

ASYMPTOTIC APPROXIMATIONS FOR STATIONARY DISTRIBUTIONS OF MANY-SERVER QUEUES WITH ABANDONMENT

WEINING KANG AND KAVITA RAMANAN

ABSTRACT. A many-server queueing system is considered in which customers arrive according to a renewal process, and have service and patience times that are drawn from two independent sequences of independent, identically distributed random variables. Customers enter service in the order of arrival and are assumed to abandon the queue if the waiting time in queue exceeds the patience time. The state $Y^{(N)}$ of the system with N servers is represented by a four-component process that consists of the backward recurrence time of the arrival process, a pair of measure-valued processes, one that keeps track of the waiting times of customers in queue and the other that keeps track of the amounts of time customers present in the system have been in service, and a real-valued process that represents the total number of customers in the system. Under general assumptions, it is first shown that $Y^{(N)}$ is a Feller process, admits a stationary distribution and is ergodic. The main result shows that when the associated fluid limit has a unique invariant state then any sequence $\{Y^{(N)}/N\}_{N \in \mathbb{N}}$ of stationary distributions of the scaled processes converges, as $N \rightarrow \infty$, to this state. In addition, a simple example is given to illustrate that, both in the presence and absence of abandonments, the $N \rightarrow \infty$ and $t \rightarrow \infty$ limits cannot always be interchanged. The stationary behavior of many-server systems is of interest for performance analysis of computer data systems and call centers.

CONTENTS

1. Introduction	2
2. Description of Model and State Dynamics	5
3. Assumptions and Main Results	9
4. Stationary Distribution of the N -Server Queue	12
5. Fluid limit	19
6. The Limit of Scaled Stationary Distributions	24
7. Concluding Remarks	28
Appendix A. Proof of Theorem 7.1	32
References	35

Date: March 15, 2010.

2000 Mathematics Subject Classification. Primary: 60K25, 68M20, 90B22; Secondary: 60F99.

Key words and phrases. Multi-server queues, stationary distribution, ergodicity, measure-valued processes, abandonment, reneging, interchange of limits, mean-field limits, call centers.

Partially supported by the National Science Foundation under Grants CMMI-0728064, CMMI-0928154.

1. INTRODUCTION

1.1. Description. An N -server queueing system is considered in which customers arrive according to a renewal process, have independent and identically distributed (i.i.d.) service requirements that are drawn from a general distribution with finite mean and also carry i.i.d. patience times that are drawn from another general distribution. Customers enter service, in the order of arrival, as soon as an idle server is available, service is non-preemptive and customers abandon the queue if the time spent waiting in queue reaches the patience time. This system is also sometimes referred to as the GI/G/N+G model. In this work, it is assumed that the sequences of service requirements and patience times are mutually independent, and that the interarrival, service and patience time distributions have densities.

The state of the N -server system is represented by a four component process $Y^{(N)}$, consisting of the backward recurrence time process associated with the renewal arrival process, a measure-valued process that keeps track of the amounts of time customers currently in service have been in service, another measure-valued process that encodes the times elapsed since customers have entered the system (for all customers for which this time has not yet exceeded their patience times), and a real-valued process that keeps track of the total number of customers in the system. This infinite-dimensional state representation was shown in Lemma B.1 of [14] to lead to a Markovian description of the dynamics (with respect to a suitable filtration). In addition, a fluid limit for this model was also established in [14], i.e., under suitable assumptions, it was shown that almost surely, $\bar{Y}^{(N)} = Y^{(N)}/N$ converges, as $N \rightarrow \infty$, to the fluid limit \bar{Y} , which is characterized as the unique solution to a set of coupled integral equations (see Definition 5.1).

The present work focuses on obtaining first-order approximations to the stationary distribution of $Y^{(N)}$, which is of fundamental interest for the performance analysis of many-server queues. In particular, this work addresses several questions that were raised in Whitt [21]. It is first shown that for each N , $Y^{(N)}$ is a Feller, strong Markov process and has a stationary distribution. Under an additional assumption (see Assumption 6), the ergodicity of each $Y^{(N)}$ is also established. The main result, Theorem 3.3, shows that if the fluid limit has a unique invariant state, then any sequence of scaled stationary distributions converges, as $N \rightarrow \infty$, to this unique invariant state. More generally, this work seeks to illustrate how an infinite-dimensional Markovian representation of a stochastic network can facilitate the (first-order) characterization of the associated stationary distributions. Moreover, examples are presented in Section 7 to illustrate some subtleties in the dynamics and to show that the $t \rightarrow \infty$ and $N \rightarrow \infty$ limits cannot in general be interchanged.

1.2. Motivation and Context. The study of many-server queueing systems with abandonment is motivated by applications to telephone call centers and (more generally) customer contact centers. The incorporation of customer abandonment captures the effect of customers' impatience, which has a substantial impact on the performance of the system. For example, customer abandonment can stabilize a system even when it is overloaded. A considerable body of work has been devoted to the study of various steady-state or stationary performance measures of many-server queues, both with and without abandonment. In the absence of abandonment, when the interarrival times and service times are exponential, an explicit expression for the steady state queue length can be found in [4], while the classical work of Kiefer

and Wolfowitz [16] (see also Foss [7]) establishes the convergence of waiting time processes (or vectors) in discrete time when the i.i.d. interarrival and service times are generally distributed. The case of continuous time is dealt with in [2]. For a many-server queue with stationary renewal arrivals, deterministic service times and no abandonments, Jelenkovic, Mandelbaum and Momčilović [13] showed that on the diffusive scale, the scaled stationary waiting times converge in distribution to the supremum of a Gaussian random walk with negative drift. For a many-server queue with stationary renewal arrivals, a finitely supported, lattice-valued service time distribution and no abandonments, Gamarnik and Momčilović [8] characterized the limiting scaled stationary queue length distribution in terms of the stationary distribution of an explicitly constructed Markov chain and obtained an explicit expression for the exponential decay rate of the moment generating function of the limiting stationary queue length.

For many-server queues with abandonment whose interarrival, service and abandonment distributions are exponential, Garnett, Mandelbaum and Reiman [10] provide exact calculations of various steady state performance measures and their approximations in the diffusive scale, both in the case of finite waiting rooms (M/M/N/B+M) and infinite waiting rooms (M/M/N+M). In the case of Poisson arrivals, exponential service distribution and general abandonment distribution (M/M/N+G), explicit formulae for the steady state distributions of the queue length and virtual waiting time were obtained by Baccelli and Hebuterne [3] (see Sections IV and V.2 therein), while several other steady state performance measures and their asymptotic approximations, in the limit as the arrival rates and servers go to infinity, were derived by Mandelbaum and Zeltyn [18].

In the previously mentioned works on characterization of stationary distributions of many-server queues, either the interarrival times and service times are assumed to be exponential or the service times are discrete and there is no abandonment. However, statistical analysis of real call centers has shown that both service times and patience times are typically not exponentially distributed [5, 18], thus providing strong motivation for this work. In general, it is difficult to derive explicit expressions for the stationary distributions of many-server queues, especially in the realistic case when service times are non-exponential and there is abandonment. This is also the case for many other classes of stochastic networks. To resolve this issue, a common approach that is taken is to identify the long time limit of the fluid or diffusion approximations, which is often more tractable, and then use this limit as an approximation for the stationary distribution of the original system. Such an approach relies on the premise that the long time behavior of the fluid limit can be characterized, and also requires an argument that justifies the interchange of (the $N \rightarrow \infty$ and $t \rightarrow \infty$) limits (see, for example, [9] for an interchange of limits result in the context of generalized Jackson networks). However, we show that this approach may not always be appropriate for stochastic network models. Indeed, for the case of many-server queues with non-exponential service distributions, the long-time behavior of the fluid is subtle and difficult to characterize, in large part due to the complexity in the dynamics introduced by the coupling of the measure-valued component of the fluid limit with the positive real-valued component by the non-idling condition. Furthermore, as the example we construct in Section 7 demonstrates, in general the order of the $N \rightarrow \infty$ and $t \rightarrow \infty$ limits cannot be interchanged. Instead, we take a different approach to showing convergence, which uses

a representation formula for the dynamics of the measure-valued state processes (see Proposition 2.1).

1.3. Outline. The outline of the paper is as follows. A precise mathematical description of the model is provided in Section 2. Section 3 introduces the basic assumptions and states the main result. The Feller property and the existence of stationary distributions of the state descriptor are proved in Section 4. The fluid equations and the invariant manifold are described in Section 5, and the asymptotics of the stationary distributions is established in Section 6. Finally, positive Harris recurrence and ergodicity of the state descriptor, the long time behavior of the fluid limit and an example that shows that the “interchange of limits” property does not always hold are discussed in Section 7. In the remainder of this section, we introduce some common notation used in the paper.

1.4. Notation and Terminology. The following notation will be used throughout the paper. \mathbb{Z} is the set of integers, \mathbb{N} is the set of positive integers, \mathbb{R} is the set of real numbers, \mathbb{Z}_+ is the set of non-negative integers and \mathbb{R}_+ the set of non-negative real numbers. For $a, b \in \mathbb{R}$, $a \vee b$ denotes the maximum of a and b , $a \wedge b$ the minimum of a and b and the short-hand a^+ is used for $a \vee 0$. $\mathbb{1}_B$ denotes the indicator function of the set B (that is, $\mathbb{1}_B(x) = 1$ if $x \in B$ and $\mathbb{1}_B(x) = 0$ otherwise).

1.4.1. Function and Measure Spaces. Given any metric space E , $\mathcal{C}_b(E)$ and $\mathcal{C}_c(E)$ are, respectively, the space of bounded, continuous functions and the space of continuous real-valued functions with compact support defined on E . Given a nondecreasing real function f on $[0, \infty)$, f^{-1} denotes the inverse function of f in the sense that $f^{-1}(y) = \inf\{x \geq 0 : f(x) \geq y\}$. The support of a function φ is denoted by $\text{supp}(\varphi)$.

The space of Radon measures on a metric space E , endowed with the Borel σ -algebra, is denoted by $\mathcal{M}(E)$, while $\mathcal{M}_F(E)$ is the subspace of finite measures in $\mathcal{M}(E)$. Recall that a Radon measure is one that assigns finite measure to every relatively compact subset of \mathbb{R}_+ . By identifying a Radon measure $\mu \in \mathcal{M}(E)$ with the mapping on $\mathcal{C}_c(E)$ defined by

$$\varphi \mapsto \int_E \varphi(x) \mu(dx),$$

one can equivalently define a Radon measure on E as a linear mapping from $\mathcal{C}_c(E)$ into \mathbb{R} such that for every compact set $\mathcal{K} \subset E$, there exists $L_{\mathcal{K}} < \infty$ such that

$$\left| \int_E \varphi(x) \mu(dx) \right| \leq L_{\mathcal{K}} \|\varphi\|_{\infty} \quad \forall \varphi \in \mathcal{C}_c(E) \text{ with } \text{supp}(\varphi) \subset \mathcal{K}.$$

The space $\mathcal{M}_F(E)$ is equipped with the weak topology, i.e., a sequence of measures $\{\mu_n\}$ in $\mathcal{M}_F(E)$ is said to converge to μ in the weak topology (denoted $\mu_n \xrightarrow{w} \mu$) if and only if for every $\varphi \in \mathcal{C}_b(E)$,

$$(1.1) \quad \int_E \varphi(x) \mu_n(dx) \rightarrow \int_E \varphi(x) \mu(dx) \quad \text{as } n \rightarrow \infty.$$

As is well-known, $\mathcal{M}_F(E)$, endowed with the weak topology, is a Polish space. The symbol δ_x will be used to denote the measure with unit mass at the point x and, by some abuse of notation, we will use $\mathbf{0}$ to denote the identically zero Radon measure on E . When E is an interval, say $[0, H)$, for notational conciseness, we will often write $\mathcal{M}[0, H)$ instead of $\mathcal{M}([0, H))$. For any Borel measurable function

$f : [0, H) \rightarrow \mathbb{R}$ that is integrable with respect to $\xi \in \mathcal{M}[0, H)$, we often use the short-hand notation

$$\langle f, \xi \rangle \doteq \int_{[0, H)} f(x) \xi(dx).$$

Also, for ease of notation, given $\xi \in \mathcal{M}[0, H)$ and an interval $(a, b) \subset [0, M)$, we will use $\xi(a, b)$ to denote $\xi((a, b))$.

1.4.2. Measure-valued Stochastic Processes. Given a Polish space \mathcal{H} , we denote by $\mathcal{D}_{\mathcal{H}}[0, T]$ (respectively, $\mathcal{D}_{\mathcal{H}}[0, \infty)$) the space of \mathcal{H} -valued, càdlàg functions on $[0, T]$ (respectively, $[0, \infty)$), and we endow this space with the usual Skorokhod J_1 -topology [20]. Then $\mathcal{D}_{\mathcal{H}}[0, T]$ and $\mathcal{D}_{\mathcal{H}}[0, \infty)$ are also Polish spaces (see [20]). In this work, we will be interested in \mathcal{H} -valued stochastic processes, where $\mathcal{H} = \mathcal{M}_F[0, H)$ for some $H \leq \infty$. These are random elements that are defined on a probability space $(\Omega, \mathcal{F}, \mathbb{P})$ and take values in $\mathcal{D}_{\mathcal{H}}[0, \infty)$, equipped with the Borel σ -algebra (generated by open sets under the Skorokhod J_1 -topology). A sequence $\{X_n\}$ of càdlàg, \mathcal{H} -valued processes, with X_n defined on the probability space $(\Omega_n, \mathcal{F}_n, \mathbb{P}_n)$, is said to converge in distribution to a càdlàg \mathcal{H} -valued process X defined on $(\Omega, \mathcal{F}, \mathbb{P})$ if, for every bounded, continuous functional $F : \mathcal{D}_{\mathcal{H}}[0, \infty) \rightarrow \mathbb{R}$, we have

$$\lim_{n \rightarrow \infty} \mathbb{E}_n [F(X_n)] = \mathbb{E} [F(X)],$$

where \mathbb{E}_n and \mathbb{E} are the expectation operators with respect to the probability measures \mathbb{P}_n and \mathbb{P} , respectively. Convergence in distribution of X_n to X will be denoted by $X_n \Rightarrow X$. Let $\mathcal{I}_{\mathbb{R}_+}[0, \infty)$ be the subset of non-decreasing functions $f \in \mathcal{D}_{\mathbb{R}_+}[0, \infty)$ with $f(0) = 0$.

2. DESCRIPTION OF MODEL AND STATE DYNAMICS

In Section 2.1 we describe the basic model, which is sometimes referred to as the GI/G/N+G model. In Section 2.2 we introduce the state descriptor, some auxiliary processes, and describe the state dynamics, and in Section 2.3, we obtain a convenient representation formula for expectations of linear functionals of the measure-valued components of the state process. In Section 2.4, we introduce a filtration, with respect to which the state descriptor is an adapted, strong Markov process. This model was also considered in [14], where functional strong law of large numbers limit for the state descriptor was established as the number of servers and the mean arrival rate both tend to infinity.

2.1. Model Description and Primitive Data. Consider a queueing system with N identical servers, in which arriving customers are served in a non-idling, First-Come-First-Serve (FCFS) manner, i.e., a newly arriving customer immediately enters service if there are any idle servers or, if all servers are busy, then the customer joins the back of the queue, and the customer at the head of the queue (if one is present) enters service as soon as a server becomes free.

It is assumed that customers are impatient, and that a customer reneges from the queue as soon as the amount of time he or she has waited in the queue reaches his or her patience time. Customers do not renege once they have entered service, and service is non-preemptive. The patience times of customers are given by an i.i.d. sequence, $\{r_i, i \in \mathbb{Z}\}$, with common cumulative distribution function G^r on $[0, \infty]$, while the service requirements of customers are given by another i.i.d. sequence, $\{v_i, i \in \mathbb{Z}\}$, with common cumulative distribution function G^s on $[0, \infty)$. For $i \in \mathbb{N}$,

r_i and v_i represent, respectively, the patience time and the service requirement of the i th customer to enter the system after time zero, while $\{r_i, i \in -\mathbb{N} \cup \{0\}\}$ and $\{v_i, i \in -\mathbb{N} \cup \{0\}\}$ represent, respectively, the patience times and the service requirements of customers that arrived prior to time zero (if such customers exist), ordered according to their arrival times (prior to time zero). We assume that G^s has density g^s and G^r , restricted to $[0, \infty)$, has density g^r . This implies, in particular, that $G^r(0+) = G^s(0+) = 0$. Let

$$\begin{aligned} H^r &\doteq \sup\{x \in [0, \infty) : G^r(x) < 1\}, \\ H^s &\doteq \sup\{x \in [0, \infty) : G^s(x) < 1\}. \end{aligned}$$

The superscript (N) will be used to refer to quantities associated with the system with N servers.

Let $E^{(N)}$ denote the cumulative arrival process associated with the system with N servers, with $E^{(N)}(t)$ representing the total number of customers that arrive into the system in the time interval $[0, t]$. We assume that $E^{(N)}$ is a renewal process with a common interarrival distribution function $F^{(N)}$, which has finite mean. Let $\lambda^{(N)}$ be the inverse of the mean of $F^{(N)}$, i.e.,

$$\lambda^{(N)} \int_0^\infty x F^{(N)}(dx) = 1.$$

The number $\lambda^{(N)}$ represents the long-run average arrival rate of customers to the system with N servers. We assume $E^{(N)}$, the sequence of service requirements $\{v_j, j \in \mathbb{Z}\}$, and the sequence of patience times $\{r_j, j \in \mathbb{Z}\}$ are mutually independent. Let $\alpha_E^{(N)}$ be a càdlàg, real-valued process defined by $\alpha_E^{(N)}(s) = s$ if $E^{(N)}(s) = 0$ and, if $E^{(N)}(s) > 0$, then

$$\alpha_E^{(N)}(s) \doteq s - \sup\left\{u < s : E^{(N)}(u) < E^{(N)}(s)\right\},$$

which denotes the time elapsed since the last arrival. Then $\alpha_E^{(N)}$ is simply the backward recurrence time process, which completely determines the cumulative arrival process $E^{(N)}$. Let $\mathcal{E}_0^{(N)}$ be an a.s. \mathbb{Z}_+ -valued random variable that represents the number of customers that entered the system prior to time zero. This random variable does not play an important role in the analysis. It is used, for bookkeeping purposes, to keep track of the indices of customers.

2.2. State Descriptor. A Markovian description of the state of the system with N servers would require one to keep track of the residual or elapsed patience times and the residual or elapsed service times of each customer present in the queue or in service. In order to do this in a succinct manner, with a common state space for all N -server systems, we use the representation introduced in [14]. The state of the N -server system consists of the backward recurrence time $\alpha_E^{(N)}$ of the renewal arrival process, a non-negative real-valued process $X^{(N)}$, which represents the total number of customers in system with N servers (including those in service and those in queue) and a pair of measure-valued processes, the ‘‘age measure’’ process, $\nu^{(N)}$, which encodes the amounts of time that customers currently receiving service have been in service and the ‘‘potential queue measure’’ process, $\eta^{(N)}$, which keeps track not only of the waiting times of customers in queue, but also of the potential waiting times (defined to be the times since entry into system) of all customers (irrespective of whether they have already entered service and possibly departed the system),

for whom the potential waiting time has not exceeded the patience time. Thus, the state of the system, denoted by $Y^{(N)}$, takes the form

$$(2.2) \quad Y^{(N)} = (\alpha_E^{(N)}, X^{(N)}, \nu^{(N)}, \eta^{(N)}).$$

Note that $X^{(N)}$ and $\eta^{(N)}$, together, yield the waiting times of customers currently in queue. Indeed, for $t \in [0, \infty)$, let $Q^{(N)}(t)$ be the number of customers waiting in queue at time t . Due to the non-idling condition, the queue length process is then given by

$$Q^{(N)}(t) = [X^{(N)}(t) - N]^+.$$

Since it is clear that

$$(2.3) \quad X^{(N)} = \langle \mathbf{1}, \nu^{(N)} \rangle + Q^{(N)},$$

the non-idling condition is equivalent to

$$(2.4) \quad N - \langle \mathbf{1}, \nu^{(N)} \rangle = [N - X^{(N)}]^+.$$

Moreover, since the head-of-the-line customer is the customer in queue with the longest waiting time, the quantity

$$(2.5) \quad \chi^{(N)}(t) \doteq \inf \left\{ x > 0 : \eta_t^{(N)}[0, x] \geq Q^{(N)}(t) \right\} = \left(F^{\eta_t^{(N)}} \right)^{-1} (Q^{(N)}(t))$$

represents the waiting time of the head-of-the-line customer in the queue at time t . Since this is an FCFS system, any mass in $\eta_t^{(N)}$ that lies to the right of $\chi^{(N)}(t)$ represents a customer that has already entered service by time t . Therefore, the queue length process $Q^{(N)}$ admits the following alternative representation in terms of $\chi^{(N)}$ and $\eta^{(N)}$:

$$(2.6) \quad Q^{(N)}(t) = \eta_t^{(N)}[0, \chi^{(N)}(t)].$$

The following auxiliary processes are useful for the evolution of the system and can be recovered from the state of the system $Y^{(N)}$ by using equations (2.9)–(2.11) and (2.14) in [14]:

- the cumulative reneging process $R^{(N)}$, where $R^{(N)}(t)$ is the cumulative number of customers that have reneged from the system in the time interval $[0, t]$;
- the cumulative potential reneging process $S^{(N)}$, where $S^{(N)}(t)$ represents the cumulative number of customers whose potential waiting times have reached their patience times in the interval $[0, t]$;
- the cumulative departure process $D^{(N)}$, where $D^{(N)}(t)$ is the cumulative number of customers that have departed the system after completion of service in the interval $[0, t]$;
- the process $K^{(N)}$, where $K^{(N)}(t)$ represents the cumulative number of customers that have entered service in the interval $[0, t]$.

It is easy to see the following mass balance for the number of customers in queue hold:

$$(2.7) \quad Q^{(N)}(0) + E^{(N)} = Q^{(N)} + R^{(N)} + K^{(N)}.$$

2.3. A Useful Representation Formula. We now establish representation formulas for expectations of linear functionals of the age and potential queue measure-valued processes. These are used to establish tightness of sequences of stationary distributions in Section 4.2.

Proposition 2.1. *For each bounded measurable function f on \mathbb{R}_+ and $t \geq 0$,*

$$(2.8) \quad \mathbb{E} \left[\langle f, \eta_t^{(N)} \rangle \right] = \mathbb{E} \left[\int_{[0, H^r)} f(x+t) \frac{1 - G^r(x+t)}{1 - G^r(x)} \eta_0^{(N)}(dx) \right] \\ + \mathbb{E} \left[\int_0^t f(t-s)(1 - G^r(t-s)) dE^{(N)}(s) \right]$$

and

$$(2.9) \quad \mathbb{E} \left[\langle f, \nu_t^{(N)} \rangle \right] = \mathbb{E} \left[\int_{[0, H^s)} f(x+t) \frac{1 - G^s(x+t)}{1 - G^s(x)} \nu_0^{(N)}(dx) \right] \\ + \mathbb{E} \left[\int_0^t f(t-s)(1 - G^s(t-s)) dK^{(N)}(s) \right].$$

Proof. We only prove (2.8) since the proof of (2.9) is exactly analogous. Fix $\varphi \in \mathcal{C}_c(\mathbb{R}^2)$. Suppose φ has compact support in $[0, T] \times [0, m]$ for some $T < \infty$ and $m < H^r$. Then, by the analog of (5.18) in Proposition 5.1(2) of [14] and (3.36) of [14], it follows that

$$\mathbb{E} \left[\int_0^T \langle \varphi(\cdot, s) h^r(\cdot), \eta_s^{(N)} \rangle ds \right] \leq C(m, T) \|\varphi\|_\infty,$$

where

$$C(m, T) \doteq 2 \left(\int_0^m h^r(x) dx \right) \mathbb{E}[X^{(N)}(0) + E^{(N)}(T)] < \infty.$$

On the other hand, let $e^N(t) \doteq \mathbb{E}[E^{(N)}(t)]$, $t \geq 0$. Then, taking expectations in (2.28) of Theorem 2.1 of [14] and using the fact that, for any bounded, continuous function φ on \mathbb{R}_+^2 , by Proposition 5.1(2) of [14],

$$\mathbb{E} \left[S_\varphi^{(N)}(t) \right] = \mathbb{E} \left[\int_0^t \langle \varphi(\cdot, s) h^r(\cdot), \eta_s^{(N)} \rangle ds \right] \\ = \int_0^t \langle \varphi(\cdot, s) h^r(\cdot), \mathbb{E}[\eta_s^{(N)}] \rangle ds,$$

we conclude that for every $t > 0$,

$$\langle \varphi(\cdot, t), \mathbb{E}[\eta_t^{(N)}] \rangle = \langle \varphi(\cdot, 0), \mathbb{E}[\eta_0^{(N)}] \rangle + \int_0^t \langle \varphi_s(\cdot, s) + \varphi_x(\cdot, s), \mathbb{E}[\eta_s^{(N)}] \rangle ds \\ - \int_0^t \langle \varphi(\cdot, s) h^r(\cdot), \mathbb{E}[\eta_s^{(N)}] \rangle ds + \int_{[0, t]} \varphi(0, s) de^{(N)}(s).$$

Thus, we have shown that the inequality (4.1) and equation (4.2) of [14] are satisfied with $\{\bar{\pi}_s\}_{s \geq 0}$, \bar{Z} and h replaced, respectively, by $\{\mathbb{E}[\eta_s^{(N)}]\}_{s \geq 0}$, $e^{(N)}$ and h^r . The result then follows from Proposition 4.1 of [14] since (4.3) of [14] reduces to (2.8), after the appropriate substitutions are made. \square

2.4. State Space and Filtration. The total number of customers in service at time t is given by

$$\langle \mathbf{1}, \nu_t^{(N)} \rangle = \nu_t^{(N)}[0, H^s),$$

and is bounded above by the number of servers N . On the other hand, it is clear (see, e.g., (2.13) of [14]) that a.s., for every $t \in [0, \infty)$,

$$\langle \mathbf{1}, \eta_t^{(N)} \rangle = \eta_t^{(N)}[0, H^r) \leq E^{(N)}(t) + \langle \mathbf{1}, \eta_0^{(N)} \rangle \leq E^{(N)}(t) + \mathcal{E}_0^{(N)} < \infty.$$

Therefore, a.s., for every $t \in [0, \infty)$, $\nu_t^{(N)} \in \mathcal{M}_F[0, H^s)$ and $\eta_t^{(N)} \in \mathcal{M}_F[0, H^r)$.

Let $\mathcal{M}_D[0, H^s)$ be the subset of measures in $\mathcal{M}_F[0, H^s)$ that can be represented as the sum of a finite number of unit Dirac measures in $[0, H^s)$, i.e., measures that take the form $\sum_{i=1}^k \delta_{x_i}$ for some $k \in \mathbb{Z}_+$ and $x_i \in [0, H^s)$, $i = 1, \dots, k$. Analogously, let $\mathcal{M}_D[0, H^r)$ be the subset of $\mathcal{M}_F[0, H^r)$ that can be expressed as the sum of a finite number of unit Dirac measures in $[0, H^r)$. Also, define

$$(2.10) \quad \mathcal{Y}^{(N)} \doteq \left\{ \begin{array}{l} (\alpha, x, \mu, \pi) \in \mathbb{R}_+ \times \mathbb{Z}_+ \times \mathcal{M}_D[0, H^s) \times \mathcal{M}_D[0, H^r) : \\ x \leq \langle \mathbf{1}, \mu \rangle + \langle \mathbf{1}, \pi \rangle, \langle \mathbf{1}, \mu \rangle \leq N \end{array} \right\},$$

where \mathbb{R}_+ is endowed with the Euclidean topology d , \mathbb{Z}_+ is endowed with the discrete topology ρ , and $\mathcal{M}_D[0, H^s)$ and $\mathcal{M}_D[0, H^r)$ are both endowed with the topology of weak convergence. The space $\mathcal{Y}^{(N)}$ is a closed subset of $\mathbb{R}_+ \times \mathbb{Z}_+ \times \mathcal{M}_F[0, H^s) \times \mathcal{M}_F[0, H^r)$ and is endowed with the usual product topology. Since $\mathbb{R}_+ \times \mathbb{Z}_+ \times \mathcal{M}_F[0, H^s) \times \mathcal{M}_F[0, H^r)$ is a Polish space, the closed subset $\mathcal{Y}^{(N)}$ is also a Polish space. It follows from the representations of $\nu_t^{(N)}$ and $\eta_t^{(N)}$ in (2.3) and (2.8) of [14] that a.s., the state descriptor $Y^{(N)}(t)$ takes values in $\mathcal{Y}^{(N)}$ for every $t \in [0, \infty)$.

For $t \in [0, \infty)$, let $\tilde{\mathcal{F}}_t^{(N)}$ be the σ -algebra generated by

$$\left\{ \mathcal{E}_0^{(N)}, X^{(N)}(0), \alpha_E^{(N)}(s), w_j^{(N)}(s), a_j^{(N)}(s), s_j^{(N)}, j \in \{-\mathcal{E}_0^{(N)} + 1, \dots, 0\} \cup \mathbb{N}, s \in [0, t] \right\},$$

where $s^{(N)} \doteq (s_j^{(N)}, j \in \mathbb{Z})$ is the ‘‘station process’’, defined on the same probability space $(\Omega, \mathcal{F}, \mathbb{P})$. For each $t \in [0, \infty)$, if customer j has already entered service by time t , then $s_j^{(N)}(t)$ is equal to the index $i \in \{1, \dots, N\}$ of the station at which customer j receives service and $s_j^{(N)}(t) \doteq 0$ otherwise. Let $\{\mathcal{F}_t^{(N)}\}$ denote the associated right-continuous filtration, completed with respect to \mathbb{P} . It is proved in Appendix A of [14] that the state descriptor $Y^{(N)}$ and the auxiliary processes $E^{(N)}$, $Q^{(N)}$, $S^{(N)}$, $R^{(N)}$, $D^{(N)}$ and $K^{(N)}$ are càdlàg and adapted to the filtration $\{\mathcal{F}_t^{(N)}\}$. Moreover, from Lemma B.1 of [14], it follows that $Y^{(N)}$ is a strong Markov process with respect to the filtration $\{\mathcal{F}_t^{(N)}\}$.

3. ASSUMPTIONS AND MAIN RESULTS

The main focus of this paper is to obtain a ‘‘first-order’’ approximation for the stationary distribution of the N -server queue, which is accurate in the limit as the number of servers goes to infinity.

3.1. Basic Assumptions. We impose the following mild first moment assumption on the patience and service time distribution functions G^r and G^s . Without loss of generality, we can normalize the service time distribution, so that its mean equals 1.

Assumption 1. *The mean patience and service times are finite:*

$$(3.11) \quad \theta^r \doteq \int_{[0, \infty)} x g^r(x) dx = \int_{[0, \infty)} (1 - G^r(x)) dx < \infty,$$

and

$$(3.12) \quad \int_{[0, \infty)} x g^s(x) dx = \int_{[0, \infty)} (1 - G^s(x)) dx = 1.$$

Let ν_* and η_* be the probability measures defined as follows:

$$(3.13) \quad \nu_*[0, x] \doteq \int_0^x (1 - G^s(y)) dy, \quad x \in [0, H^s),$$

$$(3.14) \quad \eta_*[0, x] \doteq \int_0^x (1 - G^r(y)) dy, \quad x \in [0, H^r).$$

Note that ν_* and η_* are well-defined due to Assumption 1. For $\lambda \geq 1$, define the set B_λ as follows:

$$(3.15) \quad B_\lambda \doteq \left\{ x \in [1, \infty) : G^r \left((F^{\lambda \eta_*})^{-1}((x-1)^+) \right) = \frac{\lambda-1}{\lambda} \right\}.$$

Let

$$b_l^\lambda = \inf \{x \in [1, \infty) : x \in B_\lambda\} \quad \text{and} \quad b_r^\lambda = \sup \{x \in [1, \infty) : x \in B_\lambda\}.$$

Since the functions G^r and $F^{\lambda \eta_*}$ are continuous and non-decreasing, we have $B_\lambda = [b_l^\lambda, b_r^\lambda]$. Let \mathcal{I}_λ be the set of states defined by

$$(3.16) \quad \mathcal{I}_\lambda = \begin{cases} \{(\lambda, \lambda \nu_*, \lambda \eta_*)\} & \text{if } \lambda < 1, \\ \{(x_*, \nu_*, \lambda \eta_*) : x_* \in B_\lambda\} & \text{if } \lambda \geq 1. \end{cases}$$

We show in Theorem 5.5 that \mathcal{I}_λ describes the so-called invariant manifold for the fluid limit. Suppose that \mathcal{I}_λ satisfies the following assumption.

Assumption 2. *The set \mathcal{I}_λ has a single element.*

Note that this is a non-trivial restriction only when $\lambda \geq 1$. A sufficient condition for Assumption 2 to hold is as follows.

Lemma 3.1. *If either $\lambda < 1$, or $\lambda \in [1, \infty)$ and equation*

$$(3.17) \quad G^r(x) = \frac{\lambda-1}{\lambda}$$

has a unique solution, then Assumption 2 holds. In particular, this is true if G^r is strictly increasing.

Proof. Fix $\lambda \in [1, \infty)$. It suffices to show that the set B_λ in (3.15) consists of a single point. Since the equation in (3.17) has a unique solution and the function $(F^{\lambda \eta_*})^{-1}(\cdot)$ is strictly increasing on $[0, \lambda \theta^r)$, the equation

$$G^r \left((F^{\lambda \eta_*})^{-1}((x-1)^+) \right) = \frac{\lambda-1}{\lambda}$$

has a unique solution. Thus B_λ has a single element and the lemma follows. \square

For each $N \in \mathbb{N}$, let $\bar{Y}^{(N)} = (\bar{\alpha}_E^{(N)}, \bar{X}^{(N)}, \bar{\nu}^{(N)}, \bar{\eta}^{(N)})$ be the fluid scaled state descriptor defined as follows:

$$(3.18) \quad \bar{\alpha}_E^{(N)}(t) \doteq \alpha_E^{(N)}(t), \quad \bar{X}^{(N)}(t) \doteq \frac{X^{(N)}(t)}{N},$$

$$(3.19) \quad \bar{\nu}_t^{(N)}(B) \doteq \frac{\nu_t^{(N)}(B)}{N}, \quad \bar{\eta}_t^{(N)}(B) \doteq \frac{\eta_t^{(N)}(B)}{N},$$

for $t \in [0, \infty)$ and any Borel subset B of \mathbb{R}_+ . Analogously, for $I = E, D, K, Q, R, S$, define

$$(3.20) \quad \bar{I}^{(N)} \doteq \frac{I^{(N)}}{N}.$$

The following standard assumption is imposed on the sequences of fluid scaled external arrival processes $\{\bar{E}^{(N)}\}$ and initial conditions $(\eta_0^{(N)}, \nu_0^{(N)})$, $N \in \mathbb{N}$.

Assumption 3. *The following conditions are satisfied:*

- (1) $\bar{\lambda}^{(N)} \rightarrow \lambda$ as $N \rightarrow \infty$ for some $\lambda \in [0, \infty)$, where $\bar{\lambda}^{(N)} = \lambda^{(N)}/N$;
- (2) As $N \rightarrow \infty$, $\bar{E}^{(N)} \rightarrow \bar{E}$ in $\mathcal{D}_{\mathbb{R}_+}[0, \infty)$ \mathbb{P} -a.s., where $\bar{E}(t) = \lambda t$;
- (3) $\mathbb{E}[\langle \mathbf{1}, \eta_0^{(N)} \rangle] < \infty$ and $\mathbb{E}[\langle \mathbf{1}, \nu_0^{(N)} \rangle] < \infty$ for each $N \in \mathbb{N}$.

The following technical assumption was imposed on the hazard rate functions in [14] to establish the fluid limit theorem.

Assumption 4. *There exists $L^s < H^s$ such that h^s is either bounded or lower-semicontinuous on (L^s, H^s) , and likewise, there exists $L^r < H^r$ such that h^r is either bounded or lower-semicontinuous on (L^r, H^r) .*

We conclude with a mild assumption on the inter-arrival distribution function $F^{(N)}$.

Assumption 5. *The interarrival distribution $F^{(N)}$ has a density.*

3.2. Main Results. The first result focuses on the existence of a stationary distribution for the state process.

Theorem 3.2. *For each N , under Assumption 5, the $\{Y_t^{(N)}, \mathcal{F}_t^{(N)}\}$ is a Feller process that has a stationary distribution.*

The Feller property is proved in Proposition 4.2 and the existence of a stationary distribution is established in Theorem 4.9. In Theorem 7.1, the state process is also shown to be ergodic under an additional condition (Assumption 6), which holds, for example, when the interarrival, reneging and service densities are strictly positive and the latter two have support on $(0, \infty)$.

We now state the main result, which provides a first-order approximation for stationary distributions of N -server queues.

Theorem 3.3. *Suppose Assumptions 1–5 hold. Then given any sequence of scaled stationary distributions $Y_*^{(N)}/N = (\bar{\alpha}_{E,*}^{(N)}, \bar{X}_*^{(N)}, \bar{\nu}_*^{(N)}, \bar{\eta}_*^{(N)})$, $N \in \mathbb{N}$, the sequence of marginals $(\bar{X}_*^{(N)}, \bar{\nu}_*^{(N)}, \bar{\eta}_*^{(N)})$ converges weakly, as $N \rightarrow \infty$, to the unique element in \mathcal{I}_λ .*

The proof of Theorem 3.3, which is given in Section 6, relies on a convenient representation of the mean dynamics of the many-server process (see Proposition 2.1) to directly establish convergence to the unique invariant state. In particular, the proof does not show that the fluid limit converges, as $t \rightarrow \infty$, to the unique invariant state. Indeed, as discussed in Section 7.1, characterization of the long-time behavior of the fluid limit appears to be a non-trivial task. However, a generic example is provided in Section 7.2 to show that, when the invariant manifold has multiple states, the diagram in Figure 1 need not commute. The characterization of the stationary distribution, and investigation of the possible metastable behavior of the many-server queue in the presence of multiple states, remains a subject for future investigation.

4. STATIONARY DISTRIBUTION OF THE N -SERVER QUEUE

We now establish the existence of a stationary distribution for the Markovian state descriptor $\{Y_t^{(N)}, \mathcal{F}_t^{(N)}\}$ for the system with N servers, under Assumption 5. First, in Section 4.1, $\{Y_t^{(N)}, \mathcal{F}_t^{(N)}\}_{t \geq 0}$ is shown to be a Feller process (see Proposition 4.2). Then, in Section 4.2, the Krylov-Bogoliubov existence theorem (cf. Corollary 3.1.2 of [6]) is used to show that $\{Y_t^{(N)}, \mathcal{F}_t^{(N)}\}_{t \geq 0}$ has a stationary distribution. Finally, in Appendix A, ergodicity and positive Harris recurrence of the process $\{Y_t^{(N)}, \mathcal{F}_t^{(N)}\}_{t \geq 0}$ is established under an additional condition (Assumption 6). For conciseness, in the rest of this section, N is fixed and the dependence on N is omitted from the notation.

4.1. Feller Property. It follows from the definition of Y in (2.2) and Lemma B.1 of [14] that Y is a so-called piecewise deterministic Markov process (see [12] for a precise definition) with jump times $\{\tau_1, \tau_2, \dots\}$, where each jump time is either the arrival time of a new customer, the time of a service completion, or the end of a patience time. Note that the set of jump times also includes the times of entry into service of customers since, due to the non-idling condition, each such entry time coincides with either the arrival time of that customer or the time of service completion of another customer. Let $\tau_0 = 0$. For each $i \in \mathbb{Z}_+$, Y evolves in a deterministic fashion on $[\tau_i, \tau_{i+1})$:

$$Y(\tau_i + t) = \phi_{Y(\tau_i)}(t), \quad t \in [0, \tau_{i+1} - \tau_i),$$

where, for each $y \in \mathcal{Y}$ of the form $y = (\alpha, x, \sum_{i=1}^k \delta_{u_i}, \sum_{j=1}^l \delta_{z_j})$, $k \leq N$, we define

$$(4.21) \quad \phi_y(t) \doteq \left(\alpha + t, x, \sum_{i=1}^k \delta_{u_i+t}, \sum_{j=1}^l \delta_{z_j+t} \right), \quad t \geq 0.$$

The Markovian semigroup of Y is defined in the usual way: for each $t \geq 0$, $y \in \mathcal{Y}$ and $A \in \mathcal{B}(\mathcal{Y})$, the set of Borel subsets of \mathcal{Y} , let

$$(4.22) \quad P_t(y, A) = \mathbb{P}(Y(t) \in A | Y(0) = y).$$

Moreover, for any measurable function ψ defined on \mathcal{Y} and $t \geq 0$, let $P_t\psi$ be the function on \mathcal{Y} given by

$$(4.23) \quad P_t\psi(y) = \mathbb{E}[\psi(Y(t)) | Y(0) = y], \quad y \in \mathcal{Y}.$$

We now show that the semigroup $\{P_t, t \geq 0\}$ is Feller, i.e., we show that for any $\psi \in C_b(\mathcal{Y})$ and $t \geq 0$, $P_t\psi \in C_b(\mathcal{Y})$.

For each $m \in \mathbb{Z}_+$, let Y^m be the state descriptor of an N -server queue with initial state

$$Y^m(0) = y^m = \left(\alpha^m, x^m, \sum_{i=1}^{k^m} \delta_{u_i^m}, \sum_{j=1}^{l^m} \delta_{z_j^m} \right) \in \mathcal{Y}.$$

Suppose that $\{Y^m, m \in \mathbb{Z}_+\}$ are defined on the same probability space and y^m converges to y^0 as $m \rightarrow \infty$. The convergence of y^m to y^0 implies that $x^m = x^0$, $k^m = k^0$, $l^m = l^0$ for all sufficiently large m and, as $m \rightarrow \infty$, $\alpha^m \rightarrow \alpha^0$, $u_i^m \rightarrow u_i^0$ and $z_j^m \rightarrow z_j^0$ for each $1 \leq i \leq k^0$, $1 \leq j \leq l^0$. Without loss of generality, we may assume that $x^m = x^0$, $k^m = k^0$, $l^m = l^0$ for every $m \in \mathbb{Z}_+$. For the m th N -server system, $m \in \mathbb{Z}_+$, the time to the arrival of the first customer after time 0 has distribution function $F(\alpha^m + \cdot)/(1 - F(\alpha^m))$, the distribution of the residual patience time of the initial customer associated with the point mass $\delta_{z_j^m}$ has density $g^r(z_j^m + \cdot)/(1 - G^r(z_j^m))$ and the distribution of the residual service time of the initial customer associated with the point mass $\delta_{u_i^m}$ has density $g^s(u_i^m + \cdot)/(1 - G^s(u_i^m))$. For simplicity, henceforth, we will denote k^0, l^0, x^0 simply by k, l, x . We assume that the elements of the sequence $\{Y^m, m \in \mathbb{Z}_+\}$ are coupled so that:

- the inter-arrival times after the first arrival and the sequences of service times and patience times of customers that arrive after time 0 are identical for each N -server queue Y^m , $m \in \mathbb{Z}_+$;
- the first arrival time for the m th N -server queue converges, as $m \rightarrow \infty$, to the first arrival time for the 0th N -server queue;
- for each $j = 1, \dots, l$, the remaining patience time of the customer associated with the point mass $\delta_{z_j^m}$ converges, as $m \rightarrow \infty$, to the remaining patience time of the customer associated with the point mass $\delta_{z_j^0}$;
- for each $i = 1, \dots, k$, the remaining service time of the customer associated with the point mass $\delta_{u_i^m}$ converges, as $m \rightarrow \infty$, to the remaining service time of the customer associated with the point mass $\delta_{u_i^0}$.

Lemma 4.1. *For each $m \in \mathbb{Z}_+$ and $n \in \mathbb{N}$, let τ_n^m be the n th jump time of Y^m . Then for each $n \in \mathbb{N}$, τ_n^m converges to τ_n^0 and $Y^m(\tau_n^m)$ converges in \mathcal{Y} to $Y^0(\tau_n^0)$ a.s., as $m \rightarrow \infty$.*

Proof. We prove the lemma by an induction argument. First consider $n = 1$. For each $m \in \mathbb{Z}_+$, the first jump time τ_1^m is the minimum of the first arrival time of a new customer, the remaining patience times of initial customers with potential waiting times in the set $\{z_j^m, 1 \leq j \leq l\}$ and the remaining service times of initial customers associated with ages in the set $\{u_i^m, 1 \leq i \leq k\}$. It follows directly from the assumptions on $\{Y^m, m \in \mathbb{Z}_+\}$ that for every realization,

$$(4.24) \quad \tau_1^m \rightarrow \tau_1^0, \quad \text{as } m \rightarrow \infty.$$

Since the service time distribution function G^s and the patience time distribution function G^r have densities, with probability 1, τ_1^0 coincides with exactly one of the following in the 0th system: the first arrival time of a new customer, the remaining patience time of an initial customer with initial waiting time z_j^0 , $1 \leq j \leq l$, and the remaining service time of an initial customer with age u_i^0 , $1 \leq i \leq k$. Let us

fix a realization such that τ_1^0 equals the first arrival time of a new customer in the 0th system. The remaining two cases can be handled similarly. In this case, by the convergence of τ_1^m to τ_1^0 , the convergence of the other quantities stated above and the coupling construction, for all sufficiently large m , τ_1^m equals the first arrival time of a new customer. Hence, for all sufficiently large m , the first jump of Y^m is due to the first arrival of a new customer in the m th system. For such m , since Y^m evolves in a deterministic fashion on $[0, \tau_1^m)$ described by the continuous function ϕ introduced in (4.21), we have

$$Y^m(\tau_1^{m-}) = \left(\tau_1^m, x, \sum_{i=1}^k \delta_{u_i^m + \tau_1^m}, \sum_{j=1}^l \delta_{z_j^m + \tau_1^m} \right).$$

If $k = N$ and $x \geq k = N$, then all the servers are busy and the customer that arrives at τ_1^m will have to wait in queue. Thus

$$Y^m(\tau_1^m) = \left(0, x + 1, \sum_{i=1}^k \delta_{u_i^m + \tau_1^m}, \sum_{j=1}^l \delta_{z_j^m + \tau_1^m} + \delta_0 \right).$$

On the other hand, if $k < N$, then $x = k$ and there is at least one idle server present. Hence, the customer will join service immediately upon arrival at time τ_1^m . Thus, in this case,

$$Y^m(\tau_1^m) = \left(0, x + 1, \sum_{i=1}^k \delta_{u_i^m + \tau_1^m} + \delta_0, \sum_{j=1}^l \delta_{z_j^m + \tau_1^m} + \delta_0 \right).$$

In both cases, for the chosen realization, we have $Y^m(\tau_1^m) \rightarrow Y^0(\tau_1^0)$ as $m \rightarrow \infty$.

Now, suppose that τ_i^m converges to τ_i^0 and $Y^m(\tau_i^m)$ converges to $Y^0(\tau_i^0)$ a.s., as $m \rightarrow \infty$, for $1 \leq i \leq n$, and consider $i = n + 1$. Fix a realization such that τ_n^m converges to τ_n^0 and $Y^m(\tau_n^m)$ converges to $Y^0(\tau_n^0)$ as $m \rightarrow \infty$. By the same argument as in the case $n = 1$, we may assume, without loss of generality, that for the chosen realization and $m \in \mathbb{Z}_+$, the jump at τ_n^m for Y^m is due to the arrival of a new customer. Then, for each $m \in \mathbb{Z}_+$, $Y^m(\tau_n^m)$ has the following representation:

$$Y^m(\tau_n^m) = \left(0, x_n^m, \sum_{i=1}^{k_n^m} \delta_{u_{i,n}^m}, \sum_{j=1}^{l_n^m} \delta_{z_{j,n}^m} \right),$$

for some $k_n^m, l_n^m, x_n^m \in \mathbb{Z}_+$, $u_{i,n}^m, z_{j,n}^m \in \mathbb{R}_+$ with $x_n^m \leq k_n^m + l_n^m$, $k_n^m \leq N$. Due to the induction hypothesis and the topology of \mathcal{Y} , for all sufficiently large m , $x_n^m = x_n^0$, $k_n^m = k_n^0$ and $l_n^m = l_n^0$, and $u_{i,n}^m \rightarrow u_{i,n}^0$ and $z_{j,n}^m \rightarrow z_{j,n}^0$ as $m \rightarrow \infty$ for each $1 \leq i \leq k_n^0$ and $1 \leq j \leq l_n^0$. The argument that was used for the case $n = 1$ can be again to show that τ_{n+1}^m converges to τ_{n+1}^0 and $Y^m(\tau_{n+1}^m)$ converges to $Y^0(\tau_{n+1}^0)$ a.s., as $m \rightarrow \infty$. This completes the induction argument and hence proves the lemma. \square

Proposition 4.2. *Suppose that the interarrival distribution F has a density. Then the semigroup $\{P_t, t \geq 0\}$ is Feller.*

Proof. It is easy to see from the definition of the function $P_t\psi$ in (4.23) that when ψ is bounded, $P_t\psi$ is also bounded. To prove the lemma, it suffices to show that $P_t\psi$ is a continuous function with respect to the topology on \mathcal{Y} . Fix $t \geq 0$. Let $y^m =$

$(\alpha^m, x^m, \mu^m, \pi^m), m \in \mathbb{Z}_+$, be a sequence of points in \mathcal{Y} such that y^m converges in \mathcal{Y} , as $m \rightarrow \infty$, to y^0 , for some $y^0 = (\alpha^0, x^0, \mu^0, \pi^0) \in \mathcal{Y}$. Since \mathbb{Z}_+ is a discrete space and $x^m \rightarrow x^0$ as $m \rightarrow \infty$, it must be that for all sufficiently large m , $x^m = x^0$. Without loss of generality, we assume that $x^m = x^0$ for each $m \in \mathbb{N}$. Consider a sequence of coupled N -server queues $\{Y^m, m \in \mathbb{Z}_+\}$ such that $Y^m(0) = y^m$ for each $m \in \mathbb{Z}_+$. Then $P_t \psi(y^m) = \mathbb{E}[\psi(Y^m(t))]$. To prove the continuity of $P_t \psi$, it suffices to show that $Y^m(t) \rightarrow Y^0(t)$ a.s., as $m \rightarrow \infty$. Indeed, since $\psi \in C_b(\mathcal{Y})$, the latter convergence would imply that $\psi(Y^m(t)) \rightarrow \psi(Y^0(t))$ and hence, by the bounded convergence theorem, that $P_t \psi(y^m) \rightarrow P_t \psi(y^0)$ as $m \rightarrow \infty$, which would show that $\{P_t, t \geq 0\}$ is Feller.

It only remains to prove that $Y^m(t) \rightarrow Y^0(t)$ a.s., as $m \rightarrow \infty$. Since the interarrival distribution F , service distribution G^s and patience distribution G^r all have densities, with probability 1, t does not belong to the set $\{\tau_n^0, n \in \mathbb{N}\}$ of jump times of Y^0 . Fix a realization such that t does not belong to the set $\{\tau_n^0, n \in \mathbb{N}\}$ and such that for each $n \in \mathbb{N}$, τ_n^m converges to τ_n^0 and $Y^m(\tau_n^m)$ converges in \mathcal{Y} to $Y^0(\tau_n^0)$, as $m \rightarrow \infty$. By Lemma 4.1, this can be done on a set of probability one. Let $r \doteq \sup\{n : \tau_n^0 < t\}$. Then $\tau_r^0 < t < \tau_{r+1}^0$, and hence for all sufficiently large m , $\tau_r^m < t < \tau_{r+1}^m$. By the convergence of $Y^m(\tau_r^m)$ to $Y^0(\tau_r^0)$, as $m \rightarrow \infty$, and the definition of ϕ in (4.21), we conclude that $Y^m(t) \rightarrow Y^0(t)$, as $m \rightarrow \infty$. Thus, we have shown that $Y^m(t) \rightarrow Y^0(t)$ a.s., as $m \rightarrow \infty$. \square

4.2. Existence of Stationary Distributions. In this section, it is shown that the Feller process $\{Y_t, \mathcal{F}_t\}_{t \geq 0}$ admits a stationary distribution. To achieve this, we use the Krylov-Bogoliubov theorem (cf. Corollary 3.1.2 of [6]), which requires showing that the following family $\{L_t, t \geq 0\}$ of probability measures associated with $\{Y_t, \mathcal{F}_t\}_{t \geq 0}$ is tight. For each measurable set $B \subset \mathcal{Y}$ and $t > 0$, define

$$L_t(B) \doteq \frac{1}{t} \int_0^t \mathbb{P}(Y(s) \in B) ds.$$

Obviously, for each $t \geq 0$, L_t is a probability measure on $(\mathcal{Y}, \mathcal{B}(\mathcal{Y}))$. We now recall some useful criteria for tightness of a family of random measures.

Proposition 4.3. *A family $\{\pi_t\}_{t \geq 0}$ of $\mathcal{M}_F[0, H]$ -valued random variables is tight if and only if the following two conditions hold:*

- (1) $\sup_{t \geq 0} \mathbb{E}[\langle \mathbf{1}, \pi_t \rangle] < \infty$;
- (2) $\lim_{c \rightarrow H} \sup_t \mathbb{E}[\pi_t[c, H]] \rightarrow 0$.

Lemma 4.4. *We have $\sup_{t \geq 0} \mathbb{E}[\langle \mathbf{1}, \eta_t \rangle] < \infty$ and $\sup_{t \geq 0} \mathbb{E}[\langle \mathbf{1}, \nu_t \rangle] < \infty$.*

Proof. Let $f = \mathbf{1}$ in (2.8) and (recalling that the superscript N is being suppressed from the notation) let $e(t) \doteq \mathbb{E}[E(t)]$, $t \geq 0$. Using integration-by-parts, it follows that

$$\begin{aligned}
 \mathbb{E}[\langle \mathbf{1}, \eta_t \rangle] &\leq \mathbb{E}[\langle \mathbf{1}, \eta_0 \rangle] + \int_0^t (1 - G^r(t-s)) de(s) \\
 &= \mathbb{E}[\langle \mathbf{1}, \eta_0 \rangle] + e(t) - \int_0^t e(s) g^r(t-s) ds \\
 &= \mathbb{E}[\langle \mathbf{1}, \eta_0 \rangle] + e(t)(1 - G^r(t)) - \int_0^t (e(t) - e(t-s)) g^r(s) ds.
 \end{aligned}$$

Since E is a renewal process with rate λ , $e(t)/t \rightarrow \lambda$ as $t \rightarrow \infty$ by the key renewal theorem. Moreover, the finite mean condition (3.11) implies $t(1 - G^r(t)) \rightarrow 0$ as $t \rightarrow \infty$. Therefore, we have $\sup_{t \geq 0} e(t)(1 - G^r(t)) < \infty$. The Blackwell renewal theorem (cf. Theorem 4.3 of [1]) implies that $e(t) - e(t - s) \rightarrow s\lambda$ as $t \rightarrow \infty$, and hence that $\sup_{t \geq 0} \int_0^t (e(t) - e(t - s))g^r(s)ds < \infty$. Combining these relations with (3) of Assumption 3 and the last display, we conclude that $\sup_{t \geq 0} \mathbb{E}[\langle \mathbf{1}, \eta_t \rangle] < \infty$.

On the other hand, since each ν_t is the sum of at most N unit Dirac masses, it trivially follows that $\sup_{t \geq 0} \mathbb{E}[\langle \mathbf{1}, \nu_t \rangle] \leq N < \infty$. \square

To show that $\{\eta_t\}_{t \geq 0}$ and $\{\nu_t\}_{t \geq 0}$ satisfy the second property in Proposition 4.3, note that by choosing $f = \mathbb{1}_{[c, H^r)}$, $c > 0$, in (2.8), we obtain for $t \geq 0$,

$$(4.25) \quad \mathbb{E}[\eta_t[c, H^r)] \leq \mathbb{E} \left[\int_{[0, H^r)} \mathbb{1}_{[c, H^r)}(x+t) \frac{1 - G^r(x+t)}{1 - G^r(x)} \eta_0(dx) \right] \\ + \int_0^t \mathbb{1}_{[c, H^r)}(t-s)(1 - G^r(t-s)) de(s)$$

and, likewise, by choosing $f = \mathbb{1}_{[c, H^s)}$ in (2.9) it follows that for $t \geq 0$,

$$(4.26) \quad \mathbb{E}[\nu_t[c, H^s)] = \mathbb{E} \left[\int_{[0, H^s)} \mathbb{1}_{[c, H^s)}(x+t) \frac{1 - G^s(x+t)}{1 - G^s(x)} \nu_0(dx) \right] \\ + \mathbb{E} \left[\int_0^t \mathbb{1}_{[c, H^s)}(t-s)(1 - G^s(t-s)) dK(s) \right].$$

We now establish two supporting lemmas.

Lemma 4.5. *We have*

$$(4.27) \quad \lim_{c \rightarrow H^r} \sup_{t \geq 0} \mathbb{E} \left[\int_{[0, H^r)} \mathbb{1}_{[c, H^r)}(x+t) \frac{1 - G^r(x+t)}{1 - G^r(x)} \eta_0(dx) \right] = 0$$

and

$$(4.28) \quad \lim_{c \rightarrow H^s} \sup_{t \geq 0} \mathbb{E} \left[\int_{[0, H^s)} \mathbb{1}_{[c, H^s)}(x+t) \frac{1 - G^s(x+t)}{1 - G^s(x)} \nu_0(dx) \right] = 0.$$

Proof. When $H^r < \infty$, we have

$$\sup_{t \geq 0} \mathbb{E} \left[\int_{[0, H^r)} \mathbb{1}_{[c, H^r)}(x+t) \frac{1 - G^r(x+t)}{1 - G^r(x)} \eta_0(dx) \right] \\ \leq \sup_{t \geq 0} \mathbb{E} \left[\int_{[0, H^r)} \mathbb{1}_{[c, H^r)}(x+t) \eta_0(dx) \right] \\ = \sup_{t \in [0, c)} \mathbb{E} \left[\int_{[0, H^r)} \mathbb{1}_{[c, H^r)}(x+t) \eta_0(dx) \right] \vee \sup_{t \in [c, H^r)} \mathbb{E} \left[\int_{[0, H^r)} \mathbb{1}_{[c, H^r)}(x+t) \eta_0(dx) \right].$$

It is easy to see from (3) of Assumption 3 that

$$\lim_{c \rightarrow H^r} \sup_{t \in [c, H^r)} \mathbb{E} \left[\int_{[0, H^r)} \mathbb{1}_{[c, H^r)}(x+t) \eta_0(dx) \right] \leq \lim_{c \rightarrow H^r} \mathbb{E}[\eta_0[0, H^r - c)] = 0.$$

On the other hand, we know that

$$\sup_{t \in [0, c]} \mathbb{E} \left[\int_{[0, H^r]} \mathbb{1}_{[c, H^r]}(x+t) \eta_0(dx) \right] \leq \sup_{t \in [0, c]} \mathbb{E} [\eta_0(c-t, H^r-t)].$$

We show by contradiction that $\sup_{t \in [0, c]} \mathbb{E} [\eta_0(c-t, H^r-t)] \rightarrow 0$ as $c \rightarrow H^r$. Suppose this is not true. Then there exist $\delta > 0$ and sequences $\{c_n\}_{n \in \mathbb{N}}$ and $\{t_n\}_{n \in \mathbb{N}}$ such that $c_n \rightarrow H^r$ as $n \rightarrow \infty$, $t_n \in [0, c_n]$ for each $n \in \mathbb{N}$, and $\mathbb{E} [\eta_0(c_n - t_n, H^r - t_n)] > \delta$ for each $n \in \mathbb{N}$. Since $H^r < \infty$, $\{t_n\}_{n \in \mathbb{N}}$ is bounded and so we can take a subsequence, which we call again $\{t_n\}_{n \in \mathbb{N}}$, such that $\lim_{n \rightarrow \infty} t_n = t_* \in [0, H^r]$. In turn, this implies

$$\lim_{n \rightarrow \infty} \mathbb{E} [\eta_0(c_n - t_n, H^r - t_n)] = \lim_{n \rightarrow \infty} \mathbb{E} [\eta_0(H^r - t_*, H^r - t_*)] = 0,$$

which contradicts the initial hypothesis. Thus, $\sup_{t \in [0, c]} \mathbb{E} [\eta_0(c-t, H^r-t)] \rightarrow 0$. Together with the last three displays, this implies that (4.27) holds when $H^r < \infty$. On the other hand, when $H^r = \infty$ we have

$$\begin{aligned} & \sup_{t \geq 0} \mathbb{E} \left[\int_{[0, H^r]} \mathbb{1}_{[c, H^r]}(x+t) \frac{1-G^r(x+t)}{1-G^r(x)} \eta_0(dx) \right] \\ & \leq \sup_{t \in [0, c/2]} \mathbb{E} \left[\int_{[0, \infty)} \mathbb{1}_{[c, \infty)}(x+t) \eta_0(dx) \right] \vee \sup_{t \in [c/2, \infty)} \mathbb{E} \left[\int_{[0, \infty)} \frac{1-G^r(x+t)}{1-G^r(x)} \eta_0(dx) \right] \\ & \leq \mathbb{E} [\eta_0(c/2, \infty)] \vee \mathbb{E} \left[\int_{[0, \infty)} \frac{1-G^r(x+c/2)}{1-G^r(x)} \eta_0(dx) \right]. \end{aligned}$$

Sending $c \rightarrow \infty$ on both sides, and using the fact that $\mathbb{E}[\langle \mathbf{1}, \eta_0 \rangle] < \infty$ by Assumption 3, an application of the dominated convergence theorem shows that the right-hand side vanishes, and thus (4.27) holds in this case too. The proof of (4.28) is exactly analogous, and is thus omitted. \square

Lemma 4.6. *Let $e(t) \doteq \mathbb{E}[E(t)]$, $t \geq 0$. For $(H, G) = (H^r, G^r)$ and $(H, G) = (H^s, G^s)$, we have*

$$(4.29) \quad \lim_{c \rightarrow H} \sup_{t \geq 0} \int_0^t \mathbb{1}_{[c, H]}(t-s)(1-G(t-s)) de(s) = 0.$$

Proof. E is a (delayed) renewal process with rate λ and due to Assumption 1, the function $x \mapsto \mathbb{1}_{[c, H]}(x)(1-G(x))$ is directly Riemann integrable. Thus, by the key renewal theorem, we obtain

$$\lim_{t \rightarrow \infty} \int_0^t \mathbb{1}_{[c, H]}(t-s)(1-G(t-s)) de(s) = \frac{1}{\lambda} \int_{[0, \infty)} \mathbb{1}_{[c, H]}(x)(1-G(x)) dx.$$

Since the integrability condition imposed in Assumption 1 implies that $\int_{[0, \infty)} \mathbb{1}_{[c, H]}(x)(1-G(x)) dx \rightarrow 0$ as $c \rightarrow H$, we have the desired result. \square

Lemma 4.7. *The family $\{\eta_t\}_{t \geq 0}$ of $\mathcal{M}_F[0, H^r]$ -valued random variables and the family $\{\nu_t\}_{t \geq 0}$ of $\mathcal{M}_F[0, H^s]$ -valued random variables are tight.*

Proof. Both families satisfy the first condition of Proposition 4.3 due to Lemma 4.4. Combining (4.25) with (4.27) and Lemma 4.6, for the case $(H, G) = (H^r, G^r)$, it follows that $\{\eta_t\}_{t \geq 0}$ also satisfies the second condition of Proposition 4.3, and is thus tight.

It only remains to show that $\{\nu_t\}_{t \geq 0}$ also satisfies the second condition of Proposition 4.3, for which it suffices to show that, as $c \rightarrow H^s$, the supremum of the right-hand side of (4.26) goes to zero. Now, let $k(t) \doteq \mathbb{E}[K(t)]$ for $t \geq 0$. Applying the integration-by-parts and change of variable formulas on the second term on the right hand side of (4.26), we see that

$$(4.30) \quad \begin{aligned} & \sup_{t \geq 0} \mathbb{E} \left[\int_0^t \mathbb{1}_{[c, H^s)}(t-s)(1-G^s(t-s))dK(s) \right] \\ &= \sup_{t > c} \int_0^t \mathbb{1}_{[c, H^s)}(t-s)(1-G^s(t-s))dk(s) \\ &= \sup_{t > c} \left(k(t-c)(1-G^s(c)) - k((t-H^s)^+)(1-G((t-H^s)^+)) \right. \\ & \quad \left. - \int_c^{t \wedge H^s} k(t-s)g^s(s)ds \right) \\ &\leq \sup_{t > c} \left(k(t-c)(1-G^s(t)) + \int_c^{t \wedge H^s} (k(t-c) - k(t-s))g^s(s)ds \right). \end{aligned}$$

By taking expectations on both sides of (2.7), we have for each $t \geq 0$,

$$\mathbb{E}[Q(0)] + e(t) = \mathbb{E}[Q(t)] + \mathbb{E}[R(t)] + k(t).$$

Since Q and R are non-negative and R is increasing, it follows that

$$k(t-c) \leq e(t-c) + \mathbb{E}[Q(0)]$$

and

$$k(t-c) - k(t-s) \leq e(t-c) - e(t-s) + (\mathbb{E}[Q(t-s)] - \mathbb{E}[Q(t-c)]).$$

Combining these inequalities and carrying out another integration-by-parts, we obtain

$$(4.31) \quad \begin{aligned} & \sup_{t > c} \int_0^t \mathbb{1}_{[c, H^s)}(t-s)(1-G^s(t-s))dk(s) \\ &\leq \sup_{t > 0} \int_0^t \mathbb{1}_{[c, H^s)}(t-s)(1-G^s(t-s))de(s) + \sup_{t > c} \mathbb{E}[Q(0)](1-G^s(t)) \\ &\quad + \sup_{t > c} \int_c^{t \wedge H^s} (\mathbb{E}[Q(t-s)] - \mathbb{E}[Q(t-c)])g^s(s)ds. \end{aligned}$$

Applying Lemma 4.6, with $(H, G) = (H^s, G^s)$, we have

$$\lim_{c \rightarrow H^s} \sup_{t \geq 0} \int_0^t \mathbb{1}_{[c, H^s)}(t-s)(1-G^s(t-s))de(s) = 0.$$

Moreover,

$$\lim_{c \rightarrow H^s} \sup_{t > c} \mathbb{E}[Q(0)](1-G^s(t)) = \mathbb{E}[Q(0)] \lim_{c \rightarrow H^s} (1-G^s(c)) = 0.$$

Also, since $Q(t) \leq \langle \mathbf{1}, \eta_t \rangle$ by (2.6), we have

$$\sup_{t > c} \int_c^{t \wedge H^s} (\mathbb{E}[Q(t-s)] - \mathbb{E}[Q(t-c)])g^s(s)ds \leq 2 \sup_{t \geq 0} \mathbb{E}[\langle \mathbf{1}, \eta_t \rangle] (1-G^s(c)).$$

Since Lemma 4.4 implies $\sup_{t \geq 0} \mathbb{E}[\langle \mathbf{1}, \eta_t \rangle] < \infty$, the right-hand side of the above inequality tends to zero as $c \rightarrow H^s$. Combining the last five assertions with (4.30) and (4.31) it follows that as $c \rightarrow H^s$, the supremum over $t \geq 0$ of the second term

on the right-hand side of (4.26) vanishes to zero. On the other hand, as $c \rightarrow H^s$, the supremum over $t \geq 0$ of the first term on the right-hand side of (4.26) also vanishes to zero by (4.28). Thus, we have shown that $\sup_{t \geq 0} \mathbb{E}[\nu_t[c, H^s]] \rightarrow 0$ as $c \rightarrow H^s$, and the proof of the lemma is complete. \square

Lemma 4.8. *The family of probability measures $\{L_t\}_{t \geq 0}$ is tight.*

Proof. By Lemma 4.7, we know that for each $\delta > 0$, there exist two compact subsets $\tilde{C}_\delta \subset \mathcal{M}_F[0, H^s)$ and $\tilde{D}_\delta \subset \mathcal{M}_F[0, H^r)$ such that

$$(4.32) \quad \begin{aligned} \inf_{t \geq 0} \mathbb{P}(\nu_t \in \tilde{C}_\delta) &\geq 1 - \delta/2, \\ \inf_{t \geq 0} \mathbb{P}(\eta_t \in \tilde{D}_\delta) &\geq 1 - \delta/2. \end{aligned}$$

It follows from (2.6) and (2.3) that $X(t) \leq \langle \mathbf{1}, \nu_t \rangle + \langle \mathbf{1}, \eta_t \rangle$ for each $t \geq 0$. Together with (4.32) and the fact that the map $\mu \rightarrow \langle \mathbf{1}, \mu \rangle$ is continuous, we know that there exists $b > 0$ such that

$$(4.33) \quad \inf_{t \geq 0} \mathbb{P}(X(t) \leq b) \geq 1 - \delta.$$

On the other hand, by Theorem 4.5 of [1], it follows that $\alpha_E(t)$ converges weakly, as $t \rightarrow \infty$, to the distribution

$$(4.34) \quad F_0(t) = \lambda \int_0^t (1 - F(y)) dy.$$

Thus, there exist $T_0 > 0$ and $c > 0$ such that for all $t \geq T_0$,

$$\mathbb{P}(\alpha_E(t) \leq a) \geq \mathbb{P}(F_0 \leq a) - \delta/2 \geq 1 - \delta.$$

By choosing a large enough, we may assume, without loss of generality, that

$$\inf_{t \in [0, T_0]} \mathbb{P}(\alpha_E(t) \leq a) \geq 1 - \delta.$$

Define $C_\delta = [0, a] \times [0, b] \times \tilde{C}_\delta \times \tilde{D}_\delta$. Then the set C_δ is compact and $L_t(C_\delta) \geq 1 - \delta$ for each $t \geq 0$, which proves the lemma. \square

Since $\{Y_t, \mathcal{F}_t\}_{t \geq 0}$ is a Feller process by Proposition 4.2, the Krylov-Bogoliubov theorem and Lemma 4.8 immediately yield the following result.

Theorem 4.9. *Suppose that Assumption 5 holds. Then the state descriptor (α_E, X, ν, η) has a stationary distribution.*

5. FLUID LIMIT

In Section 5.1, we describe a deterministic dynamical system that was shown in Theorems 3.5 and 3.6 of [14] to arise as the so-called fluid limit of a many-server queue with abandonment that has service time and patience time distribution functions G^s and G^r , respectively. In Section 5.2, we identify the invariant manifold associated with the fluid limit, which is subsequently used to obtain a first order asymptotic approximation to the stationary distribution of the fluid scaled state descriptor $\bar{Y}^{(N)}$.

5.1. Fluid Equations. The state of the fluid system at time t is represented by the triplet

$$(\bar{X}(t), \bar{\nu}_t, \bar{\eta}_t) \in \mathbb{R}_+ \times \mathcal{M}_F[0, H^s] \times \mathcal{M}_F[0, H^r].$$

Here, $\bar{X}(t)$ represents the mass (or, equivalently, limiting scaled number of customers) in the system at time t , $\bar{\nu}_t[0, x]$ represents the mass of customers in service at time t who have been in service for less than x units of time, while $\bar{\eta}_t[0, x]$ represents the mass of customers in the system who, at time t , have been in the system no more than x units of time and whose patience time exceeds x . The inputs to the system are the (limiting) cumulative arrival process \bar{E} , and the initial conditions $\bar{X}(0)$, $\bar{\nu}_0$ and $\bar{\eta}_0$. Thus, $\langle \mathbf{1}, \bar{\nu}_0 \rangle$ represents the total mass of customers in service at time 0, and the fluid analog of the non-idling condition (2.4) is

$$(5.35) \quad 1 - \langle \mathbf{1}, \bar{\nu}_0 \rangle = [1 - \bar{X}(0)]^+.$$

The quantity $\langle \mathbf{1}, \bar{\eta}_0 \rangle$ represents the total mass of customers at time 0 whose residual patience times are positive. Hence, we have

$$[\bar{X}(0) - 1]^+ \leq \langle \mathbf{1}, \bar{\eta}_0 \rangle.$$

Thus, the space of possible input data for the fluid equations is given by

$$(5.36) \quad \mathcal{S}_0 \doteq \left\{ \begin{array}{l} (e, x, \nu, \eta) \in \mathcal{I}_{\mathbb{R}_+}[0, \infty) \times \mathbb{R}_+ \times \mathcal{M}_F[0, H^s] \times \mathcal{M}_F[0, H^r] : \\ 1 - \langle \mathbf{1}, \nu \rangle = [1 - x]^+, [x - 1]^+ \leq \langle \mathbf{1}, \eta \rangle \end{array} \right\},$$

where recall that $\mathcal{I}_{\mathbb{R}_+}[0, \infty)$ is the subset of non-decreasing functions $f \in \mathcal{D}_{\mathbb{R}_+}[0, \infty)$ with $f(0) = 0$. Let $F^{\bar{\eta}_t}(x)$ denote $\bar{\eta}_t[0, x]$ for each $x \in [0, H^r)$.

Definition 5.1. (Fluid Equations) Given any $(\bar{E}, \bar{X}(0), \bar{\nu}_0, \bar{\eta}_0) \in \mathcal{S}_0$, we say that the càdlàg function $(\bar{X}, \bar{\nu}, \bar{\eta})$ taking values in $\mathbb{R}_+ \times \mathcal{M}_F[0, H^s] \times \mathcal{M}_F[0, H^r]$ satisfies the associated fluid equations if for every $t \in [0, \infty)$,

$$(5.37) \quad \int_0^t \langle h^r, \bar{\eta}_s \rangle ds < \infty, \quad \int_0^t \langle h^s, \bar{\nu}_s \rangle ds < \infty,$$

for every bounded Borel measurable function f defined on \mathbb{R}_+ ,

$$(5.38) \quad \int_{[0, H^s)} f(x) \bar{\nu}_t(dx) = \int_{[0, H^s)} f(x+t) \frac{1 - G^s(x+t)}{1 - G^s(x)} \bar{\nu}_0(dx) \\ + \int_0^t f(t-s)(1 - G^s(t-s)) d\bar{K}(s);$$

and

$$(5.39) \quad \int_{[0, H^r)} f(x) \bar{\eta}_t(dx) = \int_{[0, H^r)} f(x+t) \frac{1 - G^r(x+t)}{1 - G^r(x)} \bar{\eta}_0(dx) \\ + \int_0^t f(t-s)(1 - G^r(t-s)) d\bar{E}(s);$$

where

$$(5.40) \quad \bar{K}(t) = [\bar{X}(0) - 1]^+ - [\bar{X}(t) - 1]^+ + \bar{E}(t) - \bar{R}(t);$$

$$(5.41) \quad \bar{X}(t) = \bar{X}(0) + \bar{E}(t) - \int_0^t \langle h^s, \bar{\nu}_s \rangle ds - \bar{R}(t);$$

$$(5.42) \quad \bar{R}(t) = \int_0^t \left(\int_0^{[\bar{X}(s)-1]^+} h^r \left((F^{\bar{\eta}_s})^{-1}(y) \right) dy \right) ds;$$

$$(5.43) \quad 1 - \langle \mathbf{1}, \bar{\nu}_t \rangle = [1 - \bar{X}(t)]^+;$$

$$(5.44) \quad [\bar{X}(t) - 1]^+ \leq \langle \mathbf{1}, \bar{\eta}_t \rangle.$$

Note that the fluid equations defined above are equivalent to the fluid equations in Definition 3.3 of [14] (although (3.40) and (3.42) of [14], which are the analogs of (5.38) and (5.39), were only required to be satisfied for continuous functions with compact support in [14], a standard monotone convergence argument shows that they are equivalent to (5.38) and (5.39) here). Under some mild assumptions on the input data \bar{E} , $\bar{\nu}_0$, $\bar{\eta}_0$, and the hazard rate functions h^r and h^s (stated in Assumptions 3 and 4), Theorems 3.5 and 3.6 of [14] established that there exists a unique solution to the fluid equations,

For future purposes, note that if $(\bar{X}, \bar{\nu}, \bar{\eta})$ satisfy the fluid equations for some $(\bar{E}, \bar{X}(0), \bar{\nu}_0, \bar{\eta}_0) \in \mathcal{S}_0$ and \bar{K} satisfies (5.40), then \bar{K} also satisfies

$$(5.45) \quad \bar{K}(t) = \langle \mathbf{1}, \bar{\nu}_t \rangle - \langle \mathbf{1}, \bar{\nu}_0 \rangle + \int_0^t \langle h^s, \bar{\nu}_s \rangle ds.$$

Indeed, this follows from (5.40), (5.41) and (5.43). Moreover, combining (5.45) and (5.38), with $f = \mathbf{1}$, and using an integration-by-parts argument (see Corollary 4.2 of [14]) it is easy to see that \bar{K} satisfies the renewal equation:

$$(5.46) \quad \begin{aligned} \bar{K}(t) &= \langle \mathbf{1}, \bar{\nu}_t \rangle - \langle \mathbf{1}, \bar{\nu}_0 \rangle + \int_{[0, H^s)} \frac{G^s(x+t) - G^s(x)}{1 - G^s(x)} \bar{\nu}_0(dx) \\ &\quad + \int_0^t g^s(t-s) \bar{K}(s) ds. \end{aligned}$$

Since the first two terms on the right-hand side are bounded, by the key renewal theorem (see, e.g, Theorem 4.3 in Chapter V of [1]), this implies that \bar{K} admits the representation

$$(5.47) \quad \begin{aligned} \bar{K}(t) &= \langle \mathbf{1}, \bar{\nu}_t \rangle - \langle \mathbf{1}, \bar{\nu}_0 \rangle + \int_{[0, H^s)} \frac{G^s(x+t) - G^s(x)}{1 - G^s(x)} \bar{\nu}_0(dx) \\ &\quad + \int_0^t \left(\langle \mathbf{1}, \bar{\nu}_{t-s} \rangle - \langle \mathbf{1}, \bar{\nu}_0 \rangle + \int_{[0, H^s)} \frac{G^s(x+t-s) - G^s(x)}{1 - G^s(x)} \bar{\nu}_0(dx) \right) u^s(s) ds, \end{aligned}$$

where u^s is the density of the renewal function U^s associated with G^s (u^s exists since G^s is assumed to have a density). Also, it will prove convenient to introduce the fluid queue length process \bar{Q} , defined by

$$(5.48) \quad \bar{Q}(t) \doteq [\bar{X}(t) - 1]^+, \quad t \in [0, \infty).$$

Then, the inequality in (5.44) implies

$$(5.49) \quad \bar{Q}(t) \leq \langle \mathbf{1}, \bar{\eta}_t \rangle.$$

Observe that (5.40) and (5.48), when combined, show that for every $t \in [0, \infty)$,

$$(5.50) \quad \bar{Q}(0) + \bar{E}(t) = \bar{Q}(t) + \bar{K}(t) + \bar{R}(t).$$

We can also define the fluid equations without abandonment. Let

$$(5.51) \quad \tilde{\mathcal{S}}_0 \doteq \{(e, x, \nu) \in \mathcal{I}_{\mathbb{R}_+}[0, \infty) \times \mathbb{R}_+ \times \mathcal{M}_F[0, H^s) : 1 - \langle \mathbf{1}, \nu \rangle = [1 - x]^+\}.$$

Definition 5.2. Given any $(\bar{E}, \bar{X}(0), \bar{\nu}_0) \in \tilde{\mathcal{S}}_0$, we say $(\bar{X}, \bar{\nu}) \in \mathbb{R}_+ \times \mathcal{M}_F[0, H^s]$ is a solution to the associated fluid equations in the absence of abandonment if for every $t \in [0, \infty)$, the second inequality in (5.37) holds, equations (5.38), (5.40), (5.41) and (5.43) hold with $\bar{R} \equiv 0$.

Remark 5.3. The case when customers do not renege corresponds to the case when the patience time distribution G^r has all its mass at ∞ . Formally setting $dG^r = \delta_\infty$ in Definition 5.1, we obtain the fluid limit equations in the absence of abandonment specified in Definition 5.2 (also refer to Definition 3.3 of [15]). In fact, in this case, $G^r(x) = 0$ and hence $h^r(x) = 0$ for all $x \in [0, \infty)$. From this and (5.42), we can see that $\bar{R}(t) = 0$ for all $t \geq 0$. Moreover, (5.37), (5.38), (5.41), (5.43) and (5.45) are equivalent to (3.4)–(3.8) of Definition 3.3 of [15]. At last, by letting $f = \mathbf{1}$ in (5.39), since G^r is zero on $[0, \infty)$, we have $\langle \mathbf{1}, \bar{\eta}_t \rangle = \langle \mathbf{1}, \bar{\eta}_0 \rangle + \bar{E}(t)$. On the other hand, by (5.41) and (5.36), we have

$$\begin{aligned} [\bar{X}(t) - 1]^+ &\leq [[\bar{X}(0) - 1]^+ + \bar{E}(t)]^+ \\ &\leq [\langle \mathbf{1}, \bar{\eta}_0 \rangle + \bar{E}(t)]^+. \end{aligned}$$

Thus, $[\bar{X}(t) - 1]^+ \leq \langle \mathbf{1}, \bar{\eta}_t \rangle$ shows that (5.44) holds automatically when there is no abandonment.

5.2. Invariant Manifold. We now introduce a set of states associated with the fluid equations described in Definition 5.1, which we call the *invariant manifold*. As shown in Section 6, when the invariant manifold consists of a single point, it is the limit of the scaled sequence of convergent stationary distributions $(\bar{X}_*^{(N)}, \bar{\nu}_*^{(N)}, \bar{\eta}_*^{(N)}) = \frac{1}{N}(X_*^{(N)}, \nu_*^{(N)}, \eta_*^{(N)})$.

Definition 5.4. (Invariant Manifold) Given $\lambda \in (0, \infty)$, a state $(x_0, \nu_0, \eta_0) \in \mathbb{R}_+ \times \mathcal{M}_F[0, H^s] \times \mathcal{M}_F[0, H^r]$ with $[1 - x_0]^+ = 1 - \langle \mathbf{1}, \nu_0 \rangle$ and $[x_0 - 1]^+ \leq \langle \mathbf{1}, \eta_0 \rangle$ is said to be invariant for the fluid equations described in Definition 5.1 with arrival rate λ if the solution $(\bar{X}, \bar{\nu}, \bar{\eta})$ to the fluid equations associated with the input data $(\lambda \mathbf{1}, x_0, \nu_0, \eta_0)$ satisfies $(\bar{X}(t), \bar{\nu}_t, \bar{\eta}_t) = (x_0, \nu_0, \eta_0)$ for all $t \geq 0$. The set of all invariant states for the fluid equation with rate λ will be referred to as the *invariant manifold* (associated with the fluid equations with rate λ).

Theorem 5.5. (Characterization of the Invariant Manifold) *Given $\lambda \in (0, \infty)$, the set \mathcal{I}_λ defined in (3.16) is the invariant manifold associated with the fluid equations with rate λ .*

Theorem 5.5 is a consequence of the next two lemmas. Let $\lambda \in (0, \infty)$ and (x_0, ν_0, η_0) be an invariant state according to Definition 5.4. Then the solution $(\bar{X}, \bar{\nu}, \bar{\eta})$ to the fluid equations associated with the input data $(\lambda \mathbf{1}, x_0, \nu_0, \eta_0) \in \mathcal{S}_0$ satisfies $(\bar{X}(t), \bar{\nu}_t, \bar{\eta}_t) = (x_0, \nu_0, \eta_0)$ for all $t \geq 0$. Let \bar{Q} , \bar{R} , \bar{K} be the associated auxiliary processes satisfying (5.48), (5.42), (5.40), and recall the definition of the measures ν_* and η_* given in (3.13) and (3.14), respectively.

Lemma 5.6. *If (x_0, ν_0, η_0) is an invariant state, then $\eta_0(dx) = \lambda(1 - G^r(x))dx = \lambda\eta_*(dx)$.*

Proof. On substituting the relation $\eta_t = \eta_0$, $t \geq 0$, into (5.39), we see that for every $f \in \mathcal{C}_b(\mathbb{R}_+)$ and $t \in [0, \infty)$,

$$(5.52) \quad \int_{[0, H^r)} f(x) \eta_0(dx) = \int_{[0, H^r)} f(x+t) \frac{1 - G^r(x+t)}{1 - G^r(x)} \eta_0(dx) + \lambda \int_0^t f(s)(1 - G^r(s)) ds.$$

Sending $t \rightarrow \infty$ and applying the dominated convergence theorem, we obtain

$$\int_{[0, H^r)} f(x) \eta_0(dx) = \lambda \int_0^\infty f(s)(1 - G^r(s)) ds = \int_{[0, H^r)} f(s) \lambda (1 - G^r(s)) ds.$$

It then follows that $\eta_0(dx) = \lambda \eta_*(dx)$. \square

Lemma 5.7. *If (x_0, ν_0, η_0) is an invariant state, then $\nu_0(dx) = (\lambda \wedge 1) \nu_*(dx)$, and $x_0 = \lambda$ if $\lambda < 1$ and $x_0 \in B_\lambda$ if $\lambda \geq 1$. Moreover, (x_0, ν_*, η_*) is an invariant state if $x_0 = \lambda < 1$ or $\lambda > 1$ and $x_0 \in B_\lambda$.*

Proof. Suppose (x_0, ν_0, η_0) is an invariant state. Since $\bar{X}(t) = x_0$, we have $\bar{Q}(t) = \bar{Q}(0)$ by (5.48). Since, in addition, $\bar{\eta}_t = \eta_0 = \lambda \eta_*$ by Lemma 5.6, we have

$$\int_0^{[\bar{X}(t)-1]^+} h^r (F^{\bar{\eta}_t})^{-1}(y) dy = \int_0^{[x_0-1]^+} h^r (F^{\lambda \eta_*})^{-1}(y) dy.$$

Let p denote the term on the right-hand side of the above display. Then for each $t \geq 0$, by (5.42), we have $\bar{R}(t) = pt$ and, by (5.50), we have $\bar{K}(t) = (\lambda - p)t$. Substituting $\bar{\nu}_t = \nu_0$ in (5.38), we obtain for every $f \in \mathcal{C}_b(\mathbb{R}_+)$ and $t \in [0, \infty)$,

$$(5.53) \quad \int_{[0, H^s)} f(x) \nu_0(dx) = \int_{[0, H^s)} f(x+t) \frac{1 - G^s(x+t)}{1 - G^s(x)} \nu_0(dx) + \int_0^t f(s)(1 - G^s(s))(\lambda - p) ds.$$

By letting $t \rightarrow \infty$ and applying the dominated convergence theorem, we obtain

$$\int_{[0, H^s)} f(x) \nu_0(dx) = (\lambda - p) \int_0^\infty f(s)(1 - G^s(s)) ds = (\lambda - p) \int_{[0, H^s)} f(s)(1 - G^s(s)) ds.$$

Thus $\nu_0(dx) = (\lambda - p) \nu_*(dx)$, and so $\langle \mathbf{1}, \nu_0 \rangle = \lambda - p$.

To show that $\nu_0(dx) = (\lambda \wedge 1) \nu_*(dx)$, it suffices to show that $\lambda - p = \langle \mathbf{1}, \nu_0 \rangle = \lambda \wedge 1$. If $x_0 \leq 1$, then $p = 0$ by its definition. Hence, $\nu_0(dx) = \lambda \nu_*(dx)$ and $\lambda = \langle \mathbf{1}, \nu_0 \rangle \leq 1$. Thus, in this case, $\lambda - p = \lambda \wedge 1$. On the other hand, if $x_0 > 1$, it follows from (5.43) that $\langle \mathbf{1}, \nu_0 \rangle = 1$. Since we also have $\langle \mathbf{1}, \nu_0 \rangle = \lambda - p$, it follows that $\lambda = p + 1 \geq 1$. Thus, in this case too, we have $\lambda - p = \lambda \wedge 1$. This proves the first assertion of the lemma.

For the second assertion of the lemma, we observe that when $\lambda < 1$, the equality $\lambda - p = \lambda \wedge 1$ implies $p = 0$ and $\langle \mathbf{1}, \nu_0 \rangle = \lambda < 1$. Hence (5.35) implies $x_0 = \langle \mathbf{1}, \nu_0 \rangle = \lambda$. If $\lambda \geq 1$, we have $\nu_0(dx) = \nu_*(dx)$ and the equality $\lambda - p = \lambda \wedge 1$ implies $p = \lambda - 1$. Then $x_0 \geq \langle \mathbf{1}, \nu_0 \rangle = 1$ and

$$\lambda G^r \left((F^{\lambda \eta_*})^{-1}((x_0 - 1)^+) \right) = \int_0^{(x_0-1)^+} h^r((F^{\lambda \eta_*})^{-1}(y)) dy = p = \lambda - 1.$$

Hence x_0 belongs to the set B_λ defined in (3.15). The last assertion can be verified directly by substituting the initial condition into the fluid equations. This completes the proof of the lemma. \square

6. THE LIMIT OF SCALED STATIONARY DISTRIBUTIONS

This section is devoted to the proof of Theorem 3.3. Due to Theorem 4.9, under Assumption 5, there exists a sequence of \mathcal{Y} -valued random variables $\bar{Y}_*^{(N)} = (\bar{\alpha}_{E,*}^{(N)}, \bar{X}_*^{(N)}, \bar{\nu}_*^{(N)}, \bar{\eta}_*^{(N)})$, $N \in \mathbb{N}$, where the distribution of $\bar{Y}_*^{(N)}$ is stationary for the dynamics of the fluid-scaled N -server queue with abandonment (note that we do not require the stationary distribution of the N -server system to be unique here). Let $(x_*, (\lambda \wedge 1)\nu_*, \lambda\eta_*)$ be the unique element of the invariant manifold \mathcal{I}_λ . The main result of this section is to show that, under Assumptions 1–4, as $N \rightarrow \infty$,

$$(6.54) \quad (\bar{X}_*^{(N)}, \bar{\nu}_*^{(N)}, \bar{\eta}_*^{(N)}) \Rightarrow (x_*, (\lambda \wedge 1)\nu_*, \lambda\eta_*).$$

In order to prove this result, we first show in Section 6.1 that the sequence $\{(\bar{X}_*^{(N)}, \bar{\nu}_*^{(N)}, \bar{\eta}_*^{(N)})\}$, $N \in \mathbb{N}$ is tight. Then, in Section 6.2, we prove (6.54) by showing that the weak limit of every convergent subsequence must be an invariant state and using the fact that there is a unique invariant state. For both these results, for each $N \in \mathbb{N}$, we will find it convenient to define

$$(6.55) \quad \bar{Y}^{(N)} = (\bar{E}^{(N)}, \bar{X}^{(N)}, \bar{\nu}^{(N)}, \bar{\eta}^{(N)})$$

to be the fluid-scaled process for the N -server queue with abandonment associated with the initial condition $\bar{Y}^{(N)}(0) = (\bar{E}^{(N)}, \bar{X}_*^{(N)}, \bar{\nu}_*^{(N)}, \bar{\eta}_*^{(N)})$, where $\bar{E}^{(N)}$ is the fluid scaled stationary renewal process determined from the backward recurrence time process $\alpha_E^{(N)}$ with initial distribution $\bar{\alpha}_{E,*}^{(N)}$. Let $\bar{Q}^{(N)}$, $\bar{R}^{(N)}$, $\bar{K}^{(N)}$ be the associated fluid-scaled auxiliary processes described in Section 2.2.

6.1. Tightness. Recall the criteria for tightness of measure-valued random variables in Proposition 4.3.

Lemma 6.1. *Let $c \in [0, H^r)$. For each integer $n \geq 2$, $\bar{\eta}_*^{(N)}$ and $\bar{\nu}_*^{(N)}$ satisfy the following relations:*

$$(6.56) \quad \begin{aligned} & \mathbb{E} \left[\bar{\eta}_*^{(N)}[c, H^r] \right] \\ &= \mathbb{E} \left[\int_{[0, H^r)} \frac{1 - G^r(x + nc)}{1 - G^r(x)} \bar{\eta}_*^{(N)}(dx) \right] + \mathbb{E} \left[\int_0^c \sum_{j=2}^n (1 - G^r(jc - s)) d\bar{E}^{(N)}(s) \right], \end{aligned}$$

$$(6.57) \quad \begin{aligned} & \mathbb{E} \left[\bar{\nu}_*^{(N)}[c, H^s] \right] \\ &= \mathbb{E} \left[\int_{[0, H^s)} \frac{1 - G^s(x + nc)}{1 - G^s(x)} \bar{\nu}_*^{(N)}(dx) \right] + \mathbb{E} \left[\int_0^c \sum_{j=2}^n (1 - G^s(jc - s)) d\bar{K}^{(N)}(s) \right]. \end{aligned}$$

Proof. We only prove (6.56) since (6.57) can be proved in the same way. Fix $c \in [0, H^r)$. Dividing both sides of (2.8) by N and setting $\bar{\eta}_0^{(N)} = \bar{\eta}_*^{(N)}$, we obtain for each bounded measurable function f on \mathbb{R}_+ and $t > 0$,

$$(6.58) \quad \begin{aligned} \mathbb{E} \left[\langle f, \bar{\eta}_t^{(N)} \rangle \right] &= \mathbb{E} \left[\int_{[0, H^r)} f(x+t) \frac{1 - G^r(x+t)}{1 - G^r(x)} \bar{\eta}_*^{(N)}(dx) \right] \\ &+ \mathbb{E} \left[\int_0^t f(t-s) (1 - G^r(t-s)) d\bar{E}^{(N)}(s) \right]. \end{aligned}$$

Since the initial conditions are stationary, $\bar{\eta}_t^{(N)}$ has the same distribution as $\bar{\eta}_*^{(N)}$ for every $t \geq 0$. Therefore, by substituting $f = \mathbb{1}_{[c, H^r)}$ and $t = c$ in (6.58), and noting that $\mathbb{1}_{[c, H^r)}(x + c) = 1$ for every $x \geq 0$ and $\mathbb{1}_{[c, H^r)}(c - s) = 0$ for every $s \in [0, c]$, we obtain

$$\begin{aligned} \mathbb{E} \left[\bar{\eta}_*^{(N)}[c, H^r) \right] &= \mathbb{E} \left[\bar{\eta}_c^{(N)}[c, H^r) \right] = \mathbb{E} \left[\int_{[0, H^r)} \frac{1 - G^r(x + c)}{1 - G^r(x)} \bar{\eta}_*^{(N)}(dx) \right] \\ &= \mathbb{E} \left[\int_{[0, H^r)} \frac{1 - G^r(x + c)}{1 - G^r(x)} \bar{\eta}_c^{(N)}(dx) \right]. \end{aligned}$$

Next, by choosing $f = (1 - G^r(\cdot + c))/(1 - G^r(\cdot))$ and $t = c$ in (6.58), we obtain

$$\begin{aligned} \mathbb{E} \left[\int_{[0, H^r)} \frac{1 - G^r(x + c)}{1 - G^r(x)} \bar{\eta}_c^{(N)}(dx) \right] &= \mathbb{E} \left[\int_{[0, H^r)} \frac{1 - G^r(x + 2c)}{1 - G^r(x)} \bar{\eta}_*^{(N)}(dx) \right] \\ &\quad + \mathbb{E} \left[\int_0^c (1 - G^r(2c - s)) d\bar{E}^{(N)}(s) \right]. \end{aligned}$$

By combining the last two displays, we see that

$$\begin{aligned} \mathbb{E}[\bar{\eta}_*^{(N)}[c, H^r)] &= \mathbb{E} \left[\int_{[0, H^r)} \frac{1 - G^r(x + 2c)}{1 - G^r(x)} \bar{\eta}_*^{(N)}(dx) \right] \\ &\quad + \mathbb{E} \left[\int_0^c (1 - G^r(2c - s)) d\bar{E}^{(N)}(s) \right]. \end{aligned}$$

Thus, we have shown that (6.56) holds for $n = 2$. Suppose that for some integer $m \geq 2$, (6.56) holds for $n = m$, i.e.,

$$(6.59) \quad \begin{aligned} \mathbb{E}[\bar{\eta}_*^{(N)}[c, H^r)] &= \mathbb{E} \left[\int_{[0, H^r)} \frac{1 - G^r(x + mc)}{1 - G^r(x)} \bar{\eta}_*^{(N)}(dx) \right] \\ &\quad + \mathbb{E} \left[\int_0^c \sum_{j=2}^m (1 - G^r(jc - s)) d\bar{E}^{(N)}(s) \right]. \end{aligned}$$

By choosing $f = (1 - G^r(\cdot + mc))/(1 - G^r(\cdot))$ and $t = c$ in (6.58) and using the fact that $\bar{\eta}_c^{(N)}$ has the same distribution as $\bar{\eta}_*^{(N)}$, we obtain

$$(6.60) \quad \begin{aligned} &\mathbb{E} \left[\int_{[0, H^r)} \frac{1 - G^r(x + mc)}{1 - G^r(x)} \bar{\eta}_*^{(N)}(dx) \right] \\ &= \mathbb{E} \left[\int_{[0, H^r)} \frac{1 - G^r(x + mc)}{1 - G^r(x)} \bar{\eta}_c^{(N)}(dx) \right] \\ &= \mathbb{E} \left[\int_{[0, H^r)} \frac{1 - G^r(x + (m + 1)c)}{1 - G^r(x)} \bar{\eta}_*^{(N)}(dx) \right] \\ &\quad + \mathbb{E} \left[\int_0^c (1 - G^r((m + 1)c - s)) d\bar{E}^{(N)}(s) \right]. \end{aligned}$$

This, together with (6.59), yields (6.56) with $n = m + 1$. This completes the induction argument and we have the desired result. \square

Theorem 6.2. *The sequence $\{(\bar{X}_*^{(N)}, \bar{\nu}_*^{(N)}, \bar{\eta}_*^{(N)})\}_{N \in \mathbb{N}}$ is tight.*

Proof. We first show that $\{\bar{\eta}_*^{(N)}\}_{N \in \mathbb{N}}$ is tight. Note that $\langle \mathbf{1}, \bar{\eta}^{(N)} \rangle$ can be viewed as the fluid scaled queue length process associated with an infinite-server queue with arrival process $\bar{E}^{(N)}$ and service distribution function G^r . By Little's Law (cf. [17]), we know that $\mathbb{E}[\langle \mathbf{1}, \bar{\eta}_*^{(N)} \rangle] = \bar{\lambda}^{(N)} \theta^r$, where θ^r , defined in Assumption 1, is the mean of G^r . Due to the convergence of $\bar{\lambda}^{(N)}$ to λ stated in Assumption 3, this implies

$$(6.61) \quad \sup_{N \in \mathbb{N}} \mathbb{E} \left[\langle \mathbf{1}, \bar{\eta}_*^{(N)} \rangle \right] < \infty,$$

which establishes the first criterion for tightness.

Next, for each n , the function $(1 - G^r(\cdot + nc))/(1 - G^r(\cdot))$ is bounded by 1 and converges to 0 as $n \rightarrow \infty$, an application of the dominated convergence theorem shows that

$$(6.62) \quad \lim_{n \rightarrow \infty} \mathbb{E} \left[\int_{[0, H^r)} \frac{1 - G^r(x + nc)}{1 - G^r(x)} \bar{\eta}_*^{(N)}(dx) \right] = 0.$$

Sending $n \rightarrow \infty$ on the right-hand side of (6.56), and using (6.62) and the monotone convergence theorem, we have

$$(6.63) \quad \mathbb{E}[\bar{\eta}_*^{(N)}[c, H^r)] = \mathbb{E} \left[\int_0^c \sum_{j=2}^{\infty} (1 - G^r(jc - s)) d\bar{E}^{(N)}(s) \right].$$

On the other hand, we also have the simple estimate

$$(6.64) \quad \begin{aligned} \mathbb{E} \left[\int_0^c (1 - G^r(2c - s)) d\bar{E}^{(N)}(s) \right] &\leq (1 - G^r(c)) \mathbb{E}[\bar{E}^{(N)}(c)] \\ &= c(1 - G^r(c)) \frac{\mathbb{E}[\bar{E}^{(N)}(c)]}{c}. \end{aligned}$$

By carrying out an integration by parts on $\int_0^\infty (1 - G^r(x)) dx$, it follows that

$$\int_{[0, H^r)} (1 - G^r(x)) dx = \lim_{x \in H^r} x(1 - G^r(x)) + \int_{[0, H^r)} x g^r(x) dx.$$

However, since the mean θ^r is finite by (3.11), it follows that $c(1 - G^r(c)) \rightarrow 0$ as $c \rightarrow H^r$. In addition, since $\mathbb{E}[\bar{E}^{(N)}(c)]/c \rightarrow \bar{\lambda}^{(N)}$, as $c \rightarrow \infty$, by the elementary renewal theorem and $\bar{\lambda}^{(N)} \rightarrow \lambda$ as $N \rightarrow \infty$ by Assumption 3(1), it follows that

$$(6.65) \quad \limsup_{c \rightarrow H^r} \sup_N \frac{\mathbb{E}[\bar{E}^{(N)}(c)]}{c} < \infty.$$

Thus, taking the supremum over N and then $c \rightarrow H^r$ in (6.64), we obtain

$$(6.66) \quad \lim_{c \rightarrow H^r} \sup_N \mathbb{E} \left[\int_0^c (1 - G^r(2c - s)) d\bar{E}^{(N)}(s) \right] = 0.$$

Since $1 - G^r(\cdot)$ is a decreasing function, for $s \in [0, c]$,

$$\sum_{j=3}^{\infty} c(1 - G^r(jc - s)) \leq \int_{[2c-s, H^r)} (1 - G^r(x)) dx \leq \int_{[c, H^r)} (1 - G^r(x)) dx.$$

Therefore, we have

$$\sup_N \mathbb{E} \left[\int_0^c \sum_{j=3}^{\infty} (1 - G^r(jc - s)) d\bar{E}^{(N)}(s) \right] \leq \sup_N \frac{\mathbb{E}[\bar{E}^{(N)}(c)]}{c} \int_{[c, H^r)} (1 - G^r(x)) dx,$$

which tends to zero, as $c \rightarrow H^r$, because of (6.65). By combining the last assertion with (6.63) and (6.66), we see that

$$(6.67) \quad \lim_{c \rightarrow H^r} \sup_N \mathbb{E} \left[\bar{\eta}_*^{(N)}[c, H^r] \right] = 0,$$

which establishes the second criterion for tightness. Thus, the sequence $\{\bar{\eta}_*^{(N)}\}_{N \in \mathbb{N}}$ is tight.

We next show that $\{\bar{\nu}_*^{(N)}\}_{N \in \mathbb{N}}$ is tight. The analog of (6.61) holds for $\{\bar{\nu}_*^{(N)}\}_{N \in \mathbb{N}}$ automatically since $\langle \mathbf{1}, \bar{\nu}_*^{(N)} \rangle \leq 1$ for each N . On the other hand, the analog of (6.67) can be shown to hold for $\{\bar{\nu}_*^{(N)}\}_{N \in \mathbb{N}}$ by using (6.57) and an argument similar to that used above to establish (6.67), along with the additional observation that $\mathbb{E}[\bar{K}^{(N)}(c)] \leq \mathbb{E}[\bar{E}^{(N)}(c)] + \mathbb{E}[\langle \mathbf{1}, \bar{\eta}_*^{(N)} \rangle]$ implies $\limsup_{c \rightarrow H^r} \sup_N \mathbb{E}[\bar{K}^{(N)}(c)]/c < \infty$. Thus $\{\bar{\nu}_*^{(N)}\}_{N \in \mathbb{N}}$ is tight.

Finally, we show that the sequence of \mathbb{R}_+ -valued random variables $\{\bar{X}_*^{(N)}\}_{N \in \mathbb{N}}$ is tight. Since $\bar{X}_*^{(N)} \leq 1 + \langle \mathbf{1}, \bar{\eta}_*^{(N)} \rangle$ for each N , $\sup_N \mathbb{E}[\bar{X}_*^{(N)}] \leq 1 + \sup_N \mathbb{E}[\langle \mathbf{1}, \bar{\eta}_*^{(N)} \rangle]$, which is finite due to (6.61). The tightness of $\{\bar{X}_*^{(N)}\}_{N \in \mathbb{N}}$ is a direct consequence of Markov's inequality. \square

6.2. Convergence. Let $\{(\bar{X}_*^{(N)}, \bar{\nu}_*^{(N)}, \bar{\eta}_*^{(N)})\}_{N \in \mathbb{N}}$ be any sequence of (marginals) of the scaled stationary distributions, which was shown to be tight in Section 6.1. In this section we establish the convergence (6.54) by showing that any convergent subsequence must have the invariant state as its limit, and then invoke uniqueness of the invariant state.

Lemma 6.3. *The sequence $\{\bar{\eta}_*^{(N)}\}_{N \in \mathbb{N}}$ converges weakly to $\lambda \eta_*$ as $N \rightarrow \infty$.*

Proof. Since the sequence $\{(\bar{X}_*^{(N)}, \bar{\nu}_*^{(N)}, \bar{\eta}_*^{(N)})\}_{N \in \mathbb{N}}$ is tight by Theorem 6.2, there exists a convergent subsequence which, by some abuse of notation, we denote again by $\{(\bar{X}_*^{(N)}, \bar{\nu}_*^{(N)}, \bar{\eta}_*^{(N)})\}_{N \in \mathbb{N}}$. Let $(\tilde{X}_*, \tilde{\nu}_*, \tilde{\eta}_*)$ be the corresponding limit. Due to the Skorokhod representation theorem, without loss of generality, we may assume that $(\bar{X}_*^{(N)}, \bar{\nu}_*^{(N)}, \bar{\eta}_*^{(N)})$ converges almost surely, as $N \rightarrow \infty$, to $(\tilde{X}_*, \tilde{\nu}_*, \tilde{\eta}_*)$. Together with Assumption 3, this implies that Assumption 3.1 of [14] holds. For each $N \in \mathbb{N}$, consider the fluid-scaled process $\bar{Y}^{(N)}$ in (6.55). Since Assumption 3.1 of [14] holds for $\{\bar{Y}^{(N)}\}_{N \in \mathbb{N}}$ and Assumption 4 is also satisfied, the hypotheses of Theorems 6.1 and 7.1 of [14] are satisfied. Thus, we can conclude that the sequence $\{\bar{\eta}^{(N)}\}_{N \in \mathbb{N}}$ is tight and any weak limit, denoted by $\tilde{\eta}$, satisfies (3.42) of [14] with $\bar{\eta} = \tilde{\eta}$ and $\bar{\eta}_0 = \tilde{\eta}_*$. (Note that Assumption 3.2 of [14] is not used in Theorem 7.1 of [14] to prove (3.42).) Therefore, it follows from Theorem 4.1 of [15] that for $f \in \mathcal{C}_b(\mathbb{R}_+)$ and $t > 0$,

$$(6.68) \quad \langle f, \tilde{\eta}_t \rangle = \int_{[0, H^r)} f(x+t) \frac{1 - G^r(x+t)}{1 - G^r(x)} \tilde{\eta}_*(dx) + \int_0^t f(t-s)(1 - G^r(t-s)) \lambda ds.$$

Since, for each $t > 0$, $\tilde{\eta}_t^{(N)}$ and $\tilde{\eta}_*^{(N)}$ have the same distribution, $\tilde{\eta}_t$ and $\tilde{\eta}_*$ must also have the same distribution. If we restrict the two measures $\tilde{\eta}_t$ and $\tilde{\eta}_*$ to $[0, t]$, then those two restricted measures again have the same distribution. By substituting $f = \mathbb{1}_{[0, t]}$ in (6.68), we can easily see that measure $\tilde{\eta}_t$, restricted to $[0, t]$, is the same as $\lambda \mathbb{1}_{[0, t]}(x)(1 - G^r(x))dx$. Thus, the measures $\tilde{\eta}_*$ and $\lambda\eta_*$, restricted to $[0, t]$, also coincide. Thus, $\tilde{\eta}_* = \lambda\eta_*$. This shows that $\{\tilde{\eta}_*^{(N)}\}$ converges weakly to $\lambda\eta_*$, as $N \rightarrow \infty$. \square

Proof of Theorem 3.3. Since the sequence $\{(\bar{X}_*^{(N)}, \bar{\nu}_*^{(N)}, \bar{\eta}_*^{(N)})\}_{N \in \mathbb{N}}$ is tight by Theorem 6.2, there exists a convergent subsequence which, by some abuse of notation, we denote again by $\{(\bar{X}_*^{(N)}, \bar{\nu}_*^{(N)}, \bar{\eta}_*^{(N)})\}_{N \in \mathbb{N}}$. Let $(\tilde{X}_*, \tilde{\nu}_*, \tilde{\eta}_*)$ be the corresponding limit. Due to the Skorokhod representation theorem, without loss of generality, we may assume that $(\bar{X}_*^{(N)}, \bar{\nu}_*^{(N)}, \bar{\eta}_*^{(N)})$ converges almost surely to $(\tilde{X}_*, \tilde{\nu}_*, \tilde{\eta}_*)$. For each $N \in \mathbb{N}$, consider the fluid-scaled process $\bar{Y}^{(N)}$ in (6.55). It follows from Assumption 3 that the initial condition $\{\bar{Y}^{(N)}(0)\}_{N \in \mathbb{N}}$ satisfies Assumption 3.1 of [14] holds. By Lemma 6.3 and the fact that $\bar{E}(t) = \lambda t$, the initial condition $\{\bar{Y}^{(N)}(0), N \in \mathbb{N}\}$ also satisfies Assumption 3.2 of [14]. By Assumption 4 and Theorem 3.6 of [14], it then follows that $(\bar{X}_*^{(N)}, \bar{\nu}_*^{(N)}, \bar{\eta}_*^{(N)})$ converges weakly, as $N \rightarrow \infty$, to the unique solution $(\bar{X}, \bar{\nu}, \bar{\eta})$ of the fluid equations defined in Definition 5.1, associated with the initial condition $(\tilde{X}_*, \tilde{\nu}_*, \tilde{\eta}_*)$. By the stationarity of $(\bar{X}_*^{(N)}, \bar{\nu}_*^{(N)}, \bar{\eta}_*^{(N)})$ for each N , it follows that $(\tilde{X}_*, \tilde{\nu}_*, \tilde{\eta}_*)$ is an invariant state of the fluid limit $(\bar{X}, \bar{\nu}, \bar{\eta})$. Therefore $(\tilde{X}_*, \tilde{\nu}_*, \tilde{\eta}_*)$ belongs to the invariant manifold. Since the invariant manifold has a single element by Assumption 2, the usual argument by contradiction shows that $(\tilde{X}_*, \tilde{\nu}_*, \tilde{\eta}_*)$ is in fact the limit of the sequence $\{(\bar{X}_*^{(N)}, \bar{\nu}_*^{(N)}, \bar{\eta}_*^{(N)})\}_{N \in \mathbb{N}}$. \square

7. CONCLUDING REMARKS

We can establish ergodicity of the state processes, under an additional condition. Let

$$\varrho^r \doteq \sup\{u \in [0, H^r) : g^r = 0 \text{ a.e. on } [a, a + u] \text{ for some } a \in [0, \infty)\}$$

and

$$\varrho^s \doteq \sup\{u \in [0, H^s) : g^s = 0 \text{ a.e. on } [a, a + u] \text{ for some } a \in [0, \infty)\}.$$

Assumption 6. *The following three conditions hold:*

- (1) $H^r = H^s = \infty$;
- (2) $\varrho \doteq \varrho^r \vee \varrho^s < \infty$;
- (3) For every interval $[a, b] \subset [0, \infty)$ with $b - a > 0$, $F^{(N)}(b) - F^{(N)}(a) > 0$.

Theorem 7.1. *The Markov process $\{Y_t, \mathcal{F}_t\}$ is ergodic in the sense that it has a unique stationary distribution, and the distribution of $Y(t)$ converges in total variation, as $t \rightarrow \infty$, to this unique stationary distribution.*

Theorem 7.1, whose proof is deferred to the Appendix, validates the rightward arrow on the top of the “interchange of limits” diagram presented in Figure 1. On the other hand, the fluid limit theorem (Theorem 3.6 of [14]) justifies the downward arrow on the left-hand side of Figure 1. The focus of this work has been on

$$\begin{array}{ccc}
 (\bar{X}^{(N)}(t), \bar{\nu}_t^{(N)}, \bar{\eta}_t^{(N)}) & \longrightarrow & (\bar{X}_*^{(N)}, \bar{\nu}_*^{(N)}, \bar{\eta}_*^{(N)}) \\
 \downarrow & & \downarrow \\
 (\bar{X}(t), \bar{\nu}_t, \bar{\eta}_t) & \longrightarrow & (\bar{X}_*, \bar{\nu}_*, \bar{\eta}_*)
 \end{array}$$

FIGURE 1. Interchange of Limits Diagram

understanding the convergence represented by the downward arrow on the right-hand side of Figure 1. When there is a unique invariant state, this convergence is established in Theorem 3.3. Although this question is not directly relevant to the characterization of the stationary distributions, it is natural, in this setting, to ask whether the diagram in Figure 1 commutes, namely, whether the fluid limit from any initial condition converges, as $t \rightarrow \infty$, to the unique invariant state. In Section 7.1 we briefly discuss why the study of the long-time behavior of the fluid limit is a non-trivial task. Furthermore, in Section 7.2 we provide a very simple counterexample that shows that the diagram in Figure 1 need not commute and thus the limits $N \rightarrow \infty$ and $t \rightarrow \infty$ cannot always be interchanged.

7.1. Long-time behavior of the fluid limit. The long-time behavior of the fluid limit is non-trivial even in the absence of abandonment. For example, in the absence of abandonment, it was proved in Theorem 3.9 of [15] that when the service time distribution G^s has a second moment and its hazard rate function h^s is either bounded or lower-semicontinuous on (m_0, H^s) for some $m_0 < H^s$, $\bar{\nu}_t \rightarrow (\lambda \wedge 1)\nu_*$ as $t \rightarrow \infty$. The question of whether the second moment condition on the distribution is necessary for this convergence is still unresolved. Even under the second moment assumption, the long-time behavior of the component \bar{X} of the fluid limit is not easy to describe except in the cases when the system is subcritical ($\lambda < 1$) or when the system is critical or supercritical ($\lambda \geq 1$) and the service distribution is exponential. In the former case, it follows from Theorem 3.9 of [15] that $\bar{X}(t) \rightarrow \lambda \langle \mathbf{1}, \nu_* \rangle$ as $t \rightarrow \infty$, while in the latter case, if the initial condition satisfies $\bar{X}(0) \geq 1$ and $\bar{\nu}_0 \in \mathcal{M}_F[0, \infty)$, then it is easy to see that the fluid limit is given explicitly by $\bar{X}(t) = \bar{X}(0) + (\lambda - 1)t$ and $\bar{\nu}_t(dx) = \mathbb{1}_{[0, t]}e^{-x}dx + \mathbb{1}_{(t, \infty)}(x)e^{-t}\bar{\nu}_0(dx - t)$. Therefore, at criticality ($\lambda = 1$), if $\bar{X}(0) = 1$ then $\bar{X}(t) = \bar{X}(0)$ for every $t > 0$. In particular, $\bar{X}(t) \rightarrow 1$ as $t \rightarrow \infty$. However, as the following example demonstrates, the critical fluid limit need not converge to 1 (even if starting critically loaded) when the service is non-exponential.

Example 7.2. Let the fluid arrival rate be $\bar{E}(t) = t$, $t > 0$, and let the service time distribution G^s be the Erlang distribution with density

$$g^s(x) = 4xe^{-2x}, \quad x \geq 0.$$

A simple calculation shows that $\int_0^\infty (1 - G^s(x)) dx = 1$. Let $(\bar{X}, \bar{\nu})$ be the solution to the fluid equations without abandonment (see Definition 5.2) associated with the initial condition $(\mathbf{1}, \delta_0)$. We show below that in this case, $\lim_{t \rightarrow \infty} \bar{X}(t) = 5/4$, which is bigger than $1 = \bar{X}(0)$. In fact, since $\bar{\nu}_0 = \delta_0$, a straightforward calculation shows that

$$\langle h^s, \bar{\nu}_0 \rangle = \int_0^\infty \frac{g^s(x)}{1 - G^s(x)} \bar{\nu}_0(dx) = \frac{g^s(0)}{1 - G^s(0)} = g^s(0) = 0.$$

Define

$$\kappa \doteq \inf\{t \geq 0 : \langle h^s, \bar{\nu}_t \rangle \geq 1\}.$$

Then, since h^s is continuous, $\kappa > 0$ and for $t \in [0, \kappa)$, $\langle h^s, \bar{\nu}_t \rangle < \lambda = 1$. Hence, $\langle \mathbf{1}, \bar{\nu}_t \rangle = 1$ for each $t \in [0, \kappa)$ and $d\bar{K}/dt(t) = \langle h^s, \bar{\nu}_t \rangle$. For each $t \in [0, \kappa)$,

$$\langle h^s, \bar{\nu}_t \rangle = g^s(t) + \int_0^t g^s(t-s) d\bar{K}/dt(s) ds = g^s(t) + \int_0^t g^s(t-s) \langle h^s, \bar{\nu}_s \rangle ds.$$

By the key renewal theorem, we have

$$\langle h^s, \bar{\nu}_t \rangle = u^s(t) = 1 - e^{-4t}.$$

Since $u^s(t) < 1$ for all $t \geq 0$, we must have that $\kappa = \infty$, $\langle \mathbf{1}, \bar{\nu}_t \rangle = 1$ for all $t \geq 0$, and

$$\lim_{t \rightarrow \infty} \bar{Q}(t) = \int_0^\infty (1 - u^s(t)) dt = \int_0^\infty e^{-4t} dt = 1/4,$$

which yields the convergence of $\bar{X}(t)$ to $5/4$ as $t \rightarrow \infty$.

To emphasize that this phenomenon is not an artefact of the fact that the initial condition was chosen to be singular with respect to Lebesgue measure, we show that we can modify the above example by choosing $\bar{\nu}_0$ to be absolutely continuous with respect to the Lebesgue measure. For example, for some $\alpha \in (0, \infty)$, define

$$q(x) \doteq \begin{cases} \frac{1+2x}{\alpha+\alpha^2} & \text{if } x \in [0, \alpha], \\ 0 & \text{otherwise,} \end{cases}$$

and let $\bar{\nu}_0(dx) = q(x)dx$. Then $\langle \mathbf{1}, \bar{\nu}_0 \rangle = \int_0^\alpha q(x)dx = 1$, $\langle h^s, \bar{\nu}_t \rangle = 1 - ((1-\alpha)/(\alpha+1))e^{-4t}$ for each $t \geq 0$. Hence, when $\alpha < 1$, we have $\langle h^s, \bar{\nu}_t \rangle < 1$ and $\langle \mathbf{1}, \bar{\nu}_t \rangle = 1$ for all $t \geq 0$. This implies that, when $\alpha < 1$,

$$\lim_{t \rightarrow \infty} \bar{Q}(t) = \int_0^\infty \frac{1-\alpha}{\alpha+1} e^{-4t} dt = \frac{1-\alpha}{4(\alpha+1)} > 0,$$

showing that $\lim_{t \rightarrow \infty} \bar{X}(t) > 1$.

7.2. A Counterexample (Invalidity of the Interchange of Limits). In this section we show that, even for an M/M/N queue (both with and without abandonments), an ‘‘interchange of limits’’ need not hold, i.e., the diagram presented in Figure 1 may not commute.

Consider the sequence of state processes $(X^{(N)}, \nu^{(N)})$, $N \in \mathbb{N}$, of N -server queues without abandonment, where the service time distribution G^s is exponential with rate 1. For the N -th queue, let the arrival process $E^{(N)}$ be a Poisson process with parameter $\lambda^{(N)} = N - 1$ and suppose that there exists $\bar{\nu}_0 \in \mathcal{M}_F[0, \infty)$ with $\langle \mathbf{1}, \bar{\nu}_0 \rangle \leq 1$ such that a.s., as $N \rightarrow \infty$,

$$(7.69) \quad (\bar{X}^{(N)}(0), \bar{\nu}_0^{(N)}) \rightarrow (2, \bar{\nu}_0).$$

Given the exponentiality of the service time distribution, it immediately follows that Assumption 2 of [15] is satisfied. Moreover, since (7.69) holds and $\bar{\lambda}^{(N)} = (N-1)/N \rightarrow 1$ as $N \rightarrow \infty$, it follows that Assumption 1 of [15] also holds with $\lambda = 1$. On the other hand, since $G^r(x) = 0$ for all $x \in [0, \infty)$, Assumption 2 fails to hold because in this case $B_1 = [0, \infty)$ and so the invariant manifold has uncountably many points.

Now, since Assumptions 1 and 2 of [15] ensure that the conditions of the fluid limit result, Theorem 3.7 of [15], are satisfied, it follows that, a.s., as $N \rightarrow \infty$,

$(\bar{X}^{(N)}, \bar{\nu}^{(N)})$ converges weakly to the unique solution $(\bar{X}, \bar{\nu})$ of the fluid equations associated with $(\mathbf{1}, 2, \bar{\nu}_0)$, and, using the exponentiality of the service time distribution, it is easily verified that the fluid limit is given explicitly by $\bar{X}(t) = \bar{X}(0) = 2$ and $\bar{\nu}_t(dx) = \mathbb{1}_{[0,t]}e^{-x}dx + \mathbb{1}_{(t,\infty)}(x)e^{-t}\bar{\nu}_0(d(x-t))$.

For each $N \in \mathbb{N}$, since the arrival rate, which equals $N - 1$, is less than the total service rate N , by (3.2.4) and (3.2.5) of [4], it follows that $X^{(N)}$ is ergodic and has the following stationary distribution:

$$\mathbb{P}(X_*^{(N)} = k) = \begin{cases} \frac{(N-1)^k}{k!} p_0 & \text{if } k = 0, 1, \dots, N-1, \\ \frac{(N-1)^k}{N!N^{k-N}} p_0 & \text{if } k = N, N+1, \dots, \end{cases}$$

where

$$p_0 \doteq \left\{ \sum_{i=0}^{N-1} \frac{(N-1)^i}{i!} + \frac{(N-1)^N}{(N-1)!} \right\}^{-1}.$$

Elementary calculations show that

$$\begin{aligned} \mathbb{P}(X_*^{(N)} \geq N + N/2) &= \sum_{k=N+N/2}^{\infty} \frac{(N-1)^k}{N!N^{k-N}} p_0 \\ &= \frac{N^N}{N!} p_0 \sum_{k=N+N/2}^{\infty} \left(\frac{N-1}{N} \right)^k \\ &= \frac{N^N}{N!} p_0 \left(\frac{N-1}{N} \right)^{N+N/2} N \\ &= \frac{(N-1)^N}{(N-1)!} p_0 \left(\frac{N-1}{N} \right)^{N/2} \leq \left(\frac{N-1}{N} \right)^{N/2}. \end{aligned}$$

We then have

$$\limsup_{N \rightarrow \infty} \mathbb{P}(\bar{X}_*^{(N)} \geq 3/2) = \limsup_{N \rightarrow \infty} \mathbb{P}(X_*^{(N)} \geq N + N/2) \leq e^{-2} < 1.$$

Using the distribution of $X_*^{(N)}$, it can also be shown that

$$\sup_{N \in \mathbb{N}} \mathbb{E}[\bar{X}_*^{(N)}] = \sup_{N \in \mathbb{N}} \frac{\mathbb{E}[X_*^{(N)}]}{N} \leq 3.$$

An application of Markov's inequality then shows that the sequence of $\{\bar{X}_*^{(N)}\}_{N \in \mathbb{N}}$ is tight. Thus, (7.70) clearly shows that $\bar{X}_*^{(N)}$ does not converge (even along a subsequence) to 2 as $N \rightarrow \infty$.

A minor modification of the above example shows that the interchange of limits can also fail to hold in the presence of abandonments. For the same sequence of queues described above, suppose that customers abandon the queue according to a non-trivial patience time distribution G^r satisfying Assumption 4 and having support in $(3, \infty)$. For each $N \in \mathbb{N}$, consider the marginal state process $(X^{(N)}, \nu^{(N)}, \eta^{(N)})$. Suppose that there exists $(2, \bar{\nu}_0, \bar{\eta}_0) \in \mathcal{S}_0$ such that, a.s., as $N \rightarrow \infty$,

$$(7.70) \quad (\bar{X}^{(N)}(0), \bar{\nu}_0^{(N)}, \bar{\eta}_0^{(N)}) \rightarrow (2, \bar{\nu}_0, \bar{\eta}_0).$$

Given the assumption imposed on the patience time distribution, Assumption 2 fails to hold because in this case $B_1 = [0, 3]$. By the previous argument, Assumptions 4 and 3 are satisfied. Since Assumptions 1, 3 and 4 ensure that the conditions of the fluid limit result, Theorem 3.6 of [14], are satisfied, it follows that, a.s., as $N \rightarrow \infty$, $(\bar{X}^{(N)}, \bar{\nu}^{(N)}, \bar{\eta}^{(N)})$ converges weakly to the unique solution $(\bar{X}, \bar{\nu}, \bar{\eta})$ of the fluid equations associated with $(\mathbf{1}, 2, \bar{\nu}_0, \bar{\eta}_0)$. By the exponentiality of the service time distribution, we have $\bar{X}(t) = \bar{X}(0) = 2$ and $\bar{R}(t) = 0$ for each $t \geq 0$. On the other hand, let $\bar{Y}_*^{(N)} = (\bar{\alpha}_{E,*}^{(N)}, \bar{X}_*^{(N)}, \bar{\nu}_*^{(N)}, \bar{\eta}_*^{(N)})$ be the stationary distribution associated with the fluid-scaled state process, which exists by Theorem 4.9. By a simple coupling argument, it can be shown that $X^{(N)}$ is stochastically dominated by the corresponding state $\tilde{X}^{(N)}$ of an $M/M/N$ queue without abandonment that has the same arrival process $E^{(N)}$ and the same initial condition (i.e., $\mathbb{P}(\tilde{X}^{(N)} \geq c) \geq \mathbb{P}(X^{(N)} \geq c)$ for every $c > 0$). Together with the previous discussion of the case without abandonment, this can be used to show that $\limsup_{N \rightarrow \infty} \mathbb{P}(\bar{X}_*^{(N)} \geq 3/2) < 1$ and $\{\bar{X}_*^{(N)}\}_{N \in \mathbb{N}}$ is tight. Thus, in this case too, $\bar{X}_*^{(N)} = \lim_{t \rightarrow \infty} \bar{X}_*^{(N)}(t)$ does not converge (even along a convergent subsequence) to $2 = \lim_{t \rightarrow \infty} \bar{X}(t) = \lim_{t \rightarrow \infty} \lim_{N \rightarrow \infty} \bar{X}^{(N)}(t)$, as $N \rightarrow \infty$.

APPENDIX A. PROOF OF THEOREM 7.1

By Theorem 6.1 of [19], to show that the Feller process $\{Y_t, \mathcal{F}_t\}_{t \geq 0}$ is ergodic, it suffices to establish the second assertion in Lemma A.3 and Theorem A.5 below which, respectively, show that the skeleton chain $\{Y_n\}_{n \in \mathbb{N}}$ is ψ -irreducible and that $\{Y_t, \mathcal{F}_t\}_{t \geq 0}$ is positive Harris recurrent. Let

$$\mathcal{Z} = \{(\alpha, 0, \mathbf{0}, \mathbf{0}) \in \mathcal{Y} : \alpha \in [\varrho + 1, \infty)\}.$$

For each Borel subset A of \mathcal{Z} , there exists a Borel subset Γ_A of $[\varrho + 1, \infty)$ such that

$$(1.71) \quad A = \{(\alpha, 0, \mathbf{0}, \mathbf{0}) \in \mathcal{Y} : \alpha \in \Gamma_A\}.$$

Lemma A.1. *There exists a strictly positive continuous function C on \mathcal{Y} such that for every $y = (\alpha, x, \sum_{i=1}^k \delta_{u_i}, \sum_{j=1}^l \delta_{z_j}) \in \mathcal{Y}$, every Borel subset $A \subset \mathcal{Z}$ and every $t > 2\varrho + 1$,*

$$(1.72) \quad \mathbb{P}_y(Y(t) \in A) \geq C(y) \int_{\alpha + 2\varrho + 1}^{\alpha + t} \mathbb{1}_{\Gamma_A}(\alpha + t - s)(1 - F(\alpha + t - s))dF(s).$$

Proof. At time t , if the state $Y(t)$ is in the set $A \subset \mathcal{Z}$, this means that, by time t , all customers in service at time 0 with residual service times $\{u_i, 1 \leq i \leq k\}$, all customers in queue at time 0 with residual patience times $\{z_j, 1 \leq j \leq l\}$ and those new customers that arrived in the interval $[0, t)$ have completed service (if they entered service before time t) and have run out of their patience (irrespective of whether or not they entered service). Now, we consider a subset of $\{\omega : Y(t, \omega) \in A\}$, in which (a) by time $2\varrho + 1 < t$, all the initial customers with residual patience times $\{z_j, 1 \leq j \leq l\}$ and residual service times $\{u_i, 1 \leq i \leq k\}$ have finished their services (if they entered service) and run out of their patience (irrespective of whether or nor they entered service), (b) the first new customer arrived after $2\varrho + 1$, finished service before t and ran out of his/her patience time before t , (c) the difference between t and the arrival time of that customer lies in Γ_A , and (d) the second new customer arrived after time t . Let \mathcal{Q}_a , \mathcal{Q}_{ad} and \mathcal{Q}_{bd} , respectively, be

the events that property (a) holds, properties (a)–(d) hold and properties (b)–(d) hold. Then, for $y \in \mathcal{Y}$,

$$\mathbb{P}_y(Y(t) \in A) \geq \mathbb{P}_y(\mathcal{Q}_{ad}) = \mathbb{P}_y(\mathcal{Q}_a)\mathbb{P}_y(\mathcal{Q}_{bd}|\mathcal{Q}_a),$$

and, due to the independence assumptions on the service, patience and interarrival distributions, $\mathbb{P}_y(\mathcal{Q}_{bd}|\mathcal{Q}_a)$ is greater than or equal to

$$\begin{aligned} & \int_{\alpha+2\varrho+1}^{\alpha+t} G^r(\alpha+t-s)G^s(\alpha+t-s)\mathbb{1}_{\Gamma_A}(\alpha+t-s)(1-F(\alpha+t-s))\frac{dF(s)}{1-F(\alpha)} \\ & \geq \frac{G^r(\varrho+1)G^s(\varrho+1)}{1-F(\alpha)} \int_{\alpha+2\varrho+1}^{\alpha+t} \mathbb{1}_{\Gamma_A}(\alpha+t-s)(1-F(\alpha+t-s))dF(s), \end{aligned}$$

where the last inequality holds because $\alpha+t-s \geq \varrho+1$ when $\alpha+t-s \in \Gamma_A$. Let $C(y) \doteq \frac{\mathbb{P}_y(\mathcal{Q}_a)G^r(\varrho+1)G^s(\varrho+1)}{1-F(\alpha)}$. Since, due to Assumption 6(2), $G^r(A) > 0$ and $G^s(A) > 0$ for any interval A with length bigger than ϱ , $\mathbb{P}_y(\mathcal{Q}_a)$, as a function of $y \in \mathcal{Y}$, is positive and continuous. Thus C is a positive and continuous function on \mathcal{Y} , and the lemma is proved. \square

Definition A.2. Any Markov process $\{X_t\}$ with topological state space \mathcal{X} is said to be ψ -irreducible if and only if there exists a σ -finite measure ψ on $\mathcal{B}(\mathcal{X})$, the Borel σ -algebra on \mathcal{X} such that for every $x \in \mathcal{X}$ and $B \in \mathcal{B}(\mathcal{X})$,

$$\int_0^\infty \mathbb{P}_x(X(t) \in B)dt > 0 \quad \text{if } \psi(B) > 0.$$

Let $\psi = m \times \delta_0 \times \delta_0 \times \delta_0$, where $m(A) = \bar{m}(A \cap [\varrho+1, \infty))$, where \bar{m} is Lebesgue measure. Clearly, ψ is a σ -finite measure on $(\mathcal{Y}, \mathcal{B}(\mathcal{Y}))$.

Lemma A.3. The Markov process $\{Y_t, \mathcal{F}_t\}$ is ψ -irreducible and the discrete-time Markov chain $\{Y(n)\}_{n \in \mathbb{N}}$ is ψ -irreducible.

Proof. Let $B \in \mathcal{B}(\mathcal{Y})$ be such that $\psi(B) > 0$. Then $\psi(B \cap \mathcal{Z}) > 0$ by the definition of ψ . There is a Borel measurable set $\Gamma_{B \cap \mathcal{Z}} \subset [\varrho+1, \infty)$ such that $B \cap \mathcal{Z} = \{(\alpha, 0, \mathbf{0}, \mathbf{0}) \in \mathcal{Y} : \alpha \in \Gamma_{B \cap \mathcal{Z}}\}$ and $m(\Gamma_{B \cap \mathcal{Z}}) > 0$. Fix $y \in \mathcal{Y}$. It follows from Lemma A.1 that there exists a strictly positive function C on \mathcal{Y} such that

$$\begin{aligned} & \int_0^\infty \mathbb{P}_y(Y(t) \in B \cap \mathcal{Z}) dt \\ & \geq \int_{2\varrho+1}^\infty \mathbb{P}_y(Y(t) \in B \cap \mathcal{Z}) dt \\ & \geq \int_{2\varrho+1}^\infty C(y) \left(\int_{\alpha+2\varrho+1}^{\alpha+t} \mathbb{1}_{\Gamma_{B \cap \mathcal{Z}}}(\alpha+t-s)(1-F(\alpha+t-s)) dF(s) \right) dt \\ & = C(y)(1-F(\alpha+2\varrho+1)) \int_{\Gamma_{B \cap \mathcal{Z}}} (1-F(t)) dt \\ & > 0, \end{aligned}$$

where the equality follows from Fubini's theorem and the last inequality holds because $C(y) > 0$, $m(\Gamma_{B \cap \mathcal{Z}}) > 0$ and $1-F(x) > 0$ for every $x \in [0, \infty)$ by Assumption 6(3). This establishes the first assertion. On the other hand, for $n > 2\varrho+1$,

$$\mathbb{P}_y(Y(n) \in B) \geq C(y) \int_{\alpha+2\varrho+1}^{\alpha+n} \mathbb{1}_{\Gamma_{B \cap \mathcal{Z}}}(\alpha+n-s)(1-F(\alpha+n-s)) dF(s).$$

By Assumption 6(3) and the fact that $m(\Gamma_{B \cap \mathcal{Z}}) > 0$, it follows that $\mathbb{P}_y(Y(n) \in B) > 0$ for all n sufficiently large. Hence $\{Y(n)\}_{n \in \mathbb{N}}$ is ψ -irreducible \square

For each $y \in \mathcal{Y}$, $B \in \mathcal{B}(\mathcal{Y})$ and each probability measure Π on $[0, \infty)$, let

$$\mathcal{K}_\Pi(y, B) = \int_0^\infty \mathbb{P}_y(Y(t) \in B) \Pi(dt).$$

Lemma A.4. *There exists a probability measure Π on $[0, \infty)$ and a function $T : \mathcal{Y} \times \mathcal{B}(\mathcal{Y}) \rightarrow \mathbb{R}_+$ such that*

- (1) $\mathcal{K}_\Pi(y, B) \geq T(y, B)$ for all $y \in \mathcal{Y}$ and every Borel measurable set $B \in \mathcal{B}(\mathcal{Y})$;
- (2) $T(y, \mathcal{Y}) > 0$ for all $y \in \mathcal{Y}$;
- (3) $T(\cdot, B)$ is lower-semicontinuous for every $B \in \mathcal{B}(\mathcal{Y})$.

Proof. Let C be the strictly positive, continuous function C of Lemma A.1. Let Π be a probability measure with density function $e^{-(t-2\varrho-1)}$ on $[2\varrho+1, \infty)$. For each $y \in \mathcal{Y}$ and $B \subset \mathcal{Z}$, define

$$T(y, B) \doteq C(y)e^{\alpha+2\varrho+1} \int_{\alpha+2\varrho+1}^\infty e^{-s} dF(s) \int_0^\infty (1-F(t)) \mathbb{1}_{\Gamma_B}(t) e^{-t} dt,$$

and $T(y, \mathcal{Y} \setminus \mathcal{Z}) = 0$. It is easy to see that for any Borel measurable set $B \in \mathcal{B}(\mathcal{Y})$, $T(y, B) = T(y, B \cap \mathcal{Z})$ and $T(\cdot, B)$ is continuous. Moreover, $T(y, \mathcal{Y}) = T(y, \mathcal{Z}) > 0$. Now, fix $y \in \mathcal{Y}$ and $B \in \mathcal{B}(\mathcal{Y})$. By Lemma A.1, we have

$$\begin{aligned} \mathcal{K}_\Pi(y, B) &= \int_0^\infty \mathbb{P}_y(Y(t) \in B) e^{-(t-2\varrho-1)} dt \\ &\geq \int_{2\varrho+1}^\infty \mathbb{P}_y(Y(t) \in B \cap \mathcal{Z}) e^{-(t-2\varrho-1)} dt \\ &\geq \int_{2\varrho+1}^\infty C(y) \int_{\alpha+2\varrho+1}^{\alpha+t} \mathbb{1}_{\Gamma_{B \cap \mathcal{Z}}}(\alpha+t-s) (1-F(\alpha+t-s)) dF(s) e^{-(t-2\varrho-1)} dt \\ &= C(y)e^{\alpha+2\varrho+1} \int_{\alpha+2\varrho+1}^\infty e^{-s} dF(s) \int_0^\infty (1-F(t)) \mathbb{1}_{\Gamma_{B \cap \mathcal{Z}}}(t) e^{-t} dt \\ &= T(y, B \cap \mathcal{Z}) = T(y, B). \end{aligned}$$

Thus we have proved the lemma. \square

Theorem A.5. *The Markov process Y is positive Harris recurrent.*

Proof. Lemma A.4 shows that Y is a so-called T process (cf. Section 3.2 of [19]), and Lemma A.3 shows that Y is ψ -irreducible. Now, Theorem 3.2 of [19] states that any ψ -irreducible T process Y is positive Harris recurrent if Y is bounded in probability on average, that is, for each $y \in \mathcal{Y}$ and $\varepsilon > 0$, there exists a compact set $B \in \mathcal{B}(\mathcal{Y})$ such that

$$\liminf_{t \rightarrow \infty} \frac{1}{t} \int_0^t \mathbb{P}_y(Y(s) \in B) ds \geq 1 - \varepsilon.$$

However, this is satisfied by the state process Y due to Lemma 4.8. So we have the desired result. \square

REFERENCES

- [1] ASMUSSEN, S. (2003). *Applied probability and queues*, 2nd ed. Springer-Verlag, New York.
- [2] ASMUSSEN, S. AND FOSS, S. (1993). Renovation, regeneration and coupling in multiple-server queues in continuous time. *Frontiers in Pure and Applied Probability*, H. Niemi, G. Högnas, A.N. Shiryayev and A.V. Melnikov, eds., pages 1–6.
Applied probability and queues, 2nd ed. Springer-Verlag, New York.
- [3] BACCELLI, F. and HEBUTERNE, G. (1981). On queues with impatient customers. *Performance '81*, ed. E. Gelenbe (North-Holland Publ. Cy., Amsterdam), pp 159–179.
- [4] BOCHAROV, P. P., D'APICE, C., PECHINKIN, A. V. and SALERNO, S. (2004). *Queueing theory*, Walter de Gruyter.
- [5] BROWN, L., GANS, N., MANDELBAUM, A., SAKOV, A., SHEN, H., ZELTYN, S. and ZHAO, L. (2005). Statistical analysis of a telephone call center: a queueing science perspective, *JASA* **100** No.469, 36–50.
- [6] DA PRATO, G. and ZABCZYK, J. (1996). *Ergodicity for Infinite Dimensional Systems*, Cambridge University Press.
- [7] FOSS, S. G. (1983). Conditions for ergodicity in queues with many servers and waiting, *Sib. Math. J.* **24**, 961-968.
- [8] GAMARNIK, D. and MOMČILOVIĆ, P. (2008). Steady-state analysis of a multiserver queue in the Halfin-Whitt regime, *Advances in Applied Probability* **40** No. 2, 548–577.
- [9] GAMARNIK, D. and ZEEVI, A. (2006). Validity of heavy traffic steady-state approximations in generalized Jackson networks, *Annals of Applied Probability* **16** 56–90.
- [10] GARNETT, O., MANDELBAUM, A. and REIMAN, M. I. (2002). Designing a call center with impatient customers. *Manufac. Serv. Oper. Management* **4** No. 3, 208–227.
- [11] HARRISON, J. M. and WILLIAMS, R. J. (1996). A multiclass closed queueing network with unconventional heavy traffic behavior, *Annals of Applied Probability* **6** 1-47.
- [12] JACOBSEN, M. (2006). *Point process theory and applications: Markov point and piecewise deterministic processes*, Birkhäuser Boston.
- [13] JELENKOVIĆ, P., MANDELBAUM, A. and MOMČILOVIĆ, P. (2004). Heavy traffic limits for queues with many deterministic servers, *QUESTA*, **47** No. 1-2, 53–69.
- [14] KANG, W. N. and RAMANAN, K. Fluid limits of many-server queues with renegeing, *Annals of Applied Probability*, to appear.
- [15] KASPI, H. and RAMANAN, K. Law of large numbers limits for many-server queues, *Annals of Applied Probability*, to appear.
- [16] KIEFER, J. and WOLFOWITZ, J. (1955). On the theory of queues with many servers, *Trans. Amer. Math. Soc.* **78**, 1-18.
- [17] LITTLE, J. D. C. (1961). A Proof of the Queueing Formula $L = \lambda W$, *Oper. Res.*, **9**, 383-387.
- [18] MANDELBAUM, A. and ZELTYN, S. (2005). Call centers with impatient customers: many-server asymptotics of the $M/M/N + G$ queue. *QUESTA* **51**, 361–402.
- [19] MEYN, S. P. and TWEEDIE, R. L. (1993). Stability of Markovian processes II: continuous-time processes and sampled chains, *Advances in Applied Probability*, **25** No.3, 487-517.
- [20] PARTHASARATHY, K. R. (1967). *Probability measures on metric spaces*, Academic Press.
- [21] WHITT, W. (2006). Fluid models for multiserver queues with abandonments, *Operations Research* **54** No.1, 37–54.

DEPARTMENT OF MATHEMATICS & STATISTICS, UNIVERSITY OF MARYLAND, BALTIMORE COUNTY,
1000 HILLTOP CIRCLE, BALTIMORE, MD 21250

E-mail address: wkang@umbc.edu

DIVISION OF APPLIED MATHEMATICS, BROWN UNIVERSITY, PROVIDENCE, RI 02912

E-mail address: kavita@dam.brown.edu