

SYSTEMS WITH LARGE FLEXIBLE SERVER POOLS: INSTABILITY OF “NATURAL” LOAD BALANCING

ALEXANDER L. STOLYAR AND ELENA YUDOVINA

ABSTRACT. We consider general large-scale service systems with multiple customer classes and multiple server (agent) pools; mean service times depend both on the customer class and server pool. It is assumed that the allowed activities (routing choices) form a tree (in the graph with vertices being both customer classes and server pools). We study the behavior of the system under a natural (load balancing) routing/scheduling rule, *Longest-queue freest-server* (LQFS-LB), in the many-server asymptotic regime, such that the exogenous arrival rates of the customer classes, as well as the number of agents in each pool, grow to infinity in proportion to some scaling parameter r . *Equilibrium point* of the system under LQFS-LB is the desired operating point, with server pool loads minimized and perfectly balanced.

Our main results are as follows. (a) We show that, quite surprisingly (given the tree assumption), for certain parameter ranges, the *fluid limit* of the system may be *unstable* in the vicinity of the equilibrium point; such instability may occur if the activity graph is not “too small”. (b) Using (a), we demonstrate that the sequence of stationary distributions of *diffusion-scaled* processes (measuring $O(\sqrt{r})$ deviations from the equilibrium point) may be non-tight, and in fact may escape to infinity. (c) In one special case of interest, however, we show that the sequence of stationary distributions of diffusion-scaled processes is tight, and the limit of stationary distributions is the stationary distribution of the limiting diffusion process.

1. INTRODUCTION

Large-scale service systems (such as call centers) with heterogeneous customer and server (agent) populations bring up the need for efficient dynamic control policies that match arriving (or waiting) customers and available servers. In this setting, two goals are desirable. On the one hand, customers should not be kept waiting if this is possible. On the other hand, idle time should be distributed fairly among the servers. For example, one would like to avoid the situation in which one of the server pools is fully busy while another one has significant numbers of idle agents.

Consider a general system, where the arrival rate of class i customers is Λ_i , the service rate of a class i customer by type j agent is μ_{ij} , and the server pool sizes are B_j . Another very desirable feature of a dynamic control is insensitivity to parameters Λ_i and μ_{ij} . That is, the assignment of customers to server pools should, to the maximal degree possible, depend only on the current system state, and not on prior knowledge of arrival rates or mean service times, because those parameters may not be known in advance and, moreover, they may be changing in time.

If the system objective is to minimize the maximum average load of any server pool, a “static” optimal control can be obtained by solving a linear program, called *static planning problem* (SPP), which has B_j ’s, μ_{ij} ’s and Λ_i ’s as parameters. An optimal solution to the

Date: April 24, 2012.

The research of the second author was supported by the NSF Graduate Research Fellowship.

SPP will prescribe optimal average rates Λ_{ij} at which arriving customers should be routed to the server pools. Typically (in a certain sense) the solution to SPP is unique and the *basic activities*, i.e. routing choices (ij) for which $\Lambda_{ij} > 0$, form a tree; let us assume this is the case. It is possible to design a dynamic control policy, which achieves the load balancing objective without a priori knowledge of input rates Λ_i – the *Shadow Routing* policy in [12] does just that, and in the process it “automatically identifies” the basic activity tree. Shadow Routing policy, however, does need to “know” the service rates μ_{ij} .

The key question we address in this paper is as follows. Suppose, a control policy does *not* know the service rates μ_{ij} , but “somehow” it does know the structure of the basic activity tree, and restrict routing to this tree only. (For example, all feasible activities, i.e. those (ij) ’s for which $\mu_{ij} > 0$, may form a tree simply by the structure of the system. Another example: if Shadow Routing has some estimates of μ_{ij} , this will not be sufficient for it to identify the optimal routing rates, but may very well be sufficient to correctly identify the basic activity tree.) What is an efficient load balancing policy in this case?

If routing is restricted to a tree, it is very natural to conjecture that simple policies of the type considered by Gurvich and Whitt [7], Atar-Shaki-Shwartz [2], Armony and Ward [1], which are of the “serve longest queue” and “join least loaded pool” type, should “typically be good enough”. Some of the results in these (and other) papers, in fact, prove optimal behavior of simple load balancing schemes on a *finite time interval*; which further supports the above informal conjecture. One of the main contributions of our work is to show that, surprisingly, the above conjecture is *not* correct for a general parameter setting. The key reason is that a “natural” load balancing, *even if it is done along an a priori given optimal tree*, may render the system unstable in the vicinity of equilibrium point.

The specific control rule we analyze in this paper can be seen as a special case of the Queue-and-Idleness-Ratio rule considered in [7]. Within the given (basic) activity tree, if an arriving customer sees multiple available servers, it will choose the server pool with the smallest load; while if a server sees several customers waiting in queues, it will take a customer from the longest queue. We call this rule *Longest-queue freest-server* (LQFS-LB).

We consider a many-server asymptotic regime, such that $\Lambda_i = \lambda_i r$ (or sometimes $\Lambda_i = \lambda_i r + O(\sqrt{r})$), $B_j = \beta_j r$, where λ_i and β_j are some positive constants, $r \rightarrow \infty$ is a scaling parameter, and μ_{ij} remain constant. Our key results show that the *fluid limit* of the system process (obtained via space-scaling by $1/r$) can be unstable in the vicinity of the equilibrium point. This is very counterintuitive, because it would be reasonable to expect the contrary: that a simple load balancing in a system with activity graph free of cycles would be “well behaved”.

Using the fluid limit local instability (when such occurs), we prove that the sequence of stationary distributions of *diffusion-scaled* processes (measuring $O(\sqrt{r})$ deviations from the equilibrium point) may be non-tight, and in fact may escape to infinity. This of course means, in particular, that the behavior of the diffusion limit in the vicinity of equilibrium point *on a finite time interval*, may not be relevant to the system behavior in steady state, because the system “does not spend any time” in the $O(\sqrt{r})$ -vicinity of the equilibrium point.

In addition to the instability examples, we prove that in several cases the fluid limit will be (at least locally) stable. We demonstrate that fluid limit of any underloaded system with at most two customer classes, or critically loaded system with at most four customer classes, is always locally stable. We also demonstrate local stability in the case when the service rate depends only on the customer type (but not server pool, as long as it can serve it). In the

case when the service rate depends only on the server type (but not customer type, as long as it can be served), we show more – the global stability of the fluid limit.

General results on the asymptotics of stationary distributions (most importantly – their tightness), especially in the many-server systems’ diffusion limit, are notoriously difficult to derive. (For recent results in this direction see [5, 6].) In the special case when the service rate depends only on the server type, we prove that under the LQFS-LB policy the sequence of stationary distributions of diffusion-scaled processes is tight, and the limit of stationary distributions is the stationary distribution of the limiting diffusion process.

The structure of the paper is as follows. In Section 2 we present the model, define the static planning problem and related notions, and define the LQFS-LB policy. In Section 3 we define fluid models of the system, derive their basic properties in the vicinity of an equilibrium point, and define local stability. Section 4 contains fluid model stability results in the two special cases when the service rate depends on server class only or on customer type only. Our key results on local instability of fluid models are presented in Section 5. In Section 6 we consider an underloaded system (with optimal average utilization being $1 - \epsilon < 1$), and prove possible evanescence of stationary distributions of the diffusion scaled processes. Finally, Section 7 considers the so called Halfin-Whitt asymptotic regime (where the optimal average utilization is $1 - O(1/\sqrt{r})$), and contains two main results on the asymptotics of stationary distributions of the diffusion scaled processes: (a) possible evanescence under certain parameters and (b) tightness (and “limit interchange”) result for the case when the service rate depends only on the server type.

2. MODEL

2.1. The model; Static Planning (LP) Problem. Consider the model in which there are I customer classes, or types, labeled $1, 2, \dots, I$, and J server (agent) pools, or classes, labeled $1, 2, \dots, J$ (generally, we will use the subscripts i, i' for customer classes, and j, j' for server pools). The sets of customer classes and server classes will be denoted by \mathcal{I} and \mathcal{J} respectively.

We are interested in the scaling properties of the system as it grows large. The meaning of “grows large” is as follows. We consider a sequence of systems indexed by a scaling parameter r . As r grows, the arrival rates and the sizes of the service pools, but not the speed of service, increase. Specifically, in the r th system, customers of type i enter the system as a Poisson process of rate $\lambda_i^r = r\lambda_i + o(r)$, while the j th server pool has $r\beta_j$ individual servers. (All λ_i and β_j are positive parameters.) Customers may be accepted for service immediately upon arrival, or enter a queue; there is a separate queue for each customer type. Customers do not abandon the system. When a customer of type i is accepted for service by a server in pool j , the service time is exponential of rate μ_{ij} ; the service rate depends both on the customer type and the server type, but *not* on the scaling parameter r . If customers of type i cannot be served by servers of class j , the service rate is $\mu_{ij} = 0$.

We would like to balance the proportion of busy servers across the server pools, while keeping the system operating efficiently. Let λ_{ij}^r be the average rates at which type i customers are routed to server pools j . We would like the system state to be such that λ_{ij}^r are close to $\lambda_{ij}r$, where $\{\lambda_{ij}\}$ is an optimal solution to the following *static planning problem* (SPP), which is the following linear program:

$$(1) \quad \min_{\lambda_{ij}, \rho} \rho,$$

subject to

$$(2) \quad \lambda_{ij} \geq 0, \quad \forall i, j$$

$$(3) \quad \sum_j \lambda_{ij} = \lambda_i, \quad \forall i$$

$$(4) \quad \sum_i \lambda_{ij} / (\beta_j \mu_{ij}) \leq \rho, \quad \forall j.$$

We assume that the SPP has a unique optimal solution $\{\lambda_{ij}, i \in \mathcal{I}, j \in \mathcal{J}\}, \rho$; and it is such that the *basic activities*, i.e. those pairs, or edges, (ij) for which $\lambda_{ij} > 0$, form a (connected) tree in the graph with vertices set $\mathcal{I} \cup \mathcal{J}$. The set of basic activities is denoted \mathcal{E} . These assumptions constitute the *complete resource pooling* (CRP) condition, which holds “generically”; see [12, Theorem 2.2]. For a customer type i , let $\mathcal{S}(i) = \{j : (ij) \in \mathcal{E}\}$; for a server type j , let $\mathcal{C}(j) = \{i : (ij) \in \mathcal{E}\}$.

Note that under the CRP condition, all (“server pool capacity”) constraints (4) are binding; in other words, the optimal solution to SPP minimizes and “perfectly balances” server pool loads. Optimal dual variables ν_i , $i \in \mathcal{I}$, and α_j , $j \in \mathcal{J}$, corresponding to constraints (3) and (4), respectively, are unique and all strictly positive; ν_i is interpreted as the “workload” associated with one type i customer, and α_j is interpreted as the (scaled by $1/r$) maximum rate at which server pool j can process workload. The following relations hold:

$$\begin{aligned} \alpha_j &= \max_i \nu_i \beta_j \mu_{ij} & \nu_i &= \min_j \alpha_j / (\beta_j \mu_{ij}) \\ \sum_j \alpha_j &= 1 & \sum_i \lambda_i \nu_i &= \rho \sum_j \alpha_j = \rho. \end{aligned}$$

If $\rho < 1$, the system is called *underloaded*; if $\rho = 1$, the system is called *critically loaded*. In this paper we consider both cases.

In this paper, we assume that the basic activity tree is known in advance, and restrict our attention to the basic activities only. Namely, we assume that a type i customer service in pool j is allowed only if $(ij) \in \mathcal{E}$. (Equivalently, we can a priori assume that \mathcal{E} is the set of *all* possible activities, i.e. $\mu_{ij} = 0$ when $(ij) \notin \mathcal{E}$, and \mathcal{E} is a tree. In this case CRP requires that all feasible activities are basic.)

Let $\psi_{ij}^* = \lambda_{ij} / \mu_{ij}$. Continuing our interpretation of the optimal operating point of the system, let $\Psi_{ij}^r(t)$ be the number of servers of type j serving customers of type i at time t . It is desirable to have $\Psi_{ij}^r(t) = r\psi_{ij}^* + o(r)$. Later on we will be also interested in the question of whether or not the $o(r)$ term can in fact be $O(\sqrt{r})$.

2.2. Longest-queue, freest-server load balancing algorithm (LQFS-LB). For the rest of the paper, we analyze the performance of the following intuitive load balancing algorithm.

We introduce the following notation (for the system with scaling parameter r):

- $\Psi_{ij}^r(t)$: the number of servers of type j serving customers of type i at time t ;
- $\Psi_j^r(t) = \sum_i \Psi_{ij}^r(t)$: the total number of busy servers of type j at time t ;
- $\Psi_i^r(t) = \sum_j \Psi_{ij}^r(t)$: the total number of servers serving type i customers at time t ;
- $\Xi_j^r(t) = \Psi_j^r(t) / \beta_j$: the instantaneous load of server pool j at time t ;
- $Q_i^r(t)$: the number of customers of type i waiting for service at time t ;

$X_i^r(t) = \Psi_i^r(t) + Q_i^r(t)$: the total number of customers of type i in the system at time t .

The algorithm consists of two parts: routing and scheduling. “Routing” determines where an arriving customer goes if it sees available servers of several different types. “Scheduling” determines which waiting customer a server picks if it sees customers of several different types waiting in queue.

Routing: If an arriving customer of type i sees any unoccupied servers in server classes in $\mathcal{S}(i)$, it will pick a server in the least loaded server pool, i.e. $j \in \arg \min_{j \in \mathcal{S}(i)} \Xi_j^r(t)$. (Ties are broken in an arbitrary Markovian manner.)

Scheduling: If a server of type j , upon completing a service, sees a customer of a class in $\mathcal{C}(j)$ in queue, it will pick the customer from the longest queue, i.e. $i \in \arg \max_{j \in \mathcal{C}(j)} Q_i^r$. (Ties are broken in an arbitrary Markovian manner.)

By [7, Remark 2.3], the LQFS-LB algorithm described here is a special case of the algorithm proposed by Gurvich and Whitt, with constant probabilities $p_i = \frac{1}{J}$ (queues “should” be equal), $v_j = \frac{\beta_j}{\sum \beta_j}$ (the proportion of idle servers “should” be the same in all server pools).

2.3. Basic notation. Vector $(\xi_i, i \in \mathcal{I})$, where ξ can be any symbol, is often written as (ξ_i) or $\xi_{\mathcal{I}}$; similarly, $(\xi_j, j \in \mathcal{J}) = (\xi_j) = \xi_{\mathcal{J}}$ and $(\xi_{ij}, (ij) \in \mathcal{E}) = (\xi_{ij}) = \xi_{\mathcal{E}}$. We will treat $(\xi_{ij}) = \xi_{\mathcal{E}}$ as a vector, even though its elements have two indices. Unless specified otherwise, $\sum_i \xi_{ij} = \sum_{i \in \mathcal{C}(j)} \xi_{ij}$ and $\sum_j \xi_{ij} = \sum_{j \in \mathcal{S}(i)} \xi_{ij}$. For functions (or random processes) $(\xi(t), t \geq 0)$ we often write $\xi(\cdot)$. (And similarly for functions with domain different from $[0, \infty)$.) So, for example, $(\xi_i(\cdot))$ and $\xi_{\mathcal{I}}(\cdot)$ both signify $((\xi_i(t), i \in \mathcal{I}), t \geq 0)$. The indicator function of a set A is denoted $\mathbf{1}_A$; that is, $\mathbf{1}_A(\omega) = 1$ if $\omega \in A$ and 0 otherwise.

The symbol \implies denotes convergence in distribution of either random variables in the Euclidean space \mathbb{R}^d (with appropriate dimension d), or random processes in the Skorohod space $D^d[\eta, \infty)$ of RCLL (right-continuous with left limits) functions on $[\eta, \infty)$, for some constant $\eta \geq 0$. (Unless explicitly specified otherwise, $\eta = 0$.) The symbol \xrightarrow{w} denotes the weak convergence of probability measures on \mathbb{R}^d , or its one-point compactification $\overline{\mathbb{R}^d} = \mathbb{R}^d \cup \{*\}$, where $*$ is the “point at infinity”. We always consider the Borel σ -algebras on \mathbb{R}^d and $\overline{\mathbb{R}^d}$.

Standard Euclidean norm of a vector $x \in \mathbb{R}^d$ is denoted $|x|$. The symbol \rightarrow denotes ordinary convergence in \mathbb{R}^d or $\overline{\mathbb{R}^d}$. Abbreviation *u.o.c.* means *uniform on compact sets* convergence of functions, with the argument (usually in $[0, \infty)$) which is clear from the context; *w.p.1* means convergence *with probability 1*; $\dot{f}(t)$ means $(d/dt)f(t)$. Transposition of a matrix H is denoted H^\dagger ; in matrix expressions vectors are understood as column-vectors.

3. FLUID MODEL

3.1. Definition. We now consider the behavior of fluid models associated with this system. A fluid model is a set of trajectories that w.p.1 contains any limit of fluid-scaled trajectories of the original stochastic system. (We postpone proving this relationship between the fluid models and fluid limits until Section 3.4, in order to not interrupt the main content of Section 3; for now, we just formally define fluid models.)

The term *fluid model* denotes a set of Lipschitz continuous functions

$$\{(a_i(\cdot)), (x_i(\cdot)), (q_i(\cdot)), (\psi_{ij}(\cdot)), (\rho_j(\cdot))\},$$

which satisfy the equations below. (Here $a_i(\cdot) = (a_i(t), t \geq 0)$, and similarly for other components.) The last two equations involving derivatives are to be satisfied at all regular points t , when the derivatives in question exist. The interpretation of the components is as follows: $a_i(t)$ is the total number (actually, “amount”, i.e. the number, scaled by $1/r$) of arrivals of type i customers into the system by time t , $x_i(t)$ is the number (“amount”) of customers of type i in the system at time t , $q_i(t)$ is the number (“amount”) of customers of type i waiting in queue at time t , $\psi_{ij}(t)$ is the number (“amount”) of customers of type i being served by servers of type j at time t , and $\rho_j(t)$ is the instantaneous load (proportion of busy servers, the limit of $\Xi_j^r(t)/r$) in server pool j .

$$(5a) \quad a_i(t) = \lambda_i t, \quad \forall i \in \mathcal{I}$$

$$(5b) \quad x_i(t) = q_i(t) + \sum_j \psi_{ij}(t), \quad \forall i \in \mathcal{I}$$

$$(5c) \quad x_i(t) = x_i(0) + a_i(t) - \sum_j \int_0^t \mu_{ij} \psi_{ij}(s) ds, \quad \forall i \in \mathcal{I}$$

$$(5d) \quad \rho_j(t) = \frac{1}{\beta_j} \sum_i \psi_{ij}(t), \quad \forall j \in \mathcal{J}$$

$$(5e) \quad \rho_j(t) = 1 \text{ if } q_i(t) > 0 \text{ for any } i \in \mathcal{C}(j), \quad \forall j \in \mathcal{J}$$

For any set of server types $\mathcal{J}^* \subseteq \mathcal{J}$ and any set of customer types $\mathcal{I}^* \subseteq \mathcal{I}$ such that $q_i(t) > 0$ for all $i \in \mathcal{I}^*$, and $q_i(t) > q_{i'}(t)$ whenever $i \in \mathcal{I}^*$, $i' \notin \mathcal{I}^*$ and $\mathcal{S}(i) \cap \mathcal{S}(i') \cap \mathcal{J}^* \neq \emptyset$,

$$(5fa) \quad \sum_{i \in \mathcal{I}^*} \sum_{j \in \mathcal{S}(i) \cap \mathcal{J}^*} \dot{\psi}_{ij}(t) = \sum_{j \in \cup_{i \in \mathcal{I}^*} \mathcal{S}(i) \cap \mathcal{J}^*} \sum_{i' \in \mathcal{C}(j)} \mu_{i'j} \psi_{i'j}(t) - \sum_{i \in \mathcal{I}^*} \sum_{j \in \mathcal{S}(i) \cap \mathcal{J}^*} \mu_{ij} \psi_{ij}(t)$$

For any sets of customer types $\mathcal{I}_* \subseteq \mathcal{I}$, and any set of server types $\mathcal{J}_* \subseteq \mathcal{J}$ such that $\rho_j(t) < 1$ for all $j \in \mathcal{J}_*$, and $\rho_j(t) < \rho_{j'}(t)$ whenever $j \in \mathcal{J}_*$, $j' \notin \mathcal{J}_*$, and $\mathcal{C}(j) \cap \mathcal{C}(j') \cap \mathcal{I}_* \neq \emptyset$,

$$(5fb) \quad \sum_{j \in \mathcal{J}_*} \sum_{i \in \mathcal{C}(j) \cap \mathcal{I}_*} \dot{\psi}_{ij}(t) = \sum_{i \in \cup_{j \in \mathcal{J}_*} \mathcal{C}(j) \cap \mathcal{I}_*} \lambda_i - \sum_{j \in \mathcal{J}_*} \sum_{i \in \mathcal{C}(j) \cap \mathcal{I}_*} \mu_{ij} \psi_{ij}(t)$$

The meaning of (5fa) is as follows. Consider a set of server types \mathcal{J}^* . If a set of customer types \mathcal{I}^* consists of the “longest queues for \mathcal{J}^* ” (we will make this more precise), then servers in pools $j^* \in \mathcal{J}^*$, whenever they finish serving some customer, will immediately replace her with someone from a queue in \mathcal{I}^* . In this case, the total number of customers of types \mathcal{I}^* in service by servers of types \mathcal{J}^* will be increasing at the total rate of servicing all customers by servers in \mathcal{J}^* , less the rate of servicing customers of types \mathcal{I}^* by servers in \mathcal{J}^* . The requirements that \mathcal{I}^* needs to satisfy for this to be the case are, that there be no customer types outside \mathcal{I}^* with longer queues that servers in \mathcal{J}^* can serve. For example, a one-element set $\mathcal{I}^* = \{i^*\}$ is a valid choice for a one-element set $\mathcal{J}^* = \{j^*\}$ if and only if the customer type $i^* \in \mathcal{C}(j^*)$ has the (strictly) longest queue among all of the customer types that can be served by j^* .

The second equation, (5fb), describes the fact that if a set of server pools \mathcal{J}_* consists of the “least loaded server pools available to \mathcal{I}_* ”, then servers in pools $j^* \in \mathcal{J}_*$, whenever they finish serving some customer, will immediately replace her with someone from queue i^* . For example, a one-element set $\mathcal{J}_* = \{j_*\}$ is a valid choice for a one-element set $\mathcal{I}_* = \{i_*\}$ if and only if the server pool $j_* \in \mathcal{S}(i_*)$ has the (strictly) smallest load ρ_{j_*} among all of the server pools that can serve i_* .

3.2. Behavior in the vicinity of equilibrium point. We define the *equilibrium (invariant)* point of the underloaded ($\rho < 1$) fluid model to be the state $\psi_{ij} = \psi_{ij}^*$ and $q_i = q = 0$ for all $i \in \mathcal{I}$, $j \in \mathcal{J}$. (All other components of the fluid model are also constant and uniquely defined by (ψ_{ij}^*) and q .) Clearly, $\psi_{ij}(t) \equiv \psi_{ij}^*$ and $q_i(t) \equiv q$ is indeed a stationary fluid model. Desirable system behavior would be to have $(\psi_{ij}(t)) \rightarrow (\psi_{ij}^*)$ as $t \rightarrow \infty$.

Note that if the initial system state is in the vicinity of the equilibrium point (with $\rho < 1$), then there is no queuing in the system, and we can describe the system with just the variables $(\psi_{ij}(t))$. This will be true for at least some time (depending on ρ and the initial distance to the equilibrium point), because the fluid model is Lipschitz.

The following is a “state space collapse” result for the underloaded fluid model in the neighborhood of the equilibrium point.

Theorem 3.1. *Let $\rho < 1$. There exists a sufficiently small $\epsilon > 0$, depending only on the system parameters, such that for all sufficiently small δ the following holds. There exist $T_1 = T_1(\delta)$ and $T_2 = T_2(\delta)$, $0 < T_1 < T_2$, such that if the initial system state $(\psi_{ij}(0))$ satisfies*

$$|(\psi_{ij}(0)) - (\psi_{ij}^*)| < \delta,$$

then for all $t \in [T_1, T_2]$ the system state satisfies

$$|(\psi_{ij}(t)) - (\psi_{ij}^*)| < \epsilon, \quad \rho_j(t) = \rho_{j'}(t) \text{ for all } j, j' \in \mathcal{J}.$$

Moreover, $T_1 \downarrow 0$ and $T_2 \uparrow \infty$ as $\delta \downarrow 0$. The evolution of the system on $[T_1, T_2]$ is described by a linear ODE, specified below by (10).

If the fluid system is critically loaded ($\rho = 1$), it may have queues at equilibrium, and the equilibrium is non-unique. Namely, the definition of an equilibrium (invariant) point for $\rho = 1$ is the same as for the underloaded system, except the condition on the queues becomes $q_i = q$ for some constant $q \geq 0$. In the next Theorem 3.2 we will only consider the case of positive queues ($q > 0$) for the critically loaded fluid model.

Theorem 3.2. *Let $\rho = 1$, and consider an equilibrium point with $q > 0$. There exists a sufficiently small $\epsilon > 0$, depending only on the system parameters, such that for all sufficiently small $\delta > 0$ the following holds. There exist $T_1 = T_1(\delta)$ and $T_2 = T_2(\delta)$, $0 < T_1 < T_2$, such that if the initial system state satisfies*

$$|(\psi_{ij}(0)) - (\psi_{ij}^*)| < \delta, \quad |q_i(0) - q| < \delta \text{ for all } i \in \mathcal{I},$$

then for all $t \in [T_1, T_2]$ the system state satisfies

$$\begin{aligned} |(\psi_{ij}(t)) - (\psi_{ij}^*)| < \epsilon, \quad |q_i(t) - q| < \epsilon \text{ for all } i \in \mathcal{I}, \\ q_i(t) = q_{i'}(t) \text{ for all } i, i' \in \mathcal{I}. \end{aligned}$$

Moreover, $T_1 \downarrow 0$ and $T_2 \uparrow \infty$ as $\delta \downarrow 0$. The evolution of the system on $[T_1, T_2]$ is described by a linear ODE specified below by (11).

In the rest of this section and the paper, the values associated with a stationary fluid model, “sitting” at an equilibrium point, are referred to as *nominal*. For example, ψ_{ij}^* is the nominal occupancy (of pool j by type i), λ_i is the nominal arrival rate, λ_{ij} is the nominal routing rate (along activity (ij)), $\psi_{ij}^* \mu_{ij} = \lambda_{ij}$ is the nominal service rate (of type i in pool j), $\sum_j \psi_{ij}^* \mu_{ij} = \lambda_i$ is the nominal total service rate (of type i), ρ is the nominal total occupancy (of each pool j), etc.

Proof of Theorem 3.1. Let us choose a suitably small $\epsilon > 0$ (we will specify how small later). Because the fluid model trajectories are continuous, we can always choose some $T_2 > 0$ such that, for all sufficiently small $\delta > 0$, if $|(\psi_{ij}(0)) - (\psi_{ij}^*)| < \delta$ then $|(\psi_{ij}(t)) - (\psi_{ij}^*)| < \epsilon$ for all $t \leq T_2$. We will show that $\rho_j(t) = \rho_{j'}(t)$ for all $j, j' \in \mathcal{J}$, in $[T_1, T_2]$ for some T_1 depending on δ .

Consider $\rho_*(t) = \min_j \rho_j(t)$, $\rho^*(t) = \max_j \rho_j(t)$, and assume $\rho_*(t) < \rho^*(t)$. Let $\mathcal{J}_*(t) = \{j : \rho_j(t) = \rho_*(t)\}$. As long as $\rho_*(t) < \rho^*(t)$, $\mathcal{J}_*(t)$ is of course a strict subset of \mathcal{J} . The total arrival rate to servers of type $j \in \mathcal{J}_*(t)$ is $\sum_{i \in \cup_{j \in \mathcal{J}_*(t)} \mathcal{C}(j)} \lambda_i$. By the assumption of the connectedness of the basic activity tree, this is strictly greater (by a constant) than the nominal arrival rate $\sum_{i \in \mathcal{C}(j), j \in \mathcal{J}_*(t)} \lambda_{ij}$. The total rate of departures from those servers is $\sum_{i \in \mathcal{C}(j), j \in \mathcal{J}_*(t)} \mu_{ij} \psi_{ij}(t)$. For small ϵ , the assumption $|(\psi_{ij}(t)) - (\psi_{ij}^*)| < \epsilon$ implies that this is close to the nominal departure rate, so the arrival rate exceeds the service rate by at least a constant. (This determines what “suitably small” means for ϵ in terms of the system parameters.) Consequently, as long as $\rho_*(t) < \rho^*(t)$, the minimal load $\rho_*(t)$ is increasing at a rate bounded below by a constant. Similarly, as long as $\rho_*(t) < \rho^*(t)$, the maximal load $\rho^*(t)$ is decreasing at a rate bounded below by a constant. Therefore, the difference $\rho^*(t) - \rho_*(t)$ is decreasing at a rate bounded below by a constant whenever it is positive. Thus, in finite time $T_1 = T_1(\delta)$ we will arrive at a state $\rho_*(t) = \rho^*(t)$. (Clearly, $T_1(\delta) \rightarrow 0$ as $\delta \rightarrow 0$.) Since the function $\rho^*(\cdot) - \rho_*(\cdot)$ is Lipschitz (hence absolutely continuous), bounded below by 0, and (for $t \leq T_2$) has nonpositive derivative whenever it is differentiable, the condition $\rho_*(t) = \rho^*(t)$ will continue to hold for $T_1 \leq t \leq T_2$.

It remains to derive the differential equation, and to show that T_2 can be chosen depending on δ so that $T_2 \uparrow \infty$ as $\delta \downarrow 0$.

Once we are confined to the manifold $\rho_j(t) = \rho_{j'}(t) = \rho(t)$ for all t , the system evolution is determined in terms of only I independent variables. Decreasing ϵ if necessary to ensure that there is no queueing while $|(\psi_{ij}(t)) - (\psi_{ij}^*)| < \epsilon$, we can take the I variables to be $\psi_i(t) := \sum_j \psi_{ij}(t)$. Given $(\psi_i(t))$ we know $\rho(t)$ as $(\sum_i \psi_i(t)) / (\sum_j \beta_j)$. Consequently, we know $\sum_i \psi_{ij}(t) = \rho(t) \beta_j$ and $\sum_j \psi_{ij}(t) = \psi_i(t)$. On a tree, this allows us to solve for $\psi_{ij}(t)$; the relationship will clearly be linear, i.e.

$$(7) \quad (\psi_{ij}(t)) = M(\psi_i(t))$$

for some matrix M . For future reference, we define the (“load balancing”) linear mapping M from $y \in \mathbb{R}^I$ to $z = (z_{ij}, (ij) \in \mathcal{E}) \in \mathbb{R}^{I+J-1}$ as follows: $z = My$ is the unique solution of

$$(8) \quad \eta = \frac{\sum_i y_i}{\sum_j \beta_j}; \quad \sum_i z_{ij} = \eta \beta_j, \forall j; \quad \sum_j z_{ij} = y_i, \forall i.$$

The evolution of $\psi_i(t)$ is given by

$$(9) \quad \dot{\psi}_i(t) = \lambda_i - \sum_j \mu_{ij} \psi_{ij}(t), \quad \forall i.$$

(Follows from (5c) and the fact that $q_i(t) = 0$.) Then, by the above arguments we see that this entails (in matrix form)

$$(10) \quad (\dot{\psi}_i(t)) = (\lambda_i) + A_u(\psi_i(t)),$$

where A_u is an $I \times I$ matrix, $A_u = GM$. Here, G is a $I \times (I + J - 1)$ matrix with entries $G_{i,(kj)} = -\mu_{ij}$ if $i = k$, and $G_{i,(kj)} = 0$ otherwise.

It remains to justify the claim that $T_2(\delta) \uparrow \infty$ as $\delta \downarrow 0$. This follows from the fact that, as long as $t \geq T_1$ and $|(\psi_{ij}(t)) - (\psi_{ij}^*)| < \epsilon$, the evolution of the system is described by the linear ODE above. The solutions have the general form

$$\psi_{\mathcal{I}}(t) - \psi_{\mathcal{I}}^* = \exp(A_u(t - T_1))(\psi_{\mathcal{I}}(T_1) - \psi_{\mathcal{I}}^*), \quad \psi_{\mathcal{E}}(t) - \psi_{\mathcal{E}}^* = M(\psi_{\mathcal{I}}(t) - \psi_{\mathcal{I}}^*)$$

where M and A_u are constant matrices depending on the system parameters. Therefore, if $|\psi_{\mathcal{I}}(T_1) - \psi_{\mathcal{I}}^*| \leq \delta$ is sufficiently small, then the time it takes for $\psi_{\mathcal{E}}(t)$ to escape the set $|\psi_{\mathcal{E}}(t) - \psi_{\mathcal{E}}^*| < \epsilon$ can be made arbitrarily large. Since as $\delta \downarrow 0$ we have $T_1(\delta) \downarrow 0$, and the system trajectory is Lipschitz, taking $|\psi_{\mathcal{E}}(0) - \psi_{\mathcal{E}}^*| < \delta$ for small enough δ will guarantee that $|\psi_{\mathcal{I}}(T_1) - \psi_{\mathcal{I}}^*|$ is small, and hence we can choose $T_2(\delta) \uparrow \infty$. \square

The proof of Theorem 3.2 proceeds similarly; we outline only the differences.

Proof of Theorem 3.2. First, since we assume that $\epsilon > 0$ is sufficiently small and $|q_i(t) - q| < \epsilon$, $i \in \mathcal{I}$, for all $t \leq T_2$, we clearly have $\rho_j(t) = 1$, $j \in \mathcal{J}$, for all $t \leq T_2$. The equality of queue lengths in $[T_1, T_2]$ is shown analogously to the proof of $\rho_*(t) = \rho^*(t)$ for in the underloaded case. Namely, the smallest queue must increase and the largest queue must decrease (as long as not all $q_i(t)$ are equal), because it is getting less (resp. more) service than nominal (we choose ϵ small enough for this to be true provided $|(\psi_{ij}(t)) - (\psi_{ij}^*)| < \epsilon$). Thus, in $[T_1, T_2]$ we will have $q_i(t) = q_{i'}(t)$ for all $i, i' \in \mathcal{I}$.

The linear equation is modified as follows. We have

$$\dot{x}_i(t) = \lambda_i - \sum_j \mu_{ij} \psi_{ij}(t)$$

where $x_i(t) = \psi_i(t) + q_i(t)$. Since we know that all $q_i(t)$ are equal and positive, we have $q_i(t) = q(t) = \frac{1}{I}(\sum x_k(t) - \sum \beta_j)$, and therefore

$$\dot{\psi}_i(t) = \dot{x}_i(t) - \frac{1}{I} \sum_k \dot{x}_k(t).$$

The rest of the arguments proceed as above to give

$$(11) \quad (\dot{\psi}_i(t)) = (\lambda_i - \frac{1}{I} \sum_i \lambda_i) + A_c(\psi_i(t))$$

for the appropriate matrix A_c which can be computed explicitly from the basic activity tree. (Of course, in $[T_1, T_2]$, the trajectory $(\psi_{ij}(\cdot))$ uniquely determines $(\psi_i(\cdot))$, $(x_i(\cdot))$ and $(q_i(\cdot))$.)

Just as above, the existence of the linear ODE, together with the fact that $T_1 \downarrow 0$ as $\delta \downarrow 0$, implies that $T_2 \uparrow \infty$ as $\delta \downarrow 0$. \square

To compute the matrix M , and therefore the matrices A_u and A_c , we will find the following observation useful. If $(\psi_{ij}(t)) = M(\psi_i(t))$, then the common value $\rho(t) = \rho_j(t)$, $\forall j$, is

$$\rho(t) = \sum_i \psi_i(t) / \sum_j \beta_j.$$

This allows us to find the values $(\psi_{ij}(t))$ from $(\psi_i(t))$ as follows: if i is a customer-type leaf, then $\psi_{ij}(t) = \psi_i(t)$; if j is a server type leaf, then $\psi_{ij}(t) = \rho(t)\beta_j$; we now remove the leaf and continue with the smaller tree. Inductively, for an activity i_0j_0 we find

(12)

$$\psi_{i_0j_0}(t) = \sum_{i \preceq (i_0, j_0)} \psi_i(t) - \sum_{j \preceq (i_0, j_0)} \rho(t)\beta_j = \frac{1}{\sum \beta_j} \left(\sum_{i \preceq (i_0, j_0)} \sum_{j \preceq (j_0, i_0)} \psi_i(t)\beta_j - \sum_{i \preceq (j_0, i_0)} \sum_{j \preceq (i_0, j_0)} \psi_i(t)\beta_j \right)$$

Here, the relation \preceq is defined as follows. Suppose we disconnect the basic activity tree by removing the edge (i_0, j_0) . Then for any node k (either customer type or server type) we say $k \preceq (i_0, j_0)$ if it falls in the same component as i_0 ; otherwise, $k \preceq (j_0, i_0)$.

For example, consider the network in Figure 1. For it, we obtain

$$\begin{pmatrix} \psi_{A1} \\ \psi_{A2} \\ \psi_{B2} \end{pmatrix} = \begin{pmatrix} \frac{\beta_1}{\beta_1 + \beta_2} & \frac{\beta_1}{\beta_1 + \beta_2} \\ 1 - \frac{\beta_1}{\beta_1 + \beta_2} & -\frac{\beta_1}{\beta_1 + \beta_2} \\ 0 & 1 \end{pmatrix} \begin{pmatrix} \psi_A \\ \psi_B \end{pmatrix}.$$

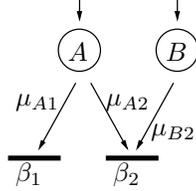


FIGURE 1. Example for calculation of the matrix M .

Since in underload we have

$$\dot{\psi}_i(t) = \lambda_i - \sum_j \mu_{ij}\psi_{ij}(t),$$

we obtain an expression for A_u , given in Lemma 3.3(i) just below.

Lemma 3.3. (i) The entries $(A_u)_{ii'}$ of the matrix A_u (for the underload case, $\rho < 1$) are as follows. The coefficient of ψ_i in $\dot{\psi}_i$ is

$$(A_u)_{ii} = -\frac{1}{\sum_j \beta_j} \sum_{j \in \mathcal{S}(i)} \mu_{ij} \sum_{j' \preceq (j, i)} \beta_{j'}.$$

The coefficient of $\psi_{i'}$ in $\dot{\psi}_i$ is

$$\begin{aligned} (A_u)_{ii'} &= \frac{1}{\sum_j \beta_j} \left[-\sum_{j \in \mathcal{S}(i), j \neq j_{ii'}} \mu_{ij} \sum_{j' \preceq (j, i)} \beta_{j'} + \mu_{ij_{ii'}} \sum_{j' \preceq (i, j_{ii'})} \beta_{j'} \right] \\ &= (A_u)_{ii} + \mu_{ij_{ii'}}. \end{aligned}$$

Here, $j_{ii'} \in \mathcal{S}(i)$ is the neighbor of i such that, after removing the edge $(i, j_{ii'})$ from the basic activity tree, nodes i and i' will be in different connected components. (Such a node is unique, since there is a unique path along the tree from i to i' .)

(ii) The matrix A_u is non-singular.

(iii) The matrix A_u depends only on (β_j) , (μ_{ij}) and the basic activity tree structure \mathcal{E} , and does not depend on (λ_i) and (ψ_{ij}^*) .

Proof. (i) In the proof of Theorem 3.1 we showed $A_u = GM$, where G is a $I \times (I + J - 1)$ matrix with entries $G_{i,(kj)} = -\delta_{ik}\mu_{ij}$, where δ_{ik} is the Kronecker's delta function, and M is the $I \times (I + J - 1)$ load-balancing matrix whose entries are determined from the expression (12). The form of the entries for A_u now follows. The equality between the two expressions for the off-diagonal entries is a consequence of the fact that, for all j' , exactly one of $j' \preceq (ij)$, $j' \preceq (ji)$ holds.

(ii) In the case $\rho < 1$, in the vicinity of the equilibrium point, the derivative $(\dot{\psi}_i) = (\lambda_i) + A_u(\psi_i)$ (which can be any real-valued I -dimensional vector, within a small neighborhood of the origin) uniquely determines (ψ_{ij}) , and then (ψ_i) as well. Indeed, we have the system of $I + J$ linear equations $\lambda_i - \sum_j \mu_{ij}\psi_{ij} = \dot{\psi}_i$, $\forall i$ and $\sum_i \psi_{ij} = \hat{\rho}\beta_j$, $\forall j$, for the $I + J$ variables $\hat{\rho}, (\psi_{ij})$. This system has unique solution, because $\hat{\rho}$ is uniquely determined by the workload derivative condition

$$\sum_i \nu_i \dot{\psi}_i = \sum_i \nu_i \lambda_i - \sum_j \hat{\rho} \alpha_j,$$

and then the values of ψ_{ij} are determined by sequentially “eliminating” leaves of the basic activity tree.

(iii) Follows from (i). □

The corresponding expression for A_c is less elegant:

Lemma 3.4. (i) *The entries $(A_c)_{ii'}$ of the matrix A_c (for the critical load case, $\rho = 1$) are as follows:*

$$(13) \quad (A_c)_{ii'} = (A_u)_{ii'} - \frac{1}{I} \sum_k (A_u)_{ki'}.$$

(ii) *The matrix A_c has rank $I - 1$. The $(I - 1)$ -dimensional subspace $L = \{y \mid \sum_i y_i = 0\}$ is invariant under the transformation A_c , i.e. $A_c L \subseteq L$. Letting π denote the matrix of the orthogonal projection (along $(1, \dots, 1)^\dagger$) onto L , we have $A_c = \pi A_u$. Restricted to L , the transformation A_c is invertible.*

(iii) *The linear transformation A_c , restricted to subspace L , depends only on (μ_{ij}) and the basic activity tree structure \mathcal{E} , and does not depend on (β_j) , (λ_i) and (ψ_{ij}^*) .*

Proof. (i) The fluid model here is such that there are always non-zero queues, which are equal across customer types. We can write

$$(14) \quad \dot{\psi}_i(t) = \dot{x}_i(t) - \frac{1}{I} \sum_k \dot{x}_k(t) = (\lambda_i - \sum_j \mu_{ij}\psi_{ij}(t)) - \frac{1}{I} \sum_k (\lambda_k - \sum_j \mu_{kj}\psi_{kj}(t)),$$

which implies (13).

(ii) First of all, it is not surprising that A_c does not have full rank: the linear ODE defining A_c is such that $\sum_i \psi_i(t) = \sum_j \beta_j$ at all times, so there are at most $(I - 1)$ degrees of freedom in the system. Also, it will be readily seen that (13) asserts precisely that $A_c = \pi A_u$. Since A_u is invertible and π has rank $I - 1$, their composition has rank $I - 1$. Since the image of A_c is contained in L , the image of A_c (as a map from \mathbb{R}^I) must be equal to all of L .

It remains to check that A_c restricted to L still has rank $I - 1$. To see this, we observe that the simple eigenvalue 0 of A_c has as its unique right eigenvector the vector $A_u^{-1}(1, 1, \dots, 1)^\dagger$. We will be done once we show that this eigenvector does not belong to L . Suppose instead that $A_u v = (1, 1, \dots, 1)^\dagger$ for some $v \in L$, $\sum_i v_i = 0$. Then, for a small $\epsilon > 0$, the state $\psi_{\mathcal{T}}^* - \epsilon v$ (with balanced pool loads, all equal to the optimal ρ) would be such that the derivatives of

all components ψ_i would be strictly negative. This is, however, impossible, because the total rate at which the system workload is served, must be zero:

$$\frac{d}{dt} \sum_i \nu_i \psi_i = \sum_i \nu_i \lambda_i - \sum_j \rho \alpha_j = 0.$$

(iii) The specific expression (13) for A_c may depend on the pool sizes (β_j) . However, A_c is a singular $I \times I$ matrix, and our claim is only about the transformation of the $(I - 1)$ -dimensional subspace L that A_c induces; this transformation does *not* depend on (β_j) , as the following argument shows.

Pick any $(ij) \in \mathcal{E}$. Modify the original system by replacing β_j by $\beta_j + \delta$ and λ_i by $\lambda_i + \delta \mu_{ij}$; this means that the nominal ψ_{ij}^* is replaced by $\psi_{ij}^* + \delta$. Then, using notation $\gamma_i(t) = \psi_i(t) - \psi_i^*$, the linear ODE

$$(15) \quad (\dot{\gamma}_i(t)) = A(\gamma_i(t)),$$

which we obtain from the ODE (14) for the original and modified systems, has exactly the same matrix A , which implies $A = A_c$. Thus, the transformation A_c must not depend on β_j .

An alternative argument is purely analytic. Recall that to compute $(A_u)_{ij}$ we used (12). In critical load, we have $\rho(t) \equiv 1$, so the (left) equation (12) for $\psi_{i_0 j_0}(t)$ simplifies to

$$(16) \quad \psi_{i_0 j_0}(t) = \sum_{i \preceq (i_0, j_0)} \psi_i(t) - \sum_{j \preceq (i_0, j_0)} \beta_j.$$

If we substitute this in the right-hand side of (14), we will obtain a different expression for $\dot{\psi}_i(t)$. While its constant term will depend on $\beta_{\mathcal{J}}$, the linear term will not, since the linear term of (16) does not depend on $\beta_{\mathcal{J}}$. That is, we have found a way of writing down a matrix for A_c which clearly does not depend on the $\beta_{\mathcal{J}}$. \square

3.3. Definition of local stability. We say that the (fluid) system is *locally stable*, if all fluid models starting in a sufficiently small neighborhood of an equilibrium point (which is unique for $\rho < 1$; and for $\rho = 1$ we consider any equilibrium point with equal queues $q > 0$) are such that, for fixed constant $C > 0$,

$$|(\psi_{ij}(t)) - (\psi_{ij}^*)| \leq \Delta_0 e^{-Ct},$$

where $\Delta_0 = |(\psi_{ij}(0)) - (\psi_{ij}^*)| + |(q_i(0)) - (q, \dots, q)^\dagger|$. Note that in the case $\rho = 1$ it is *not* required that $q_i(t) \rightarrow q$, for q associated with the chosen equilibrium point. However, local stability will guarantee convergence of queues $q_i(t) \rightarrow \bar{q}$, with some $\bar{q} > 0$ possibly different from q . Indeed, the exponentially fast convergence $\psi_{\mathcal{E}}(t) \rightarrow \psi_{\mathcal{E}}^*$ of the occupancies to the nominal, guarantees that for some fixed constant $C_1 > 0$, any i and any $s \geq t \geq 0$:

$$|x_i(s) - x_i(t)| \leq \int_t^s |\lambda_i - \sum_j \mu_{ij} \psi_{ij}(\xi)| d\xi \leq C_1 \Delta_0 e^{-Ct}.$$

Therefore, each $x_i(t)$, and then each $q_i(t)$, also converges exponentially fast. Then we can apply Theorem 3.2 to show that all $q_i(t)$ must be equal starting some time point; therefore they converge to the same value \bar{q} , which is such that that $|\bar{q} - q| \leq C_0 \Delta_0$ for some constant $C_0 > 0$ depending only on the system parameters. In other words, local stability guarantees convergence to an equilibrium point not too far from the “original” one. (We omit further detail, which are rather straightforward.)

By Theorems 3.1 and 3.2 we see that the local stability is determined by the stability of a linear ODE, which in turn is governed by the eigenvalues of the matrix A_u or A_c . We will call matrix A_u stable if all its eigenvalues have negative real part. We call matrix A_c stable if all its eigenvalues have negative real part, except one simple eigenvalue 0.¹ In this terminology, *the local stability of the system is equivalent to the stability of the matrix A in question (either A_u or A_c)*. On the other hand, if A has an eigenvalue with positive real part, the ODE has solutions diverging from equilibrium (ψ_i^*) exponentially fast; if A has (a pair of conjugate) pure imaginary eigenvalues, the ODE has oscillating, never converging solutions.

3.4. Fluid model as a fluid limit. In this section we show that the set of fluid models defined in Section 3.1 contains (in the sense specified shortly) all possible limits of “fluid scaled” processes. We consider a sequence of systems indexed by r , with the input rates being $\lambda_i^r = r\lambda_i + o(r)$, server pool sizes being $\beta_j r$, and the service rates μ_{ij} unchanged with r . Recall the notation in Section 2.2. We also add the following notation:

$A_i^r(t)$: the number of customers of type i who have entered the system by time t (a Poisson process of rate λ_i^r)

$S_{ij}^r(t)$: the number of customers of type i who have been served by servers of type j if a total time rt has been spent on these services (a Poisson process of rate $\mu_{ij}r$)

Let $\Pi_i^{(a)}(\cdot)$, $i \in \mathcal{I}$, and $\Pi_{ij}^{(s)}(\cdot)$, $(ij) \in \mathcal{E}$, be independent unit-rate Poisson processes. We can assume that, for each r ,

$$A_i^r(t) = \Pi_i^{(a)}(\lambda_i^r t), \quad S_{ij}^r(t) = \Pi_{ij}^{(s)}(\mu_{ij} r t).$$

Then, by the functional strong law of large numbers, with probability 1, uniformly on compact subsets of $[0, \infty)$,

$$(17) \quad \frac{1}{r} A_i^r(t) \rightarrow \lambda_i t, \quad \frac{1}{r} S_{ij}^r(t) \rightarrow \mu_{ij} t.$$

We consider the following scaled processes:

$$x_i^r(t) = \frac{1}{r} X_i^r(t) \quad q_i^r(t) = \frac{1}{r} Q_i^r(t) \quad \psi_{ij}^r(t) = \frac{1}{r} \Psi_{ij}^r(t) \quad \rho_j^r(t) = \frac{1}{r} \Xi_j^r(t) \quad a_i^r(t) = \frac{1}{r} A_i^r(t).$$

Theorem 3.5. *Suppose*

$$\{(x_i^r(0)), (q_i^r(0)), (\psi_{ij}^r(0)), (\rho_j^r(0))\} \rightarrow \{(x_i(0)), (q_i(0)), (\psi_{ij}(0)), (\rho_j(0))\}.$$

Then w.p.1 any subsequence of $\{r\}$ contains a further subsequence along which u.o.c.,

$$\{(a_i^r(\cdot)), (x_i^r(\cdot)), (q_i^r(\cdot)), (\psi_{ij}^r(\cdot)), (\rho_j^r(\cdot))\} \rightarrow \{(a_i(\cdot)), (x_i(\cdot)), (q_i(\cdot)), (\psi_{ij}(\cdot)), (\rho_j(\cdot))\},$$

where the limiting trajectory (in the RHS) is a fluid model.

Proof. Given property (17), the probability 1, u.o.c., convergence along a subsequence to a Lipschitz continuous set of functions easily follows. The only non-trivial properties of a fluid model that need to be verified for the limit are (5f). Let us consider a regular time t : namely, such that all the components of a limit trajectory have derivatives, and moreover the minimums and maximums over any subset of components have derivatives as well. Consider

¹A matrix A with all eigenvalues having negative real part is usually called *Hurwitz*. So, A_u stability is equivalent to A_u being Hurwitz; while A_c stability definition is slightly different, due to A_c singularity. A symmetric matrix A is Hurwitz if and only if it is negative definite, but neither A_u nor A_c is, in general, symmetric.

a sufficiently small interval $[t, t + \Delta t]$, and consider the behavior of the (fluid-scaled) pre-limit trajectory in this interval. Then, it is easy to check that the conditions (5f) on the derivatives must hold; the argument here is very standard – we omit details. \square

4. SPECIAL CASES IN WHICH FLUID MODELS ARE STABLE

In this section we analyze two special cases of the system parameters, for which we demonstrate convergence results. In Section 4.1 we consider the case when there exists a set of positive μ_j , $j \in \mathcal{J}$, such that $\mu_{ij} = \mu_j$ for $(ij) \in \mathcal{E}$ (i.e. the service rate μ_{ij} is constant across all $i \in \mathcal{C}(j)$); we show global convergence of fluid models to equilibrium. In Section 4.2 we consider the case when there exists a set of positive μ_i , $i \in \mathcal{I}$, such that $\mu_{ij} = \mu_i$ for $(ij) \in \mathcal{E}$ (i.e. the service rate μ_{ij} is constant across all $j \in \mathcal{S}(i)$); we show local stability of the fluid model (i.e. stability of A_u and A_c).

4.1. Global stability in the case $\mu_{ij} = \mu_j$, $(ij) \in \mathcal{E}$. We call the system *globally stable* if any fluid model, with arbitrary initial state, converges to an equilibrium point as $t \rightarrow \infty$. (This of course implies $\rho_j(t) \rightarrow \rho$ for all $j \in \mathcal{J}$ and $\psi_{ij}(t) \rightarrow \psi_{ij}^*$ for all $i \in \mathcal{I}$, $j \in \mathcal{J}$. Note that, in underload, the definition necessarily implies $q_i(t) \rightarrow 0$ for all $i \in \mathcal{I}$, while in critical load it requires $q_i(t) \rightarrow q$ for all $i \in \mathcal{I}$ and some $q \geq 0$.)

Theorem 4.1. *The system with $\mu_{ij} = \mu_j$, $(ij) \in \mathcal{E}$, is globally stable both for $\rho < 1$ and for $\rho = 1$. In addition, the system is locally stable as well (i.e. the matrices A_u and A_c are stable).*

Proof. Consider the underloaded system, $\rho < 1$, first. First, we show that the lowest load cannot stay too low. Suppose the minimal load $\rho_*(t) \equiv \min_j \rho_j(t)$ is smaller than ρ , and let $\mathcal{J}_*(t) \equiv \{j : \rho_j(t) = \rho_*(t)\}$. Then all customer types in $\mathcal{C}(\mathcal{J}_*(t)) \equiv \bigcup_{j \in \mathcal{J}_*(t)} \mathcal{C}(j)$ are routed to server pools in $\mathcal{J}_*(t)$, so the total arrival rate “into” $\mathcal{J}_*(t)$ is no less than nominal; on the other hand, since $\mu_{ij} = \mu_j$ and server occupancy is lower than nominal, the total departure rate “from” $\mathcal{J}_*(t)$ is smaller than nominal. This shows that if $\rho_* < \rho - \epsilon < \rho$, then $\dot{\rho}_* > \delta > 0$, where $\delta \geq c\epsilon$ for some constant $c > 0$ (depending on the system parameters). That is, if $\rho_*(t) < \rho$ then $\dot{\rho}_*(t) \geq c(\rho - \rho_*(t))$, so $\rho_*(t)$ is bounded below by a function converging exponentially fast to ρ .

Consider a fixed, sufficiently small $\epsilon > 0$; we know that $\rho_*(t) \geq \rho - \epsilon$ for all large times t . If some customer class i has a queue $q_i(t) > 0$, then all server classes $j \in \mathcal{S}(i)$ have $\rho_j(t) = 1$. It is now easy to see that the system is serving customers faster than they arrive (because $\rho < 1$ and ϵ is small). This easily implies that all $q_i(t) = 0$ after a finite time.

In the absence of queues, we can analyze $\rho^*(t) = \max_j \rho_j(t)$ similarly to the way we treated $\rho_*(t)$; namely, we show that $\rho^*(t)$ is bounded above by a function converging exponentially fast to ρ , which tells us that $\rho_j(t) \rightarrow \rho$ for all j . Once all $\rho_j(t)$ are close enough to ρ , we can use the argument essentially identical to that in the proof of Theorem 3.1 to conclude that, after a further finite time, we will have $\rho_j(t) = \rho_{j'}(t)$ for all j, j' . (The argument is even simpler, because, unlike in Theorem 3.1, where it was required that $(\psi_{ij}(t))$ were close to nominal, here it suffices that $(\rho_j(t))$ are close to nominal, because of the $\mu_{ij} = \mu_j$ assumption.) With $\rho(t) = \rho_j(t)$, $\forall j$, we then have for the total amount of “fluid” in the system:

$$(d/dt) \sum_j \beta_j \rho(t) = \sum_i \lambda_i - \sum_j \beta_j \rho(t) \mu_j.$$

This is a simple linear ODE for $\rho(t)$, which implies that (after a finite time) $\rho(t) - \rho = c_1 \exp(-c_2 t)$, with constant $c_2 > 0$ and c_1 . This in particular means that $\dot{\rho}_j(t) = \dot{\rho}(t) \rightarrow 0$. Denote by $\hat{\lambda}_{ij}(t)$ the rate at which fluid i arrives at pool j , namely

$$(18) \quad \hat{\lambda}_{ij}(t) = \mu_j \psi_{ij}(t) + \dot{\psi}_{ij}(t);$$

at any large t we have $\sum_j \hat{\lambda}_{ij}(t) = \lambda_i$. Then, for each j ,

$$\sum_i \hat{\lambda}_{ij}(t) = \sum_i \mu_j \psi_{ij}(t) + \sum_i \dot{\psi}_{ij}(t) = \beta_j \mu_j \rho_j(t) + \beta_j \dot{\rho}_j(t) \rightarrow \beta_j \mu_j \rho = \sum_i \lambda_{ij}.$$

This is only possible if each $\hat{\lambda}_{ij}(t) \rightarrow \lambda_{ij}$. But then the ODE (18) implies $\psi_{ij}(t) \rightarrow \psi_{ij}^*$.

Now, consider a critically loaded system, $\rho = 1$. Essentially same argument as above tells us that, as long as not all queues $q_i(t)$ are equal, each of the longest queues gets more service than the arrival rate into it, and so $q^*(t) = \max q_i(t)$ has strictly negative, bounded away from 0 derivative. If all $q_i(t)$ are equal and positive, then $q^*(t) = 0$. We see that $q^*(t)$ is non-increasing, and so $q^*(t) \downarrow q \geq 0$. We also have $\rho_*(t) \rightarrow \rho = 1$ exponentially fast. (Same proof as above applies.) These facts easily imply convergence to an equilibrium point. We omit further detail.

Examination of the above proof shows that it implies the following property, for both cases $\rho < 1$ and $\rho = 1$. For any fixed equilibrium point (with $q > 0$ if $\rho = 1$), there exists a sufficiently small $\epsilon > 0$ such that for all sufficiently small $\delta > 0$, any fluid model starting in the δ -neighborhood of the equilibrium point, first, never leaves the ϵ -neighborhood of the equilibrium point and, second, converges to an equilibrium point (possibly different from the “original” one, if $\rho = 1$). This property cannot hold, unless the system is locally stable (see Section 3.3). □

4.2. Local stability in the case $\mu_{ij} = \mu_i$, $(ij) \in \mathcal{E}$.

Theorem 4.2. *Assume $\rho < 1$ and $\mu_{ij} = \mu_i$ for $(ij) \in \mathcal{E}$. Then the system is locally stable (i.e. A_u is stable).*

Proof. We have

$$\dot{\psi}_i(t) = \lambda_i - \mu_i \psi_i(t)$$

and A_u is simply a diagonal matrix with entries $-\mu_i$. □

Theorem 4.3. *Assume $\rho = 1$ and $\mu_{ij} = \mu_i$ for $(ij) \in \mathcal{E}$. Then the system is locally stable (i.e., A_c is stable).*

Proof. As seen in the proof of Theorem 4.2, the matrix A_u in this case is diagonal with entries $-\mu_i$. By Lemma 3.4, A_c has off-diagonal entries $(A_c)_{ii'} = \mu_{i'}/I$ and diagonal entries $-\mu_i(1 - 1/I)$. That is, its off-diagonal entries are strictly positive. Therefore, $A_c + \eta I$ for some large enough constant $\eta > 0$ (where I is the identity matrix) is a positive matrix. By Perron-Frobenius theorem [11, Chapter 8], $A_c + \eta I$ has a real eigenvalue $p + \eta$ with the property that any other eigenvalue of $A_c + \eta I$ is smaller than $p + \eta$ in absolute value (and in particular has real part smaller than $p + \eta$). Moreover, the associated *left* eigenvector w is strictly positive, and is the unique (up to scaling) strictly positive left eigenvector of $A_c + \eta I$. Translating these statements to A_c , we get: A_c has a real eigenvalue p ; all other eigenvalues of A_c have real part smaller than p ; A_c has unique (up to scaling) strictly positive left eigenvector w ; and the eigenvalue of w is p .

Now, A_c has a positive left eigenvector with eigenvalue 0, namely $(1, 1, \dots, 1)$. Therefore, we must have $p = 0$, and we conclude that all other (i.e., non-zero) eigenvalues of A_c have real part smaller than 0, as required. \square

5. FLUID MODELS FOR GENERAL μ_{ij} : LOCAL INSTABILITY EXAMPLES.

In Sections 4.1, 4.2 we have shown that the matrices A_u and A_c are stable in the cases $\mu_{ij} = \mu_j$, $(ij) \in \mathcal{E}$ and $\mu_{ij} = \mu_i$, $(ij) \in \mathcal{E}$. Since the entries of A_u , A_c depend continuously on μ_{ij} via Lemmas 3.3, 3.4, and the eigenvalues of a matrix depend continuously on its entries, we know that the matrices will be stable for all parameter settings sufficiently close to those special cases. Therefore, there exists a non-trivial parameter domain of local stability. One might consider it to be a reasonable conjecture that local stability holds for any parameters. It turns out, however, that this conjecture is false. We will now construct examples to demonstrate that, in general, the system can be locally unstable.

Remark 5.1. In examples below, we will specify the parameters $\mu_{\mathcal{E}}$ and sometimes $\beta_{\mathcal{J}}$, but not $\lambda_{\mathcal{I}}$. It is easy to construct values of $\lambda_{\mathcal{I}}$ which will make all of the activities in \mathcal{E} basic; simply pick a strictly positive vector $\psi_{\mathcal{E}}$, such that all loads $\sum_i \psi_{ij}/\beta_j$ are equal, and set $\lambda_i = \sum_j \psi_{ij}\mu_{ij}$. Lemmas 3.3 (iii) and 3.4 (iii) guarantee that the specific values of $\lambda_{\mathcal{I}}$ do not affect the matrices A_u , A_c . In critical load, we also do not need to specify $\beta_{\mathcal{J}}$.

Local instability example 1. Consider a system with 3 customer types A, B, C and 4 server types 1 through 4, connected $1 - A - 2 - B - 3 - C - 4$. Set $\beta_1 = 0.97$ and $\beta_2 = \beta_3 = \beta_4 = 0.01$. Set $\mu_{A1} = \mu_{B2} = \mu_{C3} = 1$, and $\mu_{A2} = \mu_{B3} = \mu_{C4} = 100$. (See Figure 2.) On the other hand, we compute by Lemma 3.3

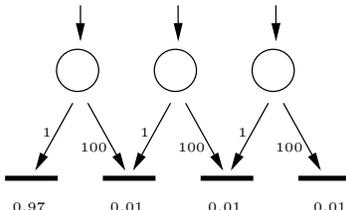


FIGURE 2. System with three customer types whose underload equilibrium is unstable

$$A_u = \begin{pmatrix} -1.99 & -0.99 & -0.99 \\ 97.02 & -2.98 & -1.98 \\ 96.03 & 96.03 & -3.97 \end{pmatrix}$$

with eigenvalues $\{-17.8, 4.45 \pm 23.4i\}$. Therefore by Theorem 3.1, the system with these parameters is described by an unstable ODE in the neighborhood of its equilibrium point.

We now show that this is a minimal instability example, in the sense made precise by the following

Lemma 5.2. *Consider an underloaded system, $\rho < 1$.*

(i) *Let $I \geq 2$. Any customer type i that is a leaf in the basic activity tree, does not affect the local stability of the system. Namely, let us modify the system by removing type i , and then modifying (if necessary) input rates λ_k of the remaining types $k \in \mathcal{I} \setminus i$ so that the basic activity tree of the modified system is $\mathcal{E} \setminus (ij)$, where (ij) is the (only) edge in \mathcal{E} adjacent to*

- i.* Then, the original system is locally stable if and only if the modified one is.
(ii) A system with two (or one) non-leaf customer types is locally stable.

Proof. (i) If type i is a leaf, the equation for $\psi_i(t)$ is simply $\dot{\psi}_i(t) = \lambda_i - \mu_{ij}\psi_i(t)$. This means (setting $i = 1$) that $(1, 0, \dots, 0)^\dagger$ is an eigenvector of A_u with eigenvalue $-\mu_{ij}$. Further, it is easy to see that: (a) the rest of the eigenvalues of A_u are those of matrix $A_u^{(-i)}$ obtained from A_u by removing the first row and first column; and (b) $A_u^{(-i)}$ is exactly the “ A_u -matrix” for the modified system.

(ii) We can assume that there are no customer type leaves. The case $I = 1$ is trivial (and is covered by Theorem 4.1), so let $I = 2$. Throughout the proof, the pool sizes β_j are fixed. From Theorem 4.1 we know that for a certain set of service rate values (namely, $\mu_{ij} = \mu_j$, $(ij) \in \mathcal{E}$), the matrix A_u is stable. Suppose that we continuously vary the parameters μ_{ij} from those initial values to the values of interest, without ever making $\mu_{ij} = 0$. If we assume that the final matrix A_u is *not* stable, then as we change μ_{ij} the (changing) matrix A_u acquires at some point two purely imaginary eigenvalues. If the eigenvalues of A_u are purely imaginary, we must have $\text{trace}(A_u) = 0$. However, as seen from the form of A_u in Lemma 3.3, the diagonal entries of A_u are always negative, and therefore $\text{trace}(A_u) < 0$. The contradiction completes the proof. \square

An argument similar to the above proof also allows us to explain how the instability example 1 was found. In degree 3, let the characteristic polynomial of A_u be $x^3 - c_2x^2 + c_1x - c_0$. A necessary and sufficient condition for all roots of the polynomial to have negative real parts is: $-c_2, c_1, -c_0 > 0$ and $c_2c_1 < c_0$ (see [3, A1.1.1]). A necessary and sufficient condition for the “boundary case” between stability and instability (i.e. the condition for a pair of conjugate purely imaginary roots) is $c_2c_1 = c_0$. Using Lemma 3.3 we can evaluate the characteristic polynomial symbolically and use the resulting expression to find parameters for which $c_2c_1 = c_0$ will hold. See online [14] for the computations.

It is possible to construct an instability example with more reasonable values of β_j , μ_{ij} , although it will be bigger. Figure 3 shows the diagram. The associated matrix A_u and its eigenvalues can also be found online [14].

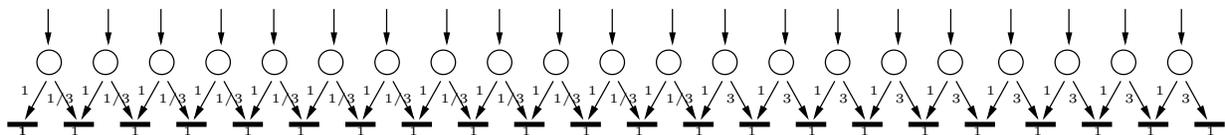


FIGURE 3. System with $\beta_j = 1$ and $\mu_{ij} \in \{1/3, 1, 3\}$ whose underload equilibrium is unstable. There are 21 customer types.

We do not have an explicit characterization of the local instability domain, beyond the necessity of $I \geq 3$.

We now analyze the critically loaded system $\rho = 1$ with queues, i.e. the stability of the matrix A_c . Recall that the transformation A_c , restricted to subspace $\{y \mid \sum_i y_i = 0\}$, and then the stability of A_c , does not depend on the values of β_j , so it suffices to specify the values μ_{ij} .

Local instability example 2. Consider the network of Figure 4, which has 5 customer types A through E and 4 server types 1 through 4, connected $A-1-B-2-C-3-D-4-E$, with the following parameters:

$$\begin{array}{cccc} \mu_{A1} = 1 & \mu_{B1} = 100 & \mu_{B2} = 1 & \mu_{C2} = 100 \\ \mu_{C3} = 1 & \mu_{D3} = 100 & \mu_{D4} = 10000 & \mu_{E4} = 100 \end{array}$$

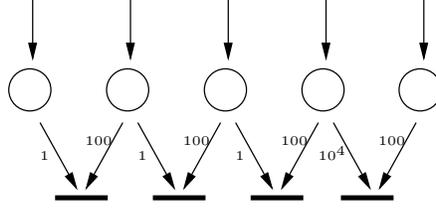


FIGURE 4. System with five customer types whose critical load equilibrium is unstable

The matrix A_c , computed from Lemma 3.4 will be given by

$$A_c = \frac{1}{20} \begin{pmatrix} 9389 & 9805 & 10201 & 10597 & -29003 \\ 10894 & 9290 & 9706 & 10102 & -29498 \\ 10399 & 10795 & 9191 & 9607 & -29993 \\ -40091 & -39695 & -39299 & -40903 & 119497 \\ 9409 & 9805 & 10201 & 10597 & -31003 \end{pmatrix}$$

and the eigenvalues of A_c are $\{0, -16.88, -2190.05, 2.565 \pm 23.23i\}$.

Again, the above example 2 is in a sense minimal:

Lemma 5.3. *Consider a critically loaded system, $\rho = 1$.*

(i) *Let $J \geq 2$. Any server type j that is a leaf in the basic activity tree does not affect the local stability of the system. Namely, let us modify the system by removing type j , and then replacing λ_i for the unique i adjacent to j by $\lambda_i - \beta_j \mu_{ij}$. Then, the original system is locally stable if and only if the modified one is.*

(ii) *Consider a system labeled S . We say that a system S' is an expansion of system S if it is obtained from S by the following modification. We pick one server type j and one customer type i adjacent to it in \mathcal{E} ; we “split” type j into two types j' and j'' ; we “connect” type i to both j' and j'' ; each of the remaining types $i' \in \mathcal{C}(j) \setminus i$ we connect to either j' or j'' (but not both); if $(i'j')$ (resp. $(i'j'')$) is a new edge, we set $\mu_{i'j'} = \mu_{i'j}$ (resp. $\mu_{i'j''} = \mu_{i'j}$). Then, S is locally stable if and only if S' is.*

(iii) *A system with four or fewer customer types is locally stable.*

Proof. (i) The argument here is a “special case” of the one used to show the independence of transformation A_c (restricted to $(I - 1)$ -dimensional invariant subspace) from (β_j) in the proof of Lemma 3.4. Namely, it is easy to check that the original and modified system share exactly same ODE (15).

(ii) Again, it is easy to see that the two systems share the same ODE (15).

(iii) We can assume that there are no server-type leaves, so that the tree \mathcal{E} has only customer-type leaves, of which it can have two, three, or four.

If it has four customer type leaves, then the tree has a total of four edges, hence five nodes, i.e. a single server pool, to which all the customer types are connected.

If the tree has three customer type leaves, then letting k be the number of edges from the fourth customer type, we have $k + 3$ total edges, so $k + 4$ nodes, of which k are server types. That is, the non-leaf customer type is connected to all of the server types. Since there are

no server type leaves, we must have $k \leq 3$; since we are assuming the fourth customer type is not a leaf, we must have $k \geq 2$; thus, $k = 2$ or $k = 3$.

The last case is of two customer type leaves. Letting k, l be the number of edges coming out of the other customer types, we have $k + l + 2$ edges. On the other hand, since each server type has at least 2 edges coming out of it, we have at most $(k + l + 2)/2$ server types, so at most $(k + l + 2)/2 + 4$ nodes. Thus, we have $(k + l + 2) + 1 \leq (k + l + 2)/2 + 4$, or $k + l + 2 \leq 6$, giving $k = l = 2$ (since they must both be ≥ 2).

We summarize the possibilities in Figure 5. Note that the bottom-left system can be obtained by a sequence of expansions from each of the top-left systems, and so this is the only system we need to consider to establish local stability for all 3- and 4-leaf cases. Thus, in total, the only two systems that need to be considered are bottom-left and right. In both of

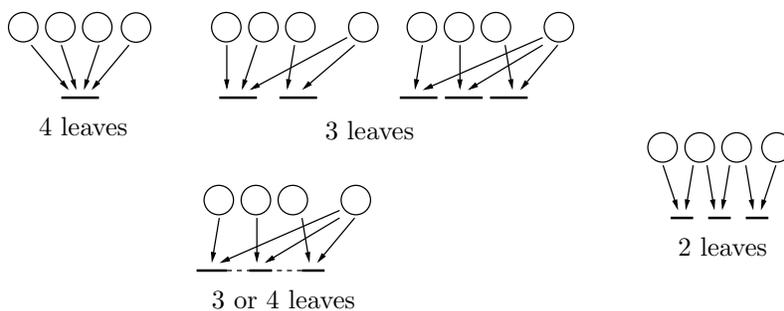


FIGURE 5. Possible arrangements of four customer types.

the resulting cases, we can use Lemma 3.4 to write out A_c and its characteristic polynomial explicitly. The characteristic polynomial will have degree 4, but one of its roots is 0, so we can reduce it to degree 3. We then symbolically verify that the cited above stability criterion [3, A1.1.1] for degree 3 polynomials, is satisfied. See online [14] for the details. \square

An argument similar to that in the above proof allows us to explain how the instability example 2 was found. We seek a condition satisfied by the coefficients of a degree 4 polynomial with two imaginary roots. Letting the polynomial be $x^4 - c_1x^3 + c_2x^2 - c_3x + c_4$, and letting the roots be $\eta_1, \eta_2, \pm iz$ (where η_1 and η_2 may be real or complex conjugates, and $z \in \mathbb{R}$), we see that $c_1 = \eta_1 + \eta_2$, $c_2 = \eta_1\eta_2 + z^2$, $c_3 = (\eta_1 + \eta_2)z^2$, and $c_4 = \eta_1\eta_2z^2$. This implies the relation $c_4c_1^2 + c_3^2 - c_1c_2c_3 = 0$, and we can find the parameters for which this is true. (The symbolic calculation will involve rather a lot of terms.) We remark that, whereas for degree 3 polynomials the condition $c_2c_1 - c_0 = 0$ is necessary and sufficient for the existence of two imaginary roots [3, A1.1.1], the condition we derive here for degree 4 polynomials is necessary, but not sufficient. (For example, the polynomial $(x - 1)^2(x + 1)^2$ has $c_1 = c_3 = 0$, so $c_4c_1^2 + c_3^2 - c_1c_2c_3 = 0$, but it has no imaginary roots.) Thus, checking the sign of the corresponding expression alone is insufficient to determine whether the system is unstable, but is a useful way of narrowing down the parameter ranges.

Finally, it is possible to construct a single system which will be unstable both for $\rho < 1$ and for $\rho = 1$ with positive queues. For the local stability of the underloaded system, the leaves of the basic activity tree corresponding to customer types are irrelevant (the corresponding occupancy on the sole available server class converges to nominal exponentially). On the other hand, for the critically loaded system, the leaves corresponding to server pools are irrelevant, since the corresponding server is fully occupied by its unique available customer

type. This observation allows us to merge the above two systems into a single one which is unstable both in underloaded and in the critically loaded case.

Consider a system with 5 customer types A through E and 5 server types 0 through 4 connected as $0 - A - 1 - B - 2 - C - 3 - D - 4 - E$, with $\mu_{A0} = 100$ and the remaining μ_{ij} as in the critically loaded case. Set $\beta_3 = 0.96$ while $\beta_0, \beta_1, \beta_2, \beta_4 = 0.01$. (See Figure 6.) By the above discussion, this system must be unstable for $\rho = 1$ and positive queues. We

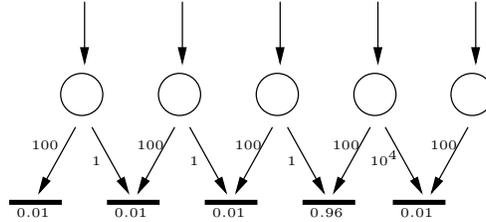


FIGURE 6. System with five customer types whose underload and critical load equilibrium points are both unstable

therefore need to consider only the first 4 customer types (E is a customer type leaf and doesn't matter) in underload. We compute

$$A_u = \begin{pmatrix} -1.99 & -0.99 & -0.99 & -0.99 \\ 97.02 & -2.98 & -1.98 & -1.98 \\ 96.03 & 96.03 & -3.97 & -2.97 \\ -99 & -99 & -99 & -199 \end{pmatrix}$$

and eigenvalues are $\{-14.6, -201.1, 3.91 \pm 18.1i\}$.

While we showed above that sufficiently small systems are at least locally stable, we will show now that, in the underload case, any sufficiently large system is locally unstable for some parameter settings.

Lemma 5.4. *In underload ($\rho < 1$), any shape of basic activity tree that includes a locally unstable system (i.e., with A_u having an eigenvalue with positive real part) as a subset will, with some set of parameters (β_j) , (μ_{ij}) , become locally unstable. In particular, any shape of basic activity tree that includes instability example 1 (Figure 2) above (for $\rho < 1$) will be locally unstable for some set of parameters β_j, μ_{ij} .*

Proof. Let U be any system whose underload ($\rho < 1$) equilibrium is locally unstable, e.g. one of the examples given above, with the associated fixed set of parameters μ_{ij}, β_j and λ_i . Let S be a system including U as a subset, namely: the activity tree of S is a superset of that of U ; the μ_{ij} and β_j in U are preserved in S ; the μ_{ij} in S are fixed. Consider a sequence of systems S^ϵ in which $\beta_j = \epsilon \rightarrow 0$ for all j not in U . For each ϵ , take λ_i^ϵ so that all of the activities are indeed basic, and such that, as $\epsilon \rightarrow 0$, $\lambda_i^\epsilon \rightarrow \lambda_i$ for i in U , and $\lambda_i^\epsilon \rightarrow 0$ for i not in U . (See Remark 5.1.) Order the ψ_i so that the customer types i in U come first. Suppose there are I customer types in U and $I+k$ customer types in S . Let A_u^ϵ be the $(I+k) \times (I+k)$ matrix associated with S^ϵ , and let A_u be the $I \times I$ matrix associated with U considered as an isolated system. Then as $\epsilon \rightarrow 0$ the top left $I \times I$ entries of A_u^ϵ converge to A_u , while the bottom left $k \times I$ entries of A_u^ϵ converge to 0. (That is, the effect of U on the stability of the rest of the system vanishes – this is due to the fact that pool size parameters β_j in U remain constant, while $\beta_j \rightarrow 0$ in the rest of the system.) Consequently, each eigenvalue of

A_u is a limit of eigenvalues of A_u^ϵ . Since A_u had an eigenvalue with positive real part, for sufficiently small ϵ the matrix A_u^ϵ will have at least one eigenvalue with positive real part as well, so the system S^ϵ will be locally unstable. \square

6. DIFFUSION SCALED PROCESS IN AN UNDERLOADED SYSTEM. POSSIBLE EVANESCENCE OF INVARIANT DISTRIBUTIONS

Above we have shown that on a fluid scale, around the equilibrium point, the system converges to a subset of its possible states, on which it evolves according to a differential equation, possibly unstable. This strongly suggests that, when the differential equation is unstable, the stochastic system is in fact “never” close to equilibrium. Our goal in this section is to demonstrate that it is the case at least on the diffusion scale. More precisely, we consider the system in underload, $\rho < 1$, and look at diffusion-scaled stationary distributions (centered at the equilibrium point and scaled down by \sqrt{r}); we show that, when the associated fluid model is locally unstable, this sequence of stationary distributions is such that the measure of any compact set vanishes.

6.1. Transient behavior of diffusion scaled process. State space collapse. In this section we cite the diffusion limit result (for the process transient behavior) that we will need from [7]. Again, we consider a sequence of systems indexed by r , with the input rates being $\lambda_i^r = r\lambda_i$, server pool sizes being $\beta_j r$, and the service rates μ_{ij} unchanged with r . (Here we drop the $o(r)$ terms in $\lambda_i^r = r\lambda_i + o(r)$, because, when $\rho < 1$, considering these terms does not make sense.) The notation for the unscaled processes is the same as in the previous section; however, we are now interested in a different – diffusion – scaling. We define

$$(19) \quad \hat{\Psi}_{ij}^r(t) = \frac{\Psi_{ij}^r(t) - r\psi_{ij}^*}{\sqrt{r}}, \quad \hat{\Psi}_i^r(t) = \sum_j \hat{\Psi}_{ij}^r(t), \quad \hat{\Psi}_j^r(t) = \sum_i \hat{\Psi}_{ij}^r(t) = \frac{\Psi_j^r(t) - \rho r \beta_j}{\sqrt{r}}$$

We will denote by M' the linear mapping from $z = (z_{ij}, (ij) \in \mathcal{E}) \in \mathbb{R}^{I+J-1}$ to $y = (y_i) \in \mathbb{R}^I$, given by $\sum_j z_{ij} = y_i$. (So, $(\hat{\Psi}_i^r(t)) \equiv M'(\hat{\Psi}_{ij}^r(t))$.) There is the obvious relation between M' and the operator M defined by (8): $M'My = y$ for any $y \in \mathbb{R}^I$. Let us define $\mathcal{M} := \{My \mid y \in \mathbb{R}^I\}$, an I -dimensional linear subspace of \mathbb{R}^{I+J-1} ; equivalently, $\mathcal{M} = \{z \in \mathbb{R}^{I+J-1} \mid z = MM'z\}$.

Theorem 6.1 (Essentially a corollary of Theorem 3.1 and Theorem 4.4 in [7]). *Let $\rho < 1$. Assume that as $r \rightarrow \infty$, $\hat{\Psi}_\mathcal{E}^r(0) \rightarrow \hat{\Psi}_\mathcal{E}(0)$ where $\hat{\Psi}_\mathcal{E}(0)$ is deterministic and finite. (Consequently, $\hat{\Psi}_\mathcal{I}^r(0) \rightarrow \hat{\Psi}_\mathcal{I}(0) = M'\hat{\Psi}_\mathcal{E}(0)$.) Then,*

$$(20) \quad \hat{\Psi}_\mathcal{I}^r(\cdot) \implies \hat{\Psi}_\mathcal{I}(\cdot) \text{ in } D^I[0, \infty),$$

and for any fixed $\eta > 0$,

$$(21) \quad \hat{\Psi}_\mathcal{E}^r(\cdot) \implies M\hat{\Psi}_\mathcal{I}(\cdot) \text{ in } D^{I+J-1}[\eta, \infty),$$

where $\hat{\Psi}_\mathcal{I}(\cdot)$ is the unique solution of the SDE

$$(22) \quad \hat{\Psi}_i(t) = \hat{\Psi}_i(0) - \sum_{j \in \mathcal{S}(i)} \mu_{ij} \int_0^t (M\hat{\Psi}_\mathcal{I}(s))_{ij} ds + \sqrt{2\lambda_i} B_i(t), \quad i \in \mathcal{I},$$

and the processes $B_i(\cdot)$ are independent standard Brownian motions.

Recalling the definition of matrix A_u (see (10)), (22) can be written as

$$(23) \quad \hat{\Psi}_{\mathcal{I}}(t) = \hat{\Psi}_{\mathcal{I}}(0) + \int_0^t A_u \hat{\Psi}_{\mathcal{I}}(s) ds + (\sqrt{2\lambda_i} B_i(t)).$$

The meaning of Theorem 6.1 is simple: the diffusion limit of the process $\hat{\Psi}_{\mathcal{I}}^r(\cdot)$ is such that, at initial time 0, it “instantly jumps” to the state $MM'\hat{\Psi}_{\mathcal{E}}(0)$ on the manifold \mathcal{M} (where $MM'\hat{\Psi}_{\mathcal{E}}(0) = \hat{\Psi}_{\mathcal{E}}(0)$ only if $\hat{\Psi}_{\mathcal{E}}(0) \in \mathcal{M}$); after this initial jump, the process stays on \mathcal{M} and evolves according to SDE (23). Theorem 6.1 is “essentially a corollary” of results in [7], because the setting in [7] is such that $\rho = 1$, while we assumed $\rho < 1$. However, our Theorem 6.1 can be proved the same way, and in a sense is easier, because when $\rho < 1$, the queues vanish in the limit (which is why the queue length process is not even present in the statement of Theorem 6.1).

6.2. Evanescence of invariant measures. In this section we show that if the matrix A_u has eigenvalues with positive real part, the stationary distribution of the (diffusion scaled) process $\hat{\Psi}_{\mathcal{E}}^r(\cdot)$ escapes to infinity as $r \rightarrow \infty$. Namely, we prove the following

Theorem 6.2. *Suppose $\rho < 1$. Consider a sequence of systems as defined in Section 6.1, and denote by μ^r the stationary distribution of the process $\hat{\Psi}_{\mathcal{E}}^r(\cdot)$, a probability measure on \mathbb{R}^{I+J-1} . Let $b_K = \{|z| \leq K\} \subset \mathbb{R}^{I+J-1}$. Suppose the matrix A_u has eigenvalues with positive real parts and no pure imaginary eigenvalues.² Then for any K , $\mu^r(b_K) \rightarrow 0$ as $r \rightarrow \infty$.*

Before we proceed with the proof, let us introduce more notation and one auxiliary result. Let $\mathcal{C}_{\mathcal{I}}$ is the submanifold of convergence (stability) of ODE $(d/dt)y = A_u y$ on \mathbb{R}^I ; namely, $\mathcal{C}_{\mathcal{I}}$ is the (real) subspace of \mathbb{R}^I spanned by the Jordan basis vectors for matrix A_u corresponding to all eigenvalues with negative real parts. Given assumptions of the theorem on A_u , the solutions to $(d/dt)y = A_u y$ converge to 0 exponentially fast if $y(0) \in \mathcal{C}_{\mathcal{I}}$, and go to infinity exponentially fast if $y(0) \in \mathbb{R}^I \setminus \mathcal{C}_{\mathcal{I}}$. Let $\mathcal{C} = M\mathcal{C}_{\mathcal{I}}$ denote the corresponding submanifold of convergence (stability) of the linear ODE $(d/dt)z = (MA_u M')z$ on $z \in \mathcal{M}$. This ODE is just the M -image of ODE $(d/dt)y = A_u y$. Therefore, a solution $z(t)$ converges to 0 exponentially fast if $z(0) \in \mathcal{C}$, and goes to infinity exponentially fast if $z(0) \in \mathcal{M} \setminus \mathcal{C}$. Let us denote $b_K(\delta_1, \delta_2) := b_K \cap \{d(z, \mathcal{M}) \leq \delta_1, d(z, \mathcal{C}) \geq \delta_2\}$, where $d(\cdot, \cdot)$ is Euclidean distance.

Lemma 6.3. *Solutions to SDE (23) have the following properties.*

(i) *For any $T > 0$ and any $\Psi_{\mathcal{I}}(0)$,*

$$\mathbb{P}\{M\hat{\Psi}_{\mathcal{I}}(T) \in \mathcal{M} \setminus \mathcal{C}\} = 1;$$

(ii) *For any $K > 0$, $\delta_2 > 0$ and $\epsilon > 0$, there exist sufficiently large T_K and $K' > K$, such that, uniformly on $M\hat{\Psi}_{\mathcal{I}}(0) \in b_K(0, \delta_2)$,*

$$\mathbb{P}\{M\hat{\Psi}_{\mathcal{I}}(T_K) \in b_{K'} \setminus b_{2K}\} \geq 1 - \epsilon.$$

Proof. Statement (i) follows from the fact that, regardless of the (deterministic) initial state $\Psi_{\mathcal{I}}(0)$, the solution to SDE (23) is such that the distribution of $\Psi_{\mathcal{I}}(T)$ is Gaussian with non-singular covariance matrix. (See [8, Section 5.6]. In our case the matrix of diffusion

²The requirement of “no pure imaginary eigenvalues” is made for convenience of differentiating between strict convergence and strict divergence. It holds for generic values of β_j, μ_{ij} : that is, any set of values β_j, μ_{ij} has a small perturbation $\tilde{\beta}_j, \tilde{\mu}_{ij}$ with for which A_u has no pure imaginary eigenvalues.

coefficients is diagonal with entries $\sqrt{2\lambda_i}$.) Therefore, the probability that $\Psi_{\mathcal{I}}(T)$ is in a subspace of lower dimension is zero.

Statement (ii) follows from the fact (again, see [8, Section 5.6]) that the expectation $m(t) = \mathbb{E}\hat{\Psi}_{\mathcal{I}}(t)$ evolves according to ODE

$$\dot{m}(t) = A_u m(t).$$

Since $d(M\hat{\Psi}_{\mathcal{I}}(0), \mathcal{C}) \geq \delta_2$ (and thus $\hat{\Psi}_{\mathcal{I}}(0)$ is also separated by a positive distance from $\mathcal{C}_{\mathcal{I}}$), we have

$$|m(t)| \geq a_1 \exp(at)$$

for some fixed $a_1, a > 0$ and all large t . (Here a_1 depends on the minimum length of the projection of $\hat{\Psi}_{\mathcal{I}}(0)$ along $\mathcal{C}_{\mathcal{I}}$ onto the (real) span of the Jordan basis vectors of A_u corresponding to eigenvalues with positive real part, and a is the smallest positive real part of an eigenvalue of A_u .) It is easy to check that if the mean of a Gaussian distribution goes to infinity, then (regardless of how the covariance matrix changes) the measure of any bounded set goes to zero. On the other hand, both $m(t)$ and the covariance matrix remain bounded for all $t \in [0, T_K]$, with any T_K ; then, for any T_K , we can always choose K' large enough so that $\mathbb{P}\{M\hat{\Psi}_{\mathcal{I}}(T_K) \in b_{K'}\}$ is arbitrarily close to 1. \square

We are now in position to give a

Proof of Theorem 6.2. We will consider measures μ^r as measures on the one-point compactification $\overline{\mathbb{R}^n} = \mathbb{R}^n \cup \{*\}$ of the space \mathbb{R}^n , where $n = I + J - 1$. In this space, any subsequence of $\{\mu^r\}$ has a further subsequence, along which $\mu^r \xrightarrow{w} \mu$ for some probability measure μ on $\overline{\mathbb{R}^n}$. We will show that the entire measure μ is concentrated on the infinity point $*$, i.e. $\mu(\mathbb{R}^n) = 0$. Suppose not, i.e. $\mu(\mathbb{R}^n) > 0$. The proof proceeds in two steps.

Step 1. We prove that $\mu(\mathbb{R}^n) = \mu(\mathcal{M} \setminus \mathcal{C})$. Indeed, let us choose any $\epsilon > 0$, and K large enough so that $\mu(b_{K/2}) > (1 - \epsilon)\mu(\mathbb{R}^n)$. Then, for all large r , $\mu^r(b_K) > (1 - \epsilon)\mu(\mathbb{R}^n)$. Choose $\delta_1 > 0$ and $T > 0$ arbitrary. From the properties of the limiting diffusion process (Lemma 6.3), we see that we can choose a sufficiently small $\delta_2 > 0$ and sufficiently large K' such that, uniformly on the initial states $\hat{\Psi}_{\mathcal{E}}^r(0) \in b_K$,

$$\liminf_{r \rightarrow \infty} \mathbb{P}\{\hat{\Psi}_{\mathcal{E}}^r(T) \in b_{K'}(\delta_1, \delta_2)\} > 1 - \epsilon.$$

This implies that for all large r ,

$$\mu^r(b_{K'}(\delta_1, \delta_2)) > (1 - \epsilon)^2 \mu(\mathbb{R}^n),$$

and then $\mu(b_{K'}(\delta_1, \delta_2)) \geq (1 - \epsilon)^2 \mu(\mathbb{R}^n)$. Since ϵ and δ_1 were arbitrary, we conclude that $\mu(\mathbb{R}^n) \leq \mu(\mathcal{M} \setminus \mathcal{C})$, and then, obviously, the equality must hold.

Step 2. We show that, for any $K > 0$, $\mu(\mathbb{R}^n \setminus b_K) = \mu(\mathbb{R}^n)$. (This is, of course, impossible when $\mu(\mathbb{R}^n) > 0$, and thus we obtain a contradiction.) It suffices to show that for any $\epsilon > 0$, we can choose a sufficiently large K , such that $\mu(\mathbb{R}^n \setminus b_K) \geq (1 - \epsilon)^2 \mu(\mathbb{R}^n)$. Let us choose (using step 1) a large K and a small $\delta_2 > 0$, such that $\mu(b_{K/2}(\delta_1/2, 2\delta_2)) > (1 - \epsilon)\mu(\mathbb{R}^n)$ for any $\delta_1 > 0$. Then, for any fixed $\delta_1 > 0$, for all large r , $\mu^r(b_K(\delta_1, \delta_2)) > (1 - \epsilon)\mu(\mathbb{R}^n)$. Now, using Lemma 6.3(ii), we can choose K' and T_K sufficiently large, and then δ_1 sufficiently small, so that, uniformly on the initial states $\hat{\Psi}_{\mathcal{E}}^r(0) \in b_K(\delta_1, \delta_2)$,

$$\liminf_{r \rightarrow \infty} \mathbb{P}\{\hat{\Psi}_{\mathcal{E}}^r(T_K) \in b_{K'} \setminus b_{2K}\} \geq 1 - \epsilon.$$

Therefore,

$$\mu^r(b_{K'} \setminus b_{2K}) > (1 - \epsilon)^2 \mu(\mathbb{R}^n)$$

for all large r , and then for the limiting measure μ we must have $\mu(\mathbb{R}^n \setminus b_K) \geq (1 - \epsilon)^2 \mu(\mathbb{R}^n)$. \square

7. DIFFUSION SCALED PROCESS IN AN CRITICALLY LOADED SYSTEM, IN HALFIN-WHITT ASYMPTOTIC REGIME.

In this section we consider the following asymptotic regime. The system is critically loaded, i.e. the optimal solution to SPP (1) is such that $\rho = 1$. As scaling parameter $r \rightarrow \infty$, assume that the server pool sizes are $r\beta_j$ (same as throughout the paper), and the input rates are $\lambda_i^r = r\lambda_i + \sqrt{r}l_i$, where the parameters (finite real numbers) $\{l_i\}$ are such that $\sum l_i\nu_i = -C < 0$. Denote by $\rho^r, \{\lambda_{ij}^r\}$ the optimal solution of SPP (1), with β_j 's and λ_i 's replaced by $r\beta_j$ and λ_i^r , respectively. (This solution is unique, as can be easily seen from the CRP condition.) Then, it is easy to check that $\rho^r = 1 + (\sum l_i\nu_i)/\sqrt{r} = 1 - C/\sqrt{r}$, which in turn easily implies that, for any r , the system process is stable with the unique stationary distribution.

We use the definitions of (19) for the diffusion scaled variables, and add to them the following ones: $\hat{X}_i^r(t) = (X_i^r(t) - \psi_i^*r)/\sqrt{r}$ for the (diffusion-scaled) number of type i customers; $\hat{Q}_i^r(t) = Q_i^r(t)/\sqrt{r}$ for the type i queue length; $\hat{Z}_j^r(t) = Z_j^r(t)/\sqrt{r}$, where $Z_j^r(t) = \Psi_j^r(t) - r\beta_j \leq 0$ is the number of idle servers of type j (with the minus sign). Note that, although the optimal average occupancy of pool j is at $\rho^r r\beta_j$, the quantity $\hat{Z}_j^r(t)$ measures the deviation from full occupancy $r\beta_j$. Our choice of signs is such that $\hat{Q}_i^r \geq 0$ while $\hat{Z}_j^r \leq 0$. We will use the vector notations, such as $\hat{X}_{\mathcal{I}}^r(t)$, as usual.

Two main results of this section are as follows: (a) it is possible for the invariant distributions to escape to infinity under certain system parameters and (b) in the special case when service rate depends on the server type only, the invariant distributions are tight.

7.1. Example of evanescence of invariant measures. Recall that π denotes the (matrix of) orthogonal projection on the subspace $L = \{y \in \mathbb{R}^I \mid \sum_i y_i = 0\}$ in \mathbb{R}^I ; this is the projection ‘‘along’’ the direction of vector $(1, \dots, 1)^\dagger$. Also recall the relation between matrices A_u and A_c :

$$A_c = \pi A_u.$$

One more notation: for $y \in \mathbb{R}^I$,

$$F[y] = \begin{cases} \pi y, & \text{if } \sum_i y_i > 0 \\ y, & \text{if } \sum_i y_i \leq 0. \end{cases}$$

Analogously to Theorem 6.1, the following fact is a corollary (this time – direct) of Theorem 3.1 and Theorem 4.4 in [7].

Theorem 7.1. *Assume that as $r \rightarrow \infty$, $\hat{X}_{\mathcal{I}}^r(0) \rightarrow \hat{X}_{\mathcal{I}}(0)$ and $\hat{\Psi}_{\mathcal{E}}^r(0) \rightarrow \hat{\Psi}_{\mathcal{E}}(0)$, where $\hat{X}_{\mathcal{I}}(0)$ and $\hat{\Psi}_{\mathcal{E}}(0)$ are deterministic and finite. Then,*

$$(24) \quad \hat{X}_{\mathcal{I}}^r(\cdot) \implies \hat{X}_{\mathcal{I}}(\cdot) \text{ in } D^I[0, \infty),$$

and for any fixed $\eta > 0$,

$$(25) \quad \hat{\Psi}_{\mathcal{E}}^r(\cdot) \implies MF[\hat{X}_{\mathcal{I}}(\cdot)] \text{ in } D^{I+J-1}[\eta, \infty),$$

where $\hat{X}_{\mathcal{I}}(\cdot)$ is the unique solution of the SDE

$$(26) \quad \hat{X}_{\mathcal{I}}(t) = \hat{X}_{\mathcal{I}}(0) + \int_0^t A_u F[\hat{X}_{\mathcal{I}}(s)] ds + (\sqrt{2\lambda_i} B_i(t)),$$

and the processes $B_i(\cdot)$ are independent standard Brownian motions.

Next we establish the following fact.

Lemma 7.2. *There exists a system and a parameter setting such that the following holds.*

(i) *Matrix A_c is unstable.*

(ii) *Matrix A_u has $(1, \dots, 1)^\dagger$ as a right eigenvector, with real non-zero eigenvalue c :*

$$(27) \quad A_u(1, \dots, 1)^\dagger = c(1, \dots, 1)^\dagger.$$

Proof. Let us start with the system in the local instability example 2 (see Figure 4) for the critical load. We will modify it as follows. We will change μ_{D3} from 100 to $100 - \epsilon$ with sufficiently small positive ϵ , so that A_c remains unstable. (The reason for this change will be explained shortly.) We will add two new server pools, 0 and 5, on the left and on the right, respectively, and set $\mu_{A0} = 100$, $\mu_{E5} = 1$; such addition of server-leaves does not change the instability of A_c . So, (i) holds.

Now, suppose all λ_i are equal, say $\lambda_i = 1$. We can choose ψ_{ij}^* such that all $\psi_i^* = \sum_j \psi_{ij}^*$ are equal, and $\sum_j \mu_{ij} \psi_{ij}^* = \lambda_i = 1$ for all i . Namely, we do the following. The reason for changing μ_{D3} from 100 to $100 - \epsilon$ is to make it possible to choose $\psi_{D3}^* > 0$ and $\psi_{D4}^* > 0$, such that $\sum_j \mu_{Dj} \psi_{Dj}^* = 1$ and $\psi_D^* = \psi_{D3}^* + \psi_{D4}^* > 1/100$. We choose $\psi_{A0}^* = 1/100 - \delta$, $\psi_{A1}^* = 100\delta$ (which guarantees $\sum_j \mu_{Aj} \psi_{Aj}^* = 1$) with $\delta > 0$ small enough so that $\psi_A^* = 1/100 + 99\delta < 1/(100 - \epsilon)$. The values of pairs $(\psi_{B1}^*, \psi_{B2}^*)$, $(\psi_{C2}^*, \psi_{C3}^*)$, $(\psi_{E4}^*, \psi_{E5}^*)$, are chosen to be equal to $(\psi_{A0}^*, \psi_{A1}^*)$. Finally, we choose $\psi_{D3}^* = (1 - \delta_1)/(100 - \epsilon)$ and $\psi_{D4}^* = \delta_1/10^4$ (which ensures $\sum_j \mu_{Dj} \psi_{Dj}^* = 1$) with $\delta_1 > 0$ satisfying

$$\psi_D^* = (1 - \delta_1)/(100 - \epsilon) + \delta_1/10^4 = 1/100 + 99\delta = \psi_A^*.$$

This completes the choice of ψ_{ij}^* .

We set $\beta_j = \sum_i \psi_{ij}^*$. We see that (ψ_{ij}^*) is the equilibrium point. It follows from the construction that (27) will hold for A_u . Indeed, if $\psi_{\mathcal{I}} - \psi_{\mathcal{I}}^* = c_1(1, \dots, 1)^\dagger$, then $\psi_{\mathcal{I}} = c_2 \psi_{\mathcal{I}}^*$, which in turn means that $\psi_{\mathcal{E}} = c_2 \psi_{\mathcal{E}}^*$; therefore, the corresponding service rates are $\sum_j \mu_{ij} \psi_{ij} = c_2 \sum_j \mu_{ij} \psi_{ij}^* = c_2 \lambda_i = c_2$ for all i ; therefore, $\psi_{\mathcal{I}} = (1 - c_2)(1, \dots, 1)^\dagger$. \square

Theorem 7.3. *Suppose we have a system with parameters satisfying Lemma 7.2, in the Halfin-Whitt regime, described in this section. Then, the sequence of stationary distributions of $\hat{X}_{\mathcal{I}}^r$ (and of $\hat{\Psi}_{\mathcal{I}}^r$) escapes to infinity: the measure of any compact set vanishes.*

Proof. Since $(1, \dots, 1)^\dagger$ is an eigenvector of A_u , for any $y \in \mathbb{R}^I$ we have:

$$\pi A_u F[y] = \pi A_u \pi y = A_c \pi y.$$

Then, taking the π -projection of equation (26), we see that $\pi \hat{X}_{\mathcal{I}}$ satisfies the following linear SDE

$$(28) \quad \pi \hat{X}_{\mathcal{I}}(t) = \pi \hat{X}_{\mathcal{I}}(0) + \int_0^t A_c \pi \hat{X}_{\mathcal{I}}(s) ds + \pi(\sqrt{2\lambda_i} B_i(t)).$$

Given instability of linear equation (28), we can repeat the argument of Section 6.2 to show that the sequence of projections of the stationary distributions of $\hat{X}_{\mathcal{I}}^r$ on L escapes to infinity. \square

7.2. Tightness of stationary distributions in the case when service rate depends on the server type only. In this section we consider a special case when there exists a set of positive rates $\{\mu_j\}$, such that $\mu_{ij} = \mu_j$ as long as $(ij) \in \mathcal{E}$. We demonstrate tightness of invariant distributions. (An analogous result holds for the underload system, $\rho < 1$, as sketched out at the end of this section.) This, in combination with the transient diffusion limit results, allows us to claim that the limit of invariant distributions is the invariant distribution of the limiting diffusion process.

Theorem 7.4. *Suppose $\mu_{ij} = \mu_j$, $(ij) \in \mathcal{E}$ and $\rho = 1$. Consider a system under the LQFS-LB rule in the asymptotic regime defined above in this section. Then, for any real*

$$\theta < \theta_0 := \frac{2 \min_i \lambda_i}{\sum_i \lambda_i + (\max_j \mu_j) \sum_j \beta_j},$$

the stationary distributions are such that

$$\limsup_r \mathbb{E} \left[\sum_i \exp(\theta \hat{Q}_i^r) + \sum_j \beta_j \exp(\theta \hat{Z}_j^r / \beta_j) \right] < \infty.$$

Proof. Note that the statement is trivial for $\theta = 0$. Also, for $\theta > 0$ each term $\exp(\theta \hat{Z}_j^r / \beta_j)$ is bounded so has finite expectation, while for $\theta < 0$ each term $\exp(\theta \hat{Q}_i^r)$ is bounded so has finite expectation.

Our method is related to that in [4]. (The exposition below is self-contained.)

Step 1: preliminary bounds. Consider the embedded Markov chain taken at the instants of (say, right after) the transitions. We will use uniformization, that is, we keep the total rate of all transitions from any state constant at $\alpha^r r = \sum_i \lambda_i^r + \sum_j r \beta_j \mu^*$, where $\mu^* = \max_j \mu_j$; note that, as $r \rightarrow \infty$, $\alpha^r \rightarrow \alpha^* = \sum_i \lambda_i + \sum_j \beta_j \mu^*$. The transitions are of three types: arrivals, departures, and virtual transitions, which do not change the state of the system. The rate of a transition due to a type i arrival is λ_i^r ; for the service completion at pool j the rate is $\mu_j(r\beta_j + Z_j^r)$ (recall $Z_j^r \leq 0$); and a virtual transition occurs at the complementary rate $\alpha^r r - \sum_i \lambda_i^r - \sum_j \mu_j(r\beta_j + Z_j^r)$. (Obviously, the probability that a transition occurring at a transition instant has a given type is the ratio of the corresponding rate and $\alpha^r r$.) The stationary distribution of the embedded Markov chain is the same as that of the original, continuous-time chain.

In the rest of the proof, $\tau \in \{0, 1, 2, \dots\}$ refers to the discrete time of the embedded Markov chain.

We will work with the following Lyapunov function

$$(29) \quad \mathcal{L}(\tau) := \sum_i \exp(\theta \hat{Q}_i^r(\tau)) + \sum_j \beta_j \exp(\theta \hat{Z}_j^r(\tau) / \beta_j).$$

Throughout, we use the bound

$$(30) \quad \exp(\theta y) \leq \exp(\theta x) \left(1 + \theta(y - x) + \frac{1}{2} \theta^2 (y - x)^2 \exp(\theta |y - x|) \right)$$

which arises from the second-order Taylor expansion of $\exp(\theta y)$.

A priori we do not know that $\mathbb{E}[\mathcal{L}(\tau)]$ exists for $\theta > 0$. Indeed, while $\hat{Z}_j^r(t)$ is bounded for any r (above by 0 and below by $-\beta_j \sqrt{r}$), the scaled queue size $\hat{Q}_i^r(t)$ is unbounded. To deal with this, we also consider the truncated Lyapunov function $\mathcal{L}^K = \min\{\mathcal{L}, K\}$.

In the equation below, let x denote the variable of interest (either \hat{Q}_i^r or \hat{Z}_j^r/β_j), and let $S(\tau)$ denote the state of the embedded Markov chain at time τ . From (30) we obtain

$$\begin{aligned} \mathbb{E}[\exp(\theta x(\tau + 1)) - \exp(\theta x(\tau)) | S(\tau)] &\leq \\ &\exp(\theta x(\tau)) (\theta \mathbb{E}[x(\tau + 1) - x(\tau) | S(\tau)] + \\ &\quad \frac{1}{2} \theta^2 \mathbb{E}[(x(\tau + 1) - x(\tau))^2 \exp(\theta |x(\tau + 1) - x(\tau)|) | S(\tau)]). \end{aligned}$$

Since for both \hat{Z}_j^r and \hat{Q}_i^r the change in a single transition is bounded by $1/\sqrt{r}$, we conclude:

$$(31) \quad \mathbb{E}[\exp(\theta \hat{Q}_i^r(\tau + 1)) - \exp(\theta \hat{Q}_i^r(\tau)) | S(\tau)] \leq \\ \exp(\theta \hat{Q}_i^r(\tau)) (\theta \mathbb{E}[\hat{Q}_i^r(\tau + 1) - \hat{Q}_i^r(\tau) | S(\tau)] + \left(\frac{1}{2} \theta^2 \exp(\theta/\sqrt{r}) \right) \frac{1}{r}),$$

$$(32) \quad \mathbb{E}[\beta_j \exp(\theta \hat{Z}_j^r(\tau + 1)/\beta_j) - \beta_j \exp(\theta \hat{Z}_j^r(\tau)/\beta_j) | S(\tau)] \leq \\ \exp(\theta \hat{Z}_j^r(\tau)/\beta_j) (\theta \mathbb{E}[\hat{Z}_j^r(\tau + 1) - \hat{Z}_j^r(\tau) | S(\tau)] + \left(\frac{1}{\beta_j} \frac{1}{2} \theta^2 \exp(\theta/\sqrt{r}) \right) \frac{1}{r}).$$

Clearly, as long as values of θ are bounded, for any fixed $C_2 > 1$ and all sufficiently (depending on C_2) large r , the second summands in (31) and (32) are upper bounded by $C_2 \frac{1}{2} \theta^2 \frac{1}{r}$ and $\frac{1}{\beta_*} C_2 \frac{1}{2} \theta^2 \frac{1}{r}$, respectively, where $\beta_* = \min_j \beta_j$. Note that the second bound is independent of j .

Next, we will obtain an upper bound on the drift

$$\mathbb{E}[\mathcal{L}(\tau + 1) - \mathcal{L}(\tau) | S(\tau)].$$

To do that, we introduce an artificial scheduling/routing rule, which acts only within one time step, and is such that the increment $\mathcal{L}(\tau + 1) - \mathcal{L}(\tau)$ under this rule is “almost” a (pathwise, w.p.1) upper bound on this increment under the actual – LQFS-LB – rule. (It is important to keep in mind that the artificial rule is *not* a rule that is applied continuously. It is limited to one time step, and its sole purpose is to derive a pathwise upper bound on the increment $\mathcal{L}(\tau + 1) - \mathcal{L}(\tau)$ within one time step.)

Step 2: Artificial scheduling/routing rule. We will use the following notation: $\mathcal{I}_+ = \mathcal{I}_+(\tau) := \{i : \hat{Q}_i^r(\tau) > 0\}$, $\mathcal{I}_0 = \mathcal{I}_0(\tau) := \{i : \hat{Q}_i^r(\tau) = 0\}$, $\mathcal{J}_- = \mathcal{J}_-(\tau) := \{j : \hat{Z}_j^r(\tau) < 0\}$, $\mathcal{J}_0 = \mathcal{J}_0(\tau) := \{j : \hat{Z}_j^r(\tau) = 0\}$.

Scheduling: Departures from servers $j \in \mathcal{J}_-$ are processed normally, i.e. reduce the corresponding $Z_j^r(\tau)$ by 1. Whenever there is a departure from a server pool $j \in \mathcal{J}_0$, the server takes up a customer of type i with probability $\lambda_{ij}^r / \sum_i \lambda_{ij}^r$, keeping $Z_j^r(\tau + 1) = 0$ and reducing $Q_i^r(\tau + 1) = Q_i^r(\tau) - 1$. However, if it happens that the chosen i is such that $Q_i^r(\tau) = 0$, i.e. $i \in \mathcal{I}_0$, then we keep $Q_i^r(\tau + 1) = Q_i^r(\tau) = 0$ and instead allow $Z_j^r(\tau + 1) = -1$.

Routing: Arrivals to customer types $i \in \mathcal{I}_+$ are processed normally, i.e. increase the corresponding $Q_i^r(\tau)$ by 1. Whenever there is an arrival to a customer type $i \in \mathcal{I}_0$, it is routed to server pool j with probability $\lambda_{ij}^r / \lambda_i^r$, keeping $Q_i^r(\tau + 1) = Q_i^r(\tau) = 0$ and increasing $Z_j^r(\tau + 1) = Z_j^r(\tau) + 1$. However, if it happens that the chosen j is such that $Z_j^r(\tau) = 0$, i.e. $j \in \mathcal{J}_0$, then we keep $Z_j^r(\tau + 1) = Z_j^r(\tau) = 0$ and instead allow $Q_i^r(\tau + 1) = 1$.

Step 3: One time-step drift under the artificial rule. For $i \in \mathcal{I}_+$,

$$\mathbb{E}[\hat{Q}_i^r(\tau+1) - \hat{Q}_i^r(\tau)|S(\tau)] = \frac{1}{\alpha^r r} \frac{1}{\sqrt{r}} \left(\lambda_i^r - \sum_j (\mu_j^r \beta_j) \frac{\lambda_{ij}^r}{\sum_k \lambda_{kj}^r} \right),$$

or, recalling that

$$(33) \quad \sum_k \lambda_{kj}^r = \mu_j \beta_j r \rho^r = \mu_j \beta_j r (1 - C/\sqrt{r}),$$

we obtain

$$(34) \quad \mathbb{E}[\hat{Q}_i^r(\tau+1) - \hat{Q}_i^r(\tau)|S(\tau)] = -\frac{C\lambda_i}{\alpha^*} \frac{1+o(1)}{r}, \quad i \in \mathcal{I}_+,$$

where $o(1)$ is a fixed function, vanishing as $r \rightarrow \infty$.

If $\hat{Q}_i^r(\tau) = 0$ (i.e. $i \in \mathcal{I}_0$), and a new type i arrival is routed to pool j with $\hat{Z}_j^r(\tau) < 0$ (i.e. $j \in \mathcal{J}_-$), then of course \hat{Q}_i^r stays at 0 and $\hat{Q}_i^r(\tau+1) - \hat{Q}_i^r(\tau) = 0$. However, if a new type i arrival has to be routed to $j \in \mathcal{J}_0$, then (by the definition of artificial rule) $\hat{Q}_i^r(\tau+1) - \hat{Q}_i^r(\tau) = \hat{Q}_i^r(\tau+1) = 1/\sqrt{r}$. Thus, we can write:

$$(35) \quad \mathbb{E}[\hat{Q}_i^r(\tau+1) - \hat{Q}_i^r(\tau)|S(\tau)] = \sum_{j \in \mathcal{J}_0} \frac{\lambda_{ij}^r}{\alpha^r r} \frac{1}{\sqrt{r}}, \quad i \in \mathcal{I}_0.$$

Note that the RHS of (35) is of order $1/\sqrt{r}$, not $1/r$. However, we will see shortly that order $1/\sqrt{r}$ terms in $\mathbb{E}[\mathcal{L}(\tau+1) - \mathcal{L}(\tau)|S(\tau)]$ cancel out, and this expected drift is in fact of order $1/r$.

The treatment of the drift of \hat{Z}_j^r is similar (and again makes use of (33)). We obtain:

$$(36) \quad \mathbb{E}[\hat{Z}_j^r(\tau+1) - \hat{Z}_j^r(\tau)|S(\tau)] = -\frac{1}{\alpha^r} \mu_j (\hat{Z}_j^r(\tau) + \beta_j C) \frac{1}{r}, \quad j \in \mathcal{J}_-,$$

$$(37) \quad \mathbb{E}[\hat{Z}_j^r(\tau+1) - \hat{Z}_j^r(\tau)|S(\tau)] = -\frac{1}{\sqrt{r}} \sum_{i \in \mathcal{I}_0} \frac{r \mu_j \beta_j}{\alpha^r r} \frac{\lambda_{ij}^r}{\sum_k \lambda_{kj}^r} = -\frac{1}{1 - C/\sqrt{r}} \sum_{i \in \mathcal{I}_0} \frac{\lambda_{ij}^r}{\alpha^r r} \frac{1}{\sqrt{r}}, \quad j \in \mathcal{J}_0.$$

We can rewrite (37) as

$$(38) \quad \mathbb{E}[\hat{Z}_j^r(\tau+1) - \hat{Z}_j^r(\tau)|S(\tau)] = -\sum_{i \in \mathcal{I}_0} \frac{\lambda_{ij}^r}{\alpha^r r} \frac{1}{\sqrt{r}} - \frac{C \sum_{i \in \mathcal{I}_0} \lambda_{ij}^r}{\alpha^*} \frac{1+o(1)}{r}, \quad j \in \mathcal{J}_0,$$

where $o(1)$ is a fixed function, vanishing as $r \rightarrow \infty$.

Note that if $\mathcal{L}(\tau) \geq K$ then $\mathcal{L}^K(\tau+1) - \mathcal{L}^K(\tau) \leq 0$, and if $\mathcal{L}(\tau) < K$ then $\mathcal{L}^K(\tau+1) - \mathcal{L}^K(\tau) \leq \mathcal{L}(\tau+1) - \mathcal{L}(\tau)$. Putting together this observation and equations (31) – (32), (34) – (38), we obtain

$$\begin{aligned}
(39a) \quad & \mathbb{E}[\mathcal{L}^K(\tau + 1) - \mathcal{L}^K(\tau) | S(\tau)] \leq \\
(39b) \quad & \mathbf{1}_{\{\mathcal{L}(\tau) \leq K\}} \left(\sum_{i \in \mathcal{I}_+} \exp(\theta \hat{Q}_i^r(\tau)) \theta \left[-\frac{C\lambda_i(1+o(1))}{\alpha^*} \right] \frac{1}{r} \right. \\
(39c) \quad & \left. + \sum_{i \in \mathcal{I}_0, j \in \mathcal{J}_0} \theta \lambda_{ij}^r \frac{1}{\alpha^r r} \frac{1}{\sqrt{r}} \right. \\
(39d) \quad & \left. + \sum_{j \in \mathcal{J}_-} \exp(\theta \hat{Z}_j^r(\tau) / \beta_j) \theta \left[-\frac{\mu_j}{\alpha^r} \right] \left[\hat{Z}_j^r(\tau) + \beta_j C \right] \frac{1}{r} \right. \\
(39e) \quad & \left. + \sum_{j \in \mathcal{J}_0, i \in \mathcal{I}_0} \theta \left[-\lambda_{ij}^r \frac{1}{\alpha^r r} \frac{1}{\sqrt{r}} - \frac{C\lambda_i(1+o(1))}{\alpha^*} \frac{1}{r} \right] \right. \\
(39f) \quad & \left. + \sum_{i \in \mathcal{I}} \exp(\theta \hat{Q}_i^r(\tau)) \left(\frac{C_2}{2} \theta^2 \right) \frac{1}{r} \right. \\
(39g) \quad & \left. + \sum_{j \in \mathcal{J}} \frac{1}{\beta_*} \exp(\theta \hat{Z}_j^r(\tau) / \beta_j) \left(\frac{C_2}{2} \theta^2 \right) \frac{1}{r} \right),
\end{aligned}$$

Note that the $O(1/\sqrt{r})$ terms in (39c) and (39e) cancel each other as promised, so there are no $O(1/\sqrt{r})$ terms in the final bound.

Step 4: One time-step drift under the LQFS-LB rule. We now explain in what sense the increment $\mathcal{L}(\tau + 1) - \mathcal{L}(\tau)$ under the artificial rule is “almost” an upper bound on this increment under LQFS-LB. To illustrate the idea, suppose first that all β_j are equal. Then, it is easy to observe that for any fixed $S(\tau)$, the increment $\mathcal{L}(\tau + 1) - \mathcal{L}(\tau)$ under the artificial rule is (with probability 1) an upper bound of this increment under LQFS-LB. Indeed, suppose first that a transition of the Markov chain is associated with a service completion in server pool j with $\hat{Z}_j^r = 0$. (If $\hat{Z}_j^r < 0$, there is no difference in what the two rules do.) The only case of interest is when the LQFS-LB “takes” a new customer for service from queue i with $\hat{Q}_i^r > 0$, while the artificial rule tries to take a customer from a different queue i' . Then $\hat{Q}_i^r \geq \hat{Q}_{i'}^r$ must hold, with $\hat{Q}_i^r > \hat{Q}_{i'}^r$ being the non-trivial case. If $\hat{Q}_{i'}^r > 0$, then the LQFS-LB will decrease the larger queue, and so the increment $\mathcal{L}(\tau + 1) - \mathcal{L}(\tau)$ under the LQFS-LB is smaller (which is true for both positive and negative θ). If $\hat{Q}_{i'}^r = 0$, then the LQFS-LB will still decrease queue $\hat{Q}_{i'}^r$, while the artificial rule will instead decrease \hat{Z}_j^r ; using convexity of $e^{\theta x}$, we verify that, again, the increment $\mathcal{L}(\tau + 1) - \mathcal{L}(\tau)$ under the LQFS-LB is smaller (for both positive and negative θ). If transition of the Markov chain is associated with a new customer arrival, we use analogous argument to show that, again, the increment $\mathcal{L}(\tau + 1) - \mathcal{L}(\tau)$ under the LQFS-LB cannot be greater than that under the artificial rule. We conclude that when all β_j are equal, the key estimate (39) of the expected drift holds, in exactly same form, for LQFS-LB rule as well.

Now consider the case of general β_j . In the event of a service completion (and then possibly taking a customer for service from one of the non-zero queues), the increment $\mathcal{L}(\tau + 1) - \mathcal{L}(\tau)$ under LQFS-LB is still no greater than under the artificial rule. (Verified similarly to the case of all β_j being equal.) The only situation when LQFS-LB can possibly cause a greater increment than the artificial rule is as follows. There is an arrival of a type i customer, which the artificial rule routes to pool j with $\hat{Z}_j^r < 0$, but the LQFS-LB will instead route it

to pool k such that $\hat{Z}_j^r/\beta_j \geq \hat{Z}_k^r/\beta_k$. Given convexity of function $e^{\theta x}$, the “worst case”, i.e. the largest increment of $\mathcal{L}(\tau + 1) - \mathcal{L}(\tau)$, occurs when \hat{Z}_k^r is such that the equality holds, $\hat{Z}_j^r/\beta_j = \hat{Z}_k^r/\beta_k$. (If $\theta > 0$ the positive increment gets larger, if we were to increase \hat{Z}_k^r ; if $\theta < 0$ the negative increment gets smaller in absolute value, if we were to increase \hat{Z}_k^r . Note also that here we allow \hat{Z}_k^r , determined by the equality, to be such that $Z_k^r = \hat{Z}_k^r \sqrt{r}$ is possibly non-integer, because we only use this value of \hat{Z}_k^r to estimate the increment of a function.) Thus, as we replace the artificial rule by LQFS-LB, in the “worst case”, the increment

$$\beta_j \exp(\theta[\hat{Z}_j^r(\tau) + r^{-1/2}]/\beta_j) - \beta_j \exp(\theta\hat{Z}_j^r(\tau)/\beta_j)$$

may need to be replaced by

$$\beta_k \exp(\theta[\hat{Z}_k^r(\tau) + r^{-1/2}]/\beta_k) - \beta_k \exp(\theta\hat{Z}_k^r(\tau)/\beta_k),$$

with $\hat{Z}_k^r(\tau)$ satisfying $\hat{Z}_j^r(\tau)/\beta_j = \hat{Z}_k^r(\tau)/\beta_k$. In this case we obtain

$$(40) \quad \beta_k \exp(\theta\hat{Z}_k^r(\tau + 1)/\beta_k) - \beta_k \exp(\theta\hat{Z}_k^r(\tau)/\beta_k) \leq \\ \exp(\theta\hat{Z}_k^r(\tau)/\beta_k) (\theta r^{-1/2} + \left(\frac{1}{\beta_k} \frac{1}{2} \theta^2 \exp(\theta/\sqrt{r}) \right) \frac{1}{r}) \leq \\ \exp(\theta\hat{Z}_j^r(\tau)/\beta_j) (\theta r^{-1/2} + \left(\frac{1}{\beta_*} \frac{1}{2} \theta^2 \exp(\theta/\sqrt{r}) \right) \frac{1}{r}),$$

This means that, *under LQFS-LB rule, the estimate (39) still holds.*

Step 5: Exponential moments estimates. Next, note that for each fixed $K > 0$ and each fixed parameter r , the values of $\exp(\theta\hat{Q}_i^r(\tau))$ are uniformly bounded over all states $S(\tau)$ satisfying condition $\mathcal{L}(\tau) \leq K$; the values of $\exp(\theta\hat{Z}_j^r(\tau)/\beta_j)$ are “automatically” uniformly bounded (for a fixed r). We take the expected values of both parts of (39) with respect to the invariant distribution. The expectation of the LHS is of course 0, and so we get rid of the factor $1/r$ from the RHS expectation. The resulting estimates we will write separately for the cases $\theta > 0$ and $\theta < 0$ (with the case $\theta = 0$ being trivial).

Case $\theta > 0$. For a fixed $\theta > 0$, the expected value of the sum of all terms not containing $\exp(\theta\hat{Q}_i^r(\tau))$ is bounded (uniformly in r). Indeed, this follows from the facts that $\hat{Z}_j^r(\tau) \leq 0$ and $0 \leq -\theta\hat{Z}_j^r(\tau) \exp(\theta\hat{Z}_j^r(\tau)/\beta_j) \leq \beta_j/e$ (because $0 \geq xe^x \geq -\frac{1}{e}$ for $x \leq 0$). Then, we obtain:

$$(41) \quad \mathbb{E} \left[\mathbf{1}_{\{\mathcal{L}(\tau) \leq K\}} \sum_{i \in \mathcal{I}_+} \exp(\theta\hat{Q}_i^r(\tau)) \left(\frac{C\lambda_i(1+o(1))}{\alpha^*} \theta - \left(\frac{C_2}{2} \theta^2 \right) \right) \right] \leq C_1$$

for some constant $C_1 = C_1(\theta) > 0$, uniformly on all sufficiently large r . Now let us fix a sufficiently small positive θ , so that all coefficients of $\exp(\theta\hat{Q}_i^r(\tau))$ are at least some $\epsilon > 0$ (for all large r). Recalling that $C_2 > 1$ can be arbitrarily close to 1, it suffices that $\theta < \theta_0 = 2(\min_i \lambda_i)/\alpha^*$. Then,

$$\mathbb{E} \left[\mathbf{1}_{\{\mathcal{L}(\tau) \leq K\}} \sum_{i \in \mathcal{I}_+} \exp(\theta\hat{Q}_i^r(\tau)) \right] \leq C_1/\epsilon,$$

from where, letting $K \rightarrow \infty$, by monotone convergence, we obtain

$$(42) \quad \mathbb{E} \left[\sum_{i \in \mathcal{I}_+} \exp(\theta \hat{Q}_i^r(\tau)) \right] \leq C_1/\epsilon < \infty,$$

uniformly on all large r .

Case $\theta < 0$. Fix arbitrary $\theta < 0$. In this case, the expected value of the sum of all terms not containing $\exp(\theta \hat{Z}_j^r(\tau))$, is bounded (uniformly on r). We can write:

$$(43) \quad \mathbb{E} \left[\mathbf{1}_{\{\mathcal{L}(\tau) \leq K\}} \sum_{j \in \mathcal{J}_-} \exp(\theta \hat{Z}_j^r(\tau)/\beta_j) \left(\theta \left[\frac{\mu_j}{\alpha^r} \right] [\hat{Z}_j^r(\tau) + \beta_j C] - \left(\frac{1}{\beta_*} \frac{C_2}{2} \theta^2 \right) \right) \right] \leq C'_1,$$

for some constant $C'_1 = C'_1(\theta) > 0$, uniformly on all sufficiently large r . Let us choose sufficiently large $K_1 > 0$, such that the condition $\hat{Z}_j^r(\tau) \leq -K_1$ implies that

$$\left(\theta \left[\frac{\mu_j}{\alpha^r} \right] [\hat{Z}_j^r(\tau) + \beta_j C] - \left(\frac{1}{\beta_*} \frac{C_2}{2} \theta^2 \right) \right) \geq \epsilon,$$

for some $\epsilon > 0$ (and all large r). Then, from (43),

$$\mathbb{E} \left[\mathbf{1}_{\{\mathcal{L}(\tau) \leq K\}} \sum_{j \in \mathcal{J}_-} \mathbf{1}_{\{\hat{Z}_j^r(\tau) \leq -K_1\}} \exp(\theta \hat{Z}_j^r(\tau)/\beta_j) \right] \leq C'_1/\epsilon,$$

from where, letting $K \rightarrow \infty$, by monotone convergence, we obtain

$$\mathbb{E} \left[\sum_{j \in \mathcal{J}_-} \mathbf{1}_{\{\hat{Z}_j^r(\tau) \leq -K_1\}} \exp(\theta \hat{Z}_j^r(\tau)/\beta_j) \right] \leq C'_1/\epsilon < \infty,$$

uniformly on all large r , which implies the required result. \square

Corollary 7.5. *The sequence of stationary distributions of the processes $((\hat{Q}_i^r(\cdot)), (\hat{Z}_j^r(\cdot)))$ has a weak limit, which is the unique stationary distribution of the limiting process $((\hat{Q}_i(\cdot)), (\hat{Z}_j(\cdot)))$, described as follows:*

$$\hat{Q}_i(t) = \max\{\hat{Y}(t)/I, 0\}, \quad \forall i, \quad \hat{Z}_j(t) = \min\left\{\frac{\beta_j}{\sum_k \beta_k} \hat{Y}(t), 0\right\}, \quad \forall j,$$

where $\hat{Y}(\cdot)$ is a one-dimensional diffusion process with constant variance parameter $2 \sum_i \lambda_i$ and piece-wise linear drift, equal at point x to

$$-\left[\sum_j \mu_j \right] [C + \min\{x, 0\}].$$

The invariant distribution density is, then, a continuous function, which is a ‘‘concatenation’’ at point 0 of exponential (for $x \geq 0$) and Gaussian (for $x \leq 0$) distribution densities.

Proof. Theorem 7.4 of course implies tightness of stationary distributions of $((\hat{Q}_i^r(\cdot)), (\hat{Z}_j^r(\cdot)))$. Then, it follows from [9, Theorem 8.5.1] (whose conditions are easily verified in our case), that as $r \rightarrow \infty$, any weak limit of the sequence of stationary distributions of the processes

$((\hat{Q}_i^r(\cdot)), (\hat{Z}_j^r(\cdot)))$ is a stationary distribution of the limit process, described in [7, Theorem 4.4], and therefore is the one-dimensional diffusion specified in the statement of the corollary. \square

Finally, we remark that a tightness result analogous to Theorem 7.4 holds for the underloaded system, $\rho < 1$, and can be proved essentially the same way.

The asymptotic regime in this case is such that $\lambda_i^r = r\lambda_i$ (there is no point in considering $O(\sqrt{r})$ terms in λ_i^r when $\rho < 1$). We denote $Z_j^r(t) = \Psi_j^r(t) - r\beta_j\rho$ (which is consistent with the definition given earlier in this section for $\rho = 1$), and keep notation $Q_i^r(t)$ for the queue length. We work with the following Lyapunov function:

$$\mathcal{L} := \sum_i \left[\exp(\theta(1-\rho)\sqrt{r} + \theta\hat{Q}_i^r) - \exp(\theta(1-\rho)\sqrt{r}) \right] + \sum_j \beta_j \exp(\theta\hat{Z}_j^r/\beta_j).$$

The same approach as in the proof of Theorem 7.4 leads to the following result: for any real θ ,

$$\limsup_r \mathbb{E} \left[\sum_j \exp(\theta\hat{Z}_j^r) \right] < \infty.$$

The limiting process for $(\hat{Z}_j^r(\cdot))$ is $(\hat{Z}_j(\cdot)) = (\frac{\beta_j}{\sum_k \beta_k} \hat{Y}(\cdot))$, with $\hat{Y}(\cdot)$ being a one-dimensional Ornstein-Uhlenbeck process, with Gaussian stationary distribution. The limit of stationary distributions of $(\hat{Z}_j^r(\cdot))$ is the stationary distribution of $(\hat{Z}_j(\cdot))$.

Acknowledgment. Authors would like to thank the referees for useful comments that helped to improve the exposition of the material.

REFERENCES

- [1] M. Armony, A. Ward. Blind Fair Routing in Large-Scale Service Systems. February 2010, preprint. http://www.stern.nyu.edu/om/faculty/armony/research/blind_fair_routing.pdf
- [2] R. Atar, Y. Shaki, A. Shwartz. A blind policy for equalizing cumulative idleness. February 2010, preprint. <http://webee.technion.ac.il/people/atar/equalization.pdf>
- [3] M. Farkas. *Dynamical Models in Biology*. Academic Press 2001.
- [4] D. Gamarnik, A. L. Stolyar. Stationary distribution of multiclass multi-server queueing system: Exponential bounds in the Halfin-Whitt regime. *Queueing Systems*, to appear.
- [5] D. Gamarnik, A. Zeevi. Validity of heavy traffic steady-state approximations in generalized Jackson networks. *The Annals of Applied Probability* vol. 16, 2006, pp.56-90.
- [6] D. Gamarnik, P. Momcilovic. Steady-state analysis of a multiserver queue in the Halfin-Whitt regime. *Advances in Applied Probability* vol. 40, 2008, pp.548-577.
- [7] I. Gurvich, W. Whitt. Queue-and-Idleness-Ratio Controls in Many-Server Service Systems. *Mathematics of OR* vol. 34 no. 2, May 2009, pp. 363-396.
- [8] I. Karatzas, S. Shreve. *Brownian Motion and Stochastic Calculus (2nd ed.)*. Springer 1996.
- [9] R. Sh. Liptser, A. N. Shiryaev. *Theory of Martingales*. Kluwer Academic Publishers 1989 (translated from Russian by K. Dzjaparidze; in Russian Nauka 1986)
- [10] A. Mandelbaum, A. L. Stolyar. Scheduling flexible servers with convex delay costs: Heavy-traffic optimality of the generalized $c\mu$ -rule. *Operations Research* vol. 52, 2004, pp. 836-855.
- [11] C. D. Meyer. *Matrix analysis and applied linear algebra, Volume 1*. SIAM 2000.
- [12] A. L. Stolyar, T. Tezcan. Shadow-routing based control of flexible multi-server pools in overload. *Operations Research* vol. 59, 2011, No.6, pp. 1427-1444.
- [13] A. L. Stolyar, T. Tezcan. Control of systems with flexible multi-server pools: A shadow routing approach. *Queueing Systems*, vol. 66, 2010, pp. 1-51.
- [14] Supporting computations. <http://www.statslab.cam.ac.uk/~ey221/LQFS-LB/>

BELL LABS, ALCATEL-LUCENT, MURRAY HILL, NJ, USA, STOLYAR@RESEARCH.BELL-LABS.COM

DEPARTMENT OF PURE MATHEMATICS AND MATHEMATICAL STATISTICS, UNIVERSITY OF CAMBRIDGE,
UK, E.YUDOVINA@STATSLAB.CAM.AC.UK