

## LARGE-SCALE JOIN-IDLE-QUEUE SYSTEM WITH GENERAL SERVICE TIMES

S. FOSS,\* *Heriot-Watt University and Novosibirsk State University*

A. L. STOLYAR,\*\* *University of Illinois at Urbana-Champaign*

### Abstract

A parallel server system with  $n$  identical servers is considered. The service time distribution has a finite mean  $1/\mu$ , but otherwise is arbitrary. Arriving customers are routed to one of the servers immediately upon arrival. The join-idle-queue routing algorithm is studied, under which an arriving customer is sent to an idle server, if such is available, and to a randomly uniformly chosen server, otherwise. We consider the asymptotic regime where  $n \rightarrow \infty$  and the customer input flow rate is  $\lambda n$ . Under the condition  $\lambda/\mu < \frac{1}{2}$ , we prove that, as  $n \rightarrow \infty$ , the sequence of (appropriately scaled) stationary distributions concentrates at the natural equilibrium point, with the fraction of occupied servers being constant at  $\lambda/\mu$ . In particular, this implies that the steady-state probability of an arriving customer waiting for service vanishes.

*Keywords:* Large-scale service system; pull-based load distribution; join-idle-queue; load balancing; fluid limit; stationary distribution; asymptotic optimality

2010 Mathematics Subject Classification: Primary 90B15  
Secondary 60K25

### 1. Introduction

We consider a parallel server system consisting of  $n$  servers, processing a single input flow of customers. The service time of any customer by any server has the same distribution with finite mean  $1/\mu$ . Each customer has to be assigned (routed) to one of the servers immediately upon arrival. (This model is sometimes referred to as the ‘supermarket’ model.) We study a join-idle-queue (JIQ) routing algorithm, under which an arriving customer is sent to an idle server, if such is available; if there are no idle servers, a customer is sent to one of the servers chosen uniformly at random.

We consider an asymptotic regime such that  $n \rightarrow \infty$  and the input rate is  $\lambda n$ , where the system load  $\lambda/\mu < 1$ . Thus, the system remains subcritically loaded. Under the additional assumption that the service time distribution has *decreasing hazard rate* (DHR), it was shown in [9] that the following property holds.

*Asymptotic optimality.* As  $n \rightarrow \infty$ , the sequence of the system stationary distributions is such that the fraction of occupied servers converges to constant  $\lambda/\mu$ ; consequently, the steady-state probability of an arriving customer being routed to a nonidle server vanishes.

The results of [9] apply to far more general systems, where servers may be nonidentical. However, the analysis in [9] does rely in an essential way on the DHR assumption on the

---

Received 20 May 2016; revision received 14 February 2017.

\* Postal address: Department of Actuarial Mathematics and Statistics, Heriot-Watt University, Edinburgh EH14 4AS, UK. Email address: s.foss@hw.ac.uk

\*\* Postal address: Department of Industrial and Enterprise Systems Engineering, University of Illinois at Urbana-Champaign, 104 S. Mathews Avenue, Office 201C, Urbana, IL 61801, USA. Email address: stolyar@illinois.edu

service times; under this assumption, the system process has the *monotonicity* property, which is a powerful tool for analysis. Informally speaking, monotonicity means that two versions of the process, such that the initial state of the first one is dominated (in the sense of some natural partial order) by that of the second one, can be coupled so that this dominance persists at all times.

When the service time distribution is general, the monotonicity under JIQ no longer holds, which requires a different approach to the analysis. In this paper we prove the following result.

*Main result.* (Theorem 2.1 in Section 2). The asymptotic optimality holds for an arbitrary service time distribution, if the system load  $\lambda/\mu < \frac{1}{2}$ .

We believe that condition  $\lambda/\mu < \frac{1}{2}$  is purely technical (required for the proof in this paper) and that our main result, in fact, holds for  $\lambda/\mu < 1$ , i.e. as long as the system is stable. This will be discussed in more detail in Section 2.1.

The key feature of the JIQ algorithm (as well as more general *pull-based* algorithms [1], [6], [9], [11]) is that it does not utilize any information about the current state of the servers besides them being idle or not. This allows for a very efficient practical implementation, requiring a very small communication overhead between the servers and the router(s) [6], [9], [11]. In fact, in the asymptotic regime that we consider, JIQ is much superior to the celebrated ‘power-of- $d$ -choices’ (or join-shortest-queue( $d$ ), or JSQ( $d$ )) algorithm [3], [4], [7], [12], in terms of both performance and communication overhead (see [9] and [11] for a detailed comparison). The JSQ( $d$ ) algorithm routes a customer to the shortest queues among the  $d \geq 1$  servers picked uniformly at random.

We note that when the service time distribution is general, there is no monotonicity under JSQ( $d$ ) (just like under JIQ in our case), and this also makes the analysis far more difficult. Specifically, the result for JSQ( $d$ ), which is a counterpart of our main result for JIQ, is Theorem 2.3 of [3], which shows the asymptotic independence of individual server states. (Our main result also implies the asymptotic independence of server states; see the formal statement in Corollary 2.1.) Theorem 2.3 of [3] imposes even stronger assumptions than ours, namely a finite second moment of the service time and load  $\lambda/\mu < \frac{1}{4}$  (for nontrivial values of  $d$ , which are  $d \geq 2$ ); our Theorem 2.1 only requires a finite first moment of the service time and load  $\lambda/\mu < \frac{1}{2}$ .

In a different asymptotic regime, the so-called Halfin–Whitt regime (when the system capacity exceeds its load by  $O(\sqrt{n})$ , as opposed to  $O(n)$ ) and with Markov assumptions (Poisson input flows and exponentially distributed service times), JIQ has been recently analyzed in [5] and [8]. In these papers the diffusion limits of the system transient behavior were studied; Markov assumptions appear to be essential for the analysis. Finally, we mention a recent paper [10], in which the author proposed and studied a version of JIQ for systems with *packing constraints* at the servers.

*Paper organization.* In Section 2 we state the formal model and main result, with Section 2.1 devoted to a discussion of the role of condition  $\lambda/\mu < \frac{1}{2}$ . A uniform stochastic upper bound on the individual server workload in steady-state is derived in Section 3. Properties of the process fluid limits are established in Section 4. Section 5 contains the proof of the main result, which relies on the above upper bound and fluid limit properties. Generalizations of the main result are presented in Section 6.

*Basic notation.* The following abbreviations are used to qualify a convergence of functions: *u.o.c.* means *uniform on compact sets*, *p.o.c.* means *convergence at points of continuity of the limit*, and *a.e.* means *almost everywhere with respect to the Lebesgue measure*. We say that a function is RCLL if it is *right-continuous with left-limits*. A scalar function  $f(t)$ ,  $t \geq 0$ ,

we will call *Lipschitz above* if there exist a constant  $L > 0$  such that  $f(t_2) - f(t_1) \leq L(t_2 - t_1)$  for any  $t_1 \leq t_2$ . The norm of a scalar function is  $\|f(\cdot)\| \doteq \sup_w |f(w)|$ . Inequalities applied to vectors (respectively, functions) are understood componentwise (respectively, for every value of the argument). We denote by ' $\xrightarrow{D}$ ' the convergence of random elements in distribution. The indicator of event or condition  $B$  is denoted by  $\mathbf{1}_B$ . Abbreviation w.l.o.g. means *without loss of generality*.

## 2. Model and main result

We consider a service system consisting of  $n$  parallel servers. The system is homogeneous in that all servers are identical, with the same customer service time distribution, given by the cumulative distribution function (CDF)  $F(\xi)$ ,  $\xi \geq 0$ . This distribution has finite mean, which w.l.o.g. can be assumed to be 1:

$$\int_0^\infty F^c(\xi) d\xi = 1,$$

where  $F^c(\xi) \doteq 1 - F(\xi)$ . Otherwise, the CDF  $F(\cdot)$  is arbitrary. The service/queueing discipline at each server is arbitrary, as long as it is work-conserving and nonidling.

Customers arrive as a Poisson process. (This assumption can be relaxed to a renewal arrival process; see Section 6.) The arrival rate is  $\lambda n$ , where  $\lambda < 1$ , so that the system load is strictly subcritical.

The routing algorithm is JIQ, which is defined as follows. (The JIQ algorithm can be viewed, in particular, as a specialization of the PULL algorithm [9], [11] to a homogeneous system with 'single router'.)

**Definition 2.1.** (*JIQ algorithm.*) An arriving customer is routed to an idle server, if there is one available. Otherwise, it is routed to server chosen uniformly at random.

We consider the sequence of systems with  $n \rightarrow \infty$ . From now on, the upper index  $n$  of a variable/quantity will indicate that it pertains to the system with  $n$  servers, or  $n$ th system. Let  $W_i^n(t)$  denote the workload, i.e. unfinished work, in queue  $i$  at time  $t$  in the  $n$ th system. Consider the following *fluid-scaled* quantities:

$$x_w^n(t) \doteq \left(\frac{1}{n}\right) \sum_i \mathbf{1}_{\{W_i^n(t) > w\}}, \quad w \geq 0. \quad (2.1)$$

That is,  $x_w^n(t)$  is the fraction of servers  $i$  with  $W_i^n(t) > w$ . Then  $x^n(t) = (x_w^n(t), w \geq 0)$  is the system state at time  $t$ ;  $\rho^n(t) \equiv x_0^n(t)$  is the fraction of busy servers (the instantaneous system load).

For any  $n$ , the state space of the process  $(x^n(t), t \geq 0)$  is a subset of a common (for all  $n$ ) state space  $\mathcal{X}$ , whose elements  $x = (x_w, w \geq 0)$  are nonincreasing RCLL functions of  $w$ , with values  $x_w \in [0, 1]$ . This state space  $\mathcal{X}$  is equipped with Skorokhod metric, topology, and corresponding Borel  $\sigma$ -algebra.

Then, for any  $n$ , process  $(x^n(t), t \geq 0)$  is Markov with state space  $\mathcal{X}$ , and sample paths being RCLL functions (with values in  $\mathcal{X}$ ), which are, in turn, elements of (another) Skorokhod space. (The Skorokhod spaces that we defined play no essential role in our analysis; we need to specify them merely to make the process well defined.)

Stability (positive Harris recurrence) of the process  $(x^n(t), t \geq 0)$ , for any  $n$ , is straightforward to verify. Indeed, as long as a server remains busy, it receives each new arrival with probability at most  $1/n$ , and, therefore, receives the new work at the average rate at most  $\lambda$ .

(We omit the details of the stability proof.) Thus, the process has unique stationary distribution. Let  $x^n(\infty)$  be a random element whose distribution is the stationary distribution of the process; in other words, this is a random system state in stationary regime.

The system *equilibrium point*  $x^* \in \mathcal{X}$  is defined as follows. Let  $\Phi^c(w)$  denote the complementary (or tail) distribution function of the steady-state residual service time; the latter is the steady-state residual time of a renewal process with renewal time distribution function  $F(\cdot)$ . We have

$$\Phi^c(w) = \int_w^\infty F^c(\xi) \, d\xi, \quad w \geq 0.$$

Then

$$x^* = (x_w^* = \lambda \Phi^c(w), w \geq 0) \in \mathcal{X}.$$

In particular, the equilibrium point is such that ‘the fraction of occupied servers’  $x_0^* = \lambda$ . Our main result is the following.

**Theorem 2.1.** *If  $\lambda < \frac{1}{2}$  then  $x^n(\infty) \xrightarrow{D} x^*$  as  $n \rightarrow \infty$ .*

The theorem shows, in particular, that if  $\lambda < \frac{1}{2}$  then, as  $n \rightarrow \infty$ , the steady-state probability of an arriving customer waiting for service (or sharing a server with other customers) vanishes. Theorem 2.1 easily generalizes to the case when:

- arrival process is renewal,
- some or all servers may have finite buffers, and
- there may be some bias in the routing when all servers are busy.

(These generalizations are described in Section 6.)

Theorem 2.1 implies the following corollary.

**Corollary 2.1.** *Assume that  $\lambda < \frac{1}{2}$ . Suppose that *JIQ* is completely symmetric with respect to the servers. Specifically, if at the time of a customer arrival there are idle servers, the customer is routed to one of them chosen uniformly at random. Then the states of individual servers in the stationary regime are asymptotically independent. Moreover, for any fixed  $m$ , the stationary distribution of  $(W_1^n, \dots, W_m^n)$  converges to that of  $(\tilde{W}_1, \dots, \tilde{W}_m)$ , with independent and identically distributed (i.i.d.) components such that  $\mathbb{P}\{\tilde{W}_1 > w\} = x_w^* = \lambda \Phi^c(w)$ ,  $w \geq 0$ .*

Indeed, by symmetry with respect to the servers, the stationary distribution of  $(W_1^n, \dots, W_m^n)$ , i.e. of the residual work on the fixed set of servers  $1, \dots, m$ , is the same as that on a set of  $m$  servers, *chosen uniformly at random*. But,  $x^n(\infty)$ , which describes the overall distribution of server workloads in the system, converges in distribution to the nonrandom point  $x^*$ . This implies Corollary 2.1.

**2.1. Discussion of condition  $\lambda < \frac{1}{2}$**

The approach we use to establish the convergence of stationary distributions in Theorem 2.1 is as follows. We find a set  $A \in \mathcal{X}$  and a fixed finite time  $T$ , such that, with high probability, for all large  $n$ ,

- (a)  $x^n(\infty) \in A$  and
- (b)  $x^n(0) \in A$  implies that  $x^n(T)$  is close to  $x^*$ .

Property (b) is key. When  $n$  is large, the trajectory  $x^n(t)$  is ‘almost deterministic’. (In fact, the problem reduces to the analysis of ‘fluid limit’ trajectories, which are the limits of  $x^n(t)$  as  $n \rightarrow \infty$ .) Then, informally speaking, property (b) reduces to

(b') trajectories  $x^n(t)$  converge to  $x^*$  as  $t \rightarrow \infty$ .

The absence of process monotonicity (described in Section 1) makes proving (b') difficult. We now describe—very informally—the key idea, which we use in our proof of convergence (b'), and which relies on the condition  $\lambda < \frac{1}{2}$ .

Suppose  $n$  is large. Consider an initial state  $x^n(0)$ , such that *the total amount of (fluid-scaled, i.e. multiplied by  $1/n$ ) unfinished work is upper bounded by  $C < \infty$* . Pick  $\alpha$  such that  $\alpha > \lambda$  and  $\alpha + \lambda < 1$ ; this is possible if and only if  $\lambda < \frac{1}{2}$ . Then, at some finite time  $\tau$ , the system must reach a state with  $\alpha n$  servers being idle. (Otherwise, if at least  $(1 - \alpha)n$  servers would continue to be busy as time goes to  $\infty$ , the unfinished work would become negative, since  $1 - \alpha > \lambda$ .) Denote by  $S_\alpha$  the set of those  $\alpha n$  servers, which are idle at time  $\tau$ . Starting time  $\tau$ , w.l.o.g., assume that all new arriving customers go to an idle server in  $S_\alpha$ , as long as there is one available. Consider the subsystem, consisting only of the servers in  $S_\alpha$ ; starting time  $\tau$  and until the (random) time when *all* servers in  $S_\alpha$  become busy, the behavior of this subsystem is obviously equivalent to that of the infinite-server system,  $M/GI/\infty$ , with idle initial state. If  $n$  is large, the behavior of  $x^n(t)$  for such an  $M/GI/\infty$  system is ‘almost deterministic’ and such that the (scaled) number of occupied servers  $x_0^n(t)$  in it is ‘almost monotone increasing, converging to  $\lambda < \alpha$ ’ and, moreover,  $x^n(t)$  ‘converges’ to  $x^*$ . But this means that after time  $\tau$  the subsystem  $S_\alpha$  will ‘always’ have idle servers, which, in turn, means that *its* state will ‘converge’ to  $x^*$  as  $t \rightarrow \infty$ . Also, after time  $\tau$ , the subsystem consisting of the servers outside  $S_\alpha$  will ‘never’ receive any new arrivals and will ‘eventually’ empty. Thus,  $x^n(t)$  for our entire system ‘converges’ to  $x^*$ .

Turning the key intuition, described above informally, into a formal proof is the subject of the rest of this paper. Set  $A \in \mathcal{X}$  is picked by using a constructed uniform in  $n$  upper bound on the stationary distribution of the workload of an individual server. The states in  $A$  are such that the total (scaled) workload is not necessarily upper bounded by a constant  $C$  (in fact, if the second moment of the service time is infinite, the steady-state total workload in the system is infinite with probability 1); however, for states in  $A$  the (scaled) workload is bounded by  $C$  on a close-to-1 fraction of servers—this suffices for the proofs. Property (b') is proved uniformly for fluid limits starting from  $A$ —from here we find that (b) holds for the pre-limit processes with high probability, uniformly for all large  $n$ .

As explained above, our proof of Theorem 2.1 relies in an essential way on condition  $\lambda < \frac{1}{2}$ . However, we believe that this condition is purely technical, and Theorem 2.1, in fact, holds for any  $\lambda < 1$ . Establishing this fact will most likely require a different proof approach, although some elements of the analysis in this paper may turn out to be useful for the proof of a more general result.

### 3. Uniform upper bound on a server workload distribution

Throughout this section, consider a fixed  $\lambda < 1$ . Consider an  $M/GI/1$  system, with arrival rate  $\lambda$  and service time distribution  $F(\cdot)$ . Let us view its workload process as regenerative with renewal points being time instants when a customer arrives into an idle system. For each  $w \geq 0$ , denote by  $x_w^{**}$  the expectation of the total time during one renewal cycle when the workload is greater than  $w$ . Clearly,  $x_w^{**}$  is nonincreasing,  $x_0^{**} = 1/(1 - \lambda)$  (the expected busy

period duration), and  $x_w^{**} \rightarrow 0, w \rightarrow \infty$ . (We will not use the exact value of  $x_0^{**}$ . Also,  $x_w^{**}$  is continuous in  $w$ , but we will not use this fact either.)

Now consider our system with any fixed  $n$ . Consider a specific server  $i$ . Consider our Markov process sampled at the ‘renewal’ instants when there is an arrival into idle server  $i$ . Time intervals between the ‘renewal’ instants are ‘renewal cycles’. Of course, such ‘renewal cycles’ are not i.i.d., the law of a cycle depends on the state of the entire system at the renewal point from which the cycle starts. However, there are uniform bounds that apply to any cycle. For a fixed  $w \geq 0$ , the expected total time within one cycle when  $W_i^n > w$  is upper bounded by  $x_w^{**}$ ; indeed, as long as the server remains busy, the probability that a new arrival will be routed to it is at most  $1/n$  (either  $1/n$  or 0). Therefore, as long as the server remains busy, the instantaneous arrival rate into it is upper bounded by  $(\lambda n)/n = \lambda$ . The mean duration of each cycle is lower bounded by the mean service time of one customer, i.e. by 1. Therefore,

$$\mathbb{P}\{W_i^n(\infty) > w\} \leq x_w^{**}, \quad w \geq 0, \tag{3.1}$$

where, recall,  $x_w^{**} \rightarrow 0, w \rightarrow \infty$ . Bound (3.1) implies the following fact.

**Lemma 3.1.** *Let  $\lambda < 1$ . Then, for any  $n$ ,  $\mathbb{E}x_w^n(\infty) \leq x_w^{**}, w \geq 0$ .*

#### 4. Fluid limits

In this section we introduce different types of the process fluid limit, which will be used later in the analysis.

We will assume that, given a fixed initial state  $x^n(0)$ , the realization of the process is determined by a common (for all  $n$ ) set of driving processes. Specifically, there is a common, rate 1, Poisson process,  $\Pi(t), t \geq 0$ ; the number of arrivals in the  $n$ th system by time  $t$  is  $\Pi(n\lambda t)$ . There is also a common sequence of i.i.d. random variables with distribution  $F(\cdot)$ , which determines the service times of arriving customers (in the order of arrivals). Let  $G^n(t, w), t \geq 0, w \geq 0$ , be the number of customer arrivals in the  $n$ th system, by time  $t$ , with the service times greater than  $w$ . Let  $g^n(t, w) = (1/n)G^n(t, w)$  and  $g(t, w) \doteq \lambda t F^c(w)$ . Then, we have the following functional strong law of large numbers (FSLLN):

$$\|g^n(t, \cdot) - g(t, \cdot)\| \rightarrow 0 \quad \text{as } n \rightarrow \infty, \text{ u.o.c. (in } t), \text{ w.p.1,} \tag{4.1}$$

where we abbreviate ‘with probability 1’ to w.p.1. Indeed, for any fixed  $t > 0$ , the total number of arrivals in  $[0, t]$ , scaled by  $1/n$ , converges to  $\lambda t$ , w.p.1; this and Glivenko–Cantelli theorem (see [2, Theorem 20.6, p. 269]) imply that  $\|g^n(t, \cdot) - g(t, \cdot)\| \rightarrow 0$ , w.p.1. But, all  $g^n(t, w)$  and  $g^n(t, w)$  are nondecreasing in  $t$ , and  $g(t, w)$  is continuous in  $t$ ; this easily implies that the convergence in (4.1) is uniform w.p.1.

The routing of arriving customers to idle servers, when such are available, is completely arbitrary w.l.o.g.; it will be specified later, in a way convenient for the analysis. The routing of arriving customers to the servers, in cases when all servers are busy, is determined by a sequence of i.i.d. random variables, uniformly distributed in  $[0, 1)$ ; these random variables are used sequentially as needed; in the  $n$ th system, a customer is routed to server  $i$  if the corresponding random variable value is in  $[(i - 1)/n, i/n)$ . (The specific construction of routing to busy servers will not be important; we need to specify it somehow, to have the process well defined.)

It will be convenient for every  $n$ , in addition to the actual system with  $n$  servers, to consider the corresponding infinite-server system; in such a system all arrivals always go to idle servers.

For a given  $n$ , for the infinite-server system the fluid-scaled quantities  $x_w^n(t)$  are still defined by (2.1), i.e. as the total number of servers with  $W_i^n > w$ , multiplied by  $1/n$ .

For every  $t \geq 0$ , let us define  $x^\uparrow(t) = (x_w^\uparrow(t), w \geq 0) \in \mathcal{X}$ ,

$$x_w^\uparrow(t) = \int_0^t F^c(w + t - \theta)\lambda \, d\theta = \int_w^{w+t} F^c(\xi)\lambda \, d\xi.$$

Clearly,  $x_w^\uparrow(t)$  is nondecreasing in  $t$ , and

$$x_w^\uparrow(t) \uparrow x_w^*, \quad t \rightarrow \infty \text{ for all } w.$$

As functions of  $w$ , all  $x_w^\uparrow(t)$  and  $x_w^*$  are nonnegative, continuous, nondecreasing, and converging to 0 as  $w \rightarrow \infty$ ; therefore, the above pointwise convergence implies uniform convergence

$$\|x^\uparrow(t) - x^*\| \rightarrow 0, \quad t \rightarrow \infty.$$

The following Lemma 4.1 is a standard fact. Informally speaking, it states that  $x^\uparrow(\cdot)$  is the ‘fluid limit’, in  $n \rightarrow \infty$ , of  $x^n(\cdot)$  for the infinite-server system, starting from an idle initial state. We state this fact in a form that is convenient for our analysis, and since it easily follows from (4.1), we provide a proof as well.

**Lemma 4.1.** *Fix arbitrary  $\lambda \geq 0$ . (Here  $\lambda \geq 1$  is allowed.) Let  $x^n(\cdot)$  be the process describing the infinite-server system, starting from an idle initial state, that is,  $x_0^n(0) = 0$ . Then, w.p.1,*

$$\|x^n(t) - x^\uparrow(t)\| \rightarrow 0, \quad \text{u.o.c.}$$

*Proof.* Fix  $t$  and  $w$ . By definition,  $x_w^n(t)$  is the scaled number of customers in the system, having the residual service time greater than  $w$ . A customer arriving at time  $\theta \in [0, t]$  counts into that number if and only if its service time is greater than  $t + w - \theta$ . Let points  $0 = t_0 < t_1 < \dots < t_\kappa = t$  partition the interval  $[0, t]$  into  $\kappa$  subintervals  $[t_k, t_{k+1})$ . (Note that, w.p.1 there are no arrivals at  $t$ .) Then

$$\begin{aligned} & \sum_{k=0}^{\kappa-1} [g^n(t_{k+1}, t + w - t_k) - g^n(t_k, t + w - t_k)] \\ & \leq x_w^n(t) \\ & \leq \sum_{k=0}^{\kappa-1} [g^n(t_{k+1}, t + w - t_{k+1}) - g^n(t_k, t + w - t_{k+1})]. \end{aligned}$$

By (4.1), w.p.1, the lower and upper bounds converge to

$$\sum_{k=0}^{\kappa-1} \lambda[t_{k+1} - t_k]F^c(t + w - t_k) \quad \text{and} \quad \sum_{k=0}^{\kappa-1} \lambda[t_{k+1} - t_k]F^c(t + w - t_{k+1}),$$

respectively. Considering a sequence of partitions with maximum subinterval size vanishing, and taking into account that  $F^c$  is nonincreasing, we obtain probability 1 convergence  $x_w^n(t) \rightarrow x_w^\uparrow(t)$ . Since  $x_w^\uparrow(t)$  and all  $x_w^n(t)$  are nonnegative, nonincreasing in  $w$ ,  $x_w^\uparrow(t)$  is continuous in  $w$ , and  $x_w^\uparrow(t) \rightarrow 0$  as  $w \rightarrow \infty$ , we obtain probability 1 convergence  $\|x^n(t) - x^\uparrow(t)\| \rightarrow 0$  for any  $t$ ; since  $x_w^\uparrow(t)$  is continuous, nondecreasing in  $t$ , this convergence is u.o.c. in  $t$ .  $\square$



Sometimes, it will be convenient to divide the set of servers into two or more subsets, and keep track of the workloads in those subsets separately. For example, suppose at time 0 the set of all servers, let us call it  $S$ , is divided (for each  $n$ ) at time 0 into two nonintersecting subsets,  $S_1$  and  $S_2$ , and these subsets do not change with time. Then, for  $\ell = 1, 2$ ,  ${}^{(\ell)}x_w^n(t)$  is the fraction of servers (out of the total number  $n$ ) which are in  $S_\ell$  and have workload  $W_i^n > w$ ,  $w \geq 0$ ;  ${}^{(\ell)}\rho^n(t) = {}^{(\ell)}x_0^n(t)$ . Of course,  $x(t) = {}^{(1)}x(t) + {}^{(2)}x(t)$ . However, we will often consider  ${}^{(\ell)}x(t)$  for only one of the subsets  $S_\ell$ .

The following fact is a corollary of Lemma 4.1.

**Lemma 4.2.** *Let  $0 \leq \lambda < 1$  and let  $\lambda < \alpha < 1$ . Consider the finite-server system. Assume that, for all  $n$ , the initial states are such that  $x_0^n(0) = \rho^n(0) \leq 1 - \alpha$ . For each  $n$ , consider the subset  $S_1 = S_1(n)$ , consisting of  $\alpha n$  servers that are initially idle. Assume w.l.o.g. that any new arrival will go to an idle server in  $S_1$ , if there is one available. Then, w.p.1, the following holds:*

$$\|{}^{(1)}x^n(t) - x^\uparrow(t)\| \rightarrow 0, \quad \text{u.o.c.}, \tag{4.2}$$

and, for any fixed  $t$ , for all sufficiently large  $n$ , all new arrivals in  $[0, t]$  will go to idle servers in  $S_1$ .

*Proof.* The behavior of the system restricted to subset  $S_1$  of servers is equivalent to that of the infinite-server system starting from idle state, as long as there are idle servers in  $S_1$ . By Lemma 4.1, w.p.1, the trajectory of the (scaled) infinite-server system converges (u.o.c.) to the trajectory  $x^\uparrow(t)$ , such that the (scaled) number of occupied servers increases and converges to  $\lambda < \alpha$ . This implies that w.p.1 the following holds for the system restricted to subset  $S_1$ : for any fixed time  $t \geq 0$ , for all sufficiently large  $n$ , subset  $S_1$  will have idle servers in the entire interval  $[0, t]$ , and then the system behavior coincides with that of the infinite-server system. This property implies (4.2), and contains the last statement of the lemma.  $\square$

Let  ${}^{(\ell)}W^n(t)$  denote the total (fluid-scaled) unfinished work at time  $t$  within a given subset  $S_\ell$  of servers:

$${}^{(\ell)}W^n(t) = \int_0^\infty {}^{(\ell)}x_w^n(t) dw.$$

The  $S_\ell = S$  case is allowed.

Denote by  ${}^{(\ell)}W^{a,n}(t)$  and  ${}^{(\ell)}W^{d,n}(t)$  the amount of (fluid-scaled) work that, respectively, arrived into and processed by subset  $S_\ell$  in the interval  $[0, t]$ . Denote by  ${}^{(\ell)}\rho^{a,n}(t)$  the (fluid-scaled) number of arrivals in  $[0, t]$  into  $S_\ell$ , that went into idle servers; such arrivals, and only they, cause  $+1/n$  jumps of  ${}^{(\ell)}\rho^n$ . Analogously, let  ${}^{(\ell)}\rho^{d,n}(t)$  denote the (fluid-scaled) number of times in  $[0, t]$  when a customer service completion occurred in  $S_\ell$ , that left a server idle; such departures, and only they, cause  $-1/n$  jumps of  ${}^{(\ell)}\rho^n$ . Functions  ${}^{(\ell)}W^{a,n}(t)$ ,  ${}^{(\ell)}W^{d,n}(t)$ ,  ${}^{(\ell)}\rho^{a,n}(t)$ , and  ${}^{(\ell)}\rho^{d,n}(t)$  are nondecreasing by definition, equal to 0 at  $t = 0$ . The following relations obviously hold for all  $t \geq 0$ :

$$\begin{aligned} {}^{(\ell)}W^n(t) &= {}^{(\ell)}W^{a,n}(t) - {}^{(\ell)}W^{d,n}(t), & {}^{(\ell)}\rho^n(t) &= {}^{(\ell)}\rho^{a,n}(t) - {}^{(\ell)}\rho^{d,n}(t), \\ {}^{(\ell)}W^{d,n}(t) &= \int_0^t {}^{(\ell)}\rho^n(\xi) d\xi. \end{aligned} \tag{4.3}$$

For future reference, let us also note the obvious fact that if there were no new arrivals into  $S_\ell$  in some time interval  $(t_1, t_2]$ , then

$${}^{(\ell)}W^n(t_2) - {}^{(\ell)}W^n(t_1) = -({}^{(\ell)}W^{d,n}(t_2) - {}^{(\ell)}W^{d,n}(t_1)) = -\int_{t_1}^{t_2} {}^{(\ell)}\rho^n(\xi) d\xi. \tag{4.4}$$



**Lemma 4.3.** *Let  $\lambda \geq 0$ . Consider the finite-server system. For each  $n$  consider a subset  $S_1 = S_1(n)$ , consisting of  $\sigma n$  servers,  $0 \leq \sigma \leq 1$ . (The  $\sigma = 1$  case is when  $S_1 = S$ .) Consider a fixed sequence (in  $n$ ) of initial states, such that  ${}^{(1)}W^n(0) \leq C < \infty$  for all  $n$ . Then, w.p.1, for any subsequence of  $n$ , there exists a further subsequence, along which the following holds:*

$${}^{(1)}W^n(t) \rightarrow {}^{(1)}W(t), \quad \text{u.o.c.,} \tag{4.5}$$

where  ${}^{(1)}W(\cdot)$  is a Lipschitz continuous function with  ${}^{(1)}W(0) \leq C$ ;

$${}^{(1)}\rho^n(t) \rightarrow {}^{(1)}\rho(t), \quad \text{p.o.c.,} \tag{4.6}$$

where  ${}^{(1)}\rho(\cdot)$  is a RCLL function, which is Lipschitz above and  ${}^{(1)}\rho(t) \in [0, \sigma]$  for all  $t$ ;

$${}^{(1)}W'(t) \leq \lambda - {}^{(1)}\rho(t), \quad \text{a.e.} \tag{4.7}$$

*Proof.* Within this proof, when we say that a function is Lipschitz continuous (respectively, Lipschitz above), we always mean that it is Lipschitz continuous (respectively, Lipschitz above) uniformly in  $n$ .

From FSLLN (4.1), we have the following fact. W.p.1, for any subsequence of  $n$ , there exists a further subsequence, along which

$${}^{(1)}\rho^{a,n}(t) \rightarrow {}^{(1)}\rho^a(t), \quad {}^{(1)}W^{a,n}(t) \rightarrow {}^{(1)}W^a(t), \quad \text{u.o.c., } n \rightarrow \infty,$$

where  ${}^{(1)}\rho^a(\cdot)$  and  ${}^{(1)}W^a(\cdot)$  are Lipschitz continuous nondecreasing, with Lipschitz constant equal to  $\lambda$ . Also, clearly, all functions  ${}^{(1)}W^{d,n}(\cdot)$  are nondecreasing Lipschitz continuous, so that we can choose a further subsequence, if necessary, along which

$${}^{(1)}W^{d,n}(t) \rightarrow {}^{(1)}W^d(t), \quad \text{u.o.c., } n \rightarrow \infty,$$

where  ${}^{(1)}W^d(\cdot)$  is Lipschitz continuous nondecreasing. This implies (4.5) with  ${}^{(1)}W(t) = {}^{(1)}W^a(t) - {}^{(1)}W^d(t)$ .

To show (4.6), observe that nondecreasing functions  ${}^{(1)}\rho^{d,n}(t)$  are uniformly bounded on any finite interval (because functions  ${}^{(1)}\rho^{a,n}(t)$  and  ${}^{(1)}\rho^n(t)$  are along the chosen subsequence). Then, we can choose a further subsequence, if necessary, such that

$${}^{(1)}\rho^{d,n}(t) \rightarrow {}^{(1)}\rho^d(t), \quad \text{p.o.c., } n \rightarrow \infty, \tag{4.8}$$

where  ${}^{(1)}\rho^d(\cdot)$  is RCLL nondecreasing. (Here we use a version of Helly's selection theorem; see [2, Theorem 25.9, p. 336].) This proves (4.6) with  ${}^{(1)}\rho(t) = {}^{(1)}\rho^a(t) - {}^{(1)}\rho^d(t)$ .

Note that the p.o.c. convergence in (4.8) implies a.e. (in  $t$ ) convergence. Then, by taking the limit in (4.3), we obtain

$${}^{(1)}W^d(t) = \int_0^t {}^{(1)}\rho(\xi) \, d\xi.$$

This and the fact that  ${}^{(1)}W^a(\cdot)$  is Lipschitz continuous with Lipschitz constant  $\lambda$ , imply (4.7). This completes the proof. □

### 5. Proof of Theorem 2.1

Here we only consider the finite systems (with  $n$  servers in the  $n$ th system). Consider a fixed  $\lambda < \frac{1}{2}$ .

By Lemma 3.1, for each  $n$ , we have  $\mathbb{E}x_w^n(\infty) \leq x_w^{**}$ , where  $x_w^{**}$  is nonincreasing and  $\lim_{w \rightarrow \infty} x_w^{**} = 0$ . Then, for any  $\delta_1 > 0$ , we can choose a sufficiently large  $b$  such that  $\mathbb{E}x_b^n(\infty) \leq \delta_1$ . This, in turn, implies that, for any  $\varepsilon > 0$  and any  $\delta > 0$ , we can pick sufficiently large  $b > 0$  such that

$$\mathbb{P}\{x_b^n(\infty) \leq \delta\} \geq 1 - \varepsilon \quad \text{for all } n. \tag{5.1}$$

For each  $n$ , consider the *stationary version of process*  $x^n(\cdot)$ ; then, for any  $t$ ,  $x^n(t)$  is equal in distribution to  $x^n(\infty)$  (by the definition of the latter). Choose  $\delta > 0$  small enough so that  $\lambda + \delta < \frac{1}{2}$ . For this  $\delta$  and arbitrarily small fixed  $\varepsilon > 0$ , choose  $b > 0$  such that (5.1) holds. Then (5.1) implies that

$$\mathbb{P}\{\text{condition (5.3) holds}\} \geq 1 - \varepsilon \quad \text{for all } n, \tag{5.2}$$

$$x^n(0) \text{ is such that there exists a subset } S_2 = S_2(n) \text{ of } (1 - \delta)n \text{ servers,} \tag{5.3}$$

each with workload at most  $b$ .

Then, to complete the proof of Theorem 2.1, it suffices to prove the following lemma.

**Lemma 5.1.** *For any  $\delta_2 > 0$  there exists  $T > 0$ , which depends on  $\varepsilon, \delta, b$ , such that, uniformly on fixed initial states  $x^n(0)$  satisfying (5.3),*

$$\mathbb{P}\{\|x^n(T) - x^*\| \leq \delta + \delta_2 \mid x^n(0)\} \rightarrow 1, \quad n \rightarrow \infty. \tag{5.4}$$

Indeed, if Lemma 5.1 holds then for  $\delta, \varepsilon, b, \delta_2, T$  chosen as specified above, and arbitrarily small  $\varepsilon_2 > 0$ , for all sufficiently large  $n$ , uniformly on  $x^n(0)$  satisfying (5.3),

$$\mathbb{P}\{\|x^n(T) - x^*\| \leq \delta + \delta_2 \mid x^n(0)\} \geq 1 - \varepsilon_2.$$

This and (5.2) imply that, for all sufficiently large  $n$ ,

$$\mathbb{P}\{\|x^n(T) - x^*\| \leq \delta + \delta_2\} \geq (1 - \varepsilon)(1 - \varepsilon_2).$$

But,  $\delta, \delta_2, \varepsilon, \varepsilon_2$  can be chosen arbitrarily small, and recall that  $x^n(T)$  is equal in distribution to  $x^n(\infty)$ . This proves Theorem 2.1.

*Proof of Lemma 5.1.* To establish (5.4) it will suffice to show that for some fixed  $T$  the following holds for any fixed sequence of initial states  $x^n(0)$  satisfying (5.3): the process can be constructed in such a way that w.p.1, for all sufficiently large  $n$ ,

$$\|x^n(T) - x^*\| \leq \delta + \delta_2. \tag{5.5}$$

Fix  $\tau > 2b/(\lambda + \delta/2)$ . Fix  $T > \tau$ . (The choice of  $T$  will be specified later.) For each  $n$ , at initial time 0, fix a subset of servers  $S_2$  as in condition (5.3); let  $S_1 = S \setminus S_2$  be the complementary subset of servers—its size is  $\delta n$ . Clearly, for each  $n$ ,

$${}^{(2)}W^n(0) \leq b, \quad {}^{(1)}\rho^n(t) \leq \delta \quad \text{for all } t.$$

Consider the Markov (stopping) time  $\tau^n$ , defined as the smallest time  $t$  in  $[0, \tau]$ , such that  ${}^{(2)}\rho^n(t) \leq \lambda + \delta/2$ ; if there is no such  $t$  then  $\tau^n = \infty$  by convention. The construction of

the process in  $[0, T]$  will be as follows: in the interval  $[0, \tau^n]$  it is driven by one set of driving processes, and in  $(\tau^n, T]$  it is driven by a different, independent set of driving processes with the same law. (However, these two sets of driving processes are common for all  $n$ .) In other words, at time  $\tau^n$  the process is ‘restarted’, with the state at  $\tau^n$  serving as initial state and with a new independent set of driving processes. By convention, if  $\tau^n = \infty$ , the process is *not* restarted.

We see that w.p.1, for all sufficiently large  $n$ ,

$$\tau^n < \tau. \tag{5.6}$$

Indeed, if we apply Lemma 4.3 to  $(^2)x^n(t)$  starting at time 0, we see that any fluid limit  $(^2)W(\cdot), (^2)\rho(\cdot)$  that can arise is such that  $(^2)W(0) \leq b$  and there exists  $t' \leq \tau/2$  such that  $(^2)\rho(t') \leq \lambda + \delta/2$ . (Otherwise,  $(^2)W(t)$  would become negative.) This implies (5.6).

Similarly, we see that w.p.1, for all sufficiently large  $n$ ,

$$(^2)W^n(\tau^n) \leq b_1 \doteq b + 2\lambda\tau. \tag{5.7}$$

Now consider any fixed sequence of  $\tau^n < \tau$  and fixed states at  $\tau^n$  satisfying (5.6) and (5.7). (Recall that starting at  $\tau^n$ , the process is controlled by a new independent set of driving processes.) Starting at time  $\tau^n$ , we keep the subset  $S_1$  as it was, but split  $S_2$  into two subsets  $S_3$  and  $S_4$  as follows:  $S_4$  will consist of  $(\frac{1}{2})n$  idle (at  $\tau^n$ ) servers (which exist by (5.6)), and  $S_3 = S_2 \setminus S_4$  will include the remaining  $[(1 - \delta) - \frac{1}{2}]n = (\frac{1}{2} - \delta)n$  servers from  $S_2$ . Clearly,  $(^3)W^n(\tau^n) = (^2)W^n(\tau^n) \leq b_1$ . To summarize, starting at  $\tau^n$ , the set of servers  $S$  is divided into three subsets,  $S_1, S_3$ , and  $S_4$ , with sizes  $\delta n, (\frac{1}{2} - \delta)n$  and  $(\frac{1}{2})n$ , respectively. Also, w.l.o.g. we assume that starting at  $\tau^n$  all new arrivals go to subset  $S_4$ , as long as there are idle servers in it. Applying Lemma 4.2, it follows that w.p.1, for all sufficiently large  $n$ , in the interval  $[\tau^n, T]$  all new arrivals go to subset  $S_4$ .

We now specify the choice of  $T$ . It has to satisfy two conditions. First, it has to be large enough so that, for any  $t \geq T - \tau, \|x^\uparrow(t) - x^*\| \leq \delta_2/3$ . Second, it has to be large enough so that

$$T - \tau > \frac{b_1}{\delta_2/3}.$$

Then applying Lemma 4.2 and (4.4), it follows that w.p.1, for all sufficiently large  $n$ ,

$$\|(^4)x^n(T) - x^*\| < \frac{\delta_2}{2}, \quad (^3)\rho^n(T) < \frac{\delta_2}{2};$$

this, in turn, implies (5.5). □

## 6. Generalizations

### 6.1. Renewal arrival process

The assumption that the arrival process is Poisson is made in order to simplify the exposition. Our main result, Theorem 2.1, and the analysis easily generalize to the case when the arrival process is renewal; in the  $n$ th system, the interarrival times are i.i.d., equal in distribution to  $A/n$ , where  $A$  is a positive random variable, and  $\mathbb{E}A = 1/\lambda$ . (Mild assumptions on the interarrival time distribution are needed to make sure that the process is positive Harris recurrent. For example, it suffices that this distribution has an absolutely continuous component.) The common process state space contains an additional scalar variable  $u$ , which is the residual interarrival time; clearly,  $u^n(\infty) \xrightarrow{D} 0$  as  $n \rightarrow \infty$ . The more general form of Theorem 2.1 is as follows: *if  $\lambda < \frac{1}{2}$  then  $(u^n, x^n)(\infty) \xrightarrow{D} (0, x^*)$ .*

The construction of the uniform stochastic upper bound on a single server workload generalizes as follows. For each  $n$ , the arrival process into a server, when it is busy, is dominated by a renewal process which is the thinned, with probability  $1/n$ , arrival process into the system. (In other words, as before, the dominating arrival process into a server, as long as the server remains busy, is such that *every* new arrival into the system goes to this server with probability  $1/n$ .) The interarrival times of this renewal process are i.i.d. with the distribution equal to that of a random variable  $A_n$ ; its mean is  $\mathbb{E}A_n = 1/\lambda$  for any  $n$ , but the distribution depends on  $n$ . However, as  $n \rightarrow \infty$ , the distribution of  $A_n$  converges to the exponential distribution. (This is a well-known property that a thinned, with probability  $1/n$  and sped up in time by factor  $n$ , renewal process converges to a Poisson process. And it is easy to check directly, since  $A_n$  is a sum of the geometrically distributed, with mean  $n$ , number of independent instances of  $A/n$ .) Then, for arbitrarily small  $\delta > 0$ , there exists a nonnegative random variable  $A^\delta$ , such that  $1/\lambda - \delta \leq \mathbb{E}A^\delta < 1/\lambda$ , and the distribution of  $A^\delta$  is dominated by that of  $A_n$  for all sufficiently large  $n$ . (For example, if  $\tilde{A}$  has an exponential distribution with mean  $1/\lambda$ , we can choose  $A^\delta = ((\tilde{A} \wedge C) - \varepsilon) \vee 0$ , where  $C > 0$  is large,  $\varepsilon > 0$  is small, and ‘ $\wedge$ ’, ‘ $\vee$ ’ denote minimum and maximum, respectively.) We fix  $\delta > 0$  such that  $1/\lambda - \delta > 1$ , and then  $\mathbb{E}A^\delta > 1$ . For all large  $n$ , the renewal arrival process with interarrival times distributed as  $A^\delta$  (and then the arrival rate  $1/\mathbb{E}A^\delta < 1$ ), dominates (pathwise, using natural coupling) the arrival process into an individual server, as long as the server remains busy. Therefore, the workload during the busy period under interarrival times  $A^\delta$  dominates that under interarrival times  $A_n$ . The rest of the construction of the uniform stochastic upper bound on a single server workload is the same. And after this bound is established, the rest of the proof of the main result remains essentially the same as well, with slight adjustments.

## 6.2. Biased routing when all servers are busy

Examination of the proof of Theorem 2.1 shows that the specific rule—uniform at random—for routing arriving customers when all servers are busy, is only used to obtain the process stability (positive Harris recurrence) and the uniform stochastic upper bound on a single server workload. In the rest of the proof, this specific rule is not used; we only use the fact that customers must go to idle servers if there are any. But, for the stability and workload upper bound, it suffices that the arrival rate into a server when it is busy is upper bounded by some  $\bar{\lambda} < 1$ , not necessarily by  $\lambda < \frac{1}{2}$ . This shows that Theorem 2.1 holds as is, even if routing when all servers are busy is biased in an arbitrary way, as long as the probability that a server receives an arrival does not exceed  $(1/n)(\bar{\lambda}/\lambda)$  for some  $\bar{\lambda} < 1$ .

## 6.3. Finite buffers

The main result, Theorem 2.1, holds as is if we allow some or all servers to have finite buffers (of the same or different sizes). If a server has a finite buffer of size  $B \geq 1$ , and already has  $B$  customers, any new customer routed to this server is blocked and leaves the system. It should be clear that our proof of Theorem 2.1 works for this more general system; additional ‘losses’ of arriving customers do not change the stochastic upper bound on a steady-state server workload; and the rest of the proof remains essentially unchanged, except for a more cumbersome state space description.

## References

- [1] BADONNEL, R. AND BURGESS, M. (2008). Dynamic pull-based load balancing for autonomic servers. In *Network Operations and Management Symposium, NOMS 2008*, IEEE, pp. 751–754.
- [2] BILLINGSLEY, P. (1995). *Probability and Measure*, 3rd edn. John Wiley, New York.

- [3] BRAMSON, M., LU, Y. AND PRABHAKAR, B. (2012). Asymptotic independence of queues under randomized load balancing. *Queueing Systems* **71**, 247–292.
- [4] BRAMSON, M., LU, Y. AND PRABHAKAR, B. (2013). Decay of tails at equilibrium for FIFO join the shortest queue networks. *Ann. Appl. Prob.* **23**, 1841–1878.
- [5] ESCHENFELDT, P. AND GAMARNIK, D. (2015). Join the shortest queue with many servers. The heavy traffic asymptotics. Preprint. Available at <https://arxiv.org/abs/1502.00999>.
- [6] LU, Y. *et al.* (2011). Join-idle-queue: a novel load balancing algorithm for dynamically scalable web services. *Performance Evaluation* **68**, 1056–1071.
- [7] MITZENMACHER, M. (2001). The power of two choices in randomized load balancing. *IEEE Trans. Parallel Distributed Systems* **12**, 1094–1104.
- [8] MUKHERJEE, D., BORST, S. C., VAN LEEUWAARDEN, J. S. H. AND WHITING, P. A. (2016). Universality of load balancing schemes on the diffusion scale. *J. Appl. Prob.* **53**, 1111–1124.
- [9] STOLYAR, A. L. (2015). Pull-based load distribution in large-scale heterogeneous service systems. *Queueing Systems* **80**, 341–361.
- [10] STOLYAR, A. L. (2017). Large-scale heterogeneous service systems with general packing constraints. *Adv. Appl. Prob.* **49**, 61–83.
- [11] STOLYAR, A. L. (2017). Pull-based load distribution among heterogeneous parallel servers: the case of multiple routers. *Queueing Systems*, **85**, 31–65.
- [12] VVEDENSKAYA, N. D., DOBRUSHIN, R. L. AND KARPELEVICH, F. I. (1996). A queueing system with a choice of the shorter of two queues—an asymptotic approach. *Problems Information Transmission* **32**, 15–27.