

FLUID LIMITS TO ANALYZE LONG-TERM FLOW RATES OF A STOCHASTIC NETWORK WITH INGRESS DISCARDING

BY JOHN MUSACCHIO^{*,†} AND JEAN WALRAND^{*,‡}

University of California, Santa Cruz[†] and University of California, Berkeley[‡]

We study a simple rate control scheme for a multiclass queuing network for which customers are partitioned into distinct flows that are queued separately at each station. The control scheme discards customers that arrive to the network ingress whenever any one of the flow's queues throughout the network holds more than a specified threshold number of customers. We prove that if the state of a corresponding fluid model tends to a set where the flow rates are equal to target rates, then there exist sufficiently high thresholds that make the long-term average flow rates of the stochastic network arbitrarily close to these target rates. The same techniques could be used to study other control schemes. To illustrate the application of our results, we analyze a network resembling a 2-input, 2-output communications network switch.

1. Introduction. We consider a multiclass queuing network whose customers are partitioned into F distinct flows. Customers of a flow $f \in \{1, \dots, F\}$ arrive according to an independent renewal process and follow a fixed, acyclic sequence of stations. The service times at each station are also independent. Each flow f has a weight $w_f \in \mathbb{R}_+$, and each of d stations is equipped with per-flow queues and serves a flow in proportion to its weight using a weighted round robin or a similar queueing discipline like weighted fair queueing or generalized head of line processor sharing.

We consider a simple scheme which we call ingress discarding for admitting customers. The ingress discarding scheme works as follows. Whenever any of a flow's queues exceed a threshold h , that flow's customers are discarded at the network ingress. There are two main objectives of the scheme: i) stability when the arrival rates in the absence of discarding would cause the utilization of some stations to exceed 1, and ii) fairness in the long-term average departure rates when the network cannot accommodate all the incoming flows. The contribution of this article is a methodology for proving

^{*}Research supported in part by NSF Grant ANI-0331659 and CNS-0953884.

AMS 2000 subject classifications: Primary 60K25, 68M20, 68M10; secondary 68K20.

Keywords and phrases: Fluid Limit, Stochastic Network

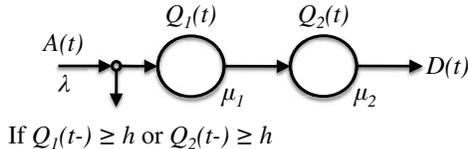


FIG 1. A queueing network with input discarding.

that the long-term average flow rates in such a network can be made arbitrarily close to those predicted by a fluid model, provided that the discarding thresholds are sufficiently high.

There are a number of applications of such a control policy. One application is for service centers such as call centers. It might be acceptable to block incoming customers but unacceptable to drop customers that have been admitted to the system, hence the appropriateness of ingress discarding. A designer of such a system might want to show that the flow rates of various types of customers are fair in some sense. This work can be used to show that if the system's fluid model achieves fair rates, then the system will achieve close to fair rates provided that the discarding thresholds are sufficiently high. Another application area is in data-packet switch design. A packet switch typically consists of several line-cards that transmit and receive the data packets, and a switch-fabric that serves as an interconnect. A design requirement might be that any packet discarding occur in the line-cards rather than in the switch fabric, since the line cards are better equipped to record statistics about the dropped packets for instance. The switch fabric can be thought of as a queueing network, and ingress discarding would be one way to fulfill the requirement that discarding only occur in the line cards. Again, this work shows that the flow rates of such a system approach those predicted by a fluid model if the discarding thresholds are made sufficiently high.

To illustrate our methodology, we consider the simple network in Figure 1. This network carries a single flow and customers arrive as a renewal process $E(t)$. There are two queues, each with i.i.d. service times with mean μ_i^{-1} in queue i ($i = 1, 2$). Designate by $Q_i(t)$ the length of queue i ($i = 1, 2$). The ingress discarding scheme discards the arrivals that occur when one of the two lengths is at least equal to threshold h . We want to show that if the thresholds are made large enough factor n , that the flow rates approach $\min\{\lambda, \mu_1, \mu_2\}$. More precisely, we want to show that for every $\epsilon > 0$ there exists some n_ϵ such that if threshold scale factor $n \geq n_\epsilon$, then the average rate of the departure process $D(t)$ exceeds $\min\{\lambda, \mu_1, \mu_2\} - \epsilon$. Note that since we scale the thresholds by a factor n , the starting value of the threshold h

is not important, so long as it is positive. Also note that we do not attempt to derive any result on the speed of convergence – how fast nh must grow to achieve rates within a smaller and smaller ϵ of the desired rates.

The analysis approach, which we believe can be extended to control strategies that change admission, service, or routing behavior when queue depths cross thresholds that can be made large, is based on deriving properties of the stochastic network using a fluid model. However for clarity of exposition, we limit our focus in this paper to the ingress discarding policy. As in work by Dai [4] we take a fluid limit by considering a sequence of larger and larger initial conditions, and scaling time and space by the size of those initial conditions. However, in order to consider stochastic networks with larger and larger thresholds, our fluid limit also considers a sequence of systems with thresholds scaled by an increasing factor n . The resulting fluid limit behaves according to a fluid model corresponding to the vector flow diagram in Figure 2. Since we scale the thresholds in our fluid limit, the thresholds appear in the fluid model with non-negligible values \bar{h} . Note that \bar{h} need not equal h since the fluid limits we consider may scale space and threshold at different rates. Also as a consequence of scaling the thresholds in taking the fluid limit, the stochastic system behaves like the fluid model (in terms of flow rates) only if the stochastic system's thresholds are sufficiently large.

First consider the case $\lambda > \mu_1 > \mu_2$. A fluid model corresponding to this case is illustrated by the vector flow diagram in the left part of Figure 2. This diagram indicates the rate of change of the vector of queue lengths as a function of its value. For instance, if the two queue lengths are between 0 and \bar{h} , then fluid enters queue 1 at rate λ and flows from that queue to queue 2 at rate μ_1 while fluid leaves queue 2 at rate μ_2 . Accordingly, the length of queue 1 increases at rate $\lambda - \mu_1$ and that of queue 2 at rate $\mu_1 - \mu_2$. The other cases can be understood similarly. The vector flow diagram shows that, irrespective of their initial values, the queue lengths converge to the pair of values $(0, \bar{h})$, which is an absorbing state for the fluid process. Moreover, when the process is close to the value $(0, \bar{h})$, the rate of the departure fluid is close to μ_2 . To conclude that the stochastic network has a departure rate close to μ_2 when h is large, one notes that the fluid process has one additional property: The time the process takes to reach the state $(0, \bar{h})$ is bounded by a linear function of the distance between the initial condition and $(0, \bar{h})$. This property, which can be seen from the vector flow diagram, can be used to show, roughly, that the stochastic system spends little time far from $(0, \bar{h})$. The intuition is that, although fluctuations occasionally move the stochastic network away from the limiting state, the system tends to follow the fluid process and get back to that state fairly quickly. This property will allow us

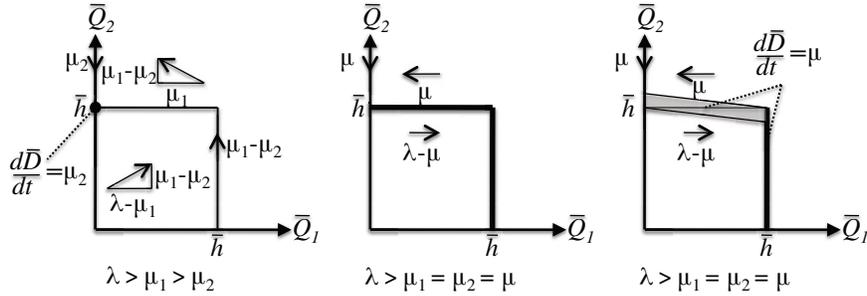


FIG 2. *The fluid process that approximates the stochastic network.*

to construct a proof that the stochastic network has a departure rate close to μ_2 most of the time.

It turns out that one needs a generalization of the above approach to cover some interesting cases. To illustrate this generalization, consider once again the network of Figure 1 but assume that $\lambda > \mu_1 = \mu_2 = \mu$. The vector flow diagram of the corresponding fluid process is shown in the middle part of Figure 2. The diagram shows that the fluid process converges to some point in the set indicated by the two thicker lines: $\{\bar{h}\} \times [0, \bar{h}] \cup [0, \bar{h}] \times \{\bar{h}\}$, depending on the initial condition. While it is true that the rate of the departure fluid is close to μ for any point close to that set, it is no longer the case that the time to reach that limiting set is bounded by a linear function of the initial distance to the set. For instance, if the initial state of the fluid process is $(\bar{h}, \bar{h} + \epsilon)$ for some arbitrary $\epsilon > 0$, the process takes at least \bar{h}/μ to reach the limiting set. To handle this situation, one considers the set shown in the right-hand part of Figure 2. That set has the following two key properties: 1) the departure flow rate is almost μ close to that set; and 2) the time to reach the set is bounded by a linear function of the initial distance to it, as can be see from the diagram. Thus, as in the previous example, one can show that the stochastic network has a departure rate close to μ most of the time.

The main technical contributions of the paper are as follows:

- A technique for scaling time, space, and threshold for finding a fluid limit for a stochastic network with threshold based ingress discarding such as in our example;
- Proof of a fluid limit for stochastic networks with thinned processes such as $\Lambda(t)$ in Figure 1;
- Proof of approximation of the rates of the stochastic network by the rates of the limiting fluid process under the two key properties indicated in our examples.

In the next subsection, we outline the key steps of our analysis. In subsection 1.2 we relate our work to other prior work, and in subsection 1.3 we review an example stochastic network with ingress discarding. Section 2 establishes the notation and initial model description, while section 3 proves the main results of the article. In section 4 we study the fluid model of a network resembling a 2 x 2 network switch and show that the fluid model has the necessary properties to employ the main results of the article. Note that Musacchio [22] shows that a more general network with ingress discarding has a fluid model with the necessary properties. Section 5 concludes the paper.

1.1. *Proof Outline.* Our goal is to show that the long-term average flow rates of the stochastic system can be made arbitrarily close to a vector of desired rates R if the discarding thresholds are made large enough. Moreover, we want to show that certain properties of the system's fluid model suffice to reach this conclusion. In this subsection we outline the arguments detailed in the rest of the paper.

The queuing network we consider has ingress discarding thresholds of nh in each queue, where $h > 0$ and $n > 0$ is a threshold scale factor that is increased to make the thresholds larger. The network is described by a Markov process $X^n = \{X^n(t), t \geq 0\}$ taking values in the state space X . The superscript emphasizes the dependence on n . The state of the Markov process includes the queue lengths, remaining service times at each queue, and the remaining time until the next exogenous arrival of each flow $f \in 1, \dots, F$. We will argue that X^n satisfies the strong Markov property.

As we discussed in the previous section, we construct fluid limits of the system by scaling time, space, and threshold scale factor in particular ways that we describe below. These fluid limits converge (in a sense also described below) to trajectories of a fluid model. The fluid model, like the original system, also has ingress discarding thresholds. However these thresholds need not equal h , since one of the fluid limits we need to consider can scale space and threshold at different rates. Therefore when referring to the system's fluid model, we need to specify \bar{h} , the discarding thresholds of each queue of the fluid model. (The queues of the fluid model have a common threshold \bar{h} , just as the queues of the original system have a common threshold nh .) The fluid model has a state space $\bar{\mathsf{X}}$ similar to that of the original system, but the queue lengths take values in \mathbb{R}^+ rather than \mathbb{Z}^+ . In what follows we adopt the notation that if $\mathcal{S} \subset \bar{\mathsf{X}}$ then the set $a\mathcal{S}$ ($a \in \mathbb{R}_+$) denotes a "scaled" set such that $\bar{x} \in a\mathcal{S}$ iff $\bar{x}/a \in \mathcal{S}$. Also let $\|\bar{x}\|_{\mathcal{S}} = \inf_{e \in \mathcal{S}} \|\bar{x} - e\|$ denote the distance between \bar{x} and the set \mathcal{S} .

Our goal is to show that if there exists a closed, bounded set $\mathcal{E} \subset \bar{X}$ and $t_0 \in \mathbb{R}^+$ such that conditions C1 and C2 below hold, then there exists a large enough n such that the stochastic network achieves long-term average departure rates arbitrarily close to R . Conditions C1 and C2 are as follows:

- C1** All trajectories of the fluid model with ingress discarding thresholds \bar{h} and initial condition $\bar{X}(0) = \bar{x}$ are absorbed by a set $\bar{h}\mathcal{E}$ in a time not more than $t_0 \|\bar{x}\|_{\bar{h}\mathcal{E}}$;
- C2** If $\bar{h} > 0$, the instantaneous departure rates of the fluid model while its state is in the set $\bar{h}\mathcal{E}$ are equal to the vector of desired rates R .

Note that C1 requires that $\bar{h}\mathcal{E}$ be an absorbing set of the fluid model with thresholds \bar{h} . For example, one can show that a minimal absorbing set of the fluid model in many cases would be, roughly, the set of states such that at least one of each flow's set of "bottleneck" queues is at its discarding threshold, and servers with a utilization below 1 have empty queues. (By "bottleneck queue", we mean a queue whose service constrains a flow's rate in the fluid model.) However, such a construction might not be sufficient to satisfy C1, particularly when flows do not have unique "bottlenecks." Recall that in the introduction we studied an example with two serially-connected queues with the same service rate. This is an example in which $\bar{h}\mathcal{E}$ needs to be made larger than the minimal absorbing set in order to satisfy C1. To see this note that even though the two line segments in the middle panel of Figure 2 constitute an absorbing set for the fluid model, if we defined \mathcal{E} so that $\bar{h}\mathcal{E}$ is equal to these two line segments (by making $\mathcal{E} = \{1\} \times [0, 1] \cup [0, 1] \times \{1\}$), condition C1 would not be met. By defining \mathcal{E} in such a way as to make $\bar{h}\mathcal{E}$ have the shape indicated by the shaded area of the right panel of Figure 2, the time it takes trajectories of the fluid model to reach $\bar{h}\mathcal{E}$ can be upper bounded by an amount proportional to the distance of the starting point of the trajectory from $\bar{h}\mathcal{E}$, thus satisfying C1.

The proof depends on two main steps.

- i) The expected flow rates associated with the process $X^n(\cdot)$, over a finite time interval of length nt_0 , and for initial conditions near a set $nh\mathcal{E}$, can be made to be arbitrarily close to R with a sufficiently large threshold scaling factor n .
- ii) The excursions of the process $X^n(\cdot)$ away from $nh\mathcal{E}$ become relatively shorter with larger threshold scaling factor n . More precisely, the first hitting time that occurs nt_0 after having started in a neighborhood of the set $nh\mathcal{E}$, can be made to be arbitrarily close to nt_0 .

In both steps we make use of the fact that a fluid limit of the process $X^n(\cdot)$ converges to a trajectory of the fluid model, but the different objectives of

the two steps require us to use different fluid limit scalings. In the first step we consider a sequence of (initial condition, scale factor) pairs $\{(x_j, n_j)\}$. To emphasize the dependence on initial condition and threshold scale factor we write $X^{\mathbf{x}_j}(\cdot)$, where the superscript $\mathbf{x}_j \triangleq (x_j, n_j)$. We require that the sequence has the properties that x_j/n_j is no more than a distance $\zeta < 1$ away from the set $h\mathcal{E}$, and $n_j \rightarrow \infty$. Otherwise, the sequence is arbitrary. We call such a sequence a *near fluid limit* sequence. (Equivalently, the near fluid limit condition has $\|x_j\|_{n_j h\mathcal{E}} < n_j \zeta$ and $n_j \rightarrow \infty$. In general it is often more intuitive to consider the distance of X/n from the set $h\mathcal{E}$ than to consider the distance of X from $nh\mathcal{E}$, so we will use whichever construction is more convenient or intuitive for the context.) We demonstrate that the sequence of scaled processes $\{\frac{1}{n_j}X^{\mathbf{x}_j}(n_j \cdot)\}$ converges along a subsequence, uniformly over compact time intervals, to a fluid model trajectory $\bar{X}(\cdot)$. The result largely follows from the fact that the process describing the cumulative time each server in the network is busy is Lipschitz continuous, and a sequence of Lipschitz continuous functions on a compact set converges along a subsequence. Consequently, the convergence to a fluid trajectory only holds on a finite time interval. The thresholds of the fluid model that $\bar{X}(\cdot)$ satisfies are of size $\bar{h} = h$. This is because we scale both space and threshold by the same amount in this fluid limit, so the two scalings cancel out. Moreover, the restrictions we put on the near fluid limit sequence ensure that the initial condition of the fluid model trajectory $\bar{X}(\cdot)$ is within a distance of ζ of $\bar{h}\mathcal{E}$. Thus, the fluid model trajectory $\bar{X}(t)$ hits $\bar{h}\mathcal{E}$ quickly (in not more than time ζt_0 by C1) and then achieves flow rates of R (by condition C2).

At this point, we have only shown convergence along a subsequence to a fluid trajectory with some desired properties. We need to show convergence along the original near fluid limit sequence in order to eventually make conclusions about the stochastic network. To that end, consider a functional \mathfrak{F} that extracts the difference between the actual flow throughput and the desired flow throughput over a compact time interval $[\zeta t_0, t_0]$ (in time scaled by n). Since $\bar{X}(t)$ hits $\bar{h}\mathcal{E}$ by time ζt_0 , the flow rates are equal to the desired rates over $[\zeta t_0, t_0]$. Consequently, $\mathfrak{F} \circ \bar{X} = 0$. This in turn allows us to argue that $\{\mathfrak{F} \circ \frac{1}{n}X^{\mathbf{x}_j}(n \cdot)\}$ converges to 0 along a subsequence. Since every near fluid limit sequence of processes (with the functional applied to them) converges along a subsequence to 0 in this way, it must be that every near fluid limit sequence also converges to 0 in this way. This fact allows us to show that the flow rates of the process $\frac{1}{n}X(n \cdot)$ can be made arbitrarily close to the desired rates, for a finite time period, from any scaled initial condition x/n near $h\mathcal{E}$, provided that n is sufficiently large. In the detailed proof the functionals we consider act on the Markov state trajectory combined with

the trajectories of some other associated processes such as the cumulative service time process. The fact that $\zeta < 1$ was chosen otherwise arbitrarily is important because it allows us to later make ζ small so that the desired rates are achieved over most of the interval $[0, t_0]$ (in scaled time).

In the second step, we again consider a sequence of (initial condition, scale factor) pairs $\{(x_j, n_j)\}$. This sequence must satisfy the properties that the distance between x_j/n and $h\mathcal{E}$ is more than a constant ζ for each j , and that $\|x_j\|_{n_j h\mathcal{E}} = n_j \|x_j/n_j\|_{h\mathcal{E}} \rightarrow \infty$. Otherwise, the sequence is arbitrary. We call such a sequence a *far fluid limit* sequence. We show that the sequence of scaled processes $\{X^{x_j}(\|x_j\|_{n_j h\mathcal{E}} \cdot) / \|x_j\|_{n_j h\mathcal{E}}\}$ converges along a subsequence of any far fluid limit sequence, uniformly over compact time intervals, to a fluid model trajectory $\bar{X}(\cdot)$ satisfying a fluid model with discarding thresholds \bar{h} . The scaled threshold sequence of the fluid limit is $\{n_j h / (n_j \|x_j/n_j\|_{h\mathcal{E}})\}$, so the choice of sequence and convergent subsequence determines a value for \bar{h} that satisfies $\bar{h} \in [0, h\zeta^{-1}]$. Also the scaling of the far fluid limit sequence ensures that the initial condition of the fluid trajectory have an initial condition that is unit distance from $\bar{h}\mathcal{E}$. This fact along with our starting assumption C1, ensure that $\|\bar{X}(t_0)\|_{\bar{h}\mathcal{E}} = 0$. The preceding two facts allow us to argue that the sequence $\{X^{x_j}(\|x_j\|_{n_j h\mathcal{E}} t_0) / \|x_j\|_{n_j h\mathcal{E}}\}$ has a distance from $n_j h\mathcal{E}$ that converges to 0 along a subsequence. Moreover since any far fluid limit sequence has a subsequence that converges to 0 in this sense, it must be that this convergence property holds for any far fluid limit sequence.

This fact is the basis for constructing an argument that

$$\mathbb{E} \|X^x(t_0 \|x\|_{nh\mathcal{E}})\|_{nh\mathcal{E}} \leq \delta \|x\|_{nh\mathcal{E}}$$

for any $\delta > 0$ provided that threshold scale factor n is sufficiently large and $\|x/n\|_{h\mathcal{E}} > \zeta$ (equivalently $\|x\|_{nh\mathcal{E}} > n\zeta$). This relation serves as a Lyapunov function which allows the construction of an argument about the recurrence time of the scaled process X/n to a neighborhood with distance ζ of $h\mathcal{E}$, and this in turn allows us to conclude (ii) above.

This recurrence time argument is adapted from [18] while the overall argument we make with the far fluid limit sequence parallels [4]. The main difference between our far fluid limit argument and that of [4] is that in [4] the fluid model and stochastic network are drawn to the origin and neighborhood of the origin respectively, whereas in our model the system is attracted to a set of states.

1.2. Relation to Prior Work. Our fluid limit proof techniques borrow heavily from work by Dai [4]. Dai shows that for networks without discard-

ing, stability of a corresponding fluid model implies positive Harris recurrence of the stochastic network. In our work we use the fluid model not only to show positive Harris recurrence of the stochastic network, but also to find its long term average flow rates. Specifically, we use two fluid limits: the far fluid limit and the near fluid limit that correspond to different sequences of initial conditions and threshold pairs.

Dai's proof considers a sequence of initial states $\{x\}$ of the Markov process describing the network, with $|x| \rightarrow \infty$, and then obtains a fluid limit by scaling time and space by $|x|$. Dai uses this result to construct a Lyapunov function to show that the expected state of the system contracts, for initial states far enough from the origin. Our far fluid limit analysis parallels this, but with the difference that our analysis focuses on the distance of the state from a set of states $h\mathcal{E}$ rather than the distance from the origin. Also, because we are interested in showing the existence of a sufficiently large threshold scaling factor n , for both the near and far fluid limits, we consider a sequence of initial condition threshold pairs $\{x, n\}$ to obtain our results rather than just a sequence of initial conditions as in [4].

Our fluid limit technique is also very similar to that found in work by Bramson [2]. In much the way we do, Bramson takes the fluid limit using a sequence of pairs, one being the initial condition and the other being a time scaling factor of both space and time. However, our results do not follow immediately from the results of Bramson because we require that the fluid model be drawn toward a set $\bar{h}\mathcal{E}$ rather than just to the origin.

Another body of work uses fluid limits to show rate stability rather than showing that the system state converges to an invariant distribution, or more precisely that the system is positive Harris recurrent. Rate stability means that the long-term average departures match the long-term average arrivals. It is a weaker concept than positive Harris recurrence because a system can be rate stable while internally the average queue lengths grow unbounded or at least fail to converge to an invariant distribution. For a treatment see [11], and examples of its application include [3] and [7]. The rate stability framework is not sufficient for our objectives because in order to show that our control policy achieves flow rates close to those predicted by a fluid model, we need to show that the vector of queue lengths settles to an invariant distribution concentrated near a particular set of lengths, as illustrated in the example of the introduction.

Another closely related work to ours is by Mandelbaum, Massey, and Reiman [17]. In [17], the authors study the fluid limit of a queueing network with state dependent routing, where the function describing the arrivals to each queue can scale with n and or \sqrt{n} , in a manner similar to the scaling of

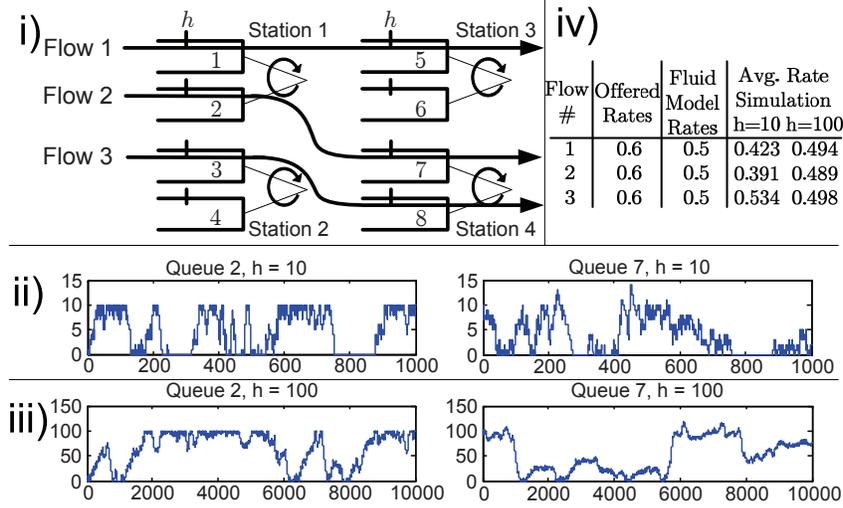


FIG 3. *i)* An illustration of the example network. *ii)* The queue lengths of queue 2 and 7 when the discarding threshold parameter h set to 10. *iii)* The lengths of queue 2 and 7 when the discarding threshold parameter h is set to 100. *iv)* The average flow rates of each of the three flows for both the $h = 10$ and $h = 100$ simulated sample paths.

our thresholds. The authors prove a functional strong law of large numbers and a functional central limit theorem in the context of their model. However, the authors assume that the network is driven by Poisson processes, rather than just the renewal assumption that we make. An earlier work by Konstantopoulos, Papadakis, and Walrand derives a functional strong law of large numbers and a functional central limit theorem for networks with state dependent service rates [16].

There are also several other works that use reflected Brownian motion models to study queueing networks with blocking [5, 13, 14]. Typically the objective of most such investigations is to approximate the distribution of the queue occupancy with a diffusion approximation. In contrast with those works, our objective is to show almost sure convergence using a strong law of large numbers scaling.

1.3. Example Network. In this subsection, we introduce an example that motivates the theory developed in this paper. The example will illustrate two important phenomena – that the long-term rates of the stochastic system get closer to those of a corresponding fluid model when discarding thresholds are raised, and that when there are not unique bottlenecks, the vector of queue depths is not attracted to a unique equilibrium point.

Our example is illustrated in Figure 3. The example is analogous to a

two-input and two-output switch. Two flows enter the network at station 1, the first input of the switch, and a third flow enters the network at station 2. We concentrate on flow 2, which shares stations 1 and 4 with flows 1 and 3 respectively. All stations are served at rate 1, have round-robin service with equal weighting to all queues, and have service times that are exponentially distributed. The arrival rate of each flow is 0.6, with Pareto interarrival distributions given by

$$P(\xi_f(j) > s) = \frac{1}{(0.6s + 1)^2}, \quad f \in \{1, 2, 3\}, \quad s \geq 0,$$

where $\xi_f(j)$ is the interarrival time preceding the j th arrival. We choose the Pareto distribution for this example to emphasize that we are interested in networks whose interarrival and service times are not necessarily memoryless.

We consider the behavior of the network's fluid model. Since stations 2 and 3 have a capacity of 1 and each carry one flow with an offered rate of 0.6, the queues of these stations should never fill. Stations 1 and 4 each carry 2 flows that offer a load of 0.6 (before considering discarding). The fluid model of the station's round robin service is that each station serves both of its queues at rate 0.5 as long as both flows are offering enough customers to be served at this rate. Consequently, when flow 1's queue at station 1 is filled below threshold, this queue grows at a rate of 0.1. However if flow 1's queue at station 1 ever went above its threshold, ingress discarding would commence, and the queue would immediately decrease. Therefore, it must be that this queue grows to its threshold, stays at this level, and then flow 1's "thinned" or post-discarding arrival process is of rate 0.5.

Similar reasoning shows that flow 3's queue at station 4 behaves in this way, and also that one of flow 2's queues must also reach the threshold and "stick" there. These steps allow us to conclude that after some time, all three flows should have rates of 0.5 in the system's fluid model. (We will verify this carefully in section 4.)

Figure 3 shows the simulated trajectories of flow 2's queues at both bottleneck stations in the stochastic network. In the $h = 10$ case, the simulation shows that queues 2 and queue 7, which both serve flow 2, are empty for over 100 time units around time 800. This empty period is significant because when flow 2's queues are empty, flow 2 misses opportunities to have its customers served by the bottleneck stations. Indeed the table included in Figure 3 shows the average rate, averaged over the last 80% of the simulation time to reduce some of the initial transient effect, is 0.391. This is substantially below the rate of 0.5 predicted by the fluid model. Most

likely, a string of long interarrival times of flow 2, caused the queues at the bottleneck stations to starve.

Raising the thresholds should reduce starvation, because larger thresholds would provide the bottleneck queues a larger backlog to smooth over fluctuations in the arrival and service processes. To test that intuition, we simulate the network with discarding thresholds of $h = 100$. Figure 3 shows the trajectories of flow 2's queues for the increased threshold. We note that neither queue spends all of the time filled to its threshold, but instead at most times at least one of the queues is near its threshold. For instance, at the beginning of the simulation, queue 7 (the second bottleneck) is chattering near the threshold while queue 2 (the first bottleneck) is below threshold. At some time before the 2000 second mark, the two queues switch these roles, and around the 6000 second mark the queues switch these roles again. We also note that flow 2 achieves an average rate of 0.489, which is much closer to the rate of 0.5 predicted by the fluid model.

2. Preliminaries. Customers of a given flow $f \in \{1, \dots, F\}$ follow the same fixed sequence of distinct stations. The service times are independent. Each flow f has a weight w_f and each station $i \in \{1, \dots, d\}$ is equipped with per-flow queues and serves each flow in proportion to its weight using a weighted round robin or a similar queueing discipline. In addition to the notion of flow, each customer also has a class $k \in \{1, \dots, K\}$ that is indicative of both the customer's flow and the station $s(k)$ it is located. Thus the class of a flow f customer changes as the customer progresses from station to station, but with the restriction that a flow f customer must always have a class in the set $K(f)$. Conversely, each class k is associated with one and only one flow $f_{(k)}$. We also adopt the numbering convention that flow f customers enter the network as class $k = f$, and thus $f \in K(f)$. The constituency matrix $C \in \{0, 1\}^{d \times K}$ records which classes are served in each station: $C_{ik} = 1$ if class k is served at station i , otherwise $C_{ik} = 0$. A customer of class k who completes service becomes a customer of class l if $P_{kl} = 1$. Thus $P \in \{0, 1\}^{K \times K}$ is a binary incidence matrix with each row containing at most one 1. Because flows follow loop-free paths, P is nilpotent.

The exogenous arrivals to the network for flow f are described by a renewal process $E_f(\cdot)$ for which the interarrival times $\{\xi_f(j), j \geq 1\}$ are i.i.d. and α_f is the mean arrival rate. Thus,

$$E_f(t) = \max\{r : U_f(0) + \xi_f(1) + \dots + \xi_f(r-1) \leq t\}, \quad t \geq 0,$$

where $U_f(t) \in \mathbb{R}_+$ is the time after time t until the next flow f customer arrives at the network ingress. We also need to assume that interarrival

times are unbounded and spread-out. More precisely, we assume that for each $k \in \{0, \dots, F\}$, there exists an integer j_k and some function $p_k(x) \geq 0$ on \mathbb{R}^+ with $\int_0^\infty p_k(x) dx > 0$, such that

$$(1) \quad \mathbb{P}[\xi_k(1) \geq x] > 0 \text{ for any } x > 0, \text{ and}$$

$$(2) \quad \mathbb{P} \left[a \leq \sum_{i=1}^{j_k} \xi_k(i) \leq b \right] \geq \int_a^b p_k(x) dx \quad \text{for any } 0 \leq a \leq b.$$

The service times $\{\eta_k(j), j \geq 1\}$ of each class k are also i.i.d. and have mean $m_k = \mu_k^{-1}$, where μ_k is the mean service rate. We also define the $K \times K$ diagonal matrix M whose k -th diagonal entry is m_k . The quantity $V_k(t) \in \mathbb{R}_+$ denotes the remaining service time of the class k customer in service, if there is one at time t , otherwise $V_k(t) = 0$. We define a service process $S_k^x(\cdot)$ as

$$S_k(t) \triangleq \max\{j : \tilde{V}_k(0) + \eta_k(1) + \dots + \eta_k(j-1) \leq t\}, \quad t \geq 0,$$

where $\tilde{V}_k(0) = V_k(0)$ if $V_k(0) > 0$, otherwise $\tilde{V}_k(0) = \eta_k(0)$ is a fresh service time with the same distribution as $\eta_k(1)$ and independent of all other service times.

In principle, our assumption that the service times are independent does not allow for service times that depend on a packet's size (taking "packets" to be "customers"). Dependence on packet size would make the service times of stations dependent on each other. To model this explicitly would require a much more complicated model. However we believe that our results in this work would still hold if this assumption were relaxed.

We define the following right-continuous processes: $A : [0, \infty) \rightarrow \mathbb{Z}_+^K$ counts the arrivals to each class k since time $t = 0$; $D : [0, \infty) \rightarrow \mathbb{Z}_+^K$ counts the departures of each class; $\Lambda : [0, \infty) \rightarrow \mathbb{Z}_+^F$, counts the exogenous arrivals of each flow that make it past the discarding point ("thinned" exogenous arrivals); $Q : [0, \infty) \rightarrow \mathbb{Z}_+^K$ is the vector process of queue depths; $T : [0, \infty) \rightarrow \mathbb{R}_+^K$ counts the total time each class k has been served since $t = 0$; and $I : [0, \infty) \rightarrow \mathbb{R}_+^d$ counts the total time each server has been idle

since $t = 0$. For each $t \geq 0$, these processes satisfy the following relations:

$$(3) \quad A(t) = P^T D(t) + \Lambda(t),$$

$$(4) \quad Q(t) = Q(0) + A(t) - D(t),$$

$$(5) \quad Q(t) \geq 0,$$

$$(6) \quad T_k(t) \text{ is nondecreasing and } T_k(0) = 0, \text{ for } k = 1 \dots K,$$

$$(7) \quad I_i(t) = t - C_i T(t) \text{ is nondecreasing and } I_i(0) = 0, \text{ for } i = 1 \dots d,$$

$$(8) \quad \int_0^\infty (CQ(t)) dI(t) = 0,$$

$$(9) \quad D_k(t) = S_k(T_k(t)) \text{ for } k = 1 \dots K.$$

Relations (3)-(5) describe the relations between the arrival, departure, and queue length processes. Statements (6)-(8) describe basic restrictions on the cumulative service time and idle time processes, with relation (8) reflecting an assumption that each station is work conserving. Equation (9) reflects that departures of class k are determined by the composition of the service time counting process $S_k(\cdot)$ and the process $T(\cdot)$.

The ingress discarding scheme drops arriving customers of flow f as they arrive whenever any queue in the set $K(f)$ exceeds a high threshold nh . Recall that n is the threshold scaling factor which we will adjust in our analysis. Conversely, when all of the queues in $K(f)$ are below a lower threshold $nh - o(n)$, flow f customers are permitted to enter the network. Note the lower threshold could be set to be the same as the upper threshold, but in some practical applications it might be beneficial to have different thresholds so that the switching between admitting and discarding is less frequent. Thus we permit this difference between upper and lower thresholds to be any function $o(n)$ that satisfies $o(n)/n \rightarrow 0$ and $o(n) \geq 0$. For instance any nonnegative constant may be used. Between these thresholds, the system has hysteresis behavior, and we define this behavior as follows. A process $H_k : [0, \infty) \rightarrow \{0, 1\}$ keeps track of whether discarding has been ‘‘turned-on’’ by each class k queue. If $Q_k(t) \geq nh$ then $H_k(t) = 1$ and if $Q_k(t) \leq nh - o(n)$ then $H_k(t) = 0$. For all t such that $Q_k(t) \in (nh - o(n), nh)$, the evolution of H_k is determined by the following rules:

- If $H_k(t) = 0$ then let $t_s = \min\{\tau \geq 0 : Q_k(t + \tau) \geq nh\}$ (note that t_s is well defined because $Q_k(\cdot)$ is right continuous). $H_k(t + \tau) = 0$ for $\tau \in [0, t_s)$ and $H_k(t_s) = 1$;
- If $H_k(t) = 1$ then let $t_s = \min\{\tau \geq 0 : Q_k(t + \tau) < nh - o(n)\}$. $H_k(t + \tau) = 1$ for $\tau \in [0, t_s)$ and $H_k(t_s) = 0$.

The flow f customers that are allowed into the network beyond the dis-

carding point depends on all the processes $H_k(\cdot)$ as

$$(10) \quad \Lambda_f(t) = \sum_{j=1}^{E_f(t)} \prod_{K(f)} (1 - H_k(\tau_j -))$$

where $\tau_j = U_f(0) + \sum_{m=1}^{j-1} \xi_f(m)$ is the time of the j th arrival to the discarding point. Here the dependence on $H_k(\tau_j -) \triangleq \lim_{t \uparrow \tau_j} H_k(t)$ rather than $H_k(\tau_j)$ is to avoid problems with causality. For instance a customer arrival that triggers discarding should not be discarded, or otherwise the customer will never arrive to the system and paradoxically the discarding will never turn on. Our modeling choice allows such a customer to enter, thus triggering discarding, which will discard future customers.

The queueing discipline of a station i serves each flow in proportion to the flow weights over long time intervals. More precisely, for some constant $c > 0$ and all $\tau > 0$,

$$(11) \quad \frac{D_k(t, t + \tau)}{w_{f(k)}} \geq \frac{D_l(t, t + \tau)}{w_{f(l)}} - c \text{ whenever } Q_k(s) > 0 \forall s \in [t, t + \tau]$$

for all $k, l \in C(i) \triangleq \{k' : C_{ik'} = 1\}$, where $D_k(t, t + \tau) \triangleq D_k(t + \tau) - D_k(t)$.

We furthermore assume that only the customer at the head of line of each queue may be served, and that the instantaneous service rate of any queue is a function of the current state. That is $\dot{T}_k(t) = f(X(t))$ for some function $f(\cdot)$ where $X(t) = [Q(t); U(t); V(t); H(t)]$.

The evolution of the queueing system depends on the particular queueing discipline. Moreover, some queueing disciplines require additional state variables. For instance, a weighted round robin scheduler visits the queues in a cyclic order, serving any customers at the head of the line. The order should be chosen so that in each cycle the number of visits of each queue is proportional to the flow weights. (Which is possible if the weights are rational multiples of each other.) Other queueing disciplines could be considered as well, though these disciplines may need additional state variables. For instance, Deficit Round Robin (DRR) requires counters for each class [23]. Also, DRR ensures that the service times given to each class are proportional rather than the number of customers served. Therefore, DRR satisfies a criterion similar to (11) except that $D(\cdot)$ is replaced by $T(\cdot)$. However, since the service times are unbounded, the criterion holds only in the limit $\tau \rightarrow \infty$, almost surely. Other disciplines require yet more complex state descriptions. For instance, Weighted Fair Queueing (WFQ) keeps track of each customer's "virtual finish time" – the time they would have departed if the

service discipline were weighted processor sharing and no more customers were to arrive [1]. To keep the presentation simple, we assume that the additional state variables required by the queueing discipline are described by a bounded vector in \mathbb{Z}_+^d . We append this to the H portion of the state description. Treatment of queueing disciplines that require more elaborate state descriptions requires some modification to the statement and proof of Theorem 1.

2.1. State Description. The dynamics of the queueing network are described by the Markov process $X = \{X(t), t \geq 0\}$. The state description contains the queue lengths $Q(t) \in \mathbb{R}_+^K$ of all the K queues in the network, as well as the residual arrival and service times $U(t) \in \mathbb{R}_+^F$ and $V(t) \in \mathbb{R}_+^K$ respectively. Recall that $U(t)$ and $V(t)$ are defined to be right-continuous. Finally the state description includes the state of the discarding hysteresis and any state variables used by the queueing discipline as described above. We assume that $H(t) \in \{0, 1\}^K \times \mathbb{Z}_+^d$. Thus the full state description is

$$X(t) = [Q(t); U(t); V(t); H(t)].$$

Let $\mathbf{X} \subset \mathbb{Z}^K \times \mathbb{R}_+^{F+K} \times \{0, 1\}^K \times \mathbb{Z}_+^d$ be the set of all states X can take. A fixed threshold scaling factor n , an initial condition $x = X(0) \in \mathbf{X}$ is sufficient to specify the statistics of the future evolution of the system.

We claim that the process X satisfies the strong Markov property, by the same argument given by Dai [4]. In turn, Dai's argument followed from Kaspi and Mandelbaum [15]. Without repeating all the details of the argument, the basic idea is that $X(\cdot)$ is a Piecewise Deterministic Markov (PDM) process – behaving deterministically between the generation of “fresh” interarrival or service time. Davis shows that a PDM process whose expected number of jumps on $[0, t]$ is finite for each t is strong Markov [9]. As we assume that the interarrival and service times have a positive and finite mean, the expected number of jumps of $X(\cdot)$ in any closed time interval is finite. Therefore $X(\cdot)$ has the strong Markov property.

The fluid model, whose defining equations will be given in Theorem 1, takes values in the state space $\bar{\mathbf{X}} \subset \mathbb{R}_+^{F+3K+d}$ since integer valued states of the original system correspond to real valued states of the fluid model.

3. Fluid Limit Analysis. In this and subsequent sections, we use the superscript $\mathbf{x} \triangleq (x, n)$ to denote the dependence on initial state x and threshold scaling factor n . As we discussed earlier, we use two different fluid limits in our analysis: the near and far fluid limits that study behavior of the stochastic network for scaled initial conditions near and far from

$h\mathcal{E}$ respectively. Recall that $\mathcal{E} \subset \bar{X}$ is a closed and bounded set. Also recall $h\mathcal{E} = \{x : x/h \in \mathcal{E}\}$. At this point we make no further assumptions on \mathcal{E} , but eventually \mathcal{E} will have to be chosen so that $\bar{h}\mathcal{E}$ is an absorbing set of the fluid model with thresholds \bar{h} to apply our final results.

For notational convenience we also define an augmented state vector process

$$\mathfrak{X}^{\mathbf{x}}(\cdot) \triangleq [X^{\mathbf{x}}(\cdot); T^{\mathbf{x}}(\cdot); \Lambda^{\mathbf{x}}(\cdot); nh]$$

which contains all the functions we want to show converge in both kinds of fluid limit.

In this section, we state Theorem 1 which shows convergence to a fluid model trajectory along a fluid limit. The convergence of the trajectory is uniformly on compact sets. More precisely, we say that $f_j(t) \rightarrow f(t)$ *uniformly on compact sets* (u.o.c.) if for each $t \geq 0$,

$$\lim_{j \rightarrow \infty} \sup_{0 \leq s \leq t} |f_j(s) - f(s)| = 0.$$

We also use the notation $\dot{f}(t) = \frac{d}{dt}f(t)$ where such a derivative exists. If a function $f(\cdot)$ is differentiable at t , we say that t is a *regular point*.

The proof, along with four lemmas used in the proof, are given in the appendix. One of these lemmas, Lemma 5, is a new result showing that the thinned arrival process converges u.o.c. to the fluid limit. In section 1.1 we previewed the two types of fluid limit, which we call the “near” and “far” fluid limits, that we will use in our analysis. In both types of fluid limit, time and space is scaled by a factor that increases. In the development that follows, that scale factor for time and space is represented by the notation a_j . Later on, we will make specific assumptions about a_j that correspond to either the near or far fluid limit. Bramson [2] takes a similar approach to defining the fluid limit. Both types of fluid limit scale the threshold no faster than time and space are scaled, and also both consider a sequence of initial conditions x_j , such that after space-scaling, the “relative initial condition” x_j/a_j is a bounded distance away from the set $\frac{n_j h}{a_j} \mathcal{E}$. More precisely, we define the following property which is common to both near and far fluid limit sequences. Thus by assuming this property in the statement of Theorem 1, the theorem applies to both near and far fluid limit sequences.

PROPERTY 1. $\{(\mathbf{x}_j, a_j)\}$ is a sequence of initial condition x_j , threshold factor n_j , and scale a_j triples for which $a_j \rightarrow \infty$. Moreover for each j ,

$n_j > 0$, $a_j > 0$, and some closed, bounded $\mathcal{E} \in \bar{X}$,

$$\frac{n_j}{a_j} \leq c_1, \quad \text{and} \quad \left\| \frac{x_j}{a_j} \right\|_{\frac{n_j h}{a_j} \mathcal{E}} \leq c_2 \text{ for some } c_1 > 0, \text{ and } c_2 > 0.$$

3.1. *Convergence to a Fluid Limit along a Subsequence.* The proof of the following theorem parallels the proof of Theorem 4.1 of Dai [4]. However, the proof of our theorem differs in that we require some specialized treatment for our fluid limit construction and for the ingress discarding feature of the network. We state the theorem here and present the proof in the appendix.

THEOREM 1. *Suppose $\{(\mathbf{x}_j, a_j)\}$ is a sequence satisfying Property 1 (on page 17). Then for almost all ω there exists a subsequence $\{(\mathbf{x}_m, a_m)\} \subseteq \{(\mathbf{x}_j, a_j)\}$ for which*

$$\frac{\mathfrak{X}^{\mathbf{x}_m}(a_m t)}{a_m} \rightarrow \bar{\mathfrak{X}}(t) \quad \text{u.o.c.}$$

for some fluid model trajectory $\bar{\mathfrak{X}}(\cdot)$ with components

$$\bar{\mathfrak{X}}(\cdot) \triangleq [\bar{X}(\cdot); \bar{T}(\cdot); \bar{\Lambda}(\cdot); \bar{h}]$$

where, in turn, the process $\bar{X}(\cdot)$ has components

$$\bar{X}(\cdot) \triangleq [\bar{Q}(\cdot); \bar{U}(\cdot); \bar{V}(\cdot); \bar{H}(\cdot)]$$

where $\bar{H}(\cdot) \equiv 0$. The process $\bar{\mathfrak{X}}(\cdot)$ may depend upon ω and the choice of subsequence $\{(\mathbf{x}_m, a_m)\}$ but must satisfy the following properties for all $t \geq 0$:

$$(12) \quad \bar{U}_f(t) = (t - \bar{U}_f(0))^+, \quad \bar{V}_k(t) = (t - \bar{V}_k(0))^+,$$

$$(13) \quad \bar{T}_k(t) \text{ is nondecreasing and starts from zero,}$$

$$(14) \quad \bar{I}_i(t) := t - C_i \bar{T}(t) \text{ is nondecreasing,}$$

$$(15) \quad \bar{D}_k(t) := \mu_{s(k)} (\bar{T}_k(t) - \bar{V}_k(0))^+,$$

$$(16) \quad \bar{A}(t) := P^\top \bar{D}(t) + \bar{\Lambda}(t),$$

$$(17) \quad \bar{Q}(t) := \bar{Q}(0) + \bar{A}(t) - \bar{D}(t),$$

$$(18) \quad \bar{Q}(t) \geq 0,$$

$$(19) \quad \int_0^\infty (C \bar{Q}(t)) d\bar{I}(t) = 0,$$

where (12), (13), and (15) hold for each flow f and class k , while (14) holds for each station i . Assignments (14), (15), (16), and (17) define $\bar{I}(t)$,

$\bar{D}(t)$, $\bar{A}(t)$, and $\bar{Q}(t)$ respectively. Also, the following hold for each flow f for regular $t \geq 0$:

$$(20) \quad \dot{\bar{\Lambda}}_f(t) = 0 \quad \text{whenever } \bar{Q}_k(t) > \bar{h} \text{ for some } k \in \mathcal{C}(f),$$

$$(21) \quad \dot{\bar{\Lambda}}_f(t) = \alpha_f 1(t \geq \bar{U}_f(0)) \quad \text{whenever } \bar{Q}_k(t) < \bar{h} \text{ for all } k \in \mathcal{C}(f),$$

$$(22) \quad \dot{\bar{\Lambda}}_f(t) \leq \alpha_f.$$

Also, for station i and for any k, l such that $\{k, l\} \in C(i)$ the following properties are satisfied for all regular $t \geq 0$:

$$(23) \quad w_k^{-1} \dot{\bar{D}}_k(t) \geq w_l^{-1} \dot{\bar{D}}_l(t) \quad \text{whenever } Q_k(t) > 0,$$

$$(24) \quad w_k^{-1} \dot{\bar{D}}_k(t) = w_l^{-1} \dot{\bar{D}}_l(t) \quad \text{whenever } Q_k(t) > 0 \text{ and } Q_l(t) > 0.$$

See the appendix for the proof. Next we state precisely the definitions of a near fluid limit sequence and far fluid limit sequence that we discussed earlier in Section 1.1. After defining these sequences, we derive two corollaries to Theorem 1 that apply to each of these types of sequence.

DEFINITION 1 (Near Fluid Limit Sequence). $\{(\mathbf{x}_j, a_j)\}$ is a **near fluid limit sequence** with respect to a closed, bounded $h\mathcal{E} \in \bar{X}$ if $a_j = n_j$, $n_j \rightarrow \infty$, and

$$\left\| \frac{x_j}{n_j} \right\|_{h\mathcal{E}} = \frac{\|x_j\|_{n_j h\mathcal{E}}}{n_j} \leq \zeta.$$

for each j and for some $\zeta > 0$.

DEFINITION 2 (Far Fluid Limit Sequence). $\{(\mathbf{x}_j, a_j)\}$ is a **far fluid limit sequence** with respect to a closed, bounded $h\mathcal{E} \in \bar{X}$ if $a_j = n_j \left\| \frac{x_j}{n_j} \right\|_{h\mathcal{E}}$, $a_j \rightarrow \infty$ and

$$\left\| \frac{x_j}{n_j} \right\|_{h\mathcal{E}} = \frac{\|x_j\|_{n_j h\mathcal{E}}}{n_j} > \zeta.$$

for each j and for some $\zeta > 0$.

As was discussed earlier, the near fluid limit sequence is defined so that the sequence of scaled initial conditions remains a bounded distance away from the set $h\mathcal{E}$ while the far fluid limit is defined so that the sequence of scaled initial conditions is bounded away from the set $h\mathcal{E}$.

COROLLARY 1. *Suppose that $\{(\mathbf{x}_j, a_j)\}$ is a near fluid limit sequence with respect to a closed, bounded $h\mathcal{E} \in \bar{X}$. Then for almost all ω there exists a subsequence $\{(\mathbf{x}_m, a_m)\} \subseteq \{(\mathbf{x}_j, a_j)\}$ for which*

$$\frac{\mathfrak{X}^{\mathbf{x}_m}(a_m t)}{a_m} \rightarrow \bar{\mathfrak{X}}(t) \quad u.o.c.$$

where $\mathfrak{X}(\cdot)$ satisfies fluid model equations (12) - (24). Moreover

$$\bar{h} = h \text{ and } \|\bar{X}(0)\|_{\bar{h}\mathcal{E}} \leq \zeta.$$

PROOF. The discarding thresholds before scaling are $n_j h$, and thus after scaling they are $n_j h/a_j = h$ for each j . Thus $\bar{h} = h$. Also $a_j \rightarrow \infty$ and $\|x_j/a_j\|_{h\mathcal{E}} \leq \zeta$ and thus the sequence $\{(\mathbf{x}_j, a_j)\}$ satisfies Property 1. By Theorem 1 there exists a subsequence $\{(\mathbf{x}_m, a_m)\}$ such that $\mathfrak{X}^{\mathbf{x}_m}(a_m t)/a_m$ converges u.o.c. to a fluid trajectory satisfying (12) - (24). By Theorem 1, the subsequence x_m/a_m converges to an initial state of the fluid trajectory $\bar{X}(0)$. Since $\|x_m/a_m\|_{h\mathcal{E}} \leq \zeta$, it must be that $\|\bar{X}(0)\|_{\bar{h}\mathcal{E}} \leq \zeta$. \square

COROLLARY 2. *Suppose that $\{(\mathbf{x}_j, a_j)\}$ is a far fluid limit sequence with respect to a closed, bounded $h\mathcal{E} \in \bar{X}$. Then for almost all ω there exists a subsequence $\{(\mathbf{x}_m, a_m)\} \subseteq \{(\mathbf{x}_j, a_j)\}$ for which*

$$\frac{\mathfrak{X}^{\mathbf{x}_m}(a_m t)}{a_m} \rightarrow \bar{\mathfrak{X}}(t) \quad u.o.c.$$

where $\mathfrak{X}(\cdot)$ satisfies fluid model equations (12) - (24). Moreover

$$\bar{h} \in [0, h/\zeta] \text{ and } \|\bar{X}(0)\|_{\bar{h}\mathcal{E}} = 1.$$

PROOF. Note that

$$\left\| \frac{x_j}{a_j} \right\|_{\frac{n_j h}{a_j} \mathcal{E}} = \frac{n_j}{a_j} \left\| \frac{x_j}{n_j} \right\|_{h\mathcal{E}} = 1$$

for each j . This combined with the fact that $a_j \rightarrow \infty$ implies that $\{(\mathbf{x}_j, a_j)\}$ satisfies Property 1. By Theorem 1 there exists a subsequence $\{(\mathbf{x}_m, a_m)\}$ such that $\mathfrak{X}^{\mathbf{x}_m}(a_m t)/a_m$ converges u.o.c. to a fluid trajectory satisfying (12) - (24). The above equation also implies that $\|\bar{X}(0)\|_{\bar{h}\mathcal{E}} = 1$. The subsequence of scaled thresholds satisfies $n_m h/a_m = h/\|x_m/n_m\|_{h\mathcal{E}} < h/\zeta$. By Theorem 1 the subsequence $n_m h/a_m$ converges, and the convergence must be to a number in the range $[0, h/\zeta]$ because of the preceding inequality relation. \square

3.2. Convergence along Subsequences to Convergence along Sequences.

In the previous section, we showed that for both near and far fluid limit sequences, we can extract a sample path dependent subsequence that converges to a fluid model trajectory. The objective of this section is to use this subsequence result to show convergence of a functional of the original sequence. In particular, we show in Lemma 1 that if a functional \mathfrak{F} of any fluid model trajectory goes to zero in a time not more than a constant times the scaled initial condition's distance from $\bar{h}\mathcal{E}$, then the value of that functional applied to the fluid limit sequence of trajectories converges almost surely. In later sections, we will invoke Lemma 1 choosing \mathfrak{F} to extract the service rates from the fluid model, and later choosing \mathfrak{F} to extract the distance from a set $h\mathcal{E}$. Lemma 1 is a generalization of an argument used by Dai in the proof of Theorem 4.2 of [4].

LEMMA 1. *Suppose that \mathfrak{F} is a functional that maps $\mathbb{R}^r \times \mathbb{R}^+$ into $\mathbb{R}^s \times \mathbb{R}^+$ where r is the dimension of $\mathfrak{X}(\cdot)$ and s is arbitrary. Also suppose that \mathfrak{F} is continuous on the topology of uniform convergence on compact sets. If the following is true:*

- *The fluid model equations (12) - (24) are such that for any trajectory $\bar{\mathfrak{X}}(\cdot)$ and $\bar{h} \geq 0$ that satisfies them, there exists some closed bounded $\mathcal{E} \in \bar{\mathfrak{X}}$ for which*

$$(25) \quad \mathfrak{F} \circ [\bar{\mathfrak{X}}(\cdot)](t) \equiv 0 \quad \forall t \geq t_0 \left\| \bar{X}(0) \right\|_{\bar{h}\mathcal{E}}.$$

Then, for any sequence $\{(\mathbf{x}_j, a_j)\}$ satisfying Property 1 where the relation $\|x_j/a_j\|_{\frac{n_j h}{a_j}\mathcal{E}} \leq c$ of Property 1 is satisfied with constant $c > 0$,

$$(26) \quad \left| \mathfrak{F} \circ \left[\frac{1}{a_j} \mathfrak{X}^{\mathbf{x}_j}(a_j \cdot) \right] (t) \right| \rightarrow 0 \quad a.s.$$

for each $t \geq ct_0$.

PROOF. By Theorem 1, for almost all sample paths ω , and for any subsequence $\{(\mathbf{x}_m, a_m)\} \subseteq \{(\mathbf{x}_j, a_j)\}$ there is a sample-path-dependent further-subsequence $\{(\mathbf{x}_{r(\omega)}, a_{r(\omega)})\} \subseteq \{(\mathbf{x}_m, a_m)\}$ for which

$$\frac{\mathfrak{X}^{\mathbf{x}_{r(\omega)}}(a_{r(\omega)}t, \omega)}{a_{r(\omega)}} \rightarrow \bar{\mathfrak{X}}(t, \omega) \quad \text{u.o.c.}$$

where $\bar{\mathfrak{X}}(t, \omega)$ satisfies (12) - (24) as well as $\|\bar{X}(0)\|_{\bar{h}\mathcal{E}} \leq c$ since each x_j/a_j has a distance from $\frac{n_j h}{a_j}\mathcal{E}$ that is no more than c by the lemma's assumption.

The notation $r(\omega)$ and $\bar{\mathfrak{X}}(t, \omega)$ emphasize that the further-subsequence and fluid trajectory depend on ω . Now fix an ω for which subsequences have convergent further subsequences as described. For the next few steps we suppress the ω arguments to simplify notation. Because \mathfrak{F} is assumed to be continuous on the topology of uniform convergence on compact sets, we have

$$\mathfrak{F} \circ \left[\frac{\mathfrak{x}^{\mathbf{x}_r}(a_r \cdot)}{a_r} \right] (t) \rightarrow \mathfrak{F} \circ [\bar{\mathfrak{X}}(\cdot)](t) \quad \text{u.o.c.}$$

Consequently

$$\left| \mathfrak{F} \circ \left[\frac{\mathfrak{x}^{\mathbf{x}_r}(a_r \cdot)}{a_r} \right] (t) \right| \rightarrow 0$$

for each $t \geq ct_0$. So for this fixed ω , any subsequence $\{(\mathbf{x}_m, a_m)\} \subseteq \{(\mathbf{x}_j, a_j)\}$ has a further subsequence $\{(\mathbf{x}_{r(\omega)}, a_{r(\omega)})\} \subseteq \{(\mathbf{x}_m, a_m)\}$ for which the above holds. Therefore the original sequence $\{(\mathbf{x}_j, a_j)\}$ converges for this fixed ω . The same argument can be used to conclude that this holds for almost all ω . Thus, we have (26). \square

3.3. Convergence to Fluid Model Rates on a Compact Time Interval.

The objective of this section is to use Lemma 1 to conclude that the rates of the stochastic system are close to those of the fluid model over a finite time interval. It will remain to show that the rates are close over the long-term.

THEOREM 2. *Suppose there exists $t_0 > 0$, a closed, bounded $\mathcal{E} \in \bar{\mathcal{X}}$, and rate vector $R \in \mathbb{R}_+^K$ such that*

$$(27) \quad M^{-1}\dot{T}(t) \equiv R \quad \forall t \geq t_0 \left\| \bar{X}(0) \right\|_{\bar{h}\mathcal{E}}$$

for any fluid model trajectory $\bar{\mathfrak{X}}(\cdot)$ and $\bar{h} > 0$ that satisfies (12) - (24). Then for any positive $\gamma < 1$ and $\zeta < 1$ there exists $L_1(\zeta, \gamma)$ such that for all $n \geq L_1$,

$$(28) \quad \inf_{\left\| \frac{x}{n} \right\|_{h\mathcal{E}} \leq \zeta} \mathbb{E} [M^{-1}T^{\mathbf{x}}(nt_0)] \geq R(1 - \zeta)(1 - \gamma)nt_0.$$

PROOF. Let $\{(\mathbf{x}_j, a_j)\}$ be a near fluid limit sequence: a sequence of threshold scale and initial condition pairs satisfying $a_j = n_j \rightarrow \infty$ and $\|x_j/n_j\|_{h\mathcal{E}} \leq \zeta$. We invoke Lemma 1 by picking \mathfrak{F} so that

$$\mathfrak{F} \circ [\bar{\mathfrak{X}}(\cdot)](t) := \bar{T}(\zeta^{-1}t) - \bar{T}(t) - MR(\zeta^{-1} - 1)t.$$

\mathfrak{F} is easily seen to be continuous on the topology of uniform convergence on compact sets. Also note that $\mathfrak{F} \circ [\bar{\mathfrak{X}}(\cdot)](t) = 0$ for all $t \geq t_0 \left\| \bar{X}(0) \right\|_{\bar{h}\mathcal{E}}$ by (27).

By Lemma 1,

$$\lim_{j \rightarrow \infty} \left| \frac{T^{\mathbf{x}_j}(n_j t_0) - T^{\mathbf{x}_j}(\zeta n_j t_0)}{n_j(1 - \zeta)t_0} - MR \right| = 0 \quad \text{a.s.},$$

where we have used the fact that $\|x_j/n_j\|_{h\mathcal{E}} \leq \zeta$ to choose the c of Lemma 1 to be ζ and selected $t = \zeta t_0$. The left hand side of the above identity is bounded from above by a constant for all j , and thus by the dominated convergence theorem [10],

$$(29) \quad \lim_{j \rightarrow \infty} \mathbb{E} \left| \frac{T^{\mathbf{x}_j}(n_j t_0) - T^{\mathbf{x}_j}(\zeta n_j t_0)}{n_j(1 - \zeta)t_0} - MR \right| = 0.$$

Also note (29) holds for any sequence $\{(\mathbf{x}_j, a_j)\}$ with $n_j = a_j \rightarrow \infty$ and $\|x_j/n_j\|_{h\mathcal{E}} \leq \zeta$, because these were the only restrictions for our initial choice of sequence.

Now pick a positive constant $\gamma < 1$. Observe that there exists a constant $L_1(\gamma, \zeta)$ such that whenever $n > L_1$,

$$\inf_{\|x/n\|_{h\mathcal{E}} \leq \zeta} \frac{\mathbb{E}[T^{\mathbf{x}}(nt_0) - T^{\mathbf{x}}(n\zeta t_0)]}{n(1 - \zeta)t_0} \geq MR(1 - \gamma)$$

for if otherwise we could construct a sequence $\{(\mathbf{x}_j, a_j)\}$ that violates (29). By the monotonicity of $T^{\mathbf{x}_j}(\cdot)$, we have (28). \square

3.4. Stochastic System Attracted to $h\mathcal{E}$. The objective of this section is to show that the scaled state of the stochastic system is attracted to $h\mathcal{E}$. In particular we show that the scaled state's expected distance from $h\mathcal{E}$ declines geometrically (roughly) for starting scaled states outside a neighborhood of $h\mathcal{E}$. Since the proof technique is similar that of Theorem 3.1 of Dai [4] we choose to provide the proof in the appendix.

THEOREM 3. *Suppose that there exists $t_0 > 0$ and a closed, bounded $\mathcal{E} \in \bar{\mathcal{X}}$ such that*

$$(30) \quad \|\bar{X}(t)\|_{\bar{h}\mathcal{E}} \equiv 0 \quad \forall t \geq t_0 \|\bar{X}(0)\|_{\bar{h}\mathcal{E}}$$

for any fluid model trajectory $\tilde{\mathbf{X}}(\cdot)$ and $\bar{h} \geq 0$ that satisfies (12) - (24). Then the following conclusions are true:

- i) For any $\zeta > 0$, and any positive $\delta < 1$ there exists $L_2(\zeta, \delta)$ such that for all $n \geq \zeta^{-1}L_2$ and all $x : \|x/n\|_{h\mathcal{E}} > \zeta$,

$$\mathbb{E} \left\| \frac{1}{n} X^{\mathbf{x}} \left(nt_0 \left\| \frac{x}{n} \right\|_{h\mathcal{E}} \right) \right\|_{h\mathcal{E}} \leq \delta \left\| \frac{x}{n} \right\|_{h\mathcal{E}}.$$

ii) For any $\zeta > 0$, and any $b > 0$ there exists $L_3(\zeta, b)$ such that for all $n \geq L_3$ and all $x : \|x/n\|_{h\mathcal{E}} \leq \zeta$,

$$\mathbb{E} \left\| \frac{1}{n} X^{\mathbf{x}}(nt_0) \right\|_{h\mathcal{E}} \leq b.$$

See the appendix for the proof.

The objective of the next lemma is to show that the results of Theorem 3 imply that the expected return time of the scaled state to the ζ ball around $h\mathcal{E}$ is small. The proof of Lemma 2 is adapted from the proof of Theorem 2.1(ii) of [19], which was for a discrete time Markov chain. Since the lemma is an adaptation of a previous result, we provide the proof in the appendix.

LEMMA 2. Suppose (1) and (2) are satisfied and for some $n > 0$, $h \geq 0$, and a closed, bounded $\mathcal{E} \in \bar{\mathcal{X}}$ we have

$$(31) \quad \mathbb{E} \left\| \frac{1}{n} X^{\mathbf{x}} \left(nt_0 \left\| \frac{x}{n} \right\|_{h\mathcal{E}} \right) \right\|_{h\mathcal{E}} \leq \delta \left\| \frac{x}{n} \right\|_{h\mathcal{E}} \quad \forall x : \|x/n\|_{h\mathcal{E}} > \zeta,$$

$$(32) \quad \mathbb{E} \left\| \frac{1}{n} X^{\mathbf{x}}(nt_0) \right\|_{h\mathcal{E}} \leq b \quad \forall x : \|x/n\|_{h\mathcal{E}} \leq \zeta.$$

Then X is positive Harris recurrent and

$$(33) \quad \sup_{x \in B} \mathbb{E}_x[\tau_B^n(nt_0)] \leq nt_0 \left[1 + \frac{\zeta + b}{1 - \delta} \right]$$

where $B \triangleq \{x : \|x/n\|_{h\mathcal{E}} \leq \zeta\}$ and $\tau_B^n(nt_0)$ is defined by

$$(34) \quad \tau_B^n(nt_0) \triangleq \inf\{t \geq nt_0 : X^n(t) \in B\}.$$

See the appendix for the proof.

3.5. *Convergence of Long Term Rates.* The objective of this section is to tie together all of the preceding results to conclude in Theorem 4 that the long-term rates of the stochastic system are close to the fluid rates for large enough n . First we pick n large enough so that the conclusions of Theorems 2, 3, and Lemma 2 apply. Theorem 2 says that the stochastic system's rates are close to the fluid rates for the first nt_0 seconds after having started with a scaled initial condition x/n in a ζ -neighborhood of $h\mathcal{E}$. To make a conclusion about the long-term, we need to show that stochastic system spends relatively little time away from the neighborhood in which Theorem 2

applies. Lemma 2 tells us that the expected first return time of X/n to a ζ -neighborhood of $h\mathcal{E}$ that happens after nt_0 seconds is no more than a constant times nt_0 . Moreover, this constant can be made arbitrarily small by picking n larger. This argument is illustrated by Figure 4. To formalize the argument we construct a sequence of stopping times that occur on the first visit of X/n to the ζ -neighborhood of $h\mathcal{E}$ that occurs at least nt_0 seconds after the last stopping time. We define random vectors ρ_i that track the cumulative service, divided by average service times, between stopping times and relate these to the desired rate vector R using Theorem 2. We use ergodicity to argue that the long term average rates exist, and that this long term limit must equal the product of the expected value of ρ_i times the lim inf of $t/N(t)$ the inverse of the arrival rate of stopping times. Due to Lemma 2, this later quantity has an upper bound of nt_0 times a constant that can be made small.

THEOREM 4. *Suppose for some $t_0 > 0$ and some closed, bounded $\mathcal{E} \in \bar{\mathcal{X}}$ both of the following are true:*

- For any fluid model trajectory $\bar{\mathfrak{X}}(\cdot)$ and $\bar{h} \geq 0$ that satisfies (12) - (24),

$$(35) \quad \|\bar{X}(t)\|_{\bar{h}\mathcal{E}} \equiv 0 \quad \forall t \geq t_0 \|\bar{X}(0)\|_{\bar{h}\mathcal{E}}.$$

- For any fluid model trajectory $\bar{\mathfrak{X}}(\cdot)$ and $\bar{h} > 0$ that satisfies (12) - (24),

$$(36) \quad M^{-1}\dot{\bar{T}}(t) \equiv R \quad \forall t \geq t_0 \|\bar{X}(0)\|_{\bar{h}\mathcal{E}}$$

where $R \in \mathbb{R}_+^K$.

Then for any $\epsilon > 0$, there exists a $n_c > 0$ such that for all $n \geq n_c$

$$\lim_{t \rightarrow \infty} \frac{D^{\mathbf{x}}(t)}{t} \geq (1 - \epsilon)R \quad a.s.$$

PROOF. We observe that equations (36) and (35) are the necessary conditions to apply Theorems 2 and 3 respectively. Therefore, we may arbitrarily pick the constants ζ , δ , and b of Theorem 3 and the constants ζ and γ of Theorem 2 (using the same ζ value in Theorems 2 as we use when we apply Theorem 3), and then fix an n satisfying

$$(37) \quad n > \max[L_1(\zeta, \gamma), \zeta^{-1}L_2(\zeta, \delta), L_3(\zeta, b)]$$

so that the conclusions of both Theorems 3 and 2 hold.

In addition, conclusions (i) and (ii) of Theorem 3 allow us to invoke Lemma 2 to conclude (33) where $\tau_B^n(nt_0)$ is defined by (34). Because the

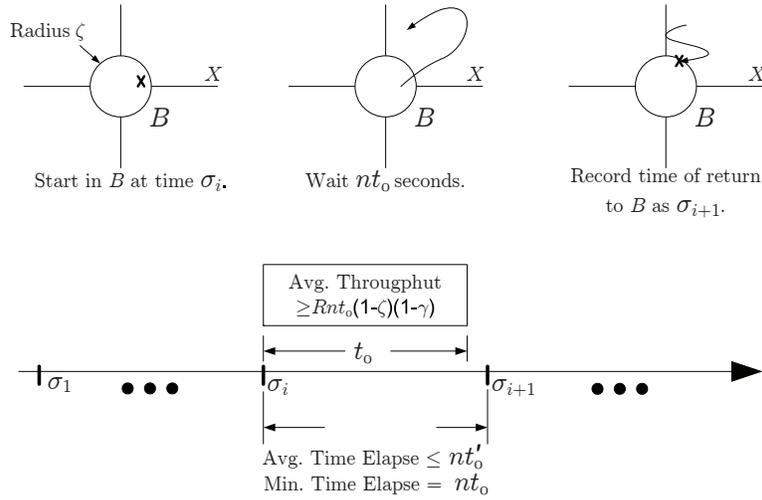


FIG 4. The top half of the figure illustrates the definition of the stopping times $\sigma_i, \sigma_{i+1}, \dots$. The bottom half illustrates the intuition behind the proof of Theorem 4 by plotting the stopping times on a time line, and showing the bound on expected throughput between such stopping times.

constants ζ, b, δ can be chosen arbitrarily, equations (33) and (37) imply that the ratio of the expected first hitting time of B (nt_0 seconds after having started in B) to nt_0 can be made to be close to 1 by choosing n large enough. We collect some of the constants in (33) in the term t'_0 defined by

$$(38) \quad nt'_0 = nt_0 \left[1 + \frac{\zeta + b}{1 - \delta} \right].$$

We have also chosen n large enough so that the following conclusion from Theorem 2 holds,

$$(39) \quad \inf_{\|x/n\|_{h\mathcal{E}} \leq \zeta} \mathbb{E}[T^x(nt_0)] \geq MR(1 - \zeta)(1 - \gamma)nt_0.$$

Define the stopping times

$$(40) \quad \sigma_0 = 0, \quad \sigma_{i+1} = \inf\{t \geq nt_0 + \sigma_i : X(t) \in B\}, \quad \forall i \geq 0.$$

Figure 4 illustrates how these stopping times are defined. Note that for any initial condition $x \in X$ (the state space of X^n) and index $i \geq 1$,

$$(41) \quad \mathbb{E}_x[\sigma_{i+1} - \sigma_i] \leq \sup_{\tilde{x} \in B} \mathbb{E}_{\tilde{x}}[\tau_B^n(nt_0)] \leq nt'_0.$$

This follows from the fact that $X^{\mathbf{x}}(\sigma_i) \in B$, the strong Markov property, the stopping time definitions (34) & (40), and expressions (33) & (38). Also, X is positive Harris recurrent by Lemma 2 and therefore, $E_x[\sigma_1] < \infty$ for any $x \in X$. We define a counting process $N(t)$ for the stopping times σ_i as $N(t) = \inf\{i : \sigma_i \leq t\}$. Because X is positive Harris recurrent, $\sigma_i < \infty$ almost surely, and therefore $N(t) \rightarrow \infty$ a.s. We now turn to bounding the expected ‘‘arrival’’ rate of the stopping times σ_i . By (41) for each i ,

$$(42) \quad \frac{E_x[\sigma_i]}{i} = \frac{\sum_{j=1}^{i-1} E_x[\sigma_{j+1} - \sigma_j] + E_x\sigma_1}{i} \leq nt'_0(1 - 1/i) + \frac{E_x\sigma_1}{i}$$

Additionally, along any sample path

$$\frac{t}{N(t)} \leq \frac{\sigma_{N(t)+1}}{N(t)+1} \frac{N(t)+1}{N(t)}.$$

Thus by taking $\liminf_{t \rightarrow \infty} E_x(\cdot)$ of both sides, and using (42) we have $\liminf_{t \rightarrow \infty} E_x \left[\frac{t}{N(t)} \right] \leq nt'_0$. Moreover, by Fatou’s Lemma

$$(43) \quad E_x \left[\liminf_{t \rightarrow \infty} \frac{t}{N(t)} \right] \leq \liminf_{t \rightarrow \infty} E_x \left[\frac{t}{N(t)} \right] \leq nt'_0.$$

We define the random vectors $\rho_i = M^{-1}[T^n(\sigma_i + \sigma_{i+1}) - T^n(\sigma_i)]$ to track the service between stopping times σ_i . Note that for $i \geq 1$ and each $x \in X$,

$$(44) \quad E_x[\rho_i] \geq \inf_{\bar{x} \in B} E_{\bar{x}}[M^{-1}T^{\bar{x},n}(nt_0)] \geq Rnt_0(1 - \zeta)(1 - \gamma).$$

This follows from the fact that $X^{\mathbf{x}}(\sigma_i) \in B$, the strong Markov property, the definition of σ_i (40), the definition of ρ_i , and relation (39). Figure 4 illustrates the fact that the throughput between stopping times σ_i and σ_{i+1} is lower-bounded according to relation (44).

By [6] the following ergodic property holds for every measurable f on X with $\pi(|f|) < \infty$,

$$\lim_{t \rightarrow \infty} \frac{1}{t} \int_0^t f(X^n(s)) ds = \pi(f) \text{ P}_x\text{-a.s. for each } x \in X$$

where π is the unique invariant distribution of X^n . Assigning the function $f(x) \triangleq M^{-1}\dot{T}^{\mathbf{x}}(0)$ to be the instantaneous service rates when the process is in state x , (Recall that we assumed the service rates are a function of the state in Section 2.) we have,

$$(45) \quad \lim_{t \rightarrow \infty} \frac{1}{t} \int_0^t f(X^{\mathbf{x}}(s)) ds = \lim_{t \rightarrow \infty} \frac{1}{t} M^{-1}\tilde{T}^{\mathbf{x}}(t) = \mathcal{R} \text{ a.s.}$$

for some constant vector \mathcal{R} .

Consider the random variable $\mathcal{N} \triangleq \liminf_{t \rightarrow \infty} \frac{t}{N(t)}$. The random variable \mathcal{N} is a P_π invariant random variable, and therefore is a constant. Moreover by (43), $\mathcal{N} \leq nt'_0$. A more detailed explanation of this argument is provided in [22].

We observe that for any sample path the following inequalities hold,

$$(46) \quad \frac{t}{N(t)} \frac{M^{-1}T^{\mathbf{x}}(t)}{t} \leq \frac{\sum_{j=0}^{N(t)} \rho_j}{N(t)} \leq \frac{t}{N(t)} \frac{\sigma_{N(t)+1}}{t} \frac{M^{-1}T^{\mathbf{x}}(\sigma_{N(t)+1})}{\sigma_{N(t)+1}}.$$

Taking the $\liminf_{t \rightarrow \infty}$ of both sides, and using (45) we have that

$$(47) \quad \liminf_{t \rightarrow \infty} \frac{\sum_{j=0}^{N(t)} \rho_j}{N(t)} = \mathcal{N}\mathcal{R} \text{ a.s.}$$

We note that $T^{\mathbf{x}}(\sigma_{N(t)+1})/\sigma_{N(t)+1} \leq \mathbf{1}$ where $\mathbf{1}$ is a column vector of 1's of appropriate dimension. This fact combined with (46) yields that for each $i > 0$,

$$\inf_{k \geq i} \frac{\sum_{j=1}^k \rho_j}{i} \leq \liminf_{t \rightarrow \infty} \frac{t}{N(t)} M^{-1}\mathbf{1} \leq nt'_0 M^{-1}\mathbf{1}.$$

Thus the random variables $\left\{ \inf_{k \geq i} i^{-1} \sum_{j=1}^k \rho_j : i > 0 \right\}$ are dominated by a constant. Consequently, $\liminf_{i \rightarrow \infty} \mathbb{E} \left[\sum_{j=1}^i \rho_j / i \right] = \mathcal{N}\mathcal{R}$ by the dominated convergence theorem. Also for each $i > 0$, $\mathbb{E} \left[\sum_{j=1}^i \rho_j / i \right] \geq Rnt_0(1-\gamma)(1-\zeta)$ by (44). Thus, $\mathcal{N}\mathcal{R} \geq Rnt_0(1-\gamma)(1-\zeta)$. Substituting (43) we have that $\mathcal{R} \geq (1-\gamma)(1-\zeta) \frac{t_0}{t'_0} R$. This implies

$$\lim_{t \rightarrow \infty} \frac{1}{t} M^{-1}T^{\mathbf{x}}(t) \geq \frac{(1-\gamma)(1-\zeta)}{1 + \frac{\zeta+b}{1-\delta}} R \text{ a.s.}$$

Recall γ , ζ , b , and δ may be chosen arbitrarily small, so long as n is chosen large enough according to (37). Thus, for any $\epsilon > 0$ there exists an n such that

$$\lim_{t \rightarrow \infty} \frac{1}{t} M^{-1}T^{\mathbf{x}}(t) \geq (1-\epsilon)R \text{ a.s.}$$

By the strong law of large numbers for renewal processes [10], $\frac{1}{t} S_k^{\mathbf{x}}(t) \rightarrow m_k$ a.s. Thus by (9), $\lim_{t \rightarrow \infty} \frac{1}{t} D^{\mathbf{x}}(t) \geq (1-\epsilon)R$ a.s. □

4. Analysis of Switch Example. In this section we apply the results of the preceding section to the example introduced in Section 1.3. Recall that this example resembles a 2-input 2-output switch and has 3 flows and is illustrated by Figure 3. As we discussed in Section 1.3, the max-min fair share rate allocation would be that all three flows achieve rates of 0.5, so we set $R = [0.5, 0.5, 0.5]^T$ to be the vector of desired rates.

To fit the framework we have developed, we must show that the fluid model with thresholds \bar{h} is drawn to a set $\bar{h}\mathcal{E}$, and that the fluid model rates while in $\bar{h}\mathcal{E}$ are R . Intuition suggests that the dynamics of the fluid model should evolve in the following way:

- One of the queues flow 2 passes through (either queue 2 or 7) reaches threshold and “chatters” there. The other queue can be anywhere at or below its threshold. By “chatters” we mean that it alternately goes a tiny amount above and below. However if the differential inclusions of the fluid model are such that: i) the queue grows whenever below threshold, ii) shrinks when above, then a fluid model trajectory would go to threshold and stay there.
- Queue 1 fills to threshold, “chatters” there, limiting flow 1’s ultimate rate.
- Queue 7 fills to threshold, “chatters” there, limiting flow 3’s ultimate rate.
- Other queues are not “bottlenecks” and should empty.

This above intuition suggests that the fluid model is drawn to the set $\bar{h}\tilde{\mathcal{E}}$ where $\tilde{\mathcal{E}}$ is given by

$$\tilde{\mathcal{E}} \triangleq \left\{ \bar{X} : \begin{array}{l} \bar{Q}_1 = \bar{Q}_8 = 1, \\ \bar{Q}_3 = \bar{Q}_4 = \bar{Q}_5 = \bar{Q}_6 = 0, \\ (\bar{Q}_2, \bar{Q}_7) \in \{[0, 1] \times 1\} \cup \{1 \times [0, 1]\}, \\ \bar{U} = 0, \bar{V} = 0, \bar{H} = 0 \end{array} \right\}.$$

As it will turn out, the most critical part of the analysis of this example’s fluid model is to show that the queues flow 2 passes through, queues 2 and 7, go to values in $\bar{h}\tilde{\mathcal{E}}$ in a time not more than a constant times their initial values. Intuition suggests that after a “settling down” period flow 1’s rate through queue 1, as well as flow 3’s rate through queue 8, settles to 0.5. After flow 1 and flow 3’s rates settle, the dynamics of $(\bar{Q}_2(t), \bar{Q}_7(t))$, the queues of flow 2, follow the relations outlined by Table 1 and illustrated by Figure 5. The entries of Table 1 are easily derived by using the observations that:

- The arrival rate to queue 2 is 0.6 when queue 2 and queue 7 are below

threshold while the arrival rate to queue 2 is 0 when one of these queues is above threshold.

- The departure rate from either queue 2 or queue 7 is 0.5 whenever the queue is nonempty or has sufficient arrivals to maintain this departure rate. (This relies on our assumption that the flow rates through queues 1 and 8 have “settled down” to 0.5).

Figure 5 is a vector flow diagram, showing the dependence of $(\dot{\bar{Q}}_2(\cdot), \dot{\bar{Q}}_7(\cdot))$ on $(\bar{Q}_2(\cdot), \bar{Q}_7(\cdot))$. It is evident from the diagram that the time to reach the set

$$\{[0, \bar{h}] \times \bar{h}\} \cup \{\bar{h} \times [0, \bar{h}]\},$$

which is the projection of $\bar{h}\tilde{\mathcal{E}}$ on to the subspace on which $(\bar{Q}_2(\cdot), \bar{Q}_7(\cdot))$ takes values, is not always less than or equal to a constant times the initial condition’s distance from this set. Consider an initial condition of $(\frac{\bar{h}}{2}, \bar{h} + \epsilon)$. This initial condition is only a distance of ϵ from $\tilde{\mathcal{E}}$, but the time it takes to reach the set $\tilde{\mathcal{E}}$ is $\bar{h} + \frac{1}{2}\epsilon$. (Note that we will use the L^1 norm throughout this section.) This is the same phenomenon we observed in the example in the introduction of the paper. There, as here, we can fix the problem by slightly enlarging the set $\tilde{\mathcal{E}}$ to a new set \mathcal{E} so that the set is reached in a time not more than a constant times the initial condition’s starting distance from the set. To this end, we define \mathcal{E} according to

$$\mathcal{E} \triangleq \left\{ \bar{X} : \left. \begin{array}{l} \bar{Q}_1 = \bar{Q}_8 = 1, \\ \bar{Q}_3 = \bar{Q}_4 = \bar{Q}_5 = \bar{Q}_6 = 0, \\ (\bar{Q}_2, \bar{Q}_7) \in \\ \left\{ (\chi, \psi) : \begin{array}{l} \chi \in [0, 1] \\ \psi \in [1 - a\chi, 1 + a(1 - \chi)] \end{array} \right\} \\ \cup \{1 \times [0, 1]\}, \\ \bar{U} = 0, \bar{V} = 0, \bar{H} = 0 \end{array} \right\} \right\}.$$

Here a is an arbitrary positive constant that should be less than 1. The projection of this set onto the subspace spanned by (Q_2, Q_7) is shown as the shaded area in Figure 5. With this definition, one can show that the set $\bar{h}\mathcal{E}$ is reached in a time not more than a constant times the initial distance from $\bar{h}\mathcal{E}$. The time to reach $\bar{h}\mathcal{E}$, along with the maximum ratio of the time to reach $\bar{h}\mathcal{E}$ divided by initial distance to $\bar{h}\mathcal{E}$ are shown in Table 1.

We are now ready to formalize the intuition we have outlined in the preceding paragraphs. We begin by stating a lemma that the system settles down so that the behavior flow 2’s queues are as described by Table 1 after a time τ_{sd} (mnemonic for “settle down”) that is in proportion to the initial condition.

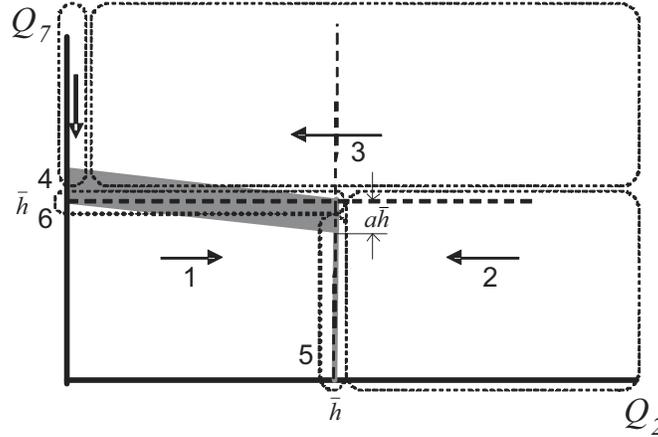


FIG 5. The evolution of $(\bar{Q}_2(t), \bar{Q}_7(t))$. The shaded area indicates the set $\bar{h}\mathcal{E}$.

LEMMA 3. *There exists a time τ_{sd} proportional to the initial condition as described by the relation*

$$\tau_{sd} = t_{01} \|\bar{X}(0)\|_{\bar{h}\mathcal{E}}$$

for some positive t_{01} such that for all regular points $t \geq \tau_{sd}$:

- The value of $(\dot{\bar{Q}}_2(t), \dot{\bar{Q}}_7(t))$ is determined by the value of $(\bar{Q}_2(t), \bar{Q}_7(t))$ as specified by Table 1.
- $\bar{Q}_3(t) = \bar{Q}_5(t) = \bar{Q}_6(t) = 0$, and $\bar{Q}_8(t) = \bar{h}$.
- The time to reach the set $\bar{h}\mathcal{E}$, as well as the maximum ratio between this time and the distance of $(\bar{Q}_2(\tau_{sd}), \bar{Q}_7(\tau_{sd}))$ from $\bar{h}\mathcal{E}|_{\bar{Q}_2, \bar{Q}_7}$ in any of the Regions 1 through 4 is as specified in Table 1. (Here $\bar{h}\mathcal{E}|_{\bar{Q}_2, \bar{Q}_7}$ denotes the projection of the set $\bar{h}\mathcal{E}$ onto the space on which (\bar{Q}_2, \bar{Q}_7) takes values.)

Lemma 3 is proved by using the relations (12) - (24) that describe the evolution of a fluid model trajectory. The proof is straightforward but slightly lengthy because it requires analysis for each entry in Table 1. We therefore omit this proof.

We now state and prove the principal result of this section.

THEOREM 5. *For any $\epsilon > 0$, there exists an $n_c > 0$ such that if the discarding thresholds of the stochastic system in Example 2 are set to nh , $n \geq n_c$, then*

$$\lim_{t \rightarrow \infty} \frac{D(t)}{t} \geq (1 - \epsilon) \frac{1}{2} \mathbf{1} \text{ a.s.}$$

	\bar{Q}_2	\bar{Q}_7	$\dot{\bar{Q}}_2$	$\dot{\bar{Q}}_7$	Time to $\bar{h}\mathcal{E}$	$\frac{\text{Time to } \bar{h}\mathcal{E}}{\ \bar{X}\ _{\bar{h}\mathcal{E}}}$
1	$[0, \bar{h})$	$[0, \bar{h})$	0.1	0	if $ \bar{Q}_7 - \bar{h} < a\bar{h}$ then $\frac{10}{a} \bar{Q}_7 - \bar{h} $ if $ \bar{Q}_7 - \bar{h} \geq a\bar{h}$ then $10 \bar{Q}_2 - \bar{h} $	$\frac{10}{a}$
2	(\bar{h}, ∞)	$[0, \bar{h})$	-0.5	0	$2 \bar{Q}_2 - \bar{h} $	2
3	$(0, \infty)$	(\bar{h}, ∞)	-0.5	0	if $ \bar{Q}_7 - \bar{h} < a\bar{h}$ then $\frac{2}{a} \bar{Q}_7 - \bar{h} $ if $ \bar{Q}_7 - \bar{h} \geq a\bar{h}$ then $2 \bar{Q}_2 - \bar{h} + 2 \bar{Q}_7 - \bar{h} - a\bar{h} $	$\frac{2}{a}$
4	0	(\bar{h}, ∞)	0	-0.5	$2 \bar{Q}_7 - \bar{h} - a $	2
5	\bar{h}	$[0, \bar{h}]$	$[-0.5, 0.1]$	$[-0.5, 0]$	0	N/A
6	$[0, \bar{h}]$	\bar{h}	$[-0.5, 0.1]$	$[-0.5, 0]$	0	N/A

TABLE 1

Dynamics of $(\bar{Q}_2(t), \bar{Q}_7(t))$, after flows 1 and 3 settle to their ultimate rates of 0.5. The rows numbers correspond to the regions labeled in the phase portrait diagram of Figure 5.

where $\mathbf{1}$ is a vector of ones of dimension K .

PROOF. By Lemma 3 the dynamics of the state variables $(\bar{Q}_2(t), \bar{Q}_7(t))$ of the fluid model trajectory evolve according to Table 1 after a time $\tau_{sd} = t_{01} \|\bar{X}(0)\|_{\bar{h}\mathcal{E}}$. From Table 1, the \bar{Q}_2 and \bar{Q}_7 components of the fluid model trajectory reach values in the set $\bar{h}\mathcal{E}$'s projection in at most an additional $\frac{10}{a} \|\bar{X}(\tau_{sd})\|_{\bar{h}\mathcal{E}}$ time units. Because the total arrival rate into the system is less than or equal to 1.8,

$$\|\bar{X}(\tau_{sd})\|_{\bar{h}\mathcal{E}} \leq (2t_{01} + 1) \|\bar{X}(0)\|_{\bar{h}\mathcal{E}}.$$

Thus after a time $t_{02} \|\bar{X}(0)\|_{\bar{h}\mathcal{E}}$, where $t_{02} = \frac{10}{a}(2t_{01} + 1) + 1$, all queues but queue 1 have been shown to reach values in the projection of the set $\bar{h}\mathcal{E}$. By Lemma 3, queue 5 is empty, so either: queue 1 is above threshold, in which case discarding is on and it will reach threshold in $2(\bar{Q}_1(t_{02} \|\bar{X}(0)\|_{\bar{h}\mathcal{E}}) - \bar{h})$ time units, or queue 1 is below threshold in which case it will reach threshold in $10(\bar{Q}_1(t_{02} \|\bar{X}(0)\|_{\bar{h}\mathcal{E}}) - \bar{h})$ time units. Once $\bar{Q}_1(t)$ reaches threshold \bar{h} , it remains there by the following reasoning. If queue 1 were to move some positive amount ϵ above \bar{h} , the discarding would have turned on before the queue grew to ϵ and prevented it from getting there. Similarly, if queue 1 were to move some positive amount ϵ below \bar{h} , the discarding would have turned off before the queue receded by ϵ , and prevented the queue from receding that much. Very loosely, we can bound the rate of growth of queue

1 before time $t_{02} \|\bar{X}(0)\|_{\bar{h}\mathcal{E}}$ by

$$\bar{Q}_1(t_{02} \|\bar{X}(0)\|_{\bar{h}\mathcal{E}}) \leq 1.6 \|\bar{X}(0)\|_{\bar{h}}.$$

Thus after a time of length $t_0 \|\bar{X}(0)\|_{\bar{h}\mathcal{E}}$, where t_0 is given by $t_0 = \frac{10}{a}(2t_{01} + 1) + 17$, all fluid model trajectories will have reached the set $\bar{h}\mathcal{E}$. The departure rates for all three flows, as well as the departure rates for each class associated with each flow, are easily seen to be 0.5 when the fluid model's state is in $\bar{h}\mathcal{E}$ and threshold $\bar{h} > 0$. Thus, by Theorem 4 we have that the asymptotic flow rates approach 0.5. \square

5. Conclusion. In this work we have shown how the analysis of the flow rates of a stochastic network with a particular flow control scheme may be reduced to an analysis of a fluid model. While we have focused on a particular flow control scheme, the same analysis could be carried out for many other control schemes. The key feature that enabled our approach was that our control scheme has a free parameter, n , which when increased makes the system look more and more like a deterministic fluid system. We have demonstrated how to use the theory developed in this paper to analyze an example network resembling a 2-input, 2-output switch.

APPENDIX

Before proving Theorem 1, we state and prove a number of lemmas. Lemma 4 is a functional form of the strong law of large numbers for renewal processes, and is taken from [4]. Lemma 5 is a new result showing that the thinned arrivals (the customers that make it beyond the discarding point) converge to a fluid limit along a subsequence. Lemma 6 is a result taken from [4] showing that the residual initial arrival and service times decline to zero at rate 1 in the fluid limit. The lemma also shows that the sequence of functions we use to take the fluid limit are uniformly integrable.

Also the lemmas will make use of fluid limits that have well defined limiting residual interarrival and service times, as defined by the following property.

PROPERTY 2. $\{(\mathbf{x}_j, a_j)\}$ is a sequence for which $\frac{U^{\mathbf{x}_j}(0)}{a_j} \rightarrow \bar{U}(0)$, $\frac{V^{\mathbf{x}_j}(0)}{a_j} \rightarrow \bar{V}(0)$, for some $\bar{U}(0) \geq 0$ and $\bar{V}(0) \geq 0$.

LEMMA 4 (Dai, Lemma 4.2 of [4]). Suppose that $\{(\mathbf{x}_j, a_j)\}$ is a sequence satisfying Properties 1 and 2. (on pages 17 and 33). Then for almost all ω ,

$$\frac{E_f^{\mathbf{x}_j}(a_j t)}{a_j} \rightarrow \alpha_f(t - \bar{U}_f(0))^+ \text{ u.o.c.}, \quad \frac{S_k^{\mathbf{x}_j}(a_j t)}{a_j} \rightarrow \mu_k(t - \bar{V}_k(0))^+ \text{ u.o.c.}$$

PROOF. See Lemma 4.2 of Dai [4]. The result is an instance of the Strong Law of Large Numbers for Renewal Processes [10]. \square

LEMMA 5 (Thinned Arrival Convergence). *Suppose that $\{(\mathbf{x}_j, a_j)\}$ is a sequence satisfying Properties 1 and 2. Then for almost all ω , there exists a subsequence $\{(\mathbf{x}_m, a_m)\} \subseteq \{(\mathbf{x}_j, a_j)\}$ such that*

$$\Lambda^{\mathbf{x}_m}(a_m t)/a_m \rightarrow \bar{\Lambda}(t) \quad \text{u.o.c.}$$

where $\bar{\Lambda}(t)$ is some Lipschitz continuous process with, for all regular $t \geq 0$,

$$(48) \quad \dot{\bar{\Lambda}}_f(t) \leq \alpha_f \quad \text{for each flow } f.$$

PROOF. By Lemma 4,

$$(49) \quad E_f^{\mathbf{x}_j}(a_j t)/a_j \rightarrow \alpha_f(t - \bar{U}_f(0))^+ \quad \text{u.o.c.}$$

for each flow f . For notational convenience in the development that follows, we define:

$$(50) \quad \bar{E}_f(t) \triangleq \alpha_f(t - \bar{U}_f(0))^+, \quad \Delta_j(t) \triangleq E^{\mathbf{x}_j}(a_j t)/a_j - \bar{E}(t).$$

Pick a compact time interval $[s_0, s_1]$. Since the number of admitted customers is not greater than the number that arrive,

$$(51) \quad \frac{1}{a_j} [\Lambda^{\mathbf{x}_j}(a_j(t + \varepsilon)) - \Lambda^{\mathbf{x}_j}(a_j t)] \leq \frac{1}{a_j} [E^{\mathbf{x}_j}(a_j(t + \varepsilon)) - E^{\mathbf{x}_j}(a_j t)]$$

for any positive $\varepsilon \leq s_1 - s_0$ and $t : s_0 \leq t \leq s_1 - \varepsilon$. Adding $-\Delta_j(t + \varepsilon)$ and $\Delta_j(t)$ to both sides and substituting (50), we have

$$\frac{\Lambda^{\mathbf{x}_j}(a_j(t + \varepsilon))}{a_j} - \Delta_j(t + \varepsilon) - \left[\frac{\Lambda^{\mathbf{x}_j}(a_j t)}{a_j} - \Delta_j(t) \right] \leq \bar{E}(t + \varepsilon) - \bar{E}(t) \leq \varepsilon \alpha.$$

Define the family of functions:

$$\mathfrak{L}_j(s_0, t) := \sup_{s \in [s_0, t]} \left[\frac{\Lambda^{\mathbf{x}_j}(a_j s)}{a_j} - \Delta_j(s) \right]$$

for $t \in [s_0, s_1]$. Because the argument of the sup function is a vector, sup is taken component-wise. Note that for any (t, ε) with $t \in [s_0, s_1 - \varepsilon]$,

$$\begin{aligned} \mathfrak{L}_j(s_0, t + \varepsilon) &= \mathfrak{L}_j(s_0, t) \vee \mathfrak{L}_j(t, t + \varepsilon) \\ \text{and } \mathfrak{L}_j(t, t + \varepsilon) &\leq \varepsilon \alpha + \mathfrak{L}_j(t, t) \leq \varepsilon \alpha + \mathfrak{L}_j(s_0, t). \end{aligned}$$

Thus $\mathfrak{L}_j(s_0, t + \varepsilon) - \mathfrak{L}_j(s_0, t) \leq \varepsilon\alpha$ and clearly $\mathfrak{L}_j(s_0, t + \varepsilon) - \mathfrak{L}_j(s_0, t) \geq 0$ because $\mathfrak{L}_j(s_0, \cdot)$ is monotone. Hence the functions $\mathfrak{L}_j(s_0, \cdot)$ are equicontinuous and individually Lipschitz continuous. Thus, by Arzela's theorem, there exists a further subsequence $\{(\mathbf{x}_m, a_m)\} \subseteq \{(\mathbf{x}_j, a_j)\}$ such that

$$\mathfrak{L}_m(s_0, t) \rightarrow \bar{\Lambda}(t)$$

uniformly on the compact set $t \in [s_0, s_1]$ for some monotone-nondecreasing, Lipschitz-continuous process $\bar{\Lambda}(t)$. But by (49), $\Delta_j(t) \rightarrow 0$ uniformly on compact sets. Because of this and the fact that $\Lambda^{\mathbf{x}_j}(a_j s)/a_j$ is monotone in s , it follows that $\mathfrak{L}_j(s_0, t)$ approaches $\Lambda^{\mathbf{x}_j}(a_j t)/a_j$ as $j \rightarrow \infty$. Thus

$$\sup_{s \in [s_0, t]} \left[\frac{\Lambda^{\mathbf{x}_j}(a_j s)}{a_j} - \Delta_j(s) \right] \rightarrow \frac{\Lambda^{\mathbf{x}_j}(a_j t)}{a_j} \rightarrow \bar{\Lambda}(t).$$

Because the choice of $[s_0, s_1]$ was arbitrary, we have $\Lambda^{\mathbf{x}_m}(a_m s)/a_m \rightarrow \bar{\Lambda}(t)$ u.o.c. Furthermore, (49) and (51) imply that $\bar{\Lambda}(t)$ satisfies (48). \square

LEMMA 6 (Lemmas 4.3 & 4.5 of Dai [4]). *Suppose that $\{(\mathbf{x}_j, a_j)\}$ is a sequence satisfying Properties 1 and 2. Then almost surely:*

$$\lim_{j \rightarrow \infty} \frac{U_f^{\mathbf{x}_j}(a_j t)}{a_j} = (\bar{U}_f(0) - t)^+ \text{ u.o.c.}, \quad \lim_{j \rightarrow \infty} \frac{V_k^{\mathbf{x}_j}(a_j t)}{a_j} = (\bar{V}_k(0) - t)^+ \text{ u.o.c.}$$

Also, for each fixed $t \geq 0$, the sets of functions:

$$\{U^{\mathbf{x}_j}(a_j t)/a_j : a_j \geq 1\}, \quad \{V^{\mathbf{x}_j}(a_j t)/a_j : a_j \geq 1\}, \\ \{Q^{\mathbf{x}_j}(a_j t)/a_j : a_j \geq 1\}$$

are uniformly integrable.

PROOF. See Lemmas 4.3 and 4.5 of Dai [4] \square

We use the following lemma later to show that because all of the systems we consider are work-conserving, the fluid limit must also be work-conserving. In the lemma below, the notation $D_{\mathbb{R}}[0, \infty)$ denotes the space of right-continuous functions on \mathbb{R}_+ having left limits on $(0, \infty)$, and endowed with the Skorohod topology [12]. $C_{\mathbb{R}}[0, \infty) \subset D_{\mathbb{R}}[0, \infty)$ is the subset of continuous paths.

LEMMA 7 (Lemma 2.4 of Dai and Williams [8]). *Let $\{(z_j, \chi_j)\}$ be a sequence in $D_{\mathbb{R}}[0, \infty) \times C_{\mathbb{R}}[0, \infty)$. Assume that χ_j is nondecreasing and (z_j, χ_j)*

converges to $(z, \chi) \in C_{\mathbb{R}}[0, \infty) \times C_{\mathbb{R}}[0, \infty)$ u.o.c. Then for any bounded continuous function f ,

$$\int_0^t f(z_j(s)) d\chi_j(s) \rightarrow \int_0^t f(z(s)) d\chi(s) \quad \text{u.o.c.}$$

PROOF. See Lemma 2.4 of Dai and Williams [8]. \square

We are now ready to prove Theorem 1.

PROOF OF THEOREM 1. Before scaling space, the discarding thresholds for each j are $n_j h$. After scaling space by a_j , the scaled thresholds are $n_j h/a_j$. Property 1 insures that n_j/a_j is upper bounded by a constant. Thus by the Bolzano-Weierstrass Theorem, there exists a subsequence $\{(\mathbf{x}_r, a_r)\} \subseteq \{(\mathbf{x}_j, a_j)\}$ for which $n_r h/a_r \rightarrow \bar{h}$ for some $\bar{h} \geq 0$.

Property 1 insures that $\|x_r/a_r\|_{n_r h/a_r \mathcal{E}}$ is upper bounded by a constant. Thus $\limsup \|x_r/a_r\|_{\bar{h} \mathcal{E}}$ is finite. Consequently, there must be some further subsequence $\{(\mathbf{x}_u, a_u)\} \subseteq \{(\mathbf{x}_r, a_r)\}$ for which $x_u/a_u \rightarrow \bar{X}(0)$ for some finite $\bar{X}(0)$.

The hysteresis variables satisfy $H^{\mathbf{x}_u}(a_u t)/a_u \rightarrow 0$ u.o.c. because $H^{\mathbf{x}_u}(a_u t)$ is bounded by a constant by its definition. This fact along with the convergence of $X^{\mathbf{x}_u}(0)/a_u \rightarrow \bar{X}(0)$ allows us to use Lemma 6 to conclude $U^{\mathbf{x}_u}(a_u t)/a_u \rightarrow \bar{U}(t)$ and $V^{\mathbf{x}_u}(a_u t)/a_u \rightarrow \bar{V}(t)$ u.o.c. where $\bar{U}(t)$ and $\bar{V}(t)$ satisfy (12).

The cumulative service time process $T^{\mathbf{x}_u}$ satisfies

$$(52) \quad [T^{\mathbf{x}_u}(a_u t) - T^{\mathbf{x}_u}(a_u s)]/a_u \leq (t - s).$$

Thus by Arzela's theorem [21], there exists a further subsequence $\{(\mathbf{x}_v, a_v)\} \subseteq \{(\mathbf{x}_u, a_u)\}$ for which $T^{\mathbf{x}_v}(a_v t)/a_v \rightarrow \bar{T}(t)$. Property (13) follows from (6). Property (7) implies $I^{\mathbf{x}_v}(a_v t)/a_v \rightarrow \bar{I}(t)$ u.o.c. where $\bar{I}(t)$ satisfies (14).

By Lemma 4, $S_k^{\mathbf{x}_v}(a_v t)/a_v \rightarrow (\mu_k t - \bar{V}_k(0))^+$ u.o.c. for each class k . This fact combined with (9) and (52) gives (15).

We have already shown that $X^{\mathbf{x}_v}(0) \rightarrow \bar{X}(0)$, therefore the $U^{\mathbf{x}_v}(0)$ and $V^{\mathbf{x}_v}(0)$ components of $X^{\mathbf{x}_v}(0)$ converge to a limiting value. This fact allows us to invoke Lemma 5 to conclude that there is a subsequence $\{(\mathbf{x}_m, a_m)\} \subseteq \{(\mathbf{x}_v, a_v)\}$ for which $\Lambda^{\mathbf{x}_m}(a_m t)/a_m \rightarrow \bar{\Lambda}(t)$ u.o.c. for some Lipschitz continuous process $\bar{\Lambda}(t)$ satisfying (22).

Lemma 4 combined with (3) gives us $A_k^{\mathbf{x}_m}(a_m t)/a_m \rightarrow \bar{A}_k(t)$ u.o.c. for each class k where $\bar{A}_k(t)$ is defined by (16). Furthermore, $\bar{A}_k(t)$ is Lipschitz continuous because it is equal to a linear combination of functions we have

already shown to be Lipschitz continuous. Thus using (4) we have that

$$(53) \quad Q^{\mathbf{x}^m}(a_mt)/a_m \rightarrow \bar{Q}(t) \quad \text{u.o.c.}$$

where $\bar{Q}(t)$ is a Lipschitz continuous function given by (17). Property (18) follows easily from (5).

The next few arguments are similar to the proof of Proposition 4.2 in [6]. Suppose that $\bar{Q}_k(t) > \bar{h}$ for some $k \in \mathcal{C}(f)$. By Lipschitz continuity of $\bar{Q}_k(t)$, there exists some small $\tau > 0$ such that $\min_{t \leq s \leq t+\tau} \bar{Q}_k(s) > \bar{h}$. By the uniformity of the queue convergence in (53) and that $n_m h/a_m \rightarrow \bar{h}$, there exists m^* such that for all $m > m^*$, $Q_k^{\mathbf{x}^m}(a_ms) > n_m h$ for all $s \in [t, t+\tau]$. Thus, by (10) one finds that $\Lambda_f^{\mathbf{x}^m}(a_ms) - \Lambda_f^{\mathbf{x}^m}(a_mt) = 0 \quad \forall s \in [t, t+\tau]$. Therefore, it follows that $\bar{\Lambda}_f(s) - \bar{\Lambda}_f(t) = 0 \quad \forall s \in [t, t+\tau]$ and consequently, $\dot{\bar{\Lambda}}_f(t) = 0$, which is (20).

Suppose that $\bar{Q}_k(t) < \bar{h}$ for all $k \in \mathcal{C}(f)$. First note that in this case $\bar{h} > 0$. By the Lipschitz continuity of $\bar{Q}_k(t)$ for each k , there exists some small $\tau > 0$ such that $\max_{k \in \mathcal{C}(f)} \max_{s \in [t, t+\tau]} \bar{Q}_k(s) < \bar{h}$. Because $n_m \rightarrow \infty$, the uniformity of the convergence in (53), and that $n_m h/a_m \rightarrow \bar{h}$, there exists m' such that for all $m > m'$, $Q_k^{\mathbf{x}^m}(a_ms) < n_m h$. Furthermore there exists a $m^* \geq m'$ such that for all $m > m^*$ and $k \in \mathcal{C}(f)$, $Q_k^{\mathbf{x}^m}(a_ms) < n_m h - o(n_m)h\varsigma$. Thus, by (10),

$$\Lambda_f^{\mathbf{x}^m}(a_ms) - \Lambda_f^{\mathbf{x}^m}(a_mt) = E_f^{\mathbf{x}^m}(a_ms) - E_f^{\mathbf{x}^m}(a_mt) \quad \forall s \in [t, t+\tau]$$

and consequently we have (21).

Suppose that for some class k , $\bar{Q}_k(t) > 0$. By the Lipschitz continuity of $\bar{Q}_k(t)$ there exists some small $\tau > 0$ such that $\min_{t \leq s \leq t+\tau} \bar{Q}_k(s) > 0$. Because of the uniformity of convergence in (53) there exists m^* such that for all $m > m^*$, $Q_k^{\mathbf{x}^m}(a_ms) > 0 \quad \forall s \in [t, t+\tau]$. By (11), for almost all ω , and all classes l we have

$$w_k^{-1}[D_k(a_ms) - D_k(a_mt)] \geq w_l^{-1}[D_l(a_ms) - D_l(a_mt)] \quad \forall s \in [t, t+\tau]$$

and thus we have (23).

If $\bar{Q}_l(t) > 0$ and $\bar{Q}_k(t) > 0$, then (23) is true as written or with the k and l and indices swapped. This implies (24).

We observe that (8) is equivalent to $\int_0^\infty f(\chi_m) dz_m = 0$ where

$$\chi_m := \frac{C_i Q^{\mathbf{x}^m}(a_mt)}{a_m}, \quad z_m := \frac{I_i^{\mathbf{x}^m}(a_mt)}{a_m}, \quad f(\cdot) := (\cdot) \wedge 1.$$

Noting that χ_m and z_m meet the required conditions for Lemma 7 we have, $\int_0^\infty [C_i \bar{Q}(t)] \wedge 1 d\bar{I}_i(t) = 0$ which is equivalent to (19). \square

PROOF OF THEOREM 3. We first prove conclusion (i). Pick any sequence of pairs $\{(\mathbf{x}_j, a_j)\}$ satisfying $a_j = n_j \|x_j/n_j\|_{h\mathcal{E}} \rightarrow \infty$ and $\|x_j/n_j\|_{h\mathcal{E}} > \zeta$ for some $\zeta > 0$ (a far fluid limit sequence). To invoke Lemma 1, we pick \mathfrak{F} while simultaneously defining the process $\bar{F}(\cdot)$ according to the expression

$$\bar{F}(t) \triangleq \mathfrak{F} \circ [\bar{X}(\cdot); \bar{T}(\cdot); \bar{\Lambda}(\cdot); \bar{h}] (t) := \|\bar{X}(t)\|_{\bar{h}\mathcal{E}} \quad \forall t \geq 0.$$

Note that $\bar{F}(\|\bar{X}(0)\|_{\bar{h}\mathcal{E}} t) = 0$ for all $t \geq t_0$ by (30), and \mathfrak{F} is easily seen to be continuous on the topology of uniform convergence on compact sets. Since $\|x_j/a_j\|_{n_j h/a_j} = 1$ as argued in Corollary 2, we can set the c of Lemma 1 to 1. Applying Lemma 1 and taking $t = t_0$ we have that

$$\frac{1}{\|x_j\|_{n_j h\mathcal{E}}} \left\| X^{\mathbf{x}_j}(\|x_j\|_{n_j h\mathcal{E}} t_0) \right\|_{n_j h\mathcal{E}} \rightarrow 0 \text{ a.s.}$$

By Lemma 6, $\frac{1}{\|x_j\|_{n_j h\mathcal{E}}} X^{\mathbf{x}_j}(\|x_j\|_{n_j h\mathcal{E}} t_0)$ is uniformly integrable. Therefore

$$(54) \quad \lim_{j \rightarrow \infty} \frac{1}{\|x_j\|_{n_j h\mathcal{E}}} \mathbb{E} \left\| X^{\mathbf{x}_j}(\|x_j\|_{n_j h\mathcal{E}} t_0) \right\|_{n_j h\mathcal{E}} = 0.$$

Using that the above holds for any far fluid limit sequence, we show by contradiction that conclusion (i) of the theorem is true. Suppose conclusion (i) were not true. Then for some $\zeta > 0$ and some positive δ , we would have that for any L_2 there would exist a pair $\mathbf{x} = (x, n)$ with $\|x\|_{nh\mathcal{E}} \geq L_2$ and $\|x/n\|_{h\mathcal{E}} > \zeta$ with $\frac{1}{\|x\|_{nh\mathcal{E}}} \mathbb{E} \|X^{\mathbf{x}}(t_0 \|x\|_{nh\mathcal{E}})\|_{h\mathcal{E}} \leq \delta$. We therefore could construct a sequence that violates (54), which is true for any far fluid limit sequence. A special case of a far fluid limit sequence is when $n > L_2 \zeta^{-1}$ and $\|x/n\|_{h\mathcal{E}} > \zeta$. Hence we have conclusion (i) of the theorem.

We now turn to showing conclusion (ii). Pick an arbitrary sequence of pairs $\{(\mathbf{x}_j, a_j)\}$ satisfying $a_j = n_j \rightarrow \infty$ and $\|x_j/n_j\|_{h\mathcal{E}} \leq \zeta$ for some constant ζ (a near fluid limit sequence). We again invoke Lemma 1 by taking \mathfrak{F} to be the same functional as before, i.e.,

$$\bar{F}(t) \triangleq \mathfrak{F} \circ [\bar{X}(\cdot); \bar{T}(\cdot); \bar{\Lambda}(\cdot); \bar{h}] (t) := \|\bar{X}(t)\|_{\bar{h}\mathcal{E}} \quad \forall t \geq 0.$$

Using Lemma 1, and the fact that $\|x_j/n_j\|_{h\mathcal{E}} \leq \zeta$ we have $\|X^{\mathbf{x}_j}(n_j t)/n_j\|_{h\mathcal{E}} \rightarrow 0$ a.s. for each $t \geq \zeta t_0$. Now take $t = t_0$, $\frac{1}{n_j} \|X^{\mathbf{x}_j}(n_j t_0)\|_{n_j h\mathcal{E}} \rightarrow 0$ a.s. By Lemma 6, $X^{\mathbf{x}_j}(n_j t_0)/n_j$ is uniformly integrable. Therefore

$$(55) \quad \lim_{j \rightarrow \infty} \mathbb{E} \left[\frac{1}{n_j} \|X^{\mathbf{x}_j}(n_j t_0)\|_{n_j h\mathcal{E}} \right] = 0.$$

We claim that the above implies conclusion (ii) is true by contradiction. Suppose (ii) were not true. Then for some choice ζ and b , we would have that for every constant L_3 , there would exist an $n \geq L_3$ and $x : \|x/n\|_{h\mathcal{E}} \leq \zeta$ satisfying $E \left\| \frac{1}{n} X^x(nt_0) \right\|_{h\mathcal{E}} > b$. This would allow us to construct a sequence that violates (55), which is a contradiction. \square

PROOF OF LEMMA 2. The argument that follows is adapted from the proof of Theorem 2.1 (ii) of Meyn and Tweedie [19]. We use the following fact taken from Theorem 14.2.2 of [18]:

Fact 1: (Meyn and Tweedie [18]) Suppose a discrete time Markov chain $\Phi = \{\Phi_k, k \in \mathbb{Z}^+\}$ is defined on a general state space X with transition kernel $P(x, A) = P_x(\Phi_1 \in A)$, where $A \in \mathfrak{B}(X)$, the Borel subsets of X . If V and f are nonnegative measurable functions satisfying

$$\int P(x, dy)V(y) \leq V(x) - f(x) + \tilde{b}1_B(x), \quad x \in X$$

then

$$E_x \left[\sum_{k=0}^{\tau_B-1} f(\Phi_k) \right] \leq V(x) + \tilde{b}$$

where $\tau_B = \inf\{k \geq 1 : \Phi_k \in B\}$.

The above fact is a form of Dynkin's formula and is shown by using the first inequality to sum bounds of the increments $E_x V(\Phi_k) - E_x V(\Phi_{k+1})$ for $k \in \{0 \dots \tau_B - 1\}$. Since $1_B(\Phi_k)$ is 1 at most once for $k \in \{0, \dots, \tau_B - 1\}$ on each sample path, \tilde{b} appears once in the final expression.

We define the set $B \triangleq \{x : \|x/n\|_{h\mathcal{E}} \leq \zeta\}$. Next, we define the following functions, the first mapping each $x \in X$ to a time $m(x)$, and the second a Lyapunov function mapping each x to a value:

$$(56) \quad m(x) \triangleq \begin{cases} n \|x/n\|_{h\mathcal{E}} t_0 & \text{if } x \notin B \\ nt_0 & \text{if } x \in B \end{cases}$$

$$(57) \quad V(x) \triangleq \frac{nt_0}{1-\delta} \|x/n\|_{h\mathcal{E}}.$$

Substituting $m(x)$ for time in relation (31), and adding a term to that relation's right hand side so that the relation holds for x both inside and outside

B , we have

$$\begin{aligned} \mathbb{E}_x \left\| \frac{1}{n} X^n(m(x)) \right\|_{h\mathcal{E}} &\leq \delta \|x/n\|_{h\mathcal{E}} + \left(\sup_{\tilde{x} \in B} \mathbb{E}_{\tilde{x}} \left\| \frac{1}{n} X^n(nt_0) \right\|_{h\mathcal{E}} \right) 1_B(x) \\ &\leq \|x/n\|_{h\mathcal{E}} + b 1_B(x) \\ &\leq \|x/n\|_{h\mathcal{E}} - \frac{1-\delta}{nt_0} m(x) + (1-\delta+b) 1_B(x) \end{aligned}$$

where the middle step follows from (32). By multiplying both sides by $nt_0/(1-\delta)$ we have

$$(58) \quad \mathbb{E}_x[V(X^n(m(x)))] \leq V(x) - m(x) + \tilde{b} 1_B(x)$$

$$(59) \quad \text{where} \quad \tilde{b} = nt_0 + \frac{nt_0}{1-\delta} b.$$

The transition kernel \mathbf{P}^t for the Markov process X^n is defined by $\mathbf{P}^t(x, A) = \mathbb{P}_x(X^n(t) \in A)$ where A is any set in $\mathfrak{B}(X)$, the Borel subsets of the state space X . We define the discrete time ‘‘embedded’’ Markov chain $\hat{\Phi} = \{\hat{\Phi}_k, k \in \mathbb{Z}_+\}$ with transition kernel $\hat{\mathbf{P}}$ given by $\hat{\mathbf{P}}(x, A) = \mathbf{P}^{m(x)}(x, A)$. Note that

$$\int \hat{\mathbf{P}}(x, dz) V(z) = \int \mathbf{P}^{m(x)}(x, dz) V(z) = \mathbb{E}_x[V(X^n(m(x)))].$$

Combining this with (58) we have

$$\int \hat{\mathbf{P}}(x, dz) V(z) \leq V(x) - m(x) + \tilde{b} 1_B(y).$$

Thus by Fact 1,

$$(60) \quad \mathbb{E}_x \left[\sum_{k=0}^{\hat{\tau}_B-1} m(\hat{\Phi}_k) \right] \leq V(x) + \tilde{b}$$

where $\hat{\tau}_B = \inf\{k \geq 1 : \hat{\Phi}_k \in B\}$. If the embedded chain hits B in $\hat{\tau}_B$ discrete steps, then the original chain must also hit B in a time less than or equal to the sum of the embedded times. Thus,

$$\inf\{t \geq 0 : X^{x,n}(t) \in B\} \leq \sum_{k=0}^{\hat{\tau}_B-1} m(\hat{\Phi}_k) \quad \mathbb{P}_x\text{-a.s.}$$

for each $x \in X$. Furthermore, whenever the initial condition $x \in B$, the first embedded time is nt_0 seconds by (56). Consequently, the time of the first

hitting of B after nt_0 seconds expire satisfies

$$\inf\{t \geq nt_0 : X^{x,n}(t) \in B\} \leq \sum_{k=0}^{\hat{\tau}_B-1} m(\hat{\Phi}_k) \quad \text{P}_x\text{-a.s.}$$

for each $x \in B$. Substituting definition (34), taking the expectation, and using (60), we have

$$\mathbb{E}_x[\tau_B^n(nt_0)] \leq V(x) + \tilde{b} \quad \text{for all } x \in B.$$

Taking the $\sup_{x \in B}$ of both sides, substituting (57) and (59) we have (33). Since B is closed and bounded, and arrivals are from an unbounded distribution (1) and spread-out (2), B is a petite set. (See [18] for a discussion of petite sets.) Therefore (33) implies X is positive Harris recurrent by Theorem 4.1 of [20]. \square

REFERENCES

- [1] BERTSEKAS, D. and GALLAGER, R. (1992). *Data Networks*. Prentice Hall, Englewood Cliffs, NJ.
- [2] BRAMSON, M. (1998). Stability of two families of queueing networks and a discussion of fluid limits. *Queueing Syst.* **28** 7-31.
- [3] CHEN, H. (1995). Fluid approximations and Stability of Multiclass Queueing Networks: Work Conserving Disciplines. *Ann. Appl. Probab.* **5** 637-665.
- [4] DAI, J. G. (1995). On Positive Harris Recurrence of Multiclass Queueing Networks: a Unified Approach via Fluid Limit Models. *Ann. Appl. Probab.* **5** 49-77.
- [5] DAI, J. G. and HARRISON, J. M. (1991). Steady-State Analysis of RBM in a Rectangle: numerical methods and a queueing application. *Ann. Appl. Probab.* **1** 16-35.
- [6] DAI, J. G. and MEYN, S. (1995). Stability and Convergence of Moments for Multiclass Queueing Networks via Fluid Limit Models. *IEEE Trans. Automat. Control* **40** 1889-1904.
- [7] DAI, J. G. and PRABHAKAR, B. (2000). The throughput of data switches with and without speedup. In *Proceedings of IEEE INFOCOM* 556-564.
- [8] DAI, J. G. and WILLIAMS, R. J. (1995). Existence and Uniqueness of Semimartingale Reflecting Brownian motions in Convex Polyhedrons. *Theory Probab. and Appl.* **40** 3-53.
- [9] DAVIS, M. H. A. (1984). Piecewise Deterministic Markov Processes: A General Class of Non-Diffusion Stochastic Models. *J. Roy. Statist. Soc.* **46** 353-388.
- [10] DURRETT, R. (2004). *Probability: Theory and Examples*, Third ed. Thomson Brooks/Cole, Belmont, CA.
- [11] EL-TAHA, M. and STIDHAM, S. (1999). *Sample-Path Analysis of Queueing Systems*. Kluwer Academic Publishers, New York.
- [12] ETHIER, S. N. and KURTZ, T. G. (1986). *Markov Processes: Characterization and Convergence*. Wiley, New York.
- [13] HARRISON, J. M. (1985). *Brownian Motion and Stochastic Flow Systems*. Wiley, New York.

- [14] HARRISON, J. M. and REIMAN, M. I. (1981). Reflected Brownian Motion on an Orthant. *Ann. Probab.* **9** 302-308.
- [15] KASPI, H. and MANDELBAUM, A. (1992). Regenerative Closed Queueing Networks. *Stochastics and Stochastics Rep.* **39** 239-258.
- [16] KONSTANTOPOULOS, T., PAPADAKIS, S. and WALRAND, J. (1994). Functional Approximation Theorems for Controlled Renewal Processes. *J. Appl. Probab.* **31** 765-776.
- [17] MANDELBAUM, A., MASSEY, W. A. and REIMAN, M. I. (1998). Strong Approximations for Markovian Service Networks. *Queueing Syst.* **30** 149-201.
- [18] MEYN, S. P. and TWEEDIE, R. L. (1993). *Markov Chains and Stochastic Stability*. Springer-Verlag, London.
- [19] MEYN, S. P. and TWEEDIE, R. L. (1994). State-Dependent Convergence Criteria of Markov Chains. *Ann. Appl. Probab.* **4** 149-168.
- [20] MEYN, S. P. and TWEEDIE, R. L. (1993). Stability of Markovian Processes III: Foster-Lyapunov Criteria for Continuous-Time Processes. *Adv. Appl. Probab.* **25** 518-548.
- [21] MUNKRES, J. (2000). *Topology*, Second ed. Prentice Hall, Upper Saddle River, NJ.
- [22] MUSACCHIO, J. (2005). Pricing and Flow Control in Communications Networks PhD thesis, Electrical Engineering and Computer Sciences, University of California, Berkeley, Berkeley, CA.
- [23] SHREEDHAR, M. and VARGHESE, G. (1996). Efficient Fair Queuing Using Deficit Round-Robin. *IEEE/ACM Tran. Netw.* **4** 375-385.

TECHNOLOGY AND INFORMATION MANAGEMENT
UNIVERSITY OF CALIFORNIA, SANTA CRUZ
1156 HIGH STREET
MS: SOE 3
SANTA CRUZ, CA 95064
E-MAIL: johnm@soe.ucsc.com

DEPARTMENT OF ELECTRICAL ENGINEERING
AND COMPUTER SCIENCES
UNIVERSITY OF CALIFORNIA, BERKELEY
CORY HALL
BERKELEY, CA 94720
E-MAIL: wlr@eecs.berkeley.edu