb{P}[f_{k,l,n} > \varepsilon] \leq (k+1)(k+n)\gamma(\varepsilon) \exp\{-l\theta(\varepsilon)\}$. Finally, for any $\varepsilon > 0$, the last inequality, (3), (5) and (7) result in, as $N \rightarrow \infty$,

$$\mathbb{P}[f_{A^N(T), \varepsilon\sqrt{N}, q_0^N} > \varepsilon] \rightarrow 0. \quad (54)$$

Next, considering whether $\{\hat{H}^N(t-s) \leq \varepsilon\}$ or $\{\hat{H}^N(t-s) > \varepsilon\}$ in (48) yields

$$\begin{aligned} \|\hat{X}_\Delta^N + \hat{I}_\Delta^N\|_T &\leq \varepsilon + \sup_{0 \leq t \leq T} \left| \int_0^t 1_{\{\hat{H}^N(t-s) > \varepsilon\}} (\hat{H}^N(t-s))^+ (F^N(t-s, ds) - F(ds)) \right| \\ &\leq \varepsilon + \|\hat{H}^N\|_T \sup_{0 \leq t \leq T} \sup_{0 \leq t \leq s} \left| (F^N(t-s, s) - F(s)) 1_{\{\hat{H}^N(t-s) > \varepsilon\}} \right| \\ &\leq \varepsilon + \|\hat{H}^N\|_T \sup_{0 \leq t \leq T} \sup_{0 \leq t \leq s} \left| F_{A^N(t-s), H^N(t-s)-N}(s) - F(s) \right| \\ &\leq \varepsilon + \|\hat{H}^N\|_T f_{A^N(T), \varepsilon\sqrt{N}, q_0^N}, \end{aligned}$$

where the third inequality is due to $F^N(t, s) = F_{A^N(t), H^N(t)-N}(s)$ on the event $\{H^N(t) > N\}$ (see (44) and (49)); the last inequality follows from (53). Now, for any $\delta > 0$ there exists $\varepsilon > 0$, small enough, so that the preceding inequality results in

$$\begin{aligned} \mathbb{P}[\|\hat{X}_\Delta^N + \hat{I}_\Delta^N\|_T > 2\delta] &\leq \mathbb{P}[\|\hat{H}^N\|_T f_{A^N(T), \varepsilon\sqrt{N}, q_0^N} > \delta] \\ &\leq \mathbb{P}[f_{A^N(T), \varepsilon\sqrt{N}, q_0^N} > \delta/c] + \mathbb{P}[\|\hat{H}^N\|_T > c], \end{aligned} \quad (55)$$

where $c > 0$ is arbitrary. Finally, taking \limsup (as $N \rightarrow \infty$) on both sides of (55) yields, due to (54),

$$\limsup_{N \rightarrow \infty} \mathbb{P}[\|\hat{X}_\Delta^N + \hat{I}_\Delta^N\|_T > \delta] \leq \limsup_{N \rightarrow \infty} \mathbb{P}[\|\hat{H}^N\|_T > c].$$

The final statement follows from the preceding by letting $c \rightarrow \infty$, Proposition 1 and Theorem 5.1 in [27]. \square

7.4. **Proof of Lemma 5.** For fixed $T > 0$ and $\Delta > 0$, the following holds:

$$\|\hat{Z}^N\|_T \leq \max_{0 \leq i \leq \lfloor T/\Delta \rfloor} |\hat{Z}^N(i\Delta)| + \max_{0 \leq i \leq \lfloor T/\Delta \rfloor} \sup_{0 \leq \delta \leq \Delta} |\hat{Z}^N(i\Delta + \delta) - \hat{Z}^N(i\Delta)|. \quad (56)$$

First, we argue that $\hat{Z}^N(t) \Rightarrow 0$, as $N \rightarrow \infty$, for any fixed $t \geq 0$. For notational purposes, it is convenient to define the random variables $z_i(t) = 1_{\{s_i > t - t_i\}} 1_{\{p_i \leq v_i\}} - \bar{F}(t - t_i)G^N(v_i)$; observe that $\mathbb{E}z_i(t) = 0$ since s_i, p_i are independent of t_i, v_i , and, hence, $\mathbb{E}[z_i(t)|t_i, v_i] = 0$. From (15) we have that $\mathbb{E}\hat{Z}^N(t) = 0$ while the second moment is given by

$$\begin{aligned} \mathbb{E}(\hat{Z}^N(t))^2 &= \frac{1}{N} \mathbb{E} \sum_{i=1}^{A^N(t)} z_i^2(t) + \frac{2}{N} \mathbb{E} \sum_{i=1}^{A^N(t)} \sum_{j=i+1}^{A^N(t)} z_i(t) z_j(t) \\ &= \frac{1}{N} \mathbb{E} \sum_{i=1}^{A^N(t)} \bar{F}(t - t_i)G^N(v_i)(1 - \bar{F}(t - t_i)G^N(v_i)); \end{aligned}$$

the expectation of the double sum equals 0 since the service requirement and patience of an arriving customer is independent of the state of the system. Then, given that F and G^N are distribution functions, it follows that, for $\varepsilon > 0$,

$$\mathbb{P}[|\hat{Z}^N(t)| > \varepsilon] \leq \frac{1}{\varepsilon^2 N} \mathbb{E} \sum_{i=1}^{A^N(t)} G^N(v_i) \rightarrow 0,$$

as $N \rightarrow \infty$, due to (3), (5), (6), and Proposition 2; thus, for fixed t , as $N \rightarrow \infty$,

$$\hat{Z}^N(t) \Rightarrow 0. \quad (57)$$

Next, we consider the second term on the right-hand side of (56). To this end, for $t > 0$ and $\delta > 0$, we have (see (14))

$$\hat{Z}^N(t + \delta) - \hat{Z}^N(t) = \frac{1}{\sqrt{N}} \sum_{i=1}^{A^N(t)} (z_i(t + \delta) - z_i(t)) + \frac{1}{\sqrt{N}} \sum_{i=A^N(t)+1}^{A^N(t+\delta)} z_i(t + \delta),$$

and upper and lower bounds follow:

$$\begin{aligned} \hat{Z}^N(t + \delta) - \hat{Z}^N(t) &\leq \frac{1}{\sqrt{N}} \sum_{i=1}^{A^N(t)} (\bar{F}(t - t_i) - \bar{F}(t + \delta - t_i))G^N(v_i) + \frac{1}{\sqrt{N}} \sum_{i=A^N(t)+1}^{A^N(t+\delta)} 1_{\{p_i \leq v_i\}} \\ &:= \hat{Z}_{(\uparrow, \delta)}^N(t); \\ \hat{Z}^N(t + \delta) - \hat{Z}^N(t) &\geq -\frac{1}{\sqrt{N}} \sum_{i=1}^{A^N(t)} (1_{\{s_i > t - t_i\}} - 1_{\{s_i > t + \delta - t_i\}})1_{\{p_i \leq v_i\}} - \frac{1}{\sqrt{N}} \sum_{i=A^N(t)+1}^{A^N(t+\delta)} G^N(v_i) \\ &:= -\hat{Z}_{(\downarrow, \delta)}^N(t). \end{aligned}$$

The nonnegativity of $\hat{Z}_{(\uparrow, \delta)}^N(t)$ and $\hat{Z}_{(\downarrow, \delta)}^N(t)$ and their monotonicity in δ imply

$$\sup_{0 \leq \delta \leq \Delta} |\hat{Z}^N(t + \delta) - \hat{Z}^N(t)| \leq \hat{Z}_{(\uparrow, \Delta)}^N(t) + \hat{Z}_{(\downarrow, \Delta)}^N(t). \quad (58)$$

For notational simplicity, introduce $A_{(c)}^N = \{A_{(c)}^N(t), t \geq 0\}$ by

$$\begin{aligned} A_{(c)}^N(t) &:= \sum_{i=1}^{A^N(t)} 1_{\{p_i \leq c/(\mu\sqrt{N})\}} \\ &= \sum_{i=1}^{A^N(t)} \left(1_{\{p_i \leq c/(\mu\sqrt{N})\}} - G^N(c/(\mu\sqrt{N})) \right) - A^N(t) G^N(c/(\mu\sqrt{N})); \end{aligned} \quad (59)$$

also, set $\tilde{t}_i = \inf\{t \geq 0 : A_{(c)}^N(t) \geq i\}$ and $\tilde{s}_i = \{s_j : \tilde{t}_i = t_j\}$. The process $A_{(c)}^N$ is the arrival process of customers with patience at most $c/(\mu\sqrt{N})$, in the N th system. Limits (3), (5) and (6) imply

$$\left\{ \frac{A^N(t)}{N} \sqrt{N} G^N(c/(\mu\sqrt{N})), t \geq 0 \right\} \rightarrow c\theta e, \quad (60)$$

a.s. u.o.c., as $N \rightarrow \infty$, while the martingale inequality [6, Corollary 1, p. 331] and (6) yield, for $T > 0$,

$$\mathbb{P} \left[\sup_{1 \leq j \leq 2\mu TN} \left| \sum_{i=1}^j \left(1_{\{p_i \leq c/(\mu\sqrt{N})\}} - G^N(c/(\mu\sqrt{N})) \right) \right| > \varepsilon \sqrt{N} \right] \leq \frac{2\mu T G^N(c/(\mu\sqrt{N}))}{\varepsilon^2} \rightarrow 0, \quad (61)$$

as $N \rightarrow \infty$. Combining (59), (60) and (61) results in

$$A_{(c)}^N / \sqrt{N} \Rightarrow c\theta e, \quad (62)$$

as $N \rightarrow \infty$. Now, for $t \leq T - \delta$, on the event $\{\|\hat{V}_{\leftarrow}^N\|_T \leq c\}$ the first term on the right-hand side of (58) can be upper bounded by using monotonicity:

$$\begin{aligned} \hat{Z}_{(\uparrow, \Delta)}^N(t) &\leq \frac{1}{\sqrt{N}} G^N(c/(\mu\sqrt{N})) \sum_{i=1}^{A^N(t)} (\bar{F}(t - t_i) - \bar{F}(t + \Delta - t_i)) + \frac{1}{\sqrt{N}} [A_{(c)}^N(t + \Delta) - A_{(c)}^N(t)] \\ &\Rightarrow c\theta \int_0^t (\bar{F}(t - s) - \bar{F}(t + \Delta - s)) ds + c\theta \Delta, \end{aligned} \quad (63)$$

as $N \rightarrow \infty$, where the limit is due to (3), (5) and (6). Similarly, on the event $\{\|\hat{V}_{\leftarrow}^N\|_T \leq c\}$, we have

$$\begin{aligned} \hat{Z}_{(\downarrow, \Delta)}^N(t) &\leq \frac{1}{\sqrt{N}} \sum_{i=1}^{A_{(c)}^N(t)} (1_{\{\tilde{s}_i > t - \tilde{t}_i\}} - 1_{\{\tilde{s}_i > t + \Delta - \tilde{t}_i\}}) + \frac{1}{\sqrt{N}} [A^N(t + \Delta) - A^N(t)] G^N(c/(\mu\sqrt{N})) \\ &\Rightarrow c\theta \int_0^t (\bar{F}(t - s) - \bar{F}(t + \Delta - s)) ds + c\theta \Delta, \end{aligned} \quad (64)$$

as $N \rightarrow \infty$, where the limit is due to (3), (5), (6) and Theorem 3 in [20]. Now, for $c > 0$, (58) implies

$$\mathbb{P} \left[\sup_{0 \leq \delta \leq \Delta} |\hat{Z}^N(t + \delta) - \hat{Z}^N(t)| > \varepsilon \right] \leq \mathbb{P}[\hat{Z}_{(\uparrow, \Delta)}^N(t) + \hat{Z}_{(\downarrow, \Delta)}^N(t) > \varepsilon, \|\hat{V}_{\leftarrow}^N\|_T \leq c] + \mathbb{P}[\|\hat{V}_{\leftarrow}^N\|_T > c].$$

Selecting Δ small enough, letting $N \rightarrow \infty$ on both sides in the preceding inequality, using (63) and (64), and then increasing $c \rightarrow \infty$ yields (for fixed t)

$$\sup_{0 \leq \delta \leq \Delta} |\hat{Z}^N(t + \delta) - \hat{Z}^N(t)| \Rightarrow 0; \quad (65)$$

the limit is also due to Proposition 5.3 in [27] and Proposition 1.

Finally, the lemma follows from (56), (57) and (65). \square

7.5. Proof of Lemma 6. In view of Lemmas 4, it is sufficient to prove $\hat{Z}_\Delta^N \Rightarrow 0$, as $N \rightarrow \infty$. Recall the definition of $A_{(c)}^N$, $\{\tilde{t}_i, i \geq 1\}$ and $\{\tilde{s}_i, i \geq 1\}$ from the proof of Lemma 5. Now, for arbitrary $T > 0$ and $\varepsilon > 0$, we have

$$\mathbb{P}[\|\hat{Z}_\Delta^N\|_T > \varepsilon] \leq \mathbb{P}[\|\hat{Z}_\Delta^N\|_T > \varepsilon, \|\hat{V}_{\leftarrow}^N\|_T \leq c] + \mathbb{P}[\|\hat{V}_{\leftarrow}^N\|_T > c].$$

On the event $\{\|\hat{V}_{\leftarrow}^N\|_T \leq c\}$, the process \hat{Z}_Δ^N , based on its definition, can be upper-bounded as follows, for all $t \in [0, T]$ and all sufficiently large N :

$$\begin{aligned} |\hat{Z}_\Delta^N(t)| &\leq \frac{1}{\sqrt{N}} \sum_{i=1}^{A_{(c)}^N(t)} (1_{\{\tilde{s}_i > t - \tilde{t}_i - \delta\}} - 1_{\{\tilde{s}_i > t - \tilde{t}_i\}}) + \frac{1}{\sqrt{N}} \int_0^t (\bar{F}(t-s-\delta) - \bar{F}(t-s)) dA_{(c)}^N(s) \\ &:= \hat{Z}_{(c,\delta)}^N(t), \end{aligned}$$

where $\delta > 0$. The preceding two inequalities render

$$\mathbb{P}[\|\hat{Z}_\Delta^N\|_T > \varepsilon] \leq \mathbb{P}[\|\hat{Z}_{(c,\delta)}^N\|_T > \varepsilon] + \mathbb{P}[\|\hat{V}_{\leftarrow}^N\|_T > c], \quad (66)$$

where $\hat{Z}_{(c,\delta)}^N = \{\hat{Z}_{(c,\delta)}^N(t), t \geq 0\}$.

Next, Theorem 3 in [20] and (62) yield, as $N \rightarrow \infty$,

$$\{\hat{Z}_{(c,\delta)}^N(t), t \geq 0\} \Rightarrow \left\{ 2c\theta \int_0^t (\bar{F}(t-s-\delta) - \bar{F}(t-s)) ds, t \geq 0 \right\}. \quad (67)$$

On the other hand, Proposition 2 implies

$$\lim_{c \rightarrow \infty} \limsup_{N \rightarrow \infty} \mathbb{P}[\|\hat{V}_{\leftarrow}^N\|_T > c] = 0. \quad (68)$$

Therefore, in view of (66), (67) and (68), given T and ε , for any $\xi > 0$, it is possible to select c and δ such that $\mathbb{P}[\|\hat{Z}_\Delta^N\|_T > \varepsilon] < \xi$ for all N large enough. The lemma holds. \square

7.6. Proof of Proposition 3. Two cases are considered separately: (i) non-deterministic service time (F), and (ii) deterministic service time (F). For $x_1, x_2 \in D[0, \infty)$ let $y_1 = \varphi(x_1)$ and $y_2 = \varphi(x_2)$.

(i) Since service time is not single-valued, there exist $\delta > 0$ and $0 < \varepsilon < 1$ such that $F(x+\delta) - F(x) < \varepsilon$, for all $x \geq 0$. Then it follows that

$$\begin{aligned} d_{L^1}^\delta(y_1, y_2) &\leq d_{L^1}^\delta(x_1, x_2) + \int_0^\delta \int_0^t |y_1(t-s) - y_2(t-s)| dF(s) dt \\ &\leq d_{L^1}^\delta(x_1, x_2) + \int_0^\delta d_{L^1}^\delta(y_1, y_2) dF(s) \\ &\leq d_{L^1}^\delta(x_1, x_2) + \varepsilon d_{L^1}^\delta(y_1, y_2), \end{aligned}$$

and, thus,

$$d_{L^1}^\delta(y_1, y_2) \leq d_{L^1}^\delta(x_1, x_2)/(1-\varepsilon). \quad (69)$$

Similarly, considering the time interval $[0, 2\delta]$ yields

$$\begin{aligned} d_{L^1}^{2\delta}(y_1, y_2) &\leq d_{L^1}^{2\delta}(x_1, x_2) + \varepsilon d_{L^1}^\delta(y_1, y_2) + \varepsilon d_{L^1}^{2\delta}(y_1, y_2) \\ &\leq d_{L^1}^{2\delta}(x_1, x_2)/(1-\varepsilon) + \varepsilon d_{L^1}^{2\delta}(y_1, y_2), \end{aligned}$$

where the second inequality is due to (69). From the preceding inequality one derives $d_{L^1}^{2\delta}(y_1, y_2) \leq d_{L^1}^{2\delta}(x_1, x_2)/(1-\varepsilon)^2$. The above argument can be applied l times iteratively to obtain $d_{L^1}^{l\delta}(y_1, y_2) \leq d_{L^1}^{l\delta}(x_1, x_2)/(1-\varepsilon)^l$. Therefore, for any T , there exists $c_T < \infty$ such that $d_{L^1}^T(y_1, y_2) \leq c_T d_{L^1}^T(x_1, x_2)$.

(ii) Let a be such that $F(a-) = 0$ and $F(a) = 1$. Then $y_i(t) = x_i(t)$, $i = 1, 2$, for $t < a$ and $d_{L^1}^T(y_1, y_2) = d_{L^1}^T(x_1, x_2)$ for $T < a$. Next, assume that $d_{L^1}^T(y_1, y_2) \leq c_T d_{L^1}^T(x_1, x_2)$ for some $T \geq a$ and $c_T < \infty$. Due to this assumption, since $y_i(t) = x_i(t) + y_i^+(t - a)$, $i = 1, 2$, for $t \geq a$, one has, for $0 < d \leq a$,

$$\begin{aligned} d_{L^1}^{T+d}(y_1, y_2) &\leq d_{L^1}^{T+d}(x_1, x_2) + d_{L^1}^T(y_1, y_2) \\ &\leq d_{L^1}^{T+d}(x_1, x_2) + c_T d_{L^1}^T(x_1, x_2) \\ &\leq (1 + c_T) d_{L^1}^{T+d}(x_1, x_2). \end{aligned}$$

The conclusion follows. \square

7.7. Proof of Lemma 8. In view of Proposition 3, it is sufficient to consider the argument of the φ operator in (21). Recall the definition of $d_{J_1}(\cdot, \cdot)$ from the proof of Proposition 3.

The non-decreasing nature of distribution functions yields

$$\begin{aligned} d_{L^1}^T(F \circ (\tau^N + V_{\leftarrow}^N), F) &\leq \int_0^T (F(t + \|V_{\leftarrow}^N\|_T) - F(t - \|e - \tau^N\|_T)) dt \\ &\leq \|V_{\leftarrow}^N\|_T + \|e - \tau^N\|_T, \end{aligned}$$

where the second inequality follows from $F(t) \leq 1$ for all t ; similarly,

$$d_{L^1}^T(F_* \circ (\tau^N + V_{\leftarrow}^N), F_*) \leq \|V_{\leftarrow}^N\|_T + \|e - \tau^N\|_T.$$

For notational simplicity, let $\hat{J}^N := (\hat{q}_0^N)^+(\bar{F} - \bar{F}_*) + \hat{q}_0^N \bar{F}_* - \sqrt{N}(1 - \rho^N)F_*$. The preceding two inequalities, jointly with (5) and (7), yield, as $N \rightarrow \infty$,

$$d_{L^1}^T(\hat{J}^N \circ (\tau^N + V_{\leftarrow}^N), \hat{J}^N) \Rightarrow 0. \quad (70)$$

The triangle inequality and the definition of $d_{L^1}^T$ (see (1)) result in

$$\begin{aligned} d_{L^1}^T((\hat{I}_{\Delta}^N + \hat{Y}_{\Delta}^N - \hat{Z}^N) \circ (\tau^N + V_{\leftarrow}^N), \hat{I}_{\Delta}^N + \hat{Y}_{\Delta}^N - \hat{Z}^N) \\ \leq d_{L^1}^T((\hat{I}_{\Delta}^N + \hat{Y}_{\Delta}^N - \hat{Z}^N) \circ (\tau^N + V_{\leftarrow}^N), 0) + d_{L^1}^T(\hat{I}_{\Delta}^N + \hat{Y}_{\Delta}^N - \hat{Z}^N, 0) \\ \leq 2T \|\hat{I}_{\Delta}^N + \hat{Y}_{\Delta}^N - \hat{Z}^N\|_{T+V_{\leftarrow}^N(T)}, \end{aligned}$$

and, thus, invoking Lemmas 5 and 6, as well as Proposition 2, yields, as $N \rightarrow \infty$,

$$d_{L^1}^T((\hat{I}_{\Delta}^N + \hat{Y}_{\Delta}^N - \hat{Z}^N) \circ (\tau^N + V_{\leftarrow}^N), \hat{I}_{\Delta}^N + \hat{Y}_{\Delta}^N - \hat{Z}^N) \Rightarrow 0. \quad (71)$$

Next, for any $\varepsilon > 0$ and $\delta > 0$, conditioning on the value of $\|\tau^N + V_{\leftarrow}^N - e\|_T$ results in

$$\begin{aligned} \mathbb{P}[d_{L^1}^T((\hat{X}^N + \hat{I}^N) \circ (\tau^N + V_{\leftarrow}^N), \hat{X}^N + \hat{I}^N) > \varepsilon] \\ \leq \mathbb{P}[d_{L^1}^T((\hat{X}^N + \hat{I}^N) \circ (\tau^N + V_{\leftarrow}^N), \hat{X}^N + \hat{I}^N) > \varepsilon, \|\tau^N + V_{\leftarrow}^N - e\|_T \leq \delta] \\ \quad + \mathbb{P}[\|\tau^N + V_{\leftarrow}^N - e\|_T > \delta] \\ \leq \mathbb{P}\left[\left\| \sup_{|s| \leq \delta} |\hat{X}^N(t+s) - \hat{X}^N(t) + \hat{I}^N(t+s) - \hat{I}^N(t)| \right\|_T > \varepsilon/T\right] \\ \quad + \mathbb{P}[\|\tau^N + V_{\leftarrow}^N - e\|_T > \delta]. \end{aligned} \quad (72)$$

Lemmas 2 and 3, the continuous mapping theorem, the continuity of the sup operator and the continuity of sample paths of \hat{X} and \hat{I} yield

$$\lim_{\delta \downarrow 0} \lim_{N \rightarrow \infty} \mathbb{P}\left[\left\| \sup_{|s| \leq \delta} |\hat{X}^N(t+s) - \hat{X}^N(t) + \hat{I}^N(t+s) - \hat{I}^N(t)| \right\|_T > \varepsilon/T\right] = 0.$$

The preceding limit, Proposition 2 and (72) imply, as $N \rightarrow \infty$,

$$d_{L^1}^T((\hat{X}^N + \hat{I}^N) \circ (\tau^N + V_{\leftarrow}^N), \hat{X}^N + \hat{I}^N) \Rightarrow 0. \quad (73)$$

Now, considering separately $s \in [0, t \wedge (\tau^N(t) + V_{\leftarrow}^N(t))]$ and $s \in (t \wedge (\tau^N(t) + V_{\leftarrow}^N(t)), t \vee (\tau^N(t) + V_{\leftarrow}^N(t))]$ we have

$$\begin{aligned} & \int_0^T \left| \int_0^{\tau^N(t) + V_{\leftarrow}^N(t)} \bar{F}(\tau^N(t) + V_{\leftarrow}^N(t) - s) \sqrt{N} G^N(V_{\leftarrow}^N(s)) d\check{A}^N(s) - \int_0^t \bar{F}(t - s) \sqrt{N} G^N(V_{\leftarrow}^N(s)) d\check{A}^N(s) \right| dt \\ & \leq \sqrt{N} G^N(\|V_{\leftarrow}^N\|_T) \int_0^T \int_0^{t \wedge (\tau^N(t) + V_{\leftarrow}^N(t))} (F(t + \|V_{\leftarrow}^N\|_T - s) - F(t - \|\tau^N - e\|_T - s)) d\check{A}^N(s) dt \\ & \quad + \sqrt{N} G^N(\|V_{\leftarrow}^N\|_T) \int_0^T (\check{A}^N(t + \|V_{\leftarrow}^N\|_T) - \check{A}^N(t - \|\tau^N - e\|_T)) dt \\ & \leq \sqrt{N} G^N(\|V_{\leftarrow}^N\|_T) (\|V_{\leftarrow}^N\|_T + \|\tau^N - e\|_T) \check{A}^N(T) \\ & \quad + \sqrt{N} G^N(\|V_{\leftarrow}^N\|_T) (\check{A}^N(T + V_{\leftarrow}^N(T)) - \check{A}^N(T)), \end{aligned} \quad (74)$$

where the last inequality is due to a change in the order of integration. The preceding inequality, (6), (3), and Proposition 2 result in, as $N \rightarrow \infty$,

$$\sqrt{N} G^N(\|V_{\leftarrow}^N\|_T) [(\|V_{\leftarrow}^N\|_T + \|\tau^N - e\|_T) \check{A}^N(T) + (\check{A}^N(T + V_{\leftarrow}^N(T)) - \check{A}^N(T))] \Rightarrow 0. \quad (75)$$

Finally, the statement of the lemma follows from (17), (70), (71), (73), (74), (75) and Proposition 3. \square

7.8. Proof of Proposition 4. The proof closely parallels the proof of Proposition 3.1 in [27]. Two cases are considered separately: (i) deterministic and (ii) non-deterministic service times. The proof of measurability is the same for the two cases and is identical to the corresponding proof in Proposition 3.1 of [27].

(i) Deterministic F . In this case $F(t) = 1_{\{t \geq a\}}$, $F_*(t) = t/a \cdot 1_{\{0 \leq t \leq a\}}$ and $\mu = 1/a$.

Existence. First consider the interval $[0, a)$ only. Let $y_0 = 0$ and

$$y_{n+1}(t) = x(t) - \theta \int_0^t y_n^+(t - s) ds \quad (76)$$

for $0 \leq t < a$ and $n \geq 1$. Then for $\delta < a$ we have

$$\begin{aligned} \|y_{n+1} - y_n\|_\delta & \leq \delta \theta \|y_n - y_{n-1}\|_\delta \\ & \leq (\delta \theta)^n \|x\|_\delta. \end{aligned}$$

The preceding will serve as a base for an induction. Assume that

$$\|y_{n+1} - y_n\|_{k\delta} \leq n^{k-1} (\delta \theta)^n \|x\|_{k\delta} \quad (77)$$

for some k ($k\delta < a$). Then, for $(k+1)\delta < a$, the inductive assumption and (76) yield

$$\begin{aligned} \|y_{n+1} - y_n\|_{(k+1)\delta} & \leq \delta \theta \sum_{i=1}^k \|y_n - y_{n-1}\|_{i\delta} + \delta \theta \|y_n - y_{n-1}\|_{(k+1)\delta} \\ & \leq (\delta \theta)^n \|x\|_{(k+1)\delta} \sum_{i=1}^k (n-1)^{i-1} + \delta \theta \|y_n - y_{n-1}\|_{(k+1)\delta} \\ & \leq n^{k-1} (\delta \theta)^n \|x\|_{(k+1)\delta} + \delta \theta \|y_n - y_{n-1}\|_{(k+1)\delta}. \end{aligned}$$

Iterating the argument from the preceding inequality results in

$$\|y_{n+1} - y_n\|_{(k+1)\delta} \leq n^k (\delta \theta)^n \|x\|_{(k+1)\delta},$$

and hence (77) holds. In view of (77), selecting $\delta < 1/\theta$ implies that $\{y_n, n \geq 0\}$ is a Cauchy sequence and there exists y such that $y_n \rightarrow y$, as $n \rightarrow \infty$. Therefore there exists a solution on the interval $[0, a)$.

Now consider the interval $[0, 2a)$. Let $y_0 = \{y_0(t) = y(t)1_{\{0 \leq t < a\}}, 0 \leq t < 2a\}$ and

$$y_{n+1}(t) = \begin{cases} y(t), & 0 \leq t < a, \\ x(t) + y(t-a) - \theta \int_{t-a}^t y_n^+(t-s) ds, & a \leq t < 2a, \end{cases}$$

where y is the solution on the interval $[0, a)$. By repeating the argument from the previous case, it is straightforward to show that there exists a solution on the interval $[0, 2a)$. Furthermore, by iterating the argument, one establishes the existence of a solution on an arbitrary interval of finite length.

Uniqueness. Let $\delta < a \wedge 1/\theta$. Suppose u and v are two solutions and consider

$$u(t) - v(t) = 1_{\{t \geq a\}}(u^+(t) - v^+(t)) - \theta \int_0^{t \wedge a} (u^+(t-s) - v^+(t-s)) ds,$$

$t \geq 0$. For $0 \leq t \leq \delta$ we have $|u(t) - v(t)| \leq \delta\theta\|u - v\|_\delta$, and, therefore, $u(t) = v(t)$ for $0 \leq t \leq \delta$. Next $|u(t) - v(t)| \leq \delta\theta\|u - v\|_\delta + \delta\theta\|u - v\|_{2\delta}$ for $\delta < t \leq 2\delta$, yielding $u(t) = v(t)$ for $0 \leq t \leq 2\delta$. Repeating this argument multiple times leads to $u(t) = v(t)$ for $0 \leq t \leq a$.

Now, assume that $u(t) = v(t)$ for $0 \leq t \leq T$, where $T \geq a$. Then, for $T < t \leq T + \delta$, we have $|u(t) - v(t)| \leq \delta\theta\|u - v\|_{T+\delta}$, resulting in $u(t) = v(t)$ for $0 \leq t \leq T + \delta$. The uniqueness follows.

Lipschitz continuity. The definition of ϕ renders, for $y = \phi(x)$ and $t < a$,

$$y(t) = x(t) - \theta \int_0^t y^+(t-s) ds$$

and, thus, $\|\phi(x_1) - \phi(x_2)\|_\delta \leq \|x_1 - x_2\|_\delta + \delta\theta\|\phi(x_1) - \phi(x_2)\|_\delta$ if $\delta < t$. By selecting $\delta > 0$ small enough such that $\delta\theta < 1$, we have

$$\|\phi(x_1) - \phi(x_2)\|_\delta \leq \|x_1 - x_2\|_\delta / (1 - \delta\theta). \quad (78)$$

Considering the interval $[0, 2\delta]$ yields

$$\|\phi(x_1) - \phi(x_2)\|_{2\delta} \leq \|x_1 - x_2\|_{2\delta} + \delta\theta\|\phi(x_1) - \phi(x_2)\|_\delta + \delta\theta\|\phi(x_1) - \phi(x_2)\|_{2\delta},$$

which, upon combining with (78), results in

$$\|\phi(x_1) - \phi(x_2)\|_{2\delta} \leq \|x_1 - x_2\|_{2\delta} / (1 - \delta\theta)^2.$$

The preceding argument can be applied repeatedly to show that ϕ is Lipschitz continuous when the interval $[0, a)$ is considered.

For $t \geq a$, $y = \phi(x)$ renders

$$y(t) = x(t) + y^+(t-a) - \theta \int_0^a y^+(t-s) ds. \quad (79)$$

When $t = a$, we obtain

$$y(a) = x(a) + x^+(0) - \theta \int_0^a y^+(s) ds,$$

and, due to the case $t < a$, it follows that there exists $c_a < \infty$ such that $\|\phi(x_1) - \phi(x_2)\|_a \leq c_a\|x_1 - x_2\|_a$. This serves as the base for the induction. Now, suppose that for some $T \geq a$ there exists $c_T < \infty$ such that $\|\phi(x_1) - \phi(x_2)\|_T \leq c_T\|x_1 - x_2\|_T$. Now, for any $\delta < \min\{a, 1/\theta\}$, from (79) we have

$$\begin{aligned} \|\phi(x_1) - \phi(x_2)\|_{T+\delta} &\leq \|x_1 - x_2\|_{T+\delta} + (1 + a\theta)\|\phi(x_1) - \phi(x_2)\|_T + \delta\theta\|\phi(x_1) - \phi(x_2)\|_{T+\delta} \\ &\leq (1 + (1 + a\theta)c_T)\|x_1 - x_2\|_{T+\delta} + \delta\theta\|\phi(x_1) - \phi(x_2)\|_{T+\delta}, \end{aligned}$$

where the second inequality is due to the inductive assumption. Hence, $\|\phi(x_1) - \phi(x_2)\|_{T+\delta} \leq c_{T+\delta}\|x_1 - x_2\|_{T+\delta}$ with $c_{T+\delta} = (1 + (1 + a\theta)c_T)/(1 - \delta\theta) < \infty$.

(ii) Non-deterministic F .

There exist $\delta > 0$ and $0 < \varepsilon < 1$ such that

$$F(t + \delta) - F(t) + \theta F_*(t + \delta)/\mu - \theta F_*(t)/\mu < \varepsilon, \tag{80}$$

for all $t \geq 0$, since F_* is absolutely continuous by definition. In view of this fact, the proof of existence, uniqueness and Lipschitz continuity is almost identical to the proof of corresponding parts in Proposition 3.1 of [27]. In particular, if $\tilde{F} := F - \theta F_*/\mu$ then

$$y(t) = x(t) + \int_0^t y^+(t - s) d\tilde{F}.$$

Note that the preceding relation can be written in terms of φ with F replaced by \tilde{F} , and, in view of (80), there exist $\delta > 0$ and $0 < \varepsilon < 1$ such that

$$\tilde{F}(t + \delta) - \tilde{F}(t) < \varepsilon, \tag{81}$$

for all $t \geq 0$. We can now apply directly the results in [27, Proposition 3.1] because the analysis of φ in [27] is based on (81). \square

ACKNOWLEDGMENTS

The authors thank an anonymous referee for careful comments that improved the presentation.

REFERENCES

- [1] F. Baccelli and P. Bremaud. *Elements of Queueing Theory*. Springer-Verlag, Berlin, 2nd edition, 2003. [7.2](#)
- [2] P. Bhattacharya and A. Ephremides. Stochastic monotonicity properties of multiserver queues with impatient customers. *J. Appl. Probab.*, 28:673–682, 1991. [3.1](#), [1](#)
- [3] P. Billingsley. *Probability and Measure*. Wiley, 3rd edition, 1995. [7.3](#)
- [4] A.A. Borovkov. On limit laws for service processes in multi-channel systems. *Siberian Math. J.*, 8:746–763, 1967. [3.2](#), [3.5](#)
- [5] S. Browne and W. Whitt. Piecewise-linear diffusion processes. In J.H. Dshalalow, editor, *Probability and Stochastic Series: Advances in Queueing: Theory, Methods and Open Problems*. CRC Press, Boca Raton, FL, 1995. [6](#)
- [6] K. L. Chung. *A Course in Probability Theory*. Academic Press, 2nd edition, 1974. [7.4](#)
- [7] J. Dai and S. He. Customer abandonment in many-server queues. Preprint. [1](#)
- [8] L. Decreusefond and P. Moyal. A functional central limit theorem for the M/GI/ ∞ queue. *Ann. Appl. Probab.*, 18(6):2156–2178, 2008. [3.2](#)
- [9] A.K. Erlang. On the rational determination of the number of circuits. In E. Brockmeyer, H.L. Halstrom, and A. Jensen, editors, *The life and works of A.K. Erlang*. The Copenhagen Telephone Company, Copenhagen, 1948. [1](#)
- [10] D. Gamarnik and P. Momčilović. Steady-state analysis of a multi-server queue in the Halfin-Whitt regime. *Adv. Appl. Probab.*, 40(2):548–577, 2008. [6](#)
- [11] N. Gans, G. Koole, and A. Mandelbaum. Telephone call centers: Tutorial, review, and research prospects. *Manufacturing and Service Operations Management*, 5(2):79–141, 2003. [1](#), [1](#)
- [12] O. Garnett, A. Mandelbaum, and M. Reiman. Designing a call center with impatient customers. *Manufacturing and Service Operations Management*, 4(3):208–227, 2002. [1](#), [6](#)
- [13] R.A. Green, P.C. Wyer, and J. Giglio. ED walkout rate correlated with ED length of stay but not with ED volume or hospital census (abstract). *Academic Emergency Medicine*, 9:514, 2002. [1](#)
- [14] S. Halfin and W. Whitt. Heavy-traffic limits for queues with many exponential servers. *Oper. Res.*, 29(3):567–588, 1981. [1](#)
- [15] D.L. Iglehart. Weak convergence of compound stochastic processes. *Stochastic Process. Appl.*, 1:11–31, 1973. [3.2](#)
- [16] D. Jagerman. Some properties of the Erlang loss function. *Bell System Techn. J.*, 53(3):525–551, 1974. [1](#)
- [17] P. Jelenković, A. Mandelbaum, and P. Momčilović. Heavy traffic limits for queues with many deterministic servers. *Queueing Syst. Theory Appl.*, 47(1-2):53–69, 2004. [3](#)

- [18] W. Kang and K. Ramanan. Fluid limits of many-server queues with reneging. Preprint. [1](#)
- [19] H. Kaspi and K. Ramanan. Law of large numbers limits for many-server queues. *Ann. Appl. Probab.*, to appear. [1](#)
- [20] E. Krichagina and A. Puhalskii. A heavy-traffic analysis of a closed queueing system with a GI/∞ service center. *Queueing Syst. Theory Appl.*, 25(1-4):235–280, 1997. [3.2](#), [2](#), [3.3](#), [3.5](#), [4.1](#), [2](#), [7.4](#), [7.5](#)
- [21] A. Mandelbaum and P. Momčilović. Queues with many servers: The virtual waiting-time process in the QED regime. *Math. Oper. Res.*, 33(3):561–586, 2008. [1](#), [2.3](#), [5](#)
- [22] A. Mandelbaum and S. Zeltyn. The impact of customers’ patience on delay and abandonment: Some empirically-driven experiments with the $M/M/N+G$ queue. *OR Spectrum*, 26(3):377–411, 2004. Special Issue on Call Centers. [1](#)
- [23] A. Mandelbaum and S. Zeltyn. The $M/M/n+G$ queue: Summary of performance measures. Available at <http://iew3.technion.ac.il/serveng/References/>, 2007. [1](#), [6](#)
- [24] A. Mandelbaum and S. Zeltyn. Staffing many-server queues with impatient customers: Constraint satisfaction in call centers. *Oper. Res.*, 57(5):1189–1205, 2009. [1](#)
- [25] A. Puhalskii. On the invariance principle for the first passage time. *Math. Oper. Res.*, 19(4):946–954, 1994. [1](#), [4.2](#)
- [26] A. Puhalskii and J. Reed. On many-server queues in heavy traffic. *Ann. Appl. Probab.*, 20(1):129–195, 2010. [1](#)
- [27] J. Reed. The $G/GI/N$ queue in the Halfin-Whitt regime. *Ann. Appl. Probab.*, 19(6):2211–2269, 2009. [1](#), [2.3](#), [3.1](#), [3.1](#), [3.5](#), [3.5](#), [3.5](#), [1](#), [3.5](#), [4](#), [4.1](#), [7.3](#), [7.3](#), [7.4](#), [7.8](#), [7.8](#), [7.8](#)
- [28] J. Reed and R. Talreja. Distribution-valued heavy-traffic limits for the $G/GI/\infty$ queue. Preprint. [3.2](#)
- [29] J. Reed and T. Tezcan. Hazard rate scaling for the $GI/M/n+GI$ queue. Preprint. [1](#)
- [30] R. Talreja and W. Whitt. Heavy-traffic limits for waiting times in many-server queues with abandonment. *Ann. Appl. Probab.*, 19(6):2137–2175, 2009. [1](#), [4.2](#)
- [31] W. Whitt. *Stochastic-Process Limits: An Introduction to Stochastic-Process Limits and Their Application to Queues*. Springer, New York, 2002. [4.1](#)
- [32] W. Whitt. Efficiency-driven heavy-traffic approximations for many-server queues with abandonments. *Management Science*, 50(10):1449–1461, 2004. [1](#)
- [33] S. Zeltyn. *Call centers with impatient customers: Exact analysis and many-server asymptotics of the $M/M/n+G$ queue*. PhD thesis, Technion, Haifa, Israel, 2005. [1](#), [5](#), [6](#)
- [34] S. Zeltyn and A. Mandelbaum. Call centers with impatient customers: Many-server asymptotics of the $M/M/n+G$ queue. *Queueing Syst. Theory Appl.*, 51(3-4):361–402, 2005. [1](#), [5](#)

FACULTY OF INDUSTRIAL ENGINEERING AND MANAGEMENT, TECHNION, HAIFA 3200, ISRAEL
E-mail address: avim@tx.technion.ac.il

DEPARTMENT OF ELECTRICAL ENGINEERING AND COMPUTER SCIENCE, UNIVERSITY OF MICHIGAN, ANN ARBOR, MI 48109
E-mail address: petar@eecs.umich.edu