

# Power-of- $d$ -Choices with Memory: Fluid Limit and Optimality

Jonatha Anselmi\* and Francois Dufour†

## Abstract

In multi-server distributed queueing systems, the access of stochastically arriving jobs to resources is often regulated by a dispatcher, also known as load balancer. A fundamental problem consists in designing a load balancing algorithm that minimizes the delays experienced by jobs. During the last two decades, the power-of- $d$ -choice algorithm, based on the idea of dispatching each job to the least loaded server out of  $d$  servers randomly sampled at the arrival of the job itself, has emerged as a breakthrough in the foundations of this area due to its versatility and appealing asymptotic properties. In this paper, we consider the power-of- $d$ -choice algorithm with the addition of a local memory that keeps track of the latest observations collected over time on the sampled servers. Then, each job is sent to a server with the lowest observation. We show that this algorithm is asymptotically optimal in the sense that the load balancer can always assign each job to an idle server in the large-system limit. This holds true if and only if the system load  $\lambda$  is less than  $1 - \frac{1}{d}$ . If this condition is not satisfied, we show that queue lengths are tightly bounded by  $\left\lceil -\frac{\log(1-\lambda)}{\log(\lambda d+1)} \right\rceil$ . This is in contrast with the classic version of the power-of- $d$ -choice algorithm, where at the fluid scale a strictly positive proportion of servers containing  $i$  jobs exists for all  $i \geq 0$ , in equilibrium. Our results quantify and highlight the importance of using memory as a means to enhance performance in randomized load balancing.

## 1 Introduction

In multi-server distributed queueing systems, the access of stochastically arriving jobs to resources, or servers, is often regulated by a central dispatcher, also known as load balancer. A fundamental problem consists in designing a load balancing algorithm able to minimize the delays experienced by jobs. In this paper, we are interested in a setting where a traffic of rate  $\lambda N$  needs to be distributed across  $N$  unit-rate parallel servers, each with its own queue, as indicated in Figure 1. The load balancer may rely on feedback information coming from the servers, which may also be stored in a local memory. Depending on the architecture, feedback information can arrive at the dispatcher through a push- or pull-based mechanism. In the former, the dispatcher initiates the communication fetching the requested information from the servers, while in the latter servers periodically send state information to the dispatcher. This type of model finds applications in computer and communication systems, hospitals and road networks, and there exists a significant and growing number of references; see, e.g., the recent works Ying et al. [25], Gardner et al. [9], Gamarnik et al. [8], Gupta and Walton [11] and the references therein. Nevertheless, it is often difficult to establish whether an algorithm is better than another because in general the answer strongly depends on the underlying architecture, application or traffic conditions. For instance, assigning jobs to servers uniformly at random or in a cyclic fashion provides a very scalable dispatching scheme as it requires neither static nor dynamic information about servers but the resulting performance is quite poor; the join-the-shortest-queue algorithm is “optimal” under some conditions, Winston [23] and Weber [22], but its applicability in large systems is debated due to the high communication overhead between the servers and the dispatcher; the join-the-idle-queue algorithm, Lu et al. [12], performs very well when  $N$  is large, see Stolyar [18], but poorly when  $\lambda$  gets close to its critical value, and in addition it requires servers to generate messages on their own; for a more complete discussion, we point the reader to the recent survey in Van der Boor et al. [20].

During the last two decades, the power-of- $d$ -choice algorithm, introduced in Mitzenmacher [15], Vvedenskaya et al. [21] and referred to as  $SQ(d)$ , has emerged as a breakthrough in the foundations of this area due

\*INRIA Bordeaux Sud Ouest, 200 av. de la Vieille Tour, 33405 Talence, France. Email: jonatha.anselmi@inria.fr

†INRIA Bordeaux Sud Ouest, 200 av. de la Vieille Tour, 33405 Talence, France. Email: francois.dufour@math.u-bordeaux.fr

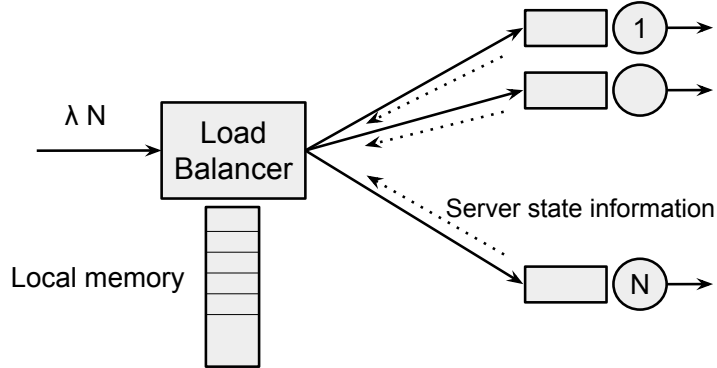


Figure 1: Architecture of the distributed system for load balancing.

to its versatility and its appealing asymptotic properties. It works as follows: upon arrival of each job,  $d \geq 2$  servers are contacted uniformly at random, their state (e.g., queue length or workload) is retrieved, and then the job is dispatched to a server in the best observed state (among the  $d$  selected). The first remarkable property is that in the large-system limit,  $N \rightarrow \infty$ , the stationary proportion of servers with at least  $i$  jobs decreases doubly exponentially in  $i$ , though it remains strictly positive for all  $i$ . This result has been generalized in Bramson et al. [3] to the case where service times are heavy-tailed rather than exponential; see also Bramson et al. [4]. In addition, it turns out that  $SQ(d)$  is *heavy-traffic optimal* in the sense that it minimizes the workload or queue-length process over all time in the diffusion limit where  $\lambda \uparrow 1$ ; see Chen and Ye [5] and Maguluri et al. [14]. In Ying et al. [25], it is also shown that the number of sampled servers can be dramatically reduced if tasks arrive in batches, which is useful to reduce the communication overhead between the load balancer and the servers. In Mitzenmacher et al. [16], the power-of- $d$ -choice algorithm is studied in the case where the load balancer is endowed with a local memory that stores the index and the state of the least loaded server out of the  $d$  sampled each time a job arrives. When the  $n$ -th job arrives, the winning server is chosen among the  $d$  servers randomly selected upon its arrival and the server associated to the observation stored in the memory. The resulting performance is better than the one achieved by  $SQ(2d)$ . In the standard memoryless case, if  $d$  is allowed to depend on  $N$  and  $d(N) \rightarrow \infty$ ,  $SQ(d)$  has been recently shown to become *fluid (or mean-field) optimal*, i.e., optimal in the large-system limit, with a diffusion limit matching the one of the celebrated join-the-shortest-queue algorithm provided that  $d(N)$  grows to infinity sufficiently fast; see Mukherjee et al. [17], Dieker and Suk [6]. At a fluid scale, optimality here is related to the ability of assigning each incoming job to an idle server. Also our work aims at achieving fluid optimality but we will consider  $d$  as a constant to keep the communication overhead at a minimum. Towards this purpose, we will show that it is enough to endow the load balancer with a local memory that keeps track of the latest observation collected on each server. This approach is also close to Mitzenmacher et al. [16], though different because in that reference the memory can only store one observation. In fact, one observation (or even a finite number of observations) is not enough to achieve fluid optimality; see Gamarnik et al. [7]. We observe that fluid optimality can also be achieved by the join-the-shortest-queue and join-the-idle-queue algorithms. However, these are not directly comparable to our algorithm because they are meant to run on a different architecture (pull-based rather than push-based).

The fact that we consider a memory with  $N$  slots has an impact on our proofs. As discussed in Mitzenmacher et al. [16], if the memory size is uniformly bounded then the observations in the local memory evolve much faster than the actual queue lengths, and in this case to establish fluid limit results one can adopt the ad-hoc proof technique developed in Luczak and Norris [13]. On the other hand, this does not apply to our case because observations and queue lengths evolve within the same timescale. Also the pull-based version of join-the-idle-queue, Lu et al. [12], requires a memory with  $N$  slots but the main difference with respect to our approach is that the information stored in the memory is always up to date, which is not the case within our algorithm.

## 1.1 Contribution.

In Algorithm 1, we provide a pseudocode for the proposed power-of- $d$ -choices algorithm with memory and  $N$  servers, referred to as SQ( $d, N$ ); some variants of such algorithm are also discussed in the Conclusions. Upon arrival of one job, the states collected from  $d$  randomly chosen servers are stored in the local array `Memory`. Then, the job is sent to a server chosen randomly (with replacement) among the ones having the lowest recorded state. Finally, the observation of the selected server is incremented by one.

---

**Algorithm 1** Power-of- $d$ -choices with memory and  $N$  servers.

---

```

1: procedure SQ( $d, N$ )
2:   Memory[ $i$ ] = 0,  $\forall i = 1, \dots, N$ ;
3:   for each job arrival do
4:     for  $i = 1, \dots, d$  do
5:       rnd_server = random( $1, \dots, N$ );
6:       Memory[rnd_server] = get_state(rnd_server);
7:     end for
8:     selected_server = random( $\arg \min_{i \in \{1, \dots, N\}} \text{Memory}[i]$ );
9:     send_job_to(selected_server);
10:    Memory[selected_server]++;
11:  end for
12: end procedure

```

---

It is intuitive that SQ( $d, N$ ) results in more balanced allocations than SQ( $d$ ). This follows by using the coupling argument developed in Theorem 3.5 of Azar et al. [1], which can be adapted to argue that at any point in time the vector of queue lengths achieved with SQ( $d, N$ ) is majorized by the vector of queue lengths achieved with SQ( $d$ ). On the other hand, it is not clear how much such improvement can be. This is the goal of the present paper.

We investigate the time-varying dynamics of SQ( $d, N$ ) by means of a continuous-time Markov chain  $X^N(t)$  that keeps track of the proportion of servers with  $i$  jobs and for which their last observation collected by the load balancer is  $j$ , for all  $i$  and  $j$ . To the best of our knowledge, this is the first paper that studies the dynamics induced by SQ( $d, N$ ). The transition rates of  $X^N(t)$  are non-Lipschitz and a satisfactory analysis of  $X^N(t)$  when  $N$  is finite seems to be out of reach. Our main contributions are as follows:

1. In Theorem 1, we let  $N \rightarrow \infty$  and identify the fluid limit of  $X^N(t)$ , an absolutely continuous function that is interpreted as a first-order approximation of the original model  $X^N(t)$ . The fluid limit is motivated by the fact that real systems are composed of many servers and that it enables a tractable analysis for the dynamics of SQ( $d, N$ ). A fluid limit is necessarily a *fluid solution*, as introduced in Definition 1. The proof of the fluid limit is the main technical part of this work and is given in Section 4. The main difficulty stands in the discontinuous structure of the drift of  $X^N(t)$ ; see Section 2.2 for further details. We obtain the fluid limit under a finite buffer assumption, though as discussed in the Conclusions we believe that this assumption can be relaxed.
2. We then study fixed points, fluid solutions that are constant over time. Theorem 2 shows that there exists a unique fixed point. The general structure of such fixed point as a function of  $\lambda$  is quite involved and implies that in equilibrium
  - a) Fluid queue lengths are uniformly and tightly bounded by  $j^* + 1$ , where

$$j^* \stackrel{\text{def}}{=} \left\lfloor -\frac{\log(1-\lambda)}{\log(\lambda d + 1)} \right\rfloor. \quad (1)$$

This is in contrast with SQ( $d$ ), where queue lengths are unbounded in the sense that a strictly positive proportion of servers containing  $i$  jobs exists for all  $i \geq 0$ , in the fluid equilibrium; see Mitzenmacher [15]. Figure 2 illustrates the behavior of the upper bound  $j^* + 1$  by varying  $\lambda$  and  $d$ , and shows that the size of the most loaded server will remain very small even when  $\lambda$  is very close to its critical value. In fact, even when  $\lambda = 0.995$  and  $d = 2$ , at the fluid scale no server will contain more than just 5 jobs.

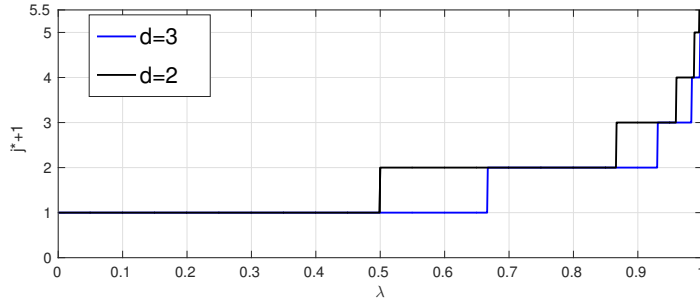


Figure 2: Plots of the maximum queue length,  $j^* + 1$ , by varying  $\lambda$  and  $d$ .

b) The load balancer memory can only contain two possible observations, namely  $j^*$  and  $j^* + 1$ .

The case of particular interest is when  $\lambda < 1 - 1/d$ , where  $j^* = 0$  and thus the load balancer memory always contains a strictly positive proportion of zeros. This means that the load balancer can always assign incoming jobs to idle servers, which is clearly the ideal situation for any incoming job. In this sense we say that  $SQ(d, N)$  is *asymptotically optimal*. When  $\lambda \geq 1 - 1/d$  the load balancer memory will never contain a strictly positive mass of zeros but it will still be able to assign a fraction of jobs to idle servers ensuring that the average number of jobs in each queue belongs to the interval  $j^* - \frac{1}{d} + \frac{1}{2} \pm \frac{1}{2}$  (Proposition 2).

3. Finally, we investigate stability properties of the unique fixed point. Theorem 3 establishes that fluid solutions converge to such point regardless of the initial condition and exponentially fast, provided that  $\lambda < 1 - 1/d$ . Thus, in this case all fluid solutions will be eventually asymptotically optimal as the load balancer memory will eventually be populated by a strictly positive mass of zeros. The proof of this result, given in Section 5.3, is based on a sort of Lyapunov argument that allows us to show that the time evolution of fluid solutions is eventually governed by the unique solution of a linear and autonomous ODE system.

In summary, the proposed algorithm  $SQ(d, N)$  has the same communication overhead of its memoryless counterpart  $SQ(d)$  but a much better performance, which is paid at the cost of endowing the controller with a memory of  $N$  slots. It is to be noted that asymptotic optimality can be obtained

## 2 Performance models

In order to describe the time varying effects of  $SQ(d, N)$  on queue lengths, we introduce a stochastic and a deterministic model. The stochastic model is meant to capture the variability of job interarrival and service times that is intrinsic in multi-server distributed queueing systems. Due to its intractability, a satisfactory analysis of such model is out of reach. In this respect, the deterministic model is convenient because it does enable analytical tractability. In this section, we also show our first result, which states that both models are connected each other: the deterministic can be interpreted as a first-order approximation of the stochastic.

In the following, we will refer to a server with  $i$  jobs and for which its last observation at the controller is  $j$  as an  $(i, j)$ -server.

### 2.1 Markov model.

First, we model the dynamics induced by  $SQ(d, N)$  as a Markov chain in continuous-time: arrivals at the load balancer are assumed to follow a Poisson process with rate  $\lambda N$ , with  $0 < \lambda < 1$ , and service times are independent, exponentially distributed random variables with unit mean. Servers process jobs according to any work-conserving discipline and each of them can contain  $I > 1$  jobs at most. A job that is sent to a server with  $I$  jobs is rejected. Each incoming job is thus assigned to one out of  $N$  queues as in Algorithm 1.

Upon each job arrival, we assume that the actions of sampling  $d$  servers and assigning the job to some queue are instantaneous and occur at the same time.

Let  $(Q^N(t), M^N(t)) = (Q_k^N(t), M_k^N(t))_{k=1}^N \in \{0, 1, \dots, I\}_+^{2N}$  be the system state at time  $t \in \mathbb{R}_+$ :  $Q_k^N(t)$  represents the number of jobs in queue  $k$  at time  $t$  and  $M_k^N(t)$  represents the last observation collected from server  $k$  by the controller at time  $t$ . To avoid unnecessary technical complication in our proofs and since the observation associated to server  $k$  is no less than the actual number of jobs in  $k$  after sampling  $k$  for the first time, for the initial condition we assume that  $Q_k^N(0) \leq M_k^N(0)$  for all  $k$ .

It is convenient to represent the system state by  $X^N(t) = (X_{i,j}^N(t) : 0 \leq i \leq j < \infty)$  where

$$X_{i,j}^N(t) \stackrel{\text{def}}{=} \frac{1}{N} \sum_{k=1}^N \mathbf{1}_{\{Q_k^N(t)=i, M_k^N(t)=j\}} \quad (2)$$

denotes the proportion of  $(i, j)$ -servers at time  $t$ . It is clear that  $X^N(t)$  is still a Markov chain with values in some finite set  $\mathcal{S}_N$  that is a subset of

$$\mathcal{S} \stackrel{\text{def}}{=} \left\{ (x_{i,j} \in \mathbb{R}_+ : 0 \leq i \leq j \leq I) : \sum_{i=0}^I \sum_{j=i}^I x_{i,j} = 1 \right\}. \quad (3)$$

The transitions and rates of the Markov chain  $X^N(t)$  that are due to server departures are easy to write because they have no impact on memory: for  $x \in \mathcal{S}_N$ , the transition  $x \mapsto x - \frac{e_{i+1,j}}{N} + \frac{e_{i,j}}{N}$  occurs with rate  $N x_{i+1,j}$  where  $e_{i,j} \stackrel{\text{def}}{=}} (\delta_{i,i'} \delta_{j,j'} \in \{0, 1\} : 0 \leq i' \leq j' \leq I)$  and  $\delta_{a,b}$  is the Kronecker delta. On the other hand, the transitions and rates of  $X^N(t)$  that are due to job arrivals are quite involved and they are omitted. However, in Section 4.1 we will show how to construct the sample paths of  $X^N(t)$ .

## 2.2 Fluid model.

For any  $x \in \mathcal{S}$ , let

$$x_{i,\cdot} \stackrel{\text{def}}{=} \sum_{j=i}^I x_{i,j} \quad \text{and} \quad x_{\cdot,j} \stackrel{\text{def}}{=} \sum_{i=0}^j x_{i,j}$$

The next definition introduces the *fluid model* for the dynamics of  $SQ(d, N)$ .

**Definition 1.** A function  $x(t) : \mathbb{R}_+ \rightarrow \mathcal{S}$  is said to be a fluid model (or fluid solution) if the following conditions are satisfied:

1.  $x(t)$  is absolutely continuous, and
2.  $\frac{dx_{i,j}(t)}{dt} = b_{i,j}(x(t))$  almost everywhere, for every  $i \geq 0$  and  $j \geq i$ ,

where  $b(x) \stackrel{\text{def}}{=}} (b_{i,j}(x) : 0 \leq i \leq j \leq I)$  is given by

$$b_{0,0}(x) = \lambda d(x_{0,\cdot} - x_{0,0}) - \lambda + \mathcal{R}_0(x) \quad (4)$$

$$\begin{aligned} b_{i,j}(x) = & x_{i+1,j} - \mathbf{1}_{\{i>0\}} x_{i,j} - \lambda d x_{i,j} - \mathcal{R}_{j-1}(x) \frac{x_{i,j}}{x_{\cdot,j}} \mathbf{1}_{\{x_{\cdot,j}>0\}} \\ & + \mathbf{1}_{\{i>0\}} \mathcal{R}_{j-2}(x) \frac{x_{i-1,j-1}}{x_{\cdot,j-1}} \mathbf{1}_{\{x_{\cdot,j-1}>0\}} + \mathbf{1}_{\{j=I, i>0\}} \mathcal{R}_{I-1}(x) \frac{x_{i-1,I}}{x_{\cdot,I}} \mathbf{1}_{\{x_{\cdot,I}>0\}}, \quad \forall i, j : i < j \end{aligned} \quad (5)$$

$$b_{1,1}(x) = -x_{1,1} + \lambda d(x_{1,\cdot} - x_{1,1}) + \lambda - \mathcal{R}_0(x) - \mathcal{R}_0(x) \frac{x_{1,1}}{x_{\cdot,1}} \mathbf{1}_{\{x_{\cdot,1}>0\}} - \mathcal{G}_1(x) \quad (6)$$

$$\begin{aligned} b_{i,i}(x) = & -x_{i,i} + \lambda d(x_{i,\cdot} - x_{i,i}) - \mathcal{R}_{i-1}(x) \frac{x_{i,i}}{x_{\cdot,i}} \mathbf{1}_{\{x_{\cdot,i}>0\}} + \mathcal{R}_{i-2}(x) \frac{x_{i-1,i-1}}{x_{\cdot,i-1}} \mathbf{1}_{\{x_{\cdot,i-1}>0\}} \\ & + \mathcal{G}_{i-1}(x) - \mathcal{G}_i(x), \quad \forall i = 2, \dots, I-1 \end{aligned} \quad (7)$$

$$b_{I,I}(x) = -x_{I,I} + \mathcal{R}_{I-2}(x) \frac{x_{I-1,I-1}}{x_{\cdot,I-1}} \mathbf{1}_{\{x_{\cdot,I-1}>0\}} + \mathcal{G}_{I-1}(x) + \mathcal{R}_{I-1}(x) \frac{x_{I-1,I}}{x_{\cdot,I}} \mathbf{1}_{\{x_{\cdot,I}>0\}} \quad (8)$$

with

$$\mathcal{R}_j(x) = 0 \vee \lambda \left( 1 - d \sum_{i=0}^j (j+1-i) x_{i,\cdot} \right) \mathbf{1}_{\{\sum_{i=0}^j x_{\cdot,i}=0\}} \quad (9)$$

$$\mathcal{G}_j(x) = \lambda d \mathbf{1}_{\{\sum_{i=0}^j x_{\cdot,i}=0, d \sum_{i=0}^j (j+1-i) x_{i,\cdot} \leq 1\}} \sum_{i=0}^j x_{i,\cdot} \quad (10)$$

and  $a \vee b \stackrel{\text{def}}{=} \max\{a, b\}$ .

The discontinuous function  $b$  will be referred to as *drift*, and to some extent it may be interpreted as the conditional expected change from state  $x$  of the Markov chain  $X^N(t)$ , though this may only be true when  $x_{0,0} > 0$ , where  $\mathcal{R}_j(x) = 0$  for all  $j$  and the formulas above become linear admitting a very intuitive explanation.

Let us provide some intuition for the drift expressions in Definition 1, and let us start with coordinates  $(0,0)$ . At the moment of each arrival at the load balancer, the states of  $d$  servers are sampled and  $k$  idle servers that the load balancer has not yet spotted are sampled with probability  $\binom{d}{k} (x_{0,\cdot} - x_{0,0})^k (1 - x_{0,\cdot} + x_{0,0})^{d-k}$ . Since  $\sum_{k=1}^d k \binom{d}{k} (x_{0,\cdot} - x_{0,0})^k (1 - x_{0,\cdot} + x_{0,0})^{d-k} = d(x_{0,\cdot} - x_{0,0})$  and arrivals occur with rate  $\lambda$ , the average rate in which  $(0, j)$ -servers,  $j \geq 1$ , are discovered is  $\lambda d(x_{0,\cdot} - x_{0,0})$  and the rationale behind the first term in (4) is justified. The dynamics that remain to specify are the ones related to the effective job assignments, that is where singularities can happen. In order to build a fluid model ‘consistent’ with the finite stochastic system  $X^N(t)$ , one should take into account the fluctuations of order  $1/N$  that appear when  $X_{0,0}^N(t) = 0$ . These bring discontinuities in the drift. Let  $z_j \stackrel{\text{def}}{=} \sum_{i \geq j} x_{i,\cdot}$  and  $R_j^N(t) \stackrel{\text{def}}{=} \sum_{i=0}^j X_{\cdot,i}^N(t)$ . We notice that  $\mathcal{R}_j(x)/\lambda$ , where  $\mathcal{R}_j(x)$  is defined in (9), will be interpreted as the proportion of time where the process  $(R_j^N(t))_{[t, t+\epsilon]}$  tends to stay on zero with the load balancer sampling  $(\cdot, j')$ -servers only, for all  $j' > j$ , in the limit where  $N \rightarrow \infty$  first and then  $\epsilon \downarrow 0$ ; this will be formalized in Section 4.3.2. Thus, the term  $\lambda - \mathcal{R}_0(x)$  represents the rate in which jobs are assigned to  $(0,0)$ -servers, which become  $(1,1)$ -servers as soon as they receive a job. This explains the drift expression in (4). The particular structure of  $\mathcal{R}_j(x)$  given in (9) will be the outcome of the stochastic analysis that will be developed in Section 4.

Let us provide some intuition also for the drift expression on coordinates  $(1,1)$  (see (6)), as it brings some additional interpretation that also applies on general coordinates. The first term says that departures from  $(1,1)$ -servers occur with rate  $x_{1,1}$  and the second one says that new  $(1,1)$ -servers are discovered with rate  $\lambda d(x_{1,\cdot} - x_{1,1})$ . This can be easily justified as done above for the first summation term of  $b_{0,0}$ . Then, we notice that the  $\lambda - \mathcal{R}_0(x)$  term has been already interpreted above and thus the dynamics that remain to specify are the ones related to job assignments at  $(1,1)$ -servers. According to  $\text{SQ}(d, N)$ , if the load balancer knows no  $(0,0)$ -server then it randomizes over the set of  $(\cdot, 1)$ -servers, and thus within this scenario  $x_{1,1}$  should decrease with rate proportional to  $\frac{x_{1,1}}{x_{\cdot,1}}$ . This is indeed the case if  $x_{\cdot,1} > 0$ . Thus,  $\frac{x_{1,1}}{x_{\cdot,1}} \mathcal{R}_0(x)$  is the rate in which jobs are assigned to  $(1,1)$ -servers when  $x_{\cdot,1} > 0$ . It remains to model the rate in which jobs are assigned to  $(1,1)$ -servers when  $x_{\cdot,1} = 0$ . Since we aim at building a deterministic model ‘consistent’ with the stochastic one, to model the rate of job assignments to  $(1,1)$ -servers when  $x_{\cdot,1} = 0$  one should take into account the fluctuations of order  $1/N$  that appear when  $R_1^N(t) = 0$ . The term  $\mathcal{G}_1(x)$  given in (10) is indeed such rate, and again it will be the outcome of the stochastic analysis developed in Section 4.

The following proposition will be proven in Section 4.

**Proposition 1.** *Fluid solutions exist.*

### 2.3 Connecting the Markov and the fluid models.

Our first result is the following connection between the stochastic and the fluid models.

**Theorem 1.** *Assume that  $X^N(0) \rightarrow x^0 \in \mathcal{S}$  almost surely. With probability one, any limit point of the stochastic process  $(X^N(t))_{t \in [0, T]}$  satisfies the conditions that define a fluid solution.*

In view of this result, proven in Section 4, a fluid solution may be interpreted as an accurate approximation of the time-dependent dynamics of the finite stochastic system  $X^N(t)$ , provided that  $N$  is sufficiently large.

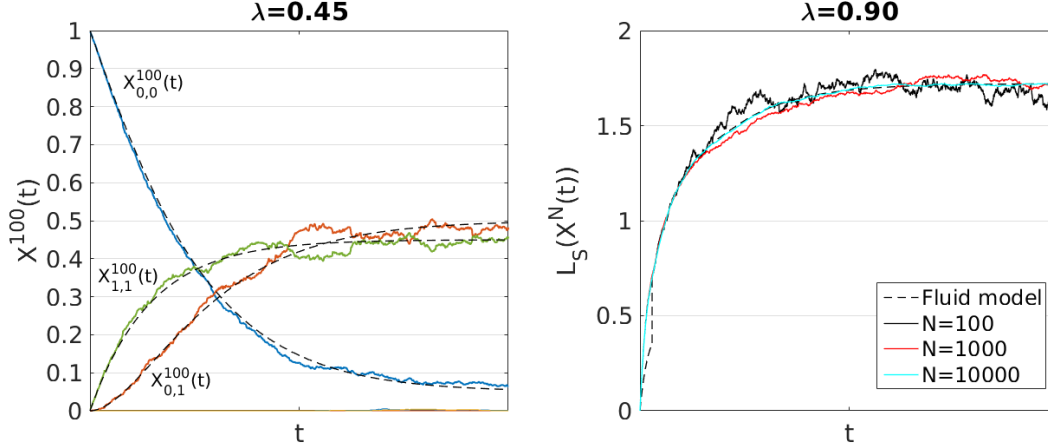


Figure 3: Numerical convergence of the stochastic model  $X^N(t)$  (continuous lines) to the fluid model  $x(t)$  (dashed line).

Given  $x \in \mathcal{S}$ , let us define the functions

$$\mathcal{L}_S(x) \stackrel{\text{def}}{=} \sum_{i=1}^I ix_i, \quad \mathcal{L}_M(x) \stackrel{\text{def}}{=} \sum_{j=1}^I jx_{\cdot,j}$$

We notice that  $\mathcal{L}_S(X^N(t))$  represents the number of jobs in the system at time  $t$  scaled by  $N$  and that  $\mathcal{L}_M(X^N(t))$  represents the number of jobs scaled by  $N$  the load balancer *believes* are in the system at time  $t$ . Since the system is symmetric with respect to the servers, the function  $\mathcal{L}_S(X^N(t))$  is also interpreted as the average number of jobs at time  $t$  in each queue.

It is clear that  $\mathcal{L}_M(x) - \mathcal{L}_S(x) = \sum_{i \geq 0} \sum_{j \geq i} (j - i)x_{i,j} \geq 0$ , which is to be expected because  $(\cdot, j)$ -servers can not contain more than  $j$  jobs by definition.

The following corollary of Theorem 1 is immediate.

**Corollary 1.** *Let  $x(t)$  be a fluid solution. Assume that  $(x(t))_{t \in [0, T]}$  is a limit point of  $(X^N(t))_{t \in [0, T]}$  with probability one. Then,  $(\mathcal{L}_S(x(t)))_{t \in [0, T]}$  and  $(\mathcal{L}_M(x(t)))_{t \in [0, T]}$  are limit points of  $(\mathcal{L}_S(X^N(t)))_{t \in [0, T]}$  and  $(\mathcal{L}_M(X^N(t)))_{t \in [0, T]}$ , respectively, with probability one.*

We complement Theorem 1 and Corollary 1 presenting some numerical simulations to support the claim that the fluid model provides a remarkably accurate approximation of the sample paths of  $X^N(t)$  even when  $N$  is finite and relatively small. Assuming  $d = 2$ , Figure 2.3 plots the time dependent dynamics of  $X^N(t)$  and  $x(t)$ . At time zero, we have chosen  $X^N(0)$  and  $x(0)$  such that  $X_{0,0}^N(0) = x_{0,0}(0) = 1$ , which means that all servers are idle and the load balancer is aware of it. Each curve on these plots is an average over ten simulations. The fluid (stochastic) model is always represented by dashed (continuous) lines. In the picture on the left ( $\lambda = 0.45$ ), we set  $N = 100$  and notice that the fluid model already captures in an accurate manner the dynamics of  $X^N(t)$ , which turn out to be concentrated more and more on just three components: namely  $(0,0)$ ,  $(0,1)$  and  $(1,1)$ . Matter of fact  $X_{0,0}^N(t) + X_{0,1}^N(t) + X_{1,1}^N(t)$  gets closer and closer to 1 when both  $N$  and  $t$  increase. In the picture on the right ( $\lambda = 0.9$ ), dynamics are distributed on several components and for convenience we have plotted  $\mathcal{L}_S(X^N(t))$  and its fluid model counterpart  $\mathcal{L}_S(x(t))$ . We notice that  $\mathcal{L}_S(x(t))$  almost overlaps the trajectory of  $\mathcal{L}_S(X^N(t))$  already when  $N = 1000$ . This size is in agreement with the magnitude of modern distributed computing such as web-server farms or data-centers, as they are often composed of (tenths of) thousands of servers.

### 3 Main results

In this section we focus on fluid solutions and investigate optimality and stability properties. First, we are interested in fixed points.

**Definition 2.** *We say that a fluid solution  $x(t)$  is a fixed point if  $b(x(t)) = 0$  for all  $t$ .*

When fluid solutions are fixed points, we drop the dependency on  $t$ .

Let us define  $j^*$  as in (1) and for simplicity let us assume that  $I > j^*$ .

The next result, proven in Section 5.1, establishes the existence and uniqueness of a fixed point and says that its mass is concentrated only on coordinates of the form  $(i, j^*)$  and  $(i, j^* + 1)$ .

**Theorem 2** (Existence and Uniqueness of Fixed Points). *There exists a unique fixed point, say  $x^*$ . It is such that  $x_{\cdot, j^*}^* + x_{\cdot, j^*+1}^* = 1$  and*

$$\lambda d x_{0, j^*}^* = (1 + \lambda d)(1 - \lambda) - \frac{1}{(1 + \lambda d)^{j^*}} \quad (11a)$$

$$x_{0, j^*}^* + x_{0, j^*+1}^* = 1 - \lambda. \quad (11b)$$

In the fixed point, our first remark is that queue lengths are bounded, by  $j^* + 1$ . As we show in our proof, an explicit expression for  $x^*$  seems to be difficult to obtain, though it can be easily computed when  $\lambda$  and  $d$  are fixed numerically. In fact, in Section 5.1 we provide an explicit expression for  $x_{i, j}^*$  when  $(i, j) \neq (j^*, j^*)$  as a function of  $x_{j^*, j^*}^*$ , and identify  $x_{j^*, j^*}^*$  by means of a polynomial equation of degree  $j^* + 1$  (see (61)).

A case of particular interest is when  $j^* = 0$ , which given (1) occurs if and only if  $\lambda < 1 - 1/d$ , where we have the following remark.

**Remark 1** (Asymptotic Optimality). *If  $\lambda < 1 - 1/d$ , then Theorem 2 implies that  $x_{0,0}^* = 1 - \lambda - 1/d$ ,  $x_{0,1}^* = 1/d$ ,  $x_{1,1}^* = \lambda$  and  $x_{i,j}^* = 0$  on the remaining coordinates. Thus, provided that dynamics converge to  $x^*$ , we have shown that a load balancer implementing  $SQ(d, N)$  is always aware of the fact that some servers are idle when  $N \rightarrow \infty$  and  $t$  is sufficiently large because  $x_{0,0}^* > 0$ . In this scenario, the load balancer can certainly assign each incoming job to one of such idle servers, and the job itself would incur zero delay. This is in fact the ideal situation for any arriving job and in this sense we say that  $SQ(d, N)$  is asymptotically optimal.*

The next proposition provides further insights on the system performance at the fixed point  $x^*$ .

**Proposition 2.** *Let  $x^*$  as in Theorem 2. Then,*

$$\mathcal{L}_M(x^*) = \mathcal{L}_S(x^*) + \frac{1}{d} \quad (12)$$

and

$$j^* - \frac{1}{d} \leq \mathcal{L}_S(x^*) \leq j^* - \frac{1}{d} + 1. \quad (13)$$

Proposition 2, proven in Section 5.2, provides simple bounds on the average number of jobs in each queue. It also says that there is a fluid mass equal to  $1/d$  that the load balancer will never spot. In other words, the samplings performed by the load balancer at each arrival will correctly build the true state of the system up to an (absolute) error of  $1/d$ .

In Remark 1, we discussed the asymptotic optimality of  $SQ(d, N)$  postulating some form of stability for fluid solutions when  $t \rightarrow \infty$ . The next result shows that fluid solutions are indeed globally stable and that convergence to  $x^*$  occurs exponentially fast, provided that  $\lambda < 1 - 1/d$ .

**Theorem 3** (Global Stability). *Let  $x(t)$  be a fluid solution. If  $\lambda < 1 - 1/d$ , then there exist  $\alpha > 0$  and  $\beta > 0$  independent of  $t$  such that*

$$\|x(t) - x^*\| \leq \alpha e^{-\beta t}, \quad \forall t \quad (14)$$

where  $\|\cdot\|$  is the Euclidean norm.



The proof of Theorem 3 is given in Section 5.3 and is based on the following ‘Lyapunov-type’ argument. When  $x_{0,0}(t) = 0$ , we first show that  $\dot{\mathcal{L}}_S(x(t)) \leq \lambda - 1 + \frac{1}{d}$ , which implies that  $\mathcal{L}_S(x(t))$  decreases with derivative bounded away from zero. However, since  $\mathcal{L}_S(x(t)) \geq 0$ ,  $x_{0,0}(t)$  must necessarily increase in finite time, and when it does we show that  $x(t)$  is uniquely determined by the unique solution of a linear ODE system of the form  $\dot{x} = A(x - x^*)$ . At this point, (14) follows by standard results of ODE theory. When  $\lambda \geq 1 - 1/d$ , a generalization of this argument is complicated by the involved structure of  $x^*$  and the fact that  $\mathcal{L}_S(x(t))$  is in general not monotone. However, we conjecture that  $x^*$  remains globally stable. This is also confirmed by the numerical simulations shown in Section 2.3.

## 4 Connection between the fluid and the Markov models

We now prove that the sequence of stochastic processes  $\{(X^N(t))_{t \in [0, T]}\}_{N=d}^{\infty}$  converges almost surely, as  $N \rightarrow \infty$ , to a fluid solution, for any  $T > 0$ . This proves Proposition 1 and Theorem 1.

Our proof is based on three steps. First, we construct the sample paths of the process  $X^N(t)$  on each pair of coordinates. This is achieved using a common coupling technique that defines the processes  $(X^N(t))_{t \in [0, T]}$  for all  $N \in \mathbb{Z}_+$  on a single probability space and in terms of a finite number of ‘‘fundamental processes’’. Then, we show that limit trajectories exist and are Lipschitz continuous with probability one. This is done by using standard arguments, e.g., Gamarnik et al. [7], Tsitsiklis and Xu [19], and Bramson [2]. Finally, we prove that any such limit trajectory must be a fluid solution, which is the main difficulty. This last step is based on technical arguments that are specific to the stochastic model under investigation.

### 4.1 Probability space and coupled construction of sample paths.

We construct a probability space where the stochastic processes  $\{(X^N(t))_{t \in [0, T]}\}_{N \geq d}$  are coupled. All the processes of interest will be a function of the following fundamental processes, all of them independent of each other:

- $\mathcal{N}_\lambda(t)$ , the Poisson processes of job arrivals, with rate  $\lambda$ , defined on  $(\Omega_A, \mathcal{A}_A, \mathbb{P}_A)$ ;
- $\mathcal{N}_1(t)$ , the Poisson processes of potential job departures, with rate 1, defined on  $(\Omega_D, \mathcal{A}_D, \mathbb{P}_D)$ ;
- $V_n^p$  for all  $p = 1, \dots, d$ ,  $(W_n)_n$ ,  $(U_n)_n$ , where the random variables  $V_n^p$ ,  $W_n$  and  $U_n$ , for all  $n$  and  $p$ , are all independent and uniformly distributed over the interval  $(0, 1]$ . These are *selection* processes:  $(V_n^p)_n$  will select the servers to sample at each arrival (see Line 5 of Algorithm 1),  $(W_n)_n$  will be used to randomize among the servers having the lowest observations (see Line 8 of Algorithm 1) and  $(U_n)_n$  will select the server that fires a departure. These  $2 + d$  processes are defined on  $(\Omega_S, \mathcal{A}_S, \mathbb{P}_S)$ ;
- $(X^N(0))_N$ , the process of the initial conditions, where each random variable  $X^N(0)$  takes values in  $\mathcal{S}$ , defined on  $(\Omega_0, \mathcal{A}_0, \mathbb{P}_0)$ .

Each process  $\{(X^N(t))_{t \in [0, T]}\}$ , with  $N \geq d$ , can be constructed on  $(\Omega, \mathcal{A}, \mathbb{P}) = (\Omega_A \times \Omega_D \times \Omega_S \times \Omega_0, \mathcal{A}_A \times \mathcal{A}_D \times \mathcal{A}_S \times \mathcal{A}_0, \mathbb{P}_A \times \mathbb{P}_D \times \mathbb{P}_S \times \mathbb{P}_0)$  by using that  $\mathcal{N}_\lambda(Nt) =_{st} \mathcal{N}_{\lambda N}(t)$ , where  $=_{st}$  denotes equality in distribution. This equality ensures that the Poisson process with rate  $\lambda N$ , which represents the arrival process associated to the  $N$ -th system, is coupled with the fundamental Poisson process  $\mathcal{N}_\lambda(t)$ . Since  $\mathcal{N}_1(Nt) =_{st} \mathcal{N}_N(t)$ , this coupling is also used for the processes of potential job departures.

Now, let  $t_n^{N, \lambda}$  and  $t_n^{N, 1}$  be the times of the  $n$ -th jump of the Poisson processes  $\mathcal{N}_\lambda(Nt)$  and  $\mathcal{N}_1(Nt)$ , respectively. Let also  $X_{i,j}^N(t^-) \stackrel{\text{def}}{=} \lim_{s \uparrow t} X_{i,j}^N(s)$  and  $X_{i,j}^N(t) \stackrel{\text{def}}{=} \sum_{j \geq i} X_{i,j}^N(t)$ . In view of the coupling discussed above, we can construct  $X_{0,0}^N(t)$  as follows

$$X_{0,0}^N(t) = X_{0,0}^N(0) + \frac{1}{N} \sum_{n=1}^{\mathcal{N}_\lambda(Nt)} \sum_{p=1}^d \mathbb{I}_{(X_{0,0}^N(t_n^{N, \lambda^-}), X_{0,0}^N(t_n^{N, \lambda^-})]}^{(V_n^p)} \quad (15a)$$

$$+ \frac{1}{N} \sum_{n=1}^{\mathcal{N}_\lambda(Nt)} \left( \mathbf{1}_{\{X_{0,0}^N(t_n^{N, \lambda^-})=0\}} \prod_{p=1}^d \mathbb{I}_{(X_{0,0}^N(t_n^{N, \lambda^-}), 1]}^{(V_n^p)} - 1 \right). \quad (15b)$$

In the above expression, the term (15a) corresponds to the action of sampling  $d$  servers and the term (15b) corresponds to the action of assigning each job to a server. At the arrival of the  $n$ -th job,  $t_n^{N,\lambda}$ , the proportion of  $(0,0)$ -servers increases by  $k/N$  if  $k$   $(0,j)$ -servers, for any  $j > 0$ , are sampled, which justifies the term in (15a), and decreases by  $1/N$  except when such proportion is zero immediately before  $t_n^{N,\lambda}$  and no idle server is sampled at time  $t_n^{N,\lambda}$ , which justifies the term in (15b).

Using the random variables  $W_n$  and  $U_n$ , an expression similar to (15) can be written for  $X_{i,j}^N(t)$  when  $(i,j) \in \{0,1,\dots,I\}^2$ . Towards this purpose, let us define

$$R_i^N(t) \stackrel{\text{def}}{=} \sum_{j=0}^i X_{i,j}^N(t), \quad S_{i,j}^N(t) \stackrel{\text{def}}{=} \sum_{i'=0}^{i-1} X_{i',j}^N(t) + \sum_{j' \geq i}^j X_{i,j'}^N(t), \quad Z_i^N(t) \stackrel{\text{def}}{=} \sum_{k \geq i} X_{k,\cdot}^N(t), \quad (16)$$

which respectively represent *i*) the proportion of servers that the controller believes have at most  $i$  jobs, *ii*) the proportion of servers with at most  $i-1$  jobs, or  $i$  jobs but with observation less than or equal to  $j$  and *iii*) the proportion of servers with at least  $i$  jobs, and

$$M_{i,j,n}^N \stackrel{\text{def}}{=} \frac{1}{N} \sum_{p=1}^d \mathbb{I}_{(S_{i,j-1}^N(t_n^{N,\lambda^-}), S_{i,j}^N(t_n^{N,\lambda^-}))}, \quad \underline{M}_{j,n}^N \stackrel{\text{def}}{=} \sum_{i=0}^j M_{i,j,n}^N, \quad \overline{M}_{j,n}^N \stackrel{\text{def}}{=} \sum_{j' \geq j} M_{j',j,n}^N.$$

We notice that  $M_{i,j,n}^N$ ,  $\underline{M}_{j,n}^N$  and  $\overline{M}_{j,n}^N$  are the scaled-by- $N$  numbers of  $(i,j)$ -,  $(\cdot,j)$ - and  $(j,\cdot)$ - servers sampled immediately before time  $t_n^{N,\lambda}$ , respectively. Furthermore, let also

$$F_{i,j,n}^N \stackrel{\text{def}}{=} \mathbb{I}_{(W_n(X_{i,j}^N(t_n^{N,\lambda^-}) + \overline{M}_{j,n}^N - \underline{M}_{j,n}^N) \in (\sum_{k=0}^{i-1} X_{k,j}^N(t_n^{N,\lambda^-}) - M_{k,j,n}^N, \sum_{k=0}^i X_{k,j}^N(t_n^{N,\lambda^-}) - M_{k,j,n}^N])} \quad (17)$$

if  $i < j$ , and

$$F_{j,j,n}^N \stackrel{\text{def}}{=} \mathbb{I}_{(W_n(X_{j,j}^N(t_n^{N,\lambda^-}) + \overline{M}_{j,n}^N - \underline{M}_{j,n}^N) \in (\sum_{k=0}^{j-1} X_{k,j}^N(t_n^{N,\lambda^-}) - M_{k,j,n}^N, X_{j,j}^N(t_n^{N,\lambda^-}) + \overline{M}_{j,n}^N - \underline{M}_{j,n}^N])} \quad (18)$$

if  $j \geq 1$ . For all  $i \leq j$ , the random variable  $F_{i,j,n}^N$  will be used to handle the randomness in Line 8 of Algorithm 1 and thus perform a job assignment to a  $(i,j)$ -server, which needs to be chosen in the set of  $(\cdot,j)$ -servers. Specifically, we will use  $F_{i,j,n}^N$ , with  $i \leq j$ , in the scenario where  $R_{i-1}^N(t_n^{N,\lambda^-}) = 0$  and  $M_{i',j',n}^N = 0$  for all  $i' < i$  and  $j' \geq i'$ , that is the case where the load balancer memory contains no observation less than  $j$  and no server containing less than  $j$  jobs is sampled immediately before  $t_n^{N,\lambda}$ . In this case, according to SQ( $d, N$ ), the  $n$ -th job must be routed to a random  $(\cdot,j)$ -server, provided that such a server exists. This randomness is captured by the uniform random variable  $W_n$  and we notice that  $N(X_{i,j}^N(t_n^{N,\lambda^-}) + \overline{M}_{j,n}^N - \underline{M}_{j,n}^N)$  is the number of  $(\cdot,j)$ -servers, or equivalently the occurrences of  $j$  in the memory of the load balancer, at the arrival of the  $n$ -th job and after having performed the associated sampling of the states of  $d$  random servers. Within these conditions, the job arriving at time  $t_n^{N,\lambda}$  is routed to an  $(i,j)$ -server if and only if  $F_{i,j,n}^N = 1$ .

Provided that  $i < j$ , the following formula constructs the process  $X^N(t)$  on coordinates  $(i,j)$

$$X_{i,j}^N(t) = X_{i,j}^N(0) + \frac{1}{N} \sum_{n=1}^{\mathcal{N}_1(Nt)} \mathbb{I}_{(U_n \in (S_{i+1,j}^N(t_n^{N,1^-}) - X_{i+1,j}^N(t_n^{N,1^-}), S_{i+1,j}^N(t_n^{N,1^-}))]} \quad (19a)$$

$$- \frac{\mathbf{1}_{\{i>0\}}}{N} \sum_{n=1}^{\mathcal{N}_1(Nt)} \mathbb{I}_{(U_n \in (S_{i,j}^N(t_n^{N,1^-}) - X_{i,j}^N(t_n^{N,1^-}), S_{i,j}^N(t_n^{N,1^-}))]} \quad (19b)$$

$$- \frac{1}{N} \sum_{n=1}^{\mathcal{N}_\lambda(Nt)} \sum_{p=1}^d \mathbb{I}_{(V_n^p \in (S_{i,j}^N(t_n^{N,\lambda^-}) - X_{i,j}^N(t_n^{N,1^-}), S_{i,j}^N(t_n^{N,\lambda^-}))]} \quad (19c)$$

$$- \frac{1}{N} \sum_{n=1}^{\mathcal{N}_\lambda(Nt)} \mathbf{1}_{\{R_{j-1}^N(t_n^{N,\lambda^-})=0\}} F_{i,j,n}^N \prod_{p=1}^d \mathbb{I}_{(1-Z_j^N(t_n^{N,\lambda^-}),1]} \quad (19d)$$

$$+ \frac{\mathbf{1}_{\{i>0\}}}{N} \sum_{n=1}^{\mathcal{N}_\lambda(Nt)} \mathbf{1}_{\{R_{j-2}^N(t_n^{N,\lambda^-})=0\}} F_{i-1,j-1,n}^N \prod_{p=1}^d \mathbb{I}_{(1-Z_{j-1}^N(t_n^{N,\lambda^-}),1]} \quad (19e)$$

$$+ \frac{\mathbf{1}_{\{j=I, i>0\}}}{N} \sum_{n=1}^{\mathcal{N}_\lambda(Nt)} \mathbf{1}_{\{R_{I-1}^N(t_n^{N,\lambda^-})=0\}} F_{i-1, I, n}^N \prod_{p=1}^d \mathbb{I}_{(1-Z_I^N(t_n^{N,\lambda^-}), 1]}^{(V_n^p)}. \quad (19f)$$

The summations in (19a) and (19b) refer, respectively, to job departures from  $(i+1, j)$ - and  $(i, j)$ -servers, the summation in (19c) refers to the case where  $k$   $(i, j)$ -servers are sampled (as soon as  $k$  of them are sampled, they become  $(i, i)$ -servers, and thus  $X_{i,j}^N$  decreases by  $k/N$ ), and the summations in (19d) and (19e) refer to the case where a job is assigned to an  $(i, j)$ -server and to an  $(i-1, j-1)$ -server, respectively. We notice that a job can be assigned at time  $t_n^{N,\lambda}$  to an  $(i, j)$ -server only if the memory contains no server with observation less than  $j-1$  immediately before  $t_n^{N,\lambda}$  (i.e.,  $R_{j-1}^N(t_n^{N,\lambda^-}) = 0$ ) and no  $(i', j')$ -server, for some  $i' < i$  and for any  $j'$ , has been sampled at time  $t_n^{N,\lambda}$ . Summation (19f) covers the boundary case where  $j = I$  and has the same intuition of term (19f).

Similarly, when  $i = j \geq 1$ , we have

$$X_{i,i}^N(t) = X_{i,i}^N(0) - \frac{1}{N} \sum_{n=1}^{\mathcal{N}_1(Nt)} \mathbb{I}_{(S_{i-1, I}^N(t_n^{N,1^-}), S_{i,i}^N(t_n^{N,1^-}))}^{(U_n)} \quad (20a)$$

$$+ \frac{1}{N} \sum_{n=1}^{\mathcal{N}_\lambda(Nt)} \sum_{p=1}^d \mathbb{I}_{(S_{i,i}^N(t_n^{N,\lambda^-}), S_{i,i}^N(t_n^{N,\lambda^-}))}^{(V_n^p)} \quad (20b)$$

$$+ \frac{\mathbf{1}_{\{i=1\}}}{N} \sum_{n=1}^{\mathcal{N}_\lambda(Nt)} \mathbf{1}_{\{X_{0,0}^N(t_n^{N,\lambda^-}) + \sum_{j=1}^I M_{0,j,n}^N > 0\}} \quad (20c)$$

$$+ \frac{\mathbf{1}_{\{i>1\}}}{N} \sum_{n=1}^{\mathcal{N}_\lambda(Nt)} \mathbf{1}_{\{R_{i-2}^N(t_n^{N,\lambda^-})=0\}} F_{i-1, i-1, n}^N \prod_{p=1}^d \mathbb{I}_{(1-Z_{i-1}^N(t_n^{N,\lambda^-}), 1]}^{(V_n^p)} \quad (20d)$$

$$- \frac{\mathbf{1}_{\{i<I\}}}{N} \sum_{n=1}^{\mathcal{N}_\lambda(Nt)} \mathbf{1}_{\{R_{i-1}^N(t_n^{N,\lambda^-})=0\}} F_{i,i,n}^N \prod_{p=1}^d \mathbb{I}_{(1-Z_i^N(t_n^{N,\lambda^-}), 1]}^{(V_n^p)} \quad (20e)$$

$$+ \frac{\mathbf{1}_{\{i=I\}}}{N} \sum_{n=1}^{\mathcal{N}_\lambda(Nt)} \mathbf{1}_{\{R_{I-1}^N(t_n^{N,\lambda^-})=0\}} F_{I-1, I, n}^N \prod_{p=1}^d \mathbb{I}_{(1-Z_I^N(t_n^{N,\lambda^-}), 1]}^{(V_n^p)}. \quad (20f)$$

The summation in (20a) refers to job departures from  $(i, i)$ -servers, the summation in (20b) refers to the sampling of  $k$  different  $(i, j)$ -servers, which become  $(i, i)$ -servers immediately after sampling. Finally, the summations in (20c) and (20e) refer to jobs assignments and have the same intuition of (19d) and (19e).

## 4.2 Limit trajectories are Lipschitz.

With respect to a set of sample paths  $\omega$  having probability one, we show that any subsequence of the sequence  $\{X^N(\omega, t)\}_N$  contains a further subsequence  $\{X^{N_k}(\omega, t)\}_k$  that converges to some Lipschitz continuous function  $x$ . This proves tightness of sample paths.

First, let us introduce the following formulas for quick reference. These can be proven in a straightforward manner using the strong law of the large numbers and the functional strong law of large numbers for the Poisson process.

**Lemma 1.** *Let  $T > 0$  and  $a, b \in [0, 1]^d$  such that  $a_k \leq b_k$  for all  $k = 1, \dots, d$ . There exists  $\mathcal{C} \subseteq \Omega$  such that  $\mathbb{P}(\mathcal{C}) = 1$  such that*

$$\begin{aligned} \lim_{N \rightarrow \infty} \sup_{t \in [0, T]} |\frac{1}{N} \mathcal{N}_\lambda(Nt, \omega) - \lambda t| &= 0, & \lim_{N \rightarrow \infty} \sup_{t \in [0, T]} |\frac{1}{N} \mathcal{N}_1(Nt, \omega) - t| &= 0, \\ \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=1}^N \sum_{p=1}^d \mathbb{I}_{(a_k, b_k]}^{(V_n^p(\omega))} &= \sum_{p=1}^d b_k - a_k, & \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=1}^N \prod_{p=1}^d \mathbb{I}_{(a_p, b_p]}^{(V_n^p(\omega))} &= \prod_{p=1}^d (b_p - a_p) \end{aligned}$$

for all  $\omega \in \mathcal{C}$ .

In the following, we will work on the set  $\mathcal{C}$  introduced in previous lemma and we will also often use that

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=1}^{\mathcal{N}_\lambda(Nt, \omega)} \prod_{p=1}^d \mathbb{I}_{(a_k, b_k]}^{(V_n^p(\omega))} = \lambda t \prod_{p=1}^d (b_k - a_k)$$

by the renewal theorem.

Let  $x^0 \in [0, 1]$ , sequences  $R_n \downarrow 0$  and  $\gamma_n \downarrow 0$ , and a constant  $L > 0$  be given. For  $n \geq 1$ , let also

$$E_N(R_N, \gamma_N, L, x^0) \stackrel{\text{def}}{=} \{x \in D[0, T] : |x(0) - x^0| \leq R_N, |x(a) - x(b)| \leq L|a - b| + \gamma_N, \forall a, b \in [0, T]\}$$

and

$$E_c(L, x^0) \stackrel{\text{def}}{=} \{x \in D[0, T] : x(0) = x^0, |x(a) - x(b)| \leq L|a - b|, \forall a, b \in [0, T]\}.$$

The next lemma says that the sample paths along any coordinates  $(i, j)$  is approximately Lipschitz continuous. The proof is omitted because follows exactly the same standard arguments used in Lemma 5.2 of Gamarnik et al. [8], which basically use the fact that the jumps of the Markov chain of interest are of the order of  $1/N$  and that the evolution of such Markov chain on a given pair of coordinates only depends on the evolution of such Markov chain on a finite number of other coordinates.

**Lemma 2.** *Fix  $T > 0$ ,  $\omega \in \mathcal{C}$ , and some  $x^0 \in \mathcal{S}$ . Suppose that  $\|X^N(\omega, 0) - x^0\| \leq \tilde{R}_N$ , for some sequence  $\tilde{R}_N \downarrow 0$ . Then, there exists sequences  $R_N \downarrow 0$  and  $\gamma_N \downarrow 0$  such that*

$$X_{i,j}^N(\omega, \cdot) \in E_N(R_N, \gamma_N, L, x^0), \quad \forall (i, j) \in \mathbb{Z}^+ : i \leq j, \forall N$$

where  $L = \lambda d + 1$ .

The next proposition says that the sample paths along any coordinates  $(i, j)$  are sufficiently close to a Lipschitz continuous function. The proof is omitted because follows exactly the same arguments used in the proof of Proposition 11 in Tsitsiklis and Xu [19]: it uses Lemma 2 and topological properties of the space  $E_c(L, x^0)$ , i.e., sequential compactness (by the Arzelà-Ascoli theorem) and closedness.

**Proposition 3.** *Fix  $T > 0$ ,  $\omega \in \mathcal{C}$ , and some  $x^0 \in \mathcal{S}$ . Suppose that  $\|X^N(\omega, 0) - x^0\| \leq \tilde{R}_N$ , for some sequence  $\tilde{R}_N \downarrow 0$ . Then, every subsequence of  $\{X^N(\omega, \cdot)\}_{N=1}^\infty$  contains a further subsequence  $\{X^{N_k}(\omega, \cdot)\}_{k=1}^\infty$  such that*

$$\lim_{k \rightarrow \infty} \sup_{t \in [0, T]} |X_{i,j}^{N_k}(\omega, t) - x_{i,j}(t)| = 0, \quad \forall i, j \geq i$$

where  $x_{i,j} \in E_c(1 + \lambda d, x^0)$  for all  $i, j \geq i$ .

Since Lipschitz continuity implies absolute continuity, we have thus obtained that limit points of  $X^N(t)$  are absolutely continuous, and it remains to show that the partial derivatives of  $x(t)$  are given by the expressions in Definition 1.

### 4.3 Limit trajectories are fluid solutions.

To conclude the proof of Theorem 1, it remains to show that any limit point is a fluid solution, i.e., it satisfies the conditions given in Definition 1. This is the main technical difficulty.

Fix  $\omega \in \mathcal{C}$  and let  $\{X^{N_k}(\omega, t)\}_{k=1}^\infty$  be a subsequence that converges to  $\bar{x}$ , i.e.

$$\lim_{k \rightarrow \infty} \sup_{t \in [0, T]} \|X^{N_k}(\omega, t) - \bar{x}(t)\| = 0. \quad (23)$$

Since  $\bar{x}_{i,j}$  must be Lipschitz continuous for all  $i$  and  $j$  by Proposition 3, it is also absolutely continuous and thus it remains to show that

$$\dot{\bar{x}}_{i,j}(t) = \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \lim_{k \rightarrow \infty} X_{i,j}^{N_k}(t + \epsilon) - X_{i,j}^{N_k}(t) = b_{i,j}(\bar{x}(t)), \quad (24)$$

whenever  $\bar{x}_{i,j}(\cdot)$  is differentiable. This will be done in the following subsections. Now, we introduce the following technical lemmas.

**Lemma 3.** Fix  $\omega \in \mathcal{C}$ ,  $\epsilon > 0$  and let (23) hold. Then, for all  $i, j$  and  $t$ ,

$$|X_{i,j}^{N_k}(u) - \bar{x}_{i,j}(t)| \leq 2L\epsilon, \quad \forall u \in [t, t + \epsilon]$$

for all  $k$  sufficiently large, where  $L = \lambda d + 1$ .

*Proof:* By Lemma 2, there exists a sequence  $\gamma_{N_k} \downarrow 0$  such that  $X_{i,j}^{N_k}(\omega, u) \in [\bar{x}_{i,j}(t) - \epsilon L - \gamma_{N_k}, \bar{x}_{i,j}(t) + \epsilon L + \gamma_{N_k}]$ , for all  $u \in [t, t + \epsilon]$ . Thus, for all  $k$  sufficiently large,  $X_{i,j}^{N_k}(\omega, u) \in [\bar{x}_{i,j}(t) - 2\epsilon L, \bar{x}_{i,j}(t) + 2\epsilon L]$ , for all  $u \in [t, t + \epsilon]$ , as desired.  $\square$

As a corollary of Lemma 3, we obtain

$$|S_{i,j}^{N_k}(u) - s_{i,j}(\bar{x}(t))| \leq C\epsilon, \quad \forall u \in [t, t + \epsilon] \quad (25)$$

for all  $k$  sufficiently large, where  $C \stackrel{\text{def}}{=} 2L(I + 1)^2$ .

**Remark 2.** In the following, we will work on any fixed trajectory  $\omega \in \mathcal{C}$  but we will write  $X^{N_k}(t)$ , instead of  $X^{N_k}(\omega, t)$ , for simplicity of notation.

**Lemma 4.** Fix  $\omega \in \mathcal{C}$  and let (23) hold. Then,

$$\lim_{\epsilon \rightarrow 0} \lim_{k \rightarrow \infty} \frac{1}{\epsilon N_k} \sum_{n=\mathcal{N}_\lambda(N_k t)+1}^{\mathcal{N}_\lambda(N_k(t+\epsilon))} \sum_{p=1}^d \mathbb{I}_{(S_{i,j}^{N_k}(t_n^{N_k, \lambda^-}) - X_{i,j}^{N_k}(t_n^{N_k, \lambda^-}), S_{i,j}^{N_k}(t_n^{N_k, \lambda^-}))}^{(V_n^p)} = \lambda d \bar{x}_{i,j}(t) \quad (26a)$$

$$\lim_{\epsilon \rightarrow 0} \lim_{k \rightarrow \infty} \frac{1}{\epsilon N_k} \sum_{n=\mathcal{N}_1(N_k t)+1}^{\mathcal{N}_1(N_k(t+\epsilon))} \mathbb{I}_{(S_{i,j}^{N_k}(t_n^{N_k, \lambda^-}) - X_{i,j}^{N_k}(t_n^{N_k, \lambda^-}), S_{i,j}^{N_k}(t_n^{N_k, \lambda^-}))}^{(U_n)} = \bar{x}_{i,j}(t). \quad (26b)$$

*Proof:* Given in the Appendix.  $\square$

#### 4.3.1 Fluid solution on coordinates (0,0).

The next lemma explicits the derivative of  $\bar{x}_{0,0}(t)$  when  $\bar{x}_{0,0}(t) > 0$ . It also implies that  $\bar{x}_{0,0}(\cdot)$  is differentiable when strictly positive.

**Lemma 5.** Fix  $\omega \in \mathcal{C}$ , let (23) hold and assume  $\bar{x}_{0,0}(t) > 0$ . Then,

$$\dot{\bar{x}}_{0,0}(t) = -\lambda + d\lambda(\bar{x}_{0,\cdot}(t) - \bar{x}_{0,0}(t)). \quad (27)$$

*Proof:* Choose  $\epsilon > 0$  small enough such that  $\bar{x}_{0,0}(t) - 2(I + 1)\epsilon L > 0$  where  $L = \lambda d + 1$ . Such  $\epsilon$  exists because  $\bar{x}_{0,0}(t) > 0$  by hypothesis. Since  $t_n^{N_k, \lambda^-} \in (t, t + \epsilon]$  when  $n \in \{\mathcal{N}_\lambda(N_k t) + 1, \dots, \mathcal{N}_\lambda(N_k(t + \epsilon))\}$ , Lemma 3 implies that for all  $k$  sufficiently large we must have

$$\mathbf{1}_{\{X_{0,0}^{N_k}(t_n^{N_k, \lambda^-}) > 0\}} = 1, \quad \forall n \in \{\mathcal{N}_\lambda(N_k t) + 1, \dots, \mathcal{N}_\lambda(N_k(t + \epsilon))\}.$$

Thus, using (15), we obtain

$$\lim_{k \rightarrow \infty} X_{0,0}^{N_k}(t + \epsilon) - X_{0,0}^{N_k}(t) = \lim_{k \rightarrow \infty} \frac{1}{N_k} \sum_{n=\mathcal{N}_\lambda(N_k t)+1}^{\mathcal{N}_\lambda(N_k(t+\epsilon))} -1 + \sum_{p=1}^d \mathbb{I}_{(X_{0,0}^{N_k}(t_n^{N_k, \lambda^-}), X_{0,\cdot}^{N_k}(t_n^{N_k, \lambda^-}))}^{(V_n^p)}.$$

A direct application of Lemma 4 concludes the proof.  $\square$

The next two lemmas give properties on the boundary where  $\bar{x}_{0,0}(t) = 0$ .

**Lemma 6.** Fix  $\omega \in \mathcal{C}$ , let (23) hold and assume  $\bar{x}_{0,0}(t) = 0$  and  $d\bar{x}_{0,\cdot}(t) > 1$ . Then,  $t$  is not a point of differentiability.

*Proof:* First of all, we notice that

$$\lim_{\epsilon \downarrow 0} \frac{1}{\epsilon} \lim_{k \rightarrow \infty} X_{0,0}^{N_k}(t + \epsilon) - X_{0,0}^{N_k}(t) \geq -\lambda + d\lambda(\bar{x}_{0,\cdot}(t) - \bar{x}_{0,0}(t)). \quad (28)$$

This holds true because

$$X_{0,0}^{N_k}(t + \epsilon) - X_{0,0}^{N_k}(t) \geq \frac{1}{N_k} \sum_{n=\mathcal{N}_\lambda(N_k t)+1}^{\mathcal{N}_\lambda(N_k(t+\epsilon))} \sum_{p=1}^d \mathbb{I}_{(X_{0,0}^{N_k}(t_n^{N_k, \lambda^-}), X_{0,\cdot}^{N_k}(t_n^{N_k, \lambda^-})]}^{(V_n^p)} - 1, \quad (29)$$

which is obvious given (15), and because the RHS of (29), once divided by  $\epsilon$ , converges to  $-\lambda + d\lambda(\bar{x}_{0,\cdot}(t) - \bar{x}_{0,0}(t))$ , by Lemmas 1 and 4, in the limit where  $k \rightarrow \infty$  first and  $\epsilon \downarrow 0$ .

Now, assume by contradiction that  $t$  is a point of differentiability. In this case, the limit in the LHS of (28) exists and must be equal to  $\dot{\bar{x}}_{0,0}(t)$ . Furthermore, if  $\bar{x}_{0,\cdot}(t) > 1/d$ , the RHS of (28) is strictly positive and thus  $\dot{\bar{x}}_{0,0}(t)$  must be strictly positive as well. On the other hand, it is not possible to have  $\dot{\bar{x}}_{0,0}(t) > 0$  and  $\bar{x}_{0,0}(t) = 0$  because the function  $\bar{x}_{0,0}$  is always non-negative. This contradicts that  $t$  is a point of differentiability of  $\bar{x}_{0,0}(\cdot)$ .  $\square$

The next lemma says that the limit trajectory  $\bar{x}_{0,0}$  remains on zero in a right neighborhood of  $t$ , provided that  $\bar{x}_{0,0}(t) = 0$  and  $0 \leq \bar{x}_{0,\cdot}(t) < 1/d$ .

**Lemma 7.** *Fix  $\omega \in \mathcal{C}$ , let (23) hold and assume  $\bar{x}_{0,0}(t) = 0$  and  $d\bar{x}_{0,\cdot}(t) < 1$ . Then,*

$$\exists \delta > 0 : \bar{x}_{0,0}(t') = 0, \quad \forall t' \in [t, t + \delta]. \quad (30)$$

*Proof:* Assume that (30) is false. Then, there exists a sequence  $t_n \downarrow t$  such that  $t_n > t_{n+1} > t$  and

$$\bar{x}_{0,0}(t_n) > 0 \text{ and } \dot{\bar{x}}_{0,0}(t_n) > 0$$

for all  $n$ . By Lemma 5, we have  $\dot{\bar{x}}_{0,0}(t_n) = -\lambda + d\lambda(\bar{x}_{0,\cdot}(t_n) - \bar{x}_{0,0}(t_n))$  and thus  $\bar{x}_{0,\cdot}(t_n) - \bar{x}_{0,0}(t_n) > \frac{1}{d}$ , for all  $n$ , and by continuity

$$\inf_n \bar{x}_{0,\cdot}(t_n) - \bar{x}_{0,0}(t_n) \geq \frac{1}{d}. \quad (31)$$

Thus, we get

$$\bar{x}_{0,\cdot}(t) - \bar{x}_{0,0}(t) = \lim_{n \rightarrow \infty} \bar{x}_{0,\cdot}(t_n) - \bar{x}_{0,0}(t_n) \geq \inf_n \bar{x}_{0,\cdot}(t_n) - \bar{x}_{0,0}(t_n) \geq \frac{1}{d}.$$

This contradicts the hypothesis.  $\square$

Summarizing,

- when  $\bar{x}_{0,0}(t) > 0$ , we have proven that  $\dot{\bar{x}}_{00}(t) = b_{0,0}(\bar{x}(t))$ ;
- when  $\bar{x}_{0,\cdot}(t) < 1/d$  and  $\bar{x}_{0,0}(t) = 0$ , we have proven that  $\bar{x}_{0,0}(t)$  remains zero on a right neighborhood, and thus if  $t$  is a point of differentiability, then  $0 = \dot{\bar{x}}_{00}(t) = b_{0,0}(\bar{x}(t))$ ;
- when  $\bar{x}_{0,\cdot}(t) > 1/d$  and  $\bar{x}_{0,0}(t) = 0$ , we have proven that  $t$  is not a point of differentiability;
- when  $\bar{x}_{0,\cdot}(t) = 1/d$  and  $\bar{x}_{0,0}(t) = 0$ , either  $t$  is not a point of differentiability or it is. In the latter case, we must have  $\dot{\bar{x}}_{00}(t) = 0$  because  $\bar{x}_{0,0}$  is a non-negative function and since also  $b_{0,0}(\bar{x}(t)) = 0$ , we have indeed  $\dot{\bar{x}}_{00}(t) = b_{0,0}(\bar{x}(t))$  as desired.

Thus,  $\dot{\bar{x}}_{0,0}(t) = b_{0,0}(\bar{x}(t))$  almost everywhere.

### 4.3.2 Fluid solution on arbitrary coordinates.

We now prove that  $\dot{\bar{x}}_{i,j}(t) = b_{i,j}(\bar{x}(t))$  almost everywhere with respect to arbitrary coordinates  $(i, j)$ . This requires a more in-depth analysis of the stochastic process  $X^N(t)$ .

Let

$$R_j(t) \stackrel{\text{def}}{=} \lim_{\epsilon \downarrow 0} \lim_{k \rightarrow \infty} \frac{1}{\epsilon N_k} \sum_{n=\mathcal{N}_\lambda(N_k t)+1}^{\mathcal{N}_\lambda(N_k(t+\epsilon))} \mathbf{1}_{\{R_j^{N_k}(t_n^{N_k, \lambda^-})=0\}} \prod_{p=1}^d \mathbb{I}_{(1-Z_{j+1}^{N_k}(t_n^{N_k, \lambda^-}), 1]}^{(V_n^p)}, \quad (32)$$

which is interpreted as the proportion of time where the process  $R_j^{N_k}(t_n^{N_k, \lambda^-})$  remains on zero in the interval  $[t, t + \epsilon]$  while the load balancer keeps sampling only  $(j', \cdot)$ -servers for all  $j' > j$  in the  $\lim_{\epsilon \downarrow 0} \lim_{k \rightarrow \infty}$  limit. In the following, we show that  $R_j(t) = \mathcal{R}_j(\bar{x}(t))$ , where  $\mathcal{R}_j$  is given in Definition 1.

The structure of  $R_0$  is easily obtained as a corollary of the analysis developed in previous section.

**Lemma 8.** *Fix  $\omega \in \mathcal{C}$ , let (23) hold, and assume that  $\bar{x}(t)$  is differentiable. Then,  $R_0(t)$  exists and is given by*

$$R_0(t) = \lambda(1 - d\bar{x}_{0,\cdot}(t)) \mathbf{1}_{\{\bar{x}_{0,0}(t)=0\}} \mathbf{1}_{\{d\bar{x}_{0,\cdot}(t) < 1\}}. \quad (33)$$

*Proof:* First, we notice that if  $\bar{x}_{0,0}(t) > 0$ , then necessarily  $R_0(t) = 0$ . In fact, if for any  $j$ ,  $\sum_{i=0}^j \bar{x}_{\cdot,i}(t) > 0$ , then we can find  $\epsilon > 0$  such that  $\sum_{i=0}^j \bar{x}_{\cdot,i}(t) - 2L(I+1)^2\epsilon > 0$ . Since  $t_n^{N_k, \lambda^-} \in (t, t + \epsilon]$  when  $n \in \{\mathcal{N}_\lambda(N_k t) + 1, \dots, \mathcal{N}_\lambda(N_k(t + \epsilon))\}$ , Lemma 3 implies that for all  $k$  sufficiently large we must have  $R_j^{N_k}(t_n^{N_k, \lambda^-}) > 0$ , and therefore  $R_j(t) = 0$  in this case.

Thus, assume that  $\bar{x}_{0,0}(t) = 0$ . In this case, since  $t$  is a point of differentiability, we necessarily have  $\dot{\bar{x}}_{0,0}(t) = 0$  and, by Lemma 6, necessarily  $d\bar{x}_{0,\cdot}(t) \leq 1$ . This gives the indicator functions in (33). Furthermore, recalling the structure of  $X_{0,0}^N(t)$  given in (15), we have

$$\begin{aligned} 0 &= \dot{\bar{x}}_{0,0}(t) \\ &= \lim_{\epsilon \downarrow 0} \frac{1}{\epsilon} \lim_{k \rightarrow \infty} X_{0,0}^{N_k}(t + \epsilon) - X_{0,0}^{N_k}(t) \\ &= \lim_{\epsilon \downarrow 0} \frac{1}{\epsilon} \lim_{k \rightarrow \infty} \frac{1}{N_k} \sum_{n=1+\mathcal{N}_\lambda(N_k t)}^{\mathcal{N}_\lambda(N_k(t+\epsilon))} \mathbf{1}_{\{X_{0,0}^{N_k}(t_n^{N_k, \lambda^-})=0\}} \prod_{p=1}^d \mathbb{I}_{(X_{0,\cdot}^{N_k}(t_n^{N_k, \lambda^-}), 1]}^{(V_n^p)} - 1 + \sum_{p=1}^d \mathbb{I}_{(X_{0,0}^{N_k}(t_n^{N_k, \lambda^-}), X_{0,\cdot}^{N_k}(t_n^{N_k, \lambda^-})]}^{(V_n^p)} \\ &= \lambda d(\bar{x}_{0,\cdot}(t) - \bar{x}_{0,0}(t)) - \lambda + R_0(t). \end{aligned}$$

In the last equality we have used Lemma 4 and the definition of  $R_0$ . This equation gives (33).  $\square$

The next lemma provides an expression for  $R_j(t)$  for all  $j$  and shows that  $R_j(t) = \mathcal{R}_j(\bar{x}(t))$ . Our proof, given in the appendix, is based on Lemma 8, which allows us to establish the existence and find the structure of  $R_j$  in an iterative manner.

**Lemma 9.** *Fix  $\omega \in \mathcal{C}$ , let (23) hold and assume that  $\bar{x}(t)$  is differentiable. Then, for all  $j$ ,  $R_j(t)$  exists and is given by*

$$R_j(t) = 0 \vee \lambda \left( 1 - d \sum_{i=0}^j (j+1-i) x_{i,\cdot}(t) \right) \mathbf{1}_{\{\sum_{i=0}^j x_{\cdot,i}(t)=0\}}. \quad (34)$$

For any  $i, j \geq i$ , let us define

$$\Gamma_{i,j}^{\epsilon,k}(t) \stackrel{\text{def}}{=} \frac{1}{\epsilon N_k} \sum_{n=\mathcal{N}_\lambda(N_k t)+1}^{\mathcal{N}_\lambda(N_k(t+\epsilon))} \mathbf{1}_{\{R_j^{N_k}(t_n^{N_k, \lambda^-})=0\}} F_{i,j+1,n}^{N_k} \prod_{p=1}^d \mathbb{I}_{(1-Z_{j+1}^{N_k}(t_n^{N_k, \lambda^-}), 1]}^{(V_n^p)} \quad (35)$$

and  $\Gamma_{i,j}(t) \stackrel{\text{def}}{=} \lim_{\epsilon \downarrow 0} \lim_{k \rightarrow \infty} \Gamma_{i,j}^{\epsilon,k}(t)$ , which is interpreted as the proportion of time where the process  $R_j^{N_k}(t_n^{N_k, \lambda^-})$  remains on zero in the interval  $[t, t + \epsilon]$  while the load balancer samples  $(j', \cdot)$ -servers only, for all  $j' > j$ , and assigns jobs to  $(i, j+1)$ -servers only when the proportion of  $(\cdot, j+1)$ -servers vanishes in the  $\lim_{\epsilon \downarrow 0} \lim_{k \rightarrow \infty}$  limit.

The next lemma, proven in the appendix, gives an expression for  $\Gamma_{i,j}(t)$  when  $\bar{x}_{\cdot,j+1}(t) > 0$  and will allow us to identify the limit behavior of terms (19d)-(19f) and (20d)-(20f).

**Lemma 10.** *Fix  $\omega \in \mathcal{C}$  and let (23) hold. Assume that  $\bar{x}(t)$  is differentiable and that  $\bar{x}_{\cdot,j+1}(t) > 0$ . Then,*

$$\Gamma_{i,j}(t) = \frac{\bar{x}_{i,j+1}(t)}{\bar{x}_{\cdot,j+1}(t)} R_j(t). \quad (36)$$

for all  $i, j$  such that  $i \leq j + 1$ .

With the lemmas above, we can identify the asymptotic behavior (in the  $\lim_{\epsilon \downarrow 0} \lim_{k \rightarrow \infty}$  limit) of each summation appearing in the expressions of  $X_{i,j}^{N_k}(t + \epsilon) - X_{i,j}^{N_k}(t)$  and  $X_{i,i}^{N_k}(t + \epsilon) - X_{i,i}^{N_k}(t)$  that are obtained using (19) and (20), respectively.

Let us first treat the case  $i < j$ .

Applying Lemma 10 in (19), when  $\bar{x}(t)$  is differentiable we obtain

$$\begin{aligned}
\dot{\bar{x}}_{i,j}(t) &= \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \lim_{k \rightarrow \infty} X_{i,j}^{N_k}(t + \epsilon) - X_{i,j}^{N_k}(t) \\
&= \bar{x}_{i+1,j}(t) - \mathbf{1}_{\{i>0\}} \bar{x}_{i,j}(t) - \lambda d \bar{x}_{i,j}(t) \\
&\quad - \frac{\bar{x}_{i,j}(t)}{\bar{x}_{\cdot,j}(t)} R_{j-1}(t) \mathbf{1}_{\{\bar{x}_{\cdot,j}(t)>0\}} \\
&\quad + \mathbf{1}_{\{i>0\}} \frac{\bar{x}_{i-1,j-1}(t)}{\bar{x}_{\cdot,j-1}(t)} R_{j-2}(t) \mathbf{1}_{\{\bar{x}_{\cdot,j-1}(t)>0\}} \\
&\quad + \mathbf{1}_{\{j=I,i>0\}} \frac{\bar{x}_{i-1,I}(t)}{\bar{x}_{\cdot,I}(t)} R_{I-1}(t) \mathbf{1}_{\{\bar{x}_{\cdot,I}(t)>0\}} \\
&\quad + \lim_{\epsilon \downarrow 0} \lim_{k \rightarrow \infty} -\Gamma_{i,j-1}^{\epsilon,k}(t) \mathbf{1}_{\{\sum_{j'=0}^j \bar{x}_{\cdot,j'}(t)=0\}} + \Gamma_{i-1,j-2}^{\epsilon,k}(t) \mathbf{1}_{\{i>0, \sum_{j'=0}^{j-1} \bar{x}_{\cdot,j'}(t)=0\}}
\end{aligned} \tag{37a}$$

where the first three terms follow by applying Lemma 4 to terms (19a), (19b) and (19c). Now, assume  $i = 0$  and  $j > 0$ . Then, if  $\sum_{j'=0}^j \bar{x}_{\cdot,j'}(t) = 0$ , then the first six terms of previous equation and the second summation term in (37a) are equal to zero, and if in addition  $t$  is a point of differentiability, then necessarily  $\dot{\bar{x}}_{i,j}(t) = 0$  (because  $\bar{x}_{i,j}(t) = 0$ ), which means that necessarily

$$\lim_{\epsilon \downarrow 0} \lim_{k \rightarrow \infty} \Gamma_{i,j-1}^{\epsilon,k}(t) \mathbf{1}_{\{\sum_{j'=0}^j \bar{x}_{\cdot,j'}(t)=0\}}$$

exists and is equal to zero. Assume  $i > 0$  and  $j > i$ . If  $\sum_{j'=0}^j \bar{x}_{\cdot,j'}(t) = 0$ , then the first six terms of previous equation again coincide with zero and if in addition  $t$  is a point of differentiability, then necessarily

$$\lim_{\epsilon \downarrow 0} \lim_{k \rightarrow \infty} -\Gamma_{i,j-1}^{\epsilon,k}(t) \mathbf{1}_{\{\sum_{j'=0}^j \bar{x}_{\cdot,j'}(t)=0\}} + \Gamma_{i-1,j-2}^{\epsilon,k}(t) \mathbf{1}_{\{i>0, \sum_{j'=0}^{j-1} \bar{x}_{\cdot,j'}(t)=0\}} \tag{38}$$

exists and is equal to zero. Furthermore, if  $\sum_{j'=0}^{j-1} \bar{x}_{\cdot,j'}(t) = 0$  and  $\bar{x}_{\cdot,j}(t) > 0$ , then (38) still exists and is equal to zero as a consequence of the fact that we have inductively shown that  $\lim_{\epsilon \downarrow 0} \lim_{k \rightarrow \infty} \Gamma_{i-1,j-2}^{\epsilon,k}(t) \mathbf{1}_{\{i>0, \bar{x}_{\cdot,j-1}(t)=0\}} = 0$ . Therefore, the limit in (37a) is always equal to zero. We have thus shown that  $\dot{\bar{x}}_{i,j} = b_{i,j}(\bar{x}_{i,j})$  when  $i < j$ .

The case  $i = j > 0$  is treated in a similar manner. Let  $G_i^{\epsilon,k}(t) \stackrel{\text{def}}{=} \mathbf{1}_{\{\sum_{j'=0}^i \bar{x}_{\cdot,j'}(t)=0\}} \Gamma_{i,i-1}^{\epsilon,k}(t)$  and  $G_i \stackrel{\text{def}}{=} G_i(t) \stackrel{\text{def}}{=} \lim_{\epsilon \downarrow 0} \lim_{k \rightarrow \infty} G_i^{\epsilon,k}(t)$ . In the following, we show that  $G_i(t) = \mathcal{G}_i(\bar{x}(t))$ , where  $\mathcal{G}_i$  is given in Definition 1.

Applying Lemma 4 to handle terms (20a) and (20b), rewriting term (20c) as

$$\begin{aligned}
\mathbf{1}_{\{X_{0,0}^N(t_n^{N,\lambda^-}) + \sum_{j=0}^I M_{0,j,n}^N > 0\}} &= 1 - \mathbf{1}_{\{X_{0,0}^N(t_n^{N,\lambda^-})=0\}} \mathbf{1}_{\{\sum_{j=0}^I M_{0,j,n}^N = 0\}} \\
&= 1 - \mathbf{1}_{\{X_{0,0}^{N_k}(t_n^{N_k,\lambda^-})=0\}} \prod_{p=1}^d \mathbb{I}_{(1-Z_1^{N_k}(t_n^{N_k,\lambda^-}), 1]}^{(V_n^p)},
\end{aligned}$$

and applying Lemma 10 to handle terms (20d), (20e) and (20f), when  $\bar{x}(t)$  is differentiable we obtain

$$\dot{\bar{x}}_{i,i}(t) = \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \lim_{k \rightarrow \infty} X_{i,i}^{N_k}(t + \epsilon) - X_{i,i}^{N_k}(t) \tag{39a}$$

$$= -\bar{x}_{i,i}(t) + \lambda d(\bar{x}_{i,\cdot}(t) - \bar{x}_{i,i}(t)) \tag{39b}$$

$$+ \mathbf{1}_{\{i=1\}} (\lambda - R_0(t) z_1^d) \tag{39c}$$

$$+ \mathbf{1}_{\{i>1\}} R_{i-2}(t) \frac{\bar{x}_{i-1,i-1}(t)}{\bar{x}_{\cdot,i-1}(t)} \mathbf{1}_{\{\bar{x}_{\cdot,i-1}(t)>0\}} \tag{39d}$$



$$- \mathbf{1}_{\{i < I\}} R_{i-1}(t) \frac{\bar{x}_{i,i}(t)}{\bar{x}_{\cdot,i}(t)} \mathbf{1}_{\{\bar{x}_{\cdot,i}(t) > 0\}} \quad (39e)$$

$$+ \mathbf{1}_{\{i=I\}} R_{I-1}(t) \frac{\bar{x}_{I-1,I}(t)}{\bar{x}_{\cdot,I}(t)} \mathbf{1}_{\{\bar{x}_{\cdot,I}(t) > 0\}} \quad (39f)$$

$$+ \lim_{\epsilon \downarrow 0} \lim_{k \rightarrow \infty} -G_i^{\epsilon,k}(t) \mathbf{1}_{\{i < I\}} + G_{i-1}^{\epsilon,k}(t) \mathbf{1}_{\{i > 1\}}. \quad (39g)$$

Now, assume that  $i = 1$ . If  $t$  is a point of differentiability and  $\bar{x}_{0,0}(t) + \bar{x}_{\cdot,1}(t) = 0$ , then we must have  $\dot{\bar{x}}_{1,1}(t) = 0$ , and thus necessarily

$$G_1(t) = (\lambda d \bar{x}_{1,\cdot}(t) + \lambda - R_0(t)) \mathbf{1}_{\{\bar{x}_{0,0}(t) + \bar{x}_{\cdot,1}(t) = 0\}} = \lambda d (\bar{x}_{1,\cdot}(t) + \bar{x}_{0,\cdot}(t)) \mathbf{1}_{\{\bar{x}_{0,0}(t) + \bar{x}_{\cdot,1}(t) = 0\}}.$$

In the last equality we have used Lemma 8 and that  $d \bar{x}_{0,\cdot}(t) \leq 1$ , which holds true because  $t$  is a point of differentiability (Lemmas 6 and 7). When  $i = 2, \dots, I-1$ , if  $t$  is a point of differentiability and  $\sum_{j'=0}^i \bar{x}_{\cdot,j'}(t) = 0$ , then necessarily  $\dot{\bar{x}}_{i,i}(t) = 0$  and proceeding in an iterative manner, we obtain

$$G_i(t) = (G_{i-1}(t) + \lambda d \bar{x}_{i,\cdot}(t)) \mathbf{1}_{\{\sum_{j'=0}^i \bar{x}_{\cdot,j'}(t) = 0\}} = \lambda d \mathbf{1}_{\{\sum_{j'=0}^i \bar{x}_{\cdot,j'}(t) = 0\}} \sum_{j'=0}^i \bar{x}_{i',\cdot}(t).$$

The following two lemmas are a generalization of Lemmas 6 and 7 and show under which conditions  $\bar{x}(t)$  is differentiable. The proofs use the same arguments in those lemmas and therefore they are omitted.

**Lemma 11.** Fix  $\omega \in \mathcal{C}$ , let (23) hold and assume  $\sum_{j'=0}^j \bar{x}_{\cdot,j'}(t) = 0$  and  $d \sum_{i=0}^j (j+1-i) \bar{x}_{i,\cdot}(t) > 1$ . Then,  $t$  is not a point of differentiability.

**Lemma 12.** Fix  $\omega \in \mathcal{C}$ , let (23) hold and assume  $\sum_{j'=0}^j \bar{x}_{\cdot,j'}(t) = 0$  and  $d \sum_{i=0}^j (j+1-i) \bar{x}_{i,\cdot}(t) < 1$ . Then,

$$\exists \delta > 0 : \sum_{j'=0}^j \bar{x}_{\cdot,j'}(t') = 0, \quad \forall t' \in [t, t + \delta]. \quad (40)$$

Now, we notice that the expressions of  $R_j(t)$  and  $G_i(t)$  obtained so far assumed that  $t$  was a point of differentiability. However, if  $\sum_{j'=0}^j \bar{x}_{\cdot,j'}(t) = 0$ , previous lemmas say that this can only be true if  $d \sum_{i=0}^j (j+1-i) \bar{x}_{i,\cdot}(t) \leq 1$ . Thus, those expressions make sense only in that case. This does not change the structure of  $R_j(t)$  obtained in (34) because

$$R_j(t) \times \mathbf{1}_{\{d \sum_{i=0}^j (j+1-i) \bar{x}_{i,\cdot}(t) \leq 1\}} = R_j(t) = \mathcal{R}_j(\bar{x}(t)).$$

but on the other hand we must have  $G_i(t) = \mathcal{G}_i(\bar{x}(t))$ , where  $\mathcal{G}_i$  is defined in (10). We have thus shown that  $\dot{\bar{x}}_{i,i}(t) = b_{i,i}(x)$ .

## 5 Proofs of Theorems 2 and 3

Let us introduce the intervals

$$\mathcal{I}_n \stackrel{\text{def}}{=} [\lambda_n^*, \lambda_{n+1}^*), \quad \forall n \geq 0$$

where  $\lambda_0^* = 0$  and  $\lambda_n^*$ , for  $n \geq 1$ , is the unique root in  $(0, 1]$  of the polynomial equation

$$(1-z)(zd+1)^n = 1.$$

The first values of  $\lambda_n^*$  are  $\lambda_1^* = 1 - \frac{1}{d}$  and  $\lambda_2^* = \frac{1}{2} - \frac{1}{d} + \sqrt{\frac{1}{4} + \frac{1}{d}}$ . We notice that  $j^*$ , defined in (1), is the unique integer such that  $\lambda \in \mathcal{I}_{j^*}$ . In fact,  $\lambda \in \mathcal{I}_n$  if and only if

$$n = -\frac{\log(1 - \lambda_n^*)}{\log(\lambda_n^* d + 1)} \leq -\frac{\log(1 - \lambda)}{\log(\lambda d + 1)}$$

and

$$n + 1 = -\frac{\log(1 - \lambda_{n+1}^*)}{\log(\lambda_{n+1}^* d + 1)} > -\frac{\log(1 - \lambda)}{\log(\lambda d + 1)},$$

which thus implies  $n = \left\lfloor -\frac{\log(1 - \lambda)}{\log(\lambda d + 1)} \right\rfloor = j^*$ .

Let  $x(t)$  be a fluid solution. Since the  $x_{i,j}(t)$ 's are absolutely continuous, both  $\mathcal{L}_S(x(t))$  and  $\mathcal{L}_M(x(t))$  are absolutely continuous as well, and thus almost everywhere differentiable. When  $t$  is a point of differentiability, it is clear that

$$\dot{\mathcal{L}}_S(x(t)) = \sum_{i=1}^I i b_{i,\cdot}(x(t)), \quad \dot{\mathcal{L}}_M(x(t)) = \sum_{j=1}^I j b_{\cdot,j}(x(t)). \quad (41)$$

The following trivial lemma gives a differentiability property of fluid solutions.

**Lemma 13.** *Let  $x(t)$  be a fluid solution. If  $\sum_{i=0}^j x_{\cdot,i}(t) = 0$ , then  $t$  is a point of differentiability if and only if*

$$d \sum_{i=0}^j (j + 1 - i) x_{i,\cdot}(t) \leq 1. \quad (42)$$

The following proposition will be crucial to prove both Theorems 2 and 3.

**Proposition 4.** *Let  $x(t)$  be a fluid solution. If  $t$  is a point of differentiability, then*

$$\dot{\mathcal{L}}_S(x(t)) = x_{0,\cdot}(t) - 1 + \lambda. \quad (43)$$

*Proof:* Let  $j^*(t) \stackrel{\text{def}}{=} \min\{j \geq 0 : x_{\cdot,j}(t) > 0\}$ . To prove (43), we consider the cases  $j^*(t) \geq 1$  and  $j^*(t) = 0$  separately. Let us drop the dependency on  $t$  for notational simplicity.

First, assume that  $j^* \geq 1$ . Using Definition 1, we obtain

$$b_{1,\cdot}(x) = x_{2,\cdot} - x_{1,\cdot} + \lambda dx_{0,\cdot} - \mathcal{R}_0(x) \frac{x_{1,1}}{x_{\cdot,1}} \mathbf{1}_{\{x_{\cdot,1} > 0\}} - \mathcal{G}_1(x) \quad (44a)$$

$$+ \sum_{j \geq 2} \mathcal{R}_{j-2}(x) \frac{x_{0,j-1}}{x_{\cdot,j-1}} \mathbf{1}_{\{x_{\cdot,j-1} > 0\}} - \mathcal{R}_{j-1}(x) \frac{x_{1,j}}{x_{\cdot,j}} \mathbf{1}_{\{x_{\cdot,j} > 0\}} + \mathbf{1}_{\{j=I\}} \mathcal{R}_{I-1}(x) \frac{x_{0,I}}{x_{\cdot,I}} \mathbf{1}_{\{x_{\cdot,I} > 0\}} \quad (44b)$$

$$b_{i,\cdot}(x) = x_{i+1,\cdot} - x_{i,\cdot} - \mathcal{R}_{i-1}(x) \frac{x_{i,i}}{x_{\cdot,i}} \mathbf{1}_{\{x_{\cdot,i} > 0\}} + \mathcal{R}_{i-2}(x) \frac{x_{i-1,i-1}}{x_{\cdot,i-1}} \mathbf{1}_{\{x_{\cdot,i-1} > 0\}} + \mathcal{G}_{i-1}(x) - \mathcal{G}_i(x) \mathbf{1}_{\{i < I\}} \quad (44c)$$

$$+ \sum_{j \geq i+1} \mathcal{R}_{j-2}(x) \frac{x_{i-1,j-1}}{x_{\cdot,j-1}} \mathbf{1}_{\{x_{\cdot,j-1} > 0\}} - \mathcal{R}_{j-1}(x) \frac{x_{i,j}}{x_{\cdot,j}} \mathbf{1}_{\{x_{\cdot,j} > 0\}} + \mathbf{1}_{\{j=I\}} \mathcal{R}_{I-1}(x) \frac{x_{i-1,I}}{x_{\cdot,I}} \mathbf{1}_{\{x_{\cdot,I} > 0\}}, \quad (44d)$$

that is,

$$b_{i,\cdot}(x) = x_{i+1,\cdot} - x_{i,\cdot} - \lambda dx_{i,\cdot} + \mathcal{R}_{j^*-1}(x) \frac{x_{i-1,j^*}}{x_{\cdot,j^*}} - \mathcal{R}_{j^*-1}(x) \frac{x_{i,j^*}}{x_{\cdot,j^*}}, \quad \forall i = 1, \dots, j^* - 1 \quad (45a)$$

$$b_{j^*,\cdot}(x) = x_{j^*+1,\cdot} - x_{j^*,\cdot} + \mathcal{G}_{j^*-1}(x) - \mathcal{R}_{j^*-1}(x) \frac{x_{j^*,j^*}}{x_{\cdot,j^*}} + \mathcal{R}_{j^*-1}(x) \frac{x_{j^*-1,j^*}}{x_{\cdot,j^*}} \quad (45b)$$

$$b_{j^*+1,\cdot}(x) = x_{j^*+2,\cdot} - x_{j^*+1,\cdot} + \mathcal{R}_{j^*-1}(x) \frac{x_{j^*,j^*}}{x_{\cdot,j^*}} \quad (45c)$$

$$b_{i,\cdot}(x) = x_{i+1,\cdot} - x_{i,\cdot}, \quad \forall i \geq j^* + 2. \quad (45d)$$

Since  $t$  is a point of differentiability, then  $d \sum_{i=0}^{j^*} (j^* + 1 - i) x_{i,\cdot} \leq 1$  (by Lemma 13) and thus

$$\mathcal{G}_{j^*-1}(x) = \lambda d \sum_{i=0}^{j^*-1} x_{i,\cdot}, \quad \mathcal{R}_{j^*-1}(x) = \lambda - \lambda d \sum_{i=0}^{j^*-1} (j^* - i) x_{i,\cdot}. \quad (46)$$

Substituting these expressions in (41), we get

$$\dot{\mathcal{L}}_S(x) = \sum_{i \geq 1} i (x_{i+1,\cdot} - x_{i,\cdot}) + \sum_{i=1}^{j^*-1} i \left( -\lambda dx_{i,\cdot} + \mathcal{R}_{j^*-1}(x) \frac{x_{i-1,j^*}}{x_{\cdot,j^*}} - \mathcal{R}_{j^*-1}(x) \frac{x_{i,j^*}}{x_{\cdot,j^*}} \right) \quad (47a)$$

$$+ j^* \left( \lambda d \sum_{i'=0}^{j^*-1} x_{i',\cdot}(t) - \mathcal{R}_{j^*-1}(x) \frac{x_{j^*,j^*}}{x_{\cdot,j^*}} + \mathcal{R}_{j^*-1}(x) \frac{x_{j^*-1,j^*}}{x_{\cdot,j^*}} \right) \quad (47b)$$

$$+ (j^* + 1) \mathcal{R}_{j^*-1}(x) \frac{x_{j^*,j^*}}{x_{\cdot,j^*}} \quad (47c)$$

$$= - \sum_{i \geq 1} x_{i,\cdot} + \lambda d \sum_{i=0}^{j^*-1} (j^* - i) x_{i,\cdot} + \mathcal{R}_{j^*-1}(x) \sum_{i=1}^{j^*-1} i \left( \frac{x_{i-1,j^*}}{x_{\cdot,j^*}} - \frac{x_{i,j^*}}{x_{\cdot,j^*}} \right) \quad (47d)$$

$$+ j^* \mathcal{R}_{j^*-1}(x) \frac{x_{j^*-1,j^*}}{x_{\cdot,j^*}} + \mathcal{R}_{j^*-1}(x) \frac{x_{j^*,j^*}}{x_{\cdot,j^*}} \quad (47e)$$

$$= x_{0,\cdot} - 1 + \lambda d \sum_{i=0}^{j^*-1} (j^* - i) x_{i,\cdot} + \mathcal{R}_{j^*-1}(x) \quad (47f)$$

$$= x_{0,\cdot} - 1 + \lambda \quad (47g)$$

as desired.

Now, assume that  $j^* = 0$ . In this case,  $x_{0,0} > 0$  and the drift  $b(x)$  given in Definition 1 takes the linear form

$$b_{0,0}(x) = -\lambda + \lambda d(x_{0,\cdot} - x_{0,0}) \quad (48a)$$

$$b_{0,1}(x) = x_{11} - \lambda d x_{0,1} \quad (48b)$$

$$b_{1,1}(x) = -x_{11} + \lambda + \lambda d(x_{1,\cdot} - x_{1,1}) \quad (48c)$$

$$b_{0,j}(x) = x_{1,j} - \lambda d x_{0,j}, \quad j > 1 \quad (48d)$$

$$b_{i,j}(x) = x_{i+1,j} - x_{i,j} - \lambda d x_{i,j}, \quad j > i, i \geq 1 \quad (48e)$$

$$b_{i,i}(x) = -x_{i,i} + \lambda d(x_{i,\cdot} - x_{i,i}), \quad i > 1. \quad (48f)$$

By taking summations in (48), we obtain

$$b_{1,\cdot}(x) = \lambda + x_{2,\cdot} - x_{1,\cdot} \quad (49a)$$

$$b_{i,\cdot}(x) = x_{i+1,\cdot} \mathbf{1}_{\{i+1 \leq I\}} - x_{i,\cdot}, \quad \forall i \geq 2. \quad (49b)$$

Thus,

$$\dot{\mathcal{L}}_S(x) = \sum_{i \geq 1} i b_{i,\cdot}(x) = \lambda + \sum_{i \geq 1} i(x_{i+1,\cdot} \mathbf{1}_{\{i+1 \leq I\}} - x_{i,\cdot}) = x_{0,\cdot} - 1 + \lambda \quad (50)$$

where in the last equality we have used the normalizing condition.  $\square$

## 5.1 Existence and uniqueness of fixed points.

Assume that  $x$  is a fixed point. Then the system of equations

$$b_{i,j}(x) = 0, \quad \forall i, j \quad (51)$$

must be satisfied.

Let  $j^* \stackrel{\text{def}}{=} j^*(x) \stackrel{\text{def}}{=} \min\{j \geq 0 : x_{\cdot,j} > 0\}$ . By Definition 1, this means that  $\mathcal{R}_j(x) = 0$  for all  $j \geq j^*$ .

If  $j^* = 0$ , then  $x_{0,0} > 0$  and the drift  $b(x)$  given in Definition 1 takes the linear form given in (48). Removing one equation from (51) and adding the normalizing condition  $\sum_{i,j} x_{i,j} = 1$ , we obtain a linear system composed of  $(I+1)(I+2)/2$  independent equations and  $(I+1)(I+2)/2$  unknowns. It is easy to see that  $x^*$ , where  $x_{0,0}^* = 1 - \lambda - 1/d$ ,  $x_{0,1}^* = 1/d$ ,  $x_{1,1}^* = \lambda$  and  $x_{i,j}^* = 0$  on the remaining coordinates, is the unique solution of such system. It is also clear that  $x^* \in \mathcal{S}$  and  $x_{0,0}^* > 0$  if and only if  $\lambda < 1 - 1/d = \lambda_1^*$ . Therefore, in the following we assume that  $j^* \geq 1$ .

If (51) holds true, then we must also have

$$b_{\cdot,j}(x) = 0, \quad \forall j. \quad (52)$$

Using Definition 1, we obtain

$$b_{\cdot,1}(x) = -\lambda dx_{0,1} - \mathcal{R}_0(x) \mathbf{1}_{\{x_{\cdot,1} > 0\}} + \lambda d(x_{1,\cdot} - x_{1,1}) + \lambda - \mathcal{R}_0(x) - \mathcal{G}_1(x) \quad (53a)$$

$$b_{\cdot,j}(x) = \lambda dx_{j,\cdot} - \lambda dx_{\cdot,j} - \mathcal{R}_{j-1}(x) \mathbf{1}_{\{x_{\cdot,j} > 0\}} + \mathcal{R}_{j-2}(x) \mathbf{1}_{\{x_{\cdot,j-1} > 0\}} + \mathcal{G}_{j-1}(x) - \mathcal{G}_j(x)$$

$$b_{\cdot,I}(x) = \mathcal{R}_{I-2}(x) \mathbf{1}_{\{x_{\cdot,I-1} > 0\}} + \mathcal{G}_{I-1}(x) - \sum_{i=0}^{I-1} \lambda dx_{i,I}. \quad (53b)$$

which can be rewritten as

$$b_{\cdot,j}(x) = 0, \quad \forall j = 0, \dots, j^* - 1 \quad (54a)$$

$$b_{\cdot,j^*}(x) = \lambda dx_{j^*,\cdot} - \lambda dx_{\cdot,j^*} - \mathcal{R}_{j^*-1}(x) + \mathcal{G}_{j^*-1}(x) \quad (54b)$$

$$b_{\cdot,j^*+1}(x) = \lambda dx_{j^*+1,\cdot} - \lambda dx_{\cdot,j^*+1} + \mathcal{R}_{j^*-1}(x) \quad (54c)$$

$$b_{\cdot,j}(x) = \lambda dx_{j,\cdot} - \lambda dx_{\cdot,j}, \quad \forall j \geq j^* + 2. \quad (54d)$$

Summing (54) for all  $j \geq j^* + 2$ , we obtain

$$0 = \sum_{j=j^*+2}^I x_{j,\cdot} - x_{\cdot,j} = - \sum_{i=0}^{j^*+1} \sum_{j \geq j^*+2} x_{i,j}$$

and since  $x$  is composed of non-negative components only, we must have  $x_{i,j} = 0$  when  $i \in \{0, 1, \dots, j^* + 1\}$  and  $j \in \{j^* + 2, \dots, I\}$ . Using this in (54) when  $j = j^* + 2$ , we obtain that necessarily  $x_{j,j'} = 0$  for all  $j' > j$ . Using again this argument when  $j = j^* + 3$ , we obtain  $x_{j,j'} = 0$  for all  $j' > j^* + 3$ , and so forth for all  $j$ . We have thus shown that the only non-zero elements of  $x$  can be on coordinates  $(i, j^*)$  and  $(i, j^* + 1)$ , for all  $i < j^* + 1$ , and  $(i, i)$  for all  $i \geq j^*$ . Substituting these properties in (5), we also obtain  $b_{j-1,j}(x) = x_{j,j}$ , provided that  $j \geq j^* + 2$ . On the other hand, (51) must hold true and therefore  $x_{j,j} = 0$  for all  $j \geq j^*$ . Thus, we have shown that the only non-zero elements of  $x$  can be on coordinates  $(i, j^*)$  and  $(i, j^* + 1)$ , for all  $i \leq j^* + 1$ . In this case, (51) simplifies to

$$0 = b_{i,j^*}(x) = x_{i+1,j^*} - \mathbf{1}_{\{i > 0\}} x_{i,j^*} - \lambda dx_{i,j^*} - \mathcal{R}_{j^*-1}(x) \frac{x_{i,j^*}}{x_{\cdot,j^*}}, \quad \forall i < j^* \quad (55a)$$

$$0 = b_{i,j^*+1}(x) = x_{i+1,j^*+1} - \mathbf{1}_{\{i > 0\}} x_{i,j^*+1} - \lambda dx_{i,j^*+1} + \mathbf{1}_{\{i > 0\}} \mathcal{R}_{j^*-1}(x) \frac{x_{i-1,j^*}}{x_{\cdot,j^*}}, \quad \forall i < j^* + 1 \quad (55b)$$

$$0 = b_{j^*,j^*}(x) = -x_{j^*,j^*} - \mathcal{R}_{j^*-1}(x) \frac{x_{j^*,j^*}}{x_{\cdot,j^*}} + \lambda d(1 - x_{j^*,j^*} - x_{j^*+1,j^*+1}) \quad (55c)$$

$$0 = b_{j^*+1,j^*+1}(x) = -x_{j^*+1,j^*+1} + \mathcal{R}_{j^*-1}(x) \frac{x_{j^*,j^*}}{x_{\cdot,j^*}} \quad (55d)$$

where  $\mathcal{R}_{j^*-1}(x) = \lambda - \lambda d \sum_{i=0}^{j^*-1} (j^* - i)(x_{i,j^*} + x_{i,j^*+1}) \geq 0$  by Lemma 13 because a fixed point is a fluid solution that is everywhere differentiable. In the remainder of the proof, we show that the system in (55) admits a unique solution that satisfies the normalizing condition  $x_{\cdot,j^*} + x_{\cdot,j^*+1} = 1$  and  $x_{\cdot,j^*} > 0$  if and only if  $\lambda \in \mathcal{I}_{j^*}$ , and that such solution satisfies as well the properties given in the statement of Theorem 2. This will conclude our proof.

Using that  $x_{0,j^*} + x_{0,j^*+1} = 1 - \lambda$  (by Proposition 4), equations (55) imply that

$$0 = b_{0,j^*}(x) + b_{1,j^*+1}(x) = y_1 - \lambda d(1 - \lambda) - (\lambda d)^2 x_{0,j^*+1}$$

$$0 = b_{i,j^*}(x) + b_{i+1,j^*+1}(x) = x_{i+1,j^*} + x_{i+2,j^*+1} - (1 + \lambda d)(x_{i,j^*} + x_{i+1,j^*+1}), \quad \forall i = 1, \dots, j^* - 1$$

$$0 = b_{j^*,j^*}(x) + b_{j^*+1,j^*+1}(x) = \lambda d(1 - x_{j^*,j^*} - x_{j^*+1,j^*+1}) - (x_{j^*,j^*} + x_{j^*+1,j^*+1}).$$

Letting  $y_i \stackrel{\text{def}}{=} x_{i,j^*} + x_{i+1,j^*+1}$ , the key observation is that

$$y_1 = \lambda d(1 - \lambda) + (\lambda d)^2 x_{0,j^*+1} \quad (57a)$$

$$y_{i+1} = (1 + \lambda d)y_i, \quad \forall i = 1, \dots, j^* - 1 \quad (57b)$$

$$y_{j^*} = \lambda d(1 - y_{j^*}), \quad (57c)$$

and we notice that the last equation is autonomous. Previous equations imply that

$$y_{j^*-i} = \frac{y_{j^*}}{(1 + \lambda d)^i}, \quad \forall i = 1, \dots, j^* - 1 \quad (58a)$$

$$y_{j^*} = \frac{\lambda d}{1 + \lambda d} \quad (58b)$$

and using (57a) we obtain the equation

$$\lambda d(1 - \lambda) + (\lambda d)^2 x_{0,j^*+1} = \frac{\lambda d}{(1 + \lambda d)^{j^*}}$$

which gives

$$\lambda d x_{0,j^*+1} = \lambda - 1 + \frac{1}{(1 + \lambda d)^{j^*}}.$$

Since it must hold true that  $x_{0,j^*} + x_{0,j^*+1} = 1 - \lambda$  (by Proposition 4), we also obtain

$$\lambda d x_{0,j^*} = (1 + \lambda d)(1 - \lambda) - \frac{1}{(1 + \lambda d)^{j^*}}.$$

Now, in order for such  $x$  to be feasible, we need that both  $x_{0,j^*}$  and  $x_{0,j^*+1}$  are non-negative. Using previous expressions, it is not difficult to see that  $x_{0,j^*} \geq 0$  and  $x_{0,j^*+1} \geq 0$  if and only if  $\lambda \in \text{cl}(\mathcal{I}_{j^*})$ , where  $\text{cl}(A)$  denotes the closure of set  $A$ .

We notice that the normalizing condition can be written as  $1 = x_{0,j^*} + x_{0,j^*+1} + x_{1,j^*+1} + \sum_{i=1}^{j^*} y_i$  and using the expressions above it is not difficult to check that it is indeed satisfied.

We now proceed with the construction of the fixed point. Substituting the properties obtained so far, we notice that

$$0 = b_{j^*+1,j^*+1}(x) - b_{j^*,j^*}(x) = 2x_{j^*,j^*} + 2\mathcal{R}_{j^*-1}(x) \frac{x_{j^*,j^*}}{x_{\cdot,j^*}} - \frac{2\lambda d}{1 + \lambda d}$$

and since necessarily  $x_{j^*,j^*} > 0$ , which is immediately implied by (55c), we obtain

$$\frac{\mathcal{R}_{j^*-1}(x)}{x_{\cdot,j^*}} = \frac{1}{x_{j^*,j^*}} \frac{\lambda d}{1 + \lambda d} - 1.$$

Using this equation in (55) and recalling that  $x_{0,j^*}$  has been already explicited, we get

$$\begin{aligned} 0 = b_{0,j^*}(x) &= x_{1,j^*} - \lambda d x_{0,j^*} - x_{0,j^*} \left( \frac{1}{x_{j^*,j^*}} \frac{\lambda d}{1 + \lambda d} - 1 \right) \\ 0 = b_{i,j^*}(x) &= x_{i+1,j^*} - x_{i,j^*} - \lambda d x_{i,j^*} - x_{i,j^*} \left( \frac{1}{x_{j^*,j^*}} \frac{\lambda d}{1 + \lambda d} - 1 \right), \quad \forall i = 1 \dots, j^* - 1 \end{aligned}$$

and thus,

$$\begin{aligned} x_{1,j^*} &= \left( \lambda d - 1 + \frac{1}{x_{j^*,j^*}} \frac{\lambda d}{1 + \lambda d} \right) x_{0,j^*} \\ x_{i+1,j^*} &= \lambda d \left( 1 + \frac{1}{x_{j^*,j^*}} \frac{1}{1 + \lambda d} \right) x_{i,j^*} = (\lambda d)^i \left( 1 + \frac{1}{x_{j^*,j^*}} \frac{1}{1 + \lambda d} \right)^i x_{1,j^*} \end{aligned}$$

for all  $i = 1, \dots, j^* - 1$ . In particular, when  $i = j^* - 1$ , the last equation allows us to identify  $x_{j^*,j^*}$  by means of the following polynomial equation

$$\begin{aligned} F(x_{j^*,j^*}) &\stackrel{\text{def}}{=} (\lambda d)^{j^*-2} \left( 1 + \frac{1}{x_{j^*,j^*}} \frac{1}{1 + \lambda d} \right)^{j^*-1} \left( \lambda d - 1 + \frac{1}{x_{j^*,j^*}} \frac{\lambda d}{1 + \lambda d} \right) \left( (1 + \lambda d)(1 - \lambda) - \frac{1}{(1 + \lambda d)^{j^*}} \right) \\ &- x_{j^*,j^*} = 0. \end{aligned} \quad (61)$$

Since  $x_{i+1,j^*+1} = y_i - x_{i,j^*}$  and the value of  $y_i$  has been already explicited for each  $i$ , to conclude the proof of existence and uniqueness of a solution of (55), it remains to show that previous equation admits a unique root in  $(0,1]$  when  $\lambda \in \mathcal{I}_{j^*}$ . This property follows easily once noted that  $\lim_{x \downarrow 0} F(x) = +\infty$ ,  $F(1) < 0$  if  $\lambda \in \mathcal{I}_{j^*}$ , and that  $F(x)$  is strictly decreasing if  $\lambda \in \mathcal{I}_{j^*}$ .

## 5.2 Bounds on fluid mass.

Whenever  $x(t)$  is differentiable, (46) and (54) imply

$$\dot{\mathcal{L}}_M = \sum_{j=1}^I j b_{\cdot,j}(x) \quad (62a)$$

$$= j^* (\lambda dx_{j^*,\cdot} - \lambda dx_{\cdot,j^*} - \mathcal{R}_{j^*-1}(x) + \mathcal{G}_{j^*-1}(x)) \quad (62b)$$

$$+ (j^* + 1) (\lambda dx_{j^*+1,\cdot} - \lambda dx_{\cdot,j^*+1} + \mathcal{R}_{j^*-1}(x)) \quad (62c)$$

$$+ \lambda d \sum_{j \geq j^*+2} j (x_{j,\cdot} - x_{\cdot,j}) \quad (62d)$$

$$= \lambda + \lambda d \sum_{i=0}^{j^*-1} i x_{i,\cdot} + \lambda d \sum_{j \geq j^*} j (x_{j,\cdot} - x_{\cdot,j}) \quad (62e)$$

$$= \lambda + \lambda d \mathcal{L}_S - \lambda d \mathcal{L}_M. \quad (62f)$$

In a fixed point  $x$ , we must have  $\dot{\mathcal{L}}_M = 0$ . This condition gives (12) and since  $x_{\cdot,j^*} + x_{\cdot,j^*+1} = 1$  (by Theorem 2) we obtain

$$\mathcal{L}_S(x) = \mathcal{L}_M(x) - \frac{1}{d} = j^* x_{\cdot,j^*} + (j^* + 1) x_{\cdot,j^*+1} - \frac{1}{d} = j^* + x_{\cdot,j^*+1} - \frac{1}{d}$$

and (13) holds true because  $x_{\cdot,j^*+1} \in [0, 1]$ .

## 5.3 Global stability.

To prove Theorem 3, we first introduce the following lemma.

**Lemma 14.** *Assume  $\lambda < 1 - \frac{1}{d}$ . Let  $x(t)$  be a fluid solution such that  $x_{0,0}(t) > 0$  for all  $t \geq 0$ . Then, (14) holds true.*

*Proof:* If  $x_{0,0} > 0$ , then the drift  $b(x)$  has the linear form given in (48). Thus, if  $x_{0,0}(t) > 0$  for all  $t \geq 0$ , then the fluid solution  $x(t)$  is uniquely determined by the ODE system

$$\dot{x}_{0,0} = -\lambda + \lambda d(x_{0,\cdot} - x_{0,0}) \quad (63a)$$

$$\dot{x}_{0,1} = x_{11} - \lambda dx_{0,1} \quad (63b)$$

$$\dot{x}_{1,1} = -x_{11} + \lambda + \lambda d(x_{1,\cdot} - x_{1,1}) \quad (63c)$$

$$\dot{x}_{0,j} = x_{1,j} - \lambda dx_{0,j}, \quad j > 1 \quad (63d)$$

$$\dot{x}_{i,j} = x_{i+1,j} - x_{i,j} - \lambda dx_{i,j}, \quad j > i, i \geq 1 \quad (63e)$$

$$\dot{x}_{i,i} = -x_{i,i} + \lambda d(x_{i,\cdot} - x_{i,i}), \quad i > 1. \quad (63f)$$

We have already shown in Section 5.1 that  $x^*$  is the unique fixed point of such linear system. The equations (63b)-(63f) do not depend on  $x_{0,0}$  and thus form an autonomous ODE system. This means that we can safely remove the equation (63a) and recall that  $x_{0,0}(t)$  can be uniquely obtained by using the normalizing condition, i.e.,  $x_{0,0}(t) = 1 - \sum_{(i,j) \neq (0,0)} x_{i,j}(t)$ . The ODE system (63b)-(63f) has the linear form  $\dot{x} = Ax + p$  where  $A$  is a triangular matrix and  $p$  is a column vector, and it is clear that the eigenvalues of  $A$  are  $-1$ ,  $-\lambda d$  and  $-(1 + \lambda d)$ . Since the eigenvalues of  $A$  are strictly negative, it follows from standard results in ODE theory that  $x(t) = x^* + e^{At}(x(0) - x^*)$ . Thus, (14) follows by the norm bound on the exponential matrix.  $\square$

For the fluid solution  $x(t)$ , either  $x_{0,0}(t) > 0$  for all  $t \geq 0$ , in which case Theorem 3 follows directly by previous lemma, or  $x_{0,0}(t_0) = 0$  for some  $t_0$ , that is the case we study in the following. Without loss of generality, let us assume  $t_0 > 0$ .

If  $t_0$  is a point of differentiability of  $x_{0,0}(\cdot)$ , then Lemma 13 implies that  $dx_{0,\cdot}(t_0) \leq 1$ . If  $t_0$  is not a point of differentiability of  $x_{0,0}(\cdot)$ , then we still have  $dx_{0,\cdot}(t_0) \leq 1$  because either there exists  $\delta$  such that

$x_{0,0}(t) = 0$  for all  $t \in [t_0 - \delta, t_0]$  and the inequality holds true again by Lemma 13 or there exists a sequence  $t_n \uparrow t_0$ ,  $n \geq 1$ , such that  $t_n < t_{n+1} < t_0$  where  $x_{0,0}(t_n) > 0$  and  $0 > \dot{x}_{0,0}(t_n) = d\lambda(x_{0,\cdot}(t_n) - x_{0,0}(t_n)) - \lambda$  (by (4)) for all  $n$ , which implies  $d\lambda x_{0,\cdot}(t_0) - \lambda = \lim_{n \rightarrow \infty} d\lambda(x_{0,\cdot}(t_n) - x_{0,0}(t_n)) - \lambda \leq 0$ .

Lemma 13 also ensures that  $dx_{0,\cdot}(t) \leq 1$  on  $[t_0, \infty)$  as long as  $x_{0,0}(t) = 0$ . Substituting  $dx_{0,\cdot}(t) \leq 1$  in (43), we obtain  $\dot{\mathcal{L}}_S(x(t)) \leq 1/d - 1 + \lambda$ . Since  $\lambda < 1 - 1/d$  by hypothesis, this means that  $x_{0,0}(t)$  cannot remain equal to zero on  $[t_0, \infty)$  because  $\mathcal{L}_S(x(t))$  would be decreasing in  $t$  with derivative bounded away from zero and necessarily  $\mathcal{L}_S(x) \geq 0$  for all  $x$ . This implies that  $t^* \stackrel{\text{def}}{=} \inf\{t \geq t_0 : x_{0,0}(t) > 0\} < \infty$  must exist. Since  $x(t)$  is continuous and  $b(x)$  is linear when  $x_{0,0} > 0$ , there exists  $\delta > 0$  such that  $x_{0,0}(t)$  is both positive and increasing on  $(t^*, t^* + \delta]$ . This means that  $0 < \dot{x}_{0,0}(t) = \lambda d(x_{0,\cdot}(t) - x_{0,0}(t)) - \lambda$  for all  $t \in (t^*, t^* + \delta]$  (by (4)), which implies  $\lim_{t \downarrow t^*} x_{0,\cdot}(t) - x_{0,0}(t) = x_{0,\cdot}(t^*) \geq 1/d$ . On the other hand, on a left neighborhood of  $t^*$ ,  $x_{0,\cdot}(t^*) = \lim_{t \uparrow t^*} x_{0,\cdot}(t) \leq 1/d$ , and thus  $x_{0,\cdot}(t^*) = 1/d$ . Summarizing, we have obtained

$$x_{0,\cdot}(t^*) = \frac{1}{d}, \quad x_{0,0}(t^*) = 0 \quad (64a)$$

$$x_{0,0}(t) > 0, \quad x_{0,\cdot}(t) - x_{0,0}(t) > 1/d, \quad \forall t \in (t^*, t^* + \delta]. \quad (64b)$$

These conditions, together with the fact that  $b(x)$  is linear when  $x_{0,0} > 0$ , imply that also the function  $w_0(t) \stackrel{\text{def}}{=} x_{0,\cdot}(t) - x_{0,0}(t)$  must be increasing on a right neighborhood of  $t^*$ . Thus,

$$\dot{w}_0(t^*) = \lim_{t \downarrow t^*} \dot{w}_0(t) \geq 0. \quad (65)$$

On  $[t^*, t^* + \delta]$ , we have shown that the fluid solution  $x(t)$  is uniquely determined by the solution of the ODE system (63) where the initial condition  $x(t^*)$  is such that (64) and (65) hold true. In the remaining part of the proof, we study (63) under these conditions and show that  $x_{0,0}(t) > 0$  for all  $t > t^*$ . This will conclude the proof in view of Lemma 14.

Without loss of generality and by means of a time shift, let us assume  $t^* = 0$ . By taking proper summations in (63), we obtain

$$\dot{w}_{0,\cdot} = -\lambda dw_0 + x_{1,\cdot}. \quad (66)$$

Now, given that  $w_0(0) = 1/d$ , (65) ensures that  $\dot{w}_0(0) = -\lambda dw_0(0) + x_{1,\cdot}(0) \geq 0$ , which means

$$x_{1,\cdot}(0) \geq \lambda. \quad (67)$$

By taking proper summations in (63), we also obtain

$$\dot{x}_{0,\cdot} = -\lambda + x_{1,\cdot}. \quad (68a)$$

$$\dot{x}_{1,\cdot} = \lambda + x_{2,\cdot} - x_{1,\cdot}. \quad (68b)$$

$$\dot{x}_{i,\cdot} = x_{i+1,\cdot} \mathbf{1}_{\{i+1 \leq I\}} - x_{i,\cdot} \quad \forall i \geq 2. \quad (68c)$$

and solving for such autonomous ODE system,

$$\begin{aligned} x_{I,\cdot}(t) &= x_{I,I} = x_{I,I}(0) e^{-t} \\ x_{I-1,\cdot}(t) &= x_{I,I}(0) e^{-t} - x_{I-1,\cdot} = (x_{I,I}(0)t + x_{I-1,\cdot}(0)) e^{-t} \\ x_{I-i,\cdot}(t) &= e^{-t} \sum_{j=0}^i x_{I-i+j,\cdot}(0) \frac{t^j}{j!} \end{aligned}$$

for all  $i = 2, \dots, I-2$ , and thus

$$x_{2,\cdot}(t) = e^{-t} \sum_{j=0}^{I-2} x_{2+j,\cdot}(0) \frac{t^j}{j!}.$$

Substituting previous equation in (68b), we obtain

$$\dot{x}_{1,\cdot} = \lambda - x_{1,\cdot} + e^{-t} \sum_{j=0}^{I-2} x_{2+j,\cdot}(0) \frac{t^j}{j!} \quad (69)$$

and solving for such ODE we obtain,

$$x_{1,\cdot}(t) = \lambda + (x_{1,\cdot}(0) - \lambda)e^{-t} + e^{-t} \sum_{j=1}^{I-1} x_{1+j,\cdot}(0) \frac{t^j}{j!}$$

Substituting previous equation in (68a), we obtain

$$\dot{x}_{0,\cdot}(t) = (x_{1,\cdot}(0) - \lambda)e^{-t} + e^{-t} \sum_{j=1}^{I-1} x_{1+j,\cdot}(0) \frac{t^j}{j!}.$$

Integrating both sides, we obtain

$$x_{0,\cdot}(t) - \frac{1}{d} = x_{1,\cdot}(0) - \lambda - (x_{1,\cdot}(0) - \lambda)e^{-t} + \sum_{j=1}^{I-1} \frac{x_{1+j,\cdot}(0)}{j!} \int_0^t e^{-s} s^j ds.$$

Similarly, substituting (69) in (66) and solving for  $w_0(t)$ , we obtain

$$w_0(t) = \frac{1}{d} + \frac{x_{1,\cdot}(0) - \lambda}{\lambda d - 1} (e^{-t} - e^{-\lambda dt}) + e^{-\lambda dt} \sum_{j=1}^{I-1} x_{1+j,\cdot}(0) \int_0^t \frac{s^j}{j!} e^{(\lambda d - 1)s} ds$$

and thus

$$\begin{aligned} x_{0,0}(t) &= x_{0,\cdot}(t) - w_0(t) \\ &= (x_{1,\cdot}(0) - \lambda)(1 - e^{-t}) - \frac{x_{1,\cdot}(0) - \lambda}{\lambda d - 1} (e^{-t} - e^{-\lambda dt}) + \sum_{j=1}^{I-1} x_{1+j,\cdot}(0) \int_0^t \frac{s^j}{j!} e^{-s} (1 - e^{-\lambda d(t-s)}) ds. \end{aligned}$$

We now use the condition (67). If  $x_{1,\cdot}(0) = \lambda$ , then  $\sum_{i \geq 2} x_{i,\cdot}(0) = 1 - 1/d - \lambda > 0$  and therefore

$$x_{0,0}(t) = \sum_{j=1}^{I-1} \frac{x_{1+j,\cdot}(0)}{j!} \int_0^t s^j e^{-s} (1 - e^{-\lambda d(t-s)}) ds \geq \left(1 - \frac{1}{d} - \lambda\right) \min_{j=1}^{I-1} \int_0^t \frac{s^j}{j!} e^{-s} (1 - e^{-\lambda d(t-s)}) ds > 0$$

for all  $t > 0$ , as desired. If  $x_{1,\cdot}(0) > \lambda$ , then

$$x_{0,0}(t) \geq (x_{1,\cdot}(0) - \lambda)(1 - e^{-t}) - \frac{x_{1,\cdot}(0) - \lambda}{\lambda d - 1} (e^{-t} - e^{-\lambda dt}).$$

Given that  $x_{0,0}(0) = 0$ , to conclude that  $x_{0,0}(t) > 0$  for all  $t > t^* = 0$  it is sufficient to show that the RHS of last equation is strictly increasing in  $t$ . This follows easily once noted that the derivative of the RHS of last equation is strictly positive if and only if

$$\frac{e^{-t} - e^{-\lambda dt}}{\lambda d - 1} > 0, \quad \forall t > 0.$$

## 6 Conclusions

In this paper, we have provided new insights on randomized load balancing: if a load balancer is endowed with a local memory storing the last observation collected on each server, the celebrated power-of- $d$ -choices algorithm can be made asymptotically optimal in the sense that arriving jobs can be always routed to idle servers. Our approach provides an algorithm that is *both* fluid ( $N \rightarrow \infty$ ) and heavy-traffic ( $\lambda \uparrow 1$ ) optimal while employing a fair control message rate that scales linearly with the system size  $N$ . This means that randomized load balancing can be made robust to orthogonal variations of both  $N$  and  $\lambda$ , which can for instance occur in presence of unexpected workload peaks or server breakdowns.

On the practical side, Algorithm 1 can be improved in several ways to enhance performance:



- Server selections can be made without replacement, instead of with replacement. In view of the results in Gast and Van Houdt [10], this may also improve the convergence speed of  $X^N$  to fluid solutions.
- Since the action of sampling  $(0, 0)$ -servers does not bring any additional information to the load balancer, server selections can be restricted to  $(\cdot, j)$ -servers, with  $j \geq 1$ .
- Upon a job arrival, if  $i$  is both the least load of the  $d$  sampled servers and the least observation contained in the memory immediately before the last sampling, then the job may be randomly assigned to one of the  $(\cdot, i)$ -server known to the load balancer immediately before the sampling. In fact, by the time of the last update, one of such servers may have decreased its load.
- At any point in time, it is clear that the observation collected on a specific server is an upper bound on the actual state of that server. This observation leads us to consider a variant of  $\text{SQ}(d, N)$  where to each server is associated a *timer* representing the age of its observation. Specifically, at the moment where a new job arrives, all timers are incremented by one except the ones associated to the  $d$  sampled servers, whose timers are set to zero. Then, the load-balancer dispatches the job to a server with the lowest stored observation and for which the timer is the largest.

It is intuitive that all the above variations of  $\text{SQ}(d, N)$  yield performance improvements. It may be less intuitive that at the fluid scale only the last variation can yield performance improvements. In view of the results presented in this paper, such improvements can only appear when  $\lambda > 1 - 1/d$ . We leave this subject as future research.

A last variant of Algorithm 1 consists in swapping Lines 4–7 and 8–10. This is meant to perform each job assignment before the  $d$  servers are sampled (the collected information will be thus used for future assignments). It can be easily shown, *mutatis mutandis*, that this yields the same fluid limit. This is not surprising: if  $x_{0,0} > 0$ , the number of zeros in the memory is proportional to  $N$  and getting  $d$  more observations does not make  $x_{0,0}$  zero.

In our analysis, we have assumed that each server has a finite buffer of size  $I$ . We conjecture that our results generalize to the case where  $I = \infty$  and a first step to prove this claim consists in adapting the proofs of Lemma 2 and Proposition 2 on coordinates  $(i, i)$  only. Provided that servers are initially empty, this conjecture is coherent with the numerical observation that the Lyapunov function  $\mathcal{L}_S(x(t))$  monotonically increases in  $t$  to its limit point, which is necessarily less than  $j^* + 1$ ; see Figure 2.3. If  $\lambda < 1 - 1/d$  and  $x_{0,0}(0) > 0$ , this can be easily proven by using Lemma 14, which ensures that the drift function  $b$  takes the linear form in (48). When  $\lambda \geq 1 - 1/d$ , a proof is complicated by the involved structure of  $\mathcal{L}_S$ .

Our model can be generalized to a setting where servers have bin-packing constraints. Specifically, each server has  $B$  units of a resource, there are  $R$  types of jobs, type- $r$  jobs requires  $b_r$  units of resource, and an arriving job is ‘blocked’ if it does not find the required amount of resource at the server. Memoryless power-of- $d$ -choice strategies have been recently applied to this type of models in Xie et al. [24], though the resulting blocking probability does not converge to zero in the fluid limit. A further direction for future research aims at evaluating whether or not a local memory at the dispatcher can still be exploited to achieve fluid optimality in this setting.

## A Proof of Lemma 4

We give a proof when  $\epsilon \downarrow 0$  as the same arguments can be applied when  $\epsilon \uparrow 0$ . Let  $s_{i,j} \stackrel{\text{def}}{=} s_{i,j}(\bar{x}) \stackrel{\text{def}}{=} \sum_{i'=0}^{i-1} \bar{x}_{i'} + \sum_{j' \geq i}^j \bar{x}_{i,j'}$ . For any  $\epsilon > 0$  and all  $k$  sufficiently large, Lemma 3 states that the inclusions

$$\begin{aligned} (s_{i,i}(t) + 2\epsilon L(I+1)^2, s_{i,\cdot}(t) - 2\epsilon L(I+1)^2) &\subseteq (S_{i,i}^{N_k}(t_n^{N_k, \lambda^-}), S_{i,\cdot}^{N_k}(t_n^{N_k, \lambda^-})) \\ &\subseteq (s_{i,i}(t) - 2\epsilon L(I+1)^2, s_{i,\cdot}(t) + 2\epsilon L(I+1)^2) \end{aligned}$$

hold true, as we recall that under the coupled construction given in Section 4.1 we have that  $t_n^{N_k, \lambda^-} \in (t, t + \epsilon]$  when  $n \in \{\mathcal{N}_\lambda(N_k t) + 1, \dots, \mathcal{N}_\lambda(N_k(t + \epsilon))\}$ . Using these inclusions and Lemma 1, we obtain

$$-\lambda\epsilon + d\lambda(\bar{x}_{i,\cdot}(t) - \bar{x}_{i,i}(t) - 4\epsilon L(I+1)^2)\epsilon \leq \lim_{k \rightarrow \infty} \frac{1}{N_k} \sum_{n=\mathcal{N}_\lambda(N_k t)+1}^{\mathcal{N}_\lambda(N_k(t+\epsilon))} \sum_{p=1}^d \mathbb{I}_{(S_{i,i}^{N_k}(t_n^{N_k, \lambda^-}), S_{i,i}^{N_k}(t_n^{N_k, \lambda^-}))}^{(V_n^p)}$$

$$\leq -\lambda\epsilon + d\lambda(\bar{x}_{i,\cdot}(t) - \bar{x}_{i,i}(t) + 4\epsilon L(I+1)^2)\epsilon$$

for all  $\epsilon > 0$ , and dividing these inequalities by  $\epsilon$  and letting  $\epsilon \downarrow 0$ , we obtain (26a).

Finally, (26b) is proven using the same argument.

## B Proof of Lemma 9

At the beginning of the proof of Lemma 8, we have already shown that  $R_j(t) = 0$  if  $\sum_{j'=0}^j \bar{x}_{\cdot,j'}(t) > 0$ . Thus, in the following we assume that  $\sum_{j'=0}^j \bar{x}_{\cdot,j'}(t) = 0$ .

First, for the indicator function in (20c), we will use that

$$\mathbf{1}_{\{X_{0,0}^N(t_n^{N,\lambda^-}) + \sum_{j'=1}^I M_{0,j',n}^N > 0\}} = 1 - \mathbf{1}_{\{X_{0,0}^N(t_n^{N,\lambda^-}) = 0\}} \mathbf{1}_{\{\sum_{j'=1}^I M_{0,j',n}^N = 0\}} \quad (72a)$$

$$= 1 - \mathbf{1}_{\{X_{0,0}^N(t_n^{N,\lambda^-}) = 0\}} \prod_{p=1}^d \mathbb{I}_{(1-Z_1^N(t_n^{N,\lambda^-}), 1]}^{(V_n^p)}. \quad (72b)$$

We also notice the following sequence of equalities (see (17) and (18) for the definition of  $F_{i,j,n}^N$ )

$$\sum_{i=0}^j F_{i,j,n}^N = \mathbb{I}_{\left( \begin{array}{c} (W_n(X_{\cdot,j}^N(t_n^{N,\lambda^-}) + \overline{M}_{j,n}^N - \underline{M}_{j,n}^N) \\ 0, X_{\cdot,j}^N(t_n^{N,\lambda^-}) + \overline{M}_{j,n}^N - \underline{M}_{j,n}^N \end{array} \right)} \quad (73a)$$

$$= \mathbf{1}_{\{X_{\cdot,j}^N(t_n^{N,\lambda^-}) + \overline{M}_{j,n}^N - \underline{M}_{j,n}^N > 0\}} = 1 - \mathbf{1}_{\{X_{\cdot,j}^N(t_n^{N,\lambda^-}) - \underline{M}_{j,n}^N = 0\}} \mathbf{1}_{\{\overline{M}_{j,n}^N = 0\}} \quad (73b)$$

$$= 1 - \mathbf{1}_{\{X_{\cdot,j}^N(t_n^{N,\lambda^-}) - \underline{M}_{j,n}^N = 0\}} \prod_{p=1}^d \mathbb{I}_{(0, 1-Z_j^N(t_n^{N,\lambda^-})] \cup (1-Z_{j+1}^N(t_n^{N,\lambda^-}), 1]}^{(V_n^p)} \quad (73c)$$

$$= 1 - \left( \mathbf{1}_{\{X_{\cdot,j}^N(t_n^{N,\lambda^-}) = 0\}} + \sum_{p=1}^d \mathbf{1}_{\{NX_{\cdot,j}^N(t_n^{N,\lambda^-}) = p\}} \mathbf{1}_{\{N\underline{M}_{j,n}^N = p\}} \right) \prod_{p=1}^d \mathbb{I}_{(0, 1-Z_j^N(t_n^{N,\lambda^-})] \cup (1-Z_{j+1}^N(t_n^{N,\lambda^-}), 1]}^{(V_n^p)}$$

$$= 1 - \mathbf{1}_{\{X_{\cdot,j}^N(t_n^{N,\lambda^-}) = 0\}} \prod_{p=1}^d \mathbb{I}_{(0, 1-Z_j^N(t_n^{N,\lambda^-})] \cup (1-Z_{j+1}^N(t_n^{N,\lambda^-}), 1]}^{(V_n^p)} \quad (73d)$$

$$- \sum_{p=1}^d \mathbf{1}_{\{NX_{\cdot,j}^N(t_n^{N,\lambda^-}) = p\}} \mathbf{1}_{\{N\underline{M}_{j,n}^N = p\}} \times \prod_{p=1}^d \mathbb{I}_{(0, 1-Z_j^N(t_n^{N,\lambda^-})] \cup (1-Z_{j+1}^N(t_n^{N,\lambda^-}), 1]}^{(V_n^p)}. \quad (73e)$$

For the summation terms in (73e), we observe that

$$\begin{aligned} \mathbf{1}_{\{N\underline{M}_{j,n}^N = p\}} &= \mathbf{1}_{\{\sum_{i=0}^j NM_{i,j,n}^N = p\}} = \mathbf{1}_{\left\{ \sum_{i=0}^j \sum_{q=1}^d \mathbb{I}_{(S_{i,j-1}^N(t_n^{N,\lambda^-}), S_{i,j}^N(t_n^{N,\lambda^-})]}^{(V_n^q)} = p \right\}} \\ &= \mathbf{1}_{\left\{ \sum_{q=1}^d \mathbb{I}_{\bigcup_{i=0}^j (S_{i,j-1}^N(t_n^{N,\lambda^-}), S_{i,j}^N(t_n^{N,\lambda^-})]}^{(V_n^q)} = p \right\}} \\ &= \sum_{I \subseteq \{1, \dots, d\}; \|I\|=p} \prod_{q \in I} \mathbb{I}_{\bigcup_{i=0}^j (S_{i,j-1}^N(t_n^{N,\lambda^-}), S_{i,j}^N(t_n^{N,\lambda^-})]}^{(V_n^q)} \times \prod_{q \notin I} \mathbb{I}_{\neg \bigcup_{i=0}^j (S_{i,j-1}^N(t_n^{N,\lambda^-}), S_{i,j}^N(t_n^{N,\lambda^-})]}^{(V_n^q)} \\ &\leq \sum_{I \subseteq \{1, \dots, d\}; \|I\|=p} \prod_{q \in I} \mathbb{I}_{\bigcup_{i=0}^j (S_{i,j-1}^N(t_n^{N,\lambda^-}), S_{i,j}^N(t_n^{N,\lambda^-})]}^{(V_n^q)} \end{aligned}$$

where  $\neg A$  denotes the complement of set  $A$  and we have defined  $S_{j,j-1}^N(t_n^{N,\lambda^-}) \stackrel{\text{def}}{=} S_{j-1,I}^N(t_n^{N,\lambda^-})$ . Thus, for any  $\epsilon > 0$ , on the interval  $[t, t + \epsilon]$ , (25) ensures that

$$0 \leq \mathbf{1}_{\{N_k \underline{M}_{j,n}^N = p\}} \leq \sum_{\substack{I \subseteq \{1, \dots, d\}; \\ \|I\|=p}} \prod_{q \in I} \mathbb{I}_{\bigcup_{i=0}^j (0 \vee s_{i,j-1} - \epsilon, s_{i,j} + \epsilon]}^{(V_n^q)} \quad (74)$$

for all  $k$  sufficiently large, where  $s_{i,j} \stackrel{\text{def}}{=} s_{i,j}(\bar{x}(t)) \stackrel{\text{def}}{=} \sum_{i'=0}^{i-1} \bar{x}_{i',\cdot}(t) + \sum_{j' \geq i}^j \bar{x}_{i,j'}(t)$  for all  $i \leq j$  and  $s_{j,j-1} = s_{j-1,I}$ .

Let us treat the cases  $j = 1$  and  $j > 1$  separately.

Assume for now  $j = 1$ . Substituting (72) in the sample path expressions (19) and (20), we obtain

$$\begin{aligned} X_{\cdot,1}^N(t) &= X_{\cdot,1}^N(0) + \frac{1}{N} \sum_{n=1}^{\mathcal{N}_\lambda(Nt)} \sum_{p=1}^d \mathbb{I}_{(S_{1,1}^N(t_n^{N,\lambda^-}), S_{1,I}^N(t_n^{N,\lambda^-}))}^{(V_n^p)} \\ &\quad + \frac{1}{N} \sum_{n=1}^{\mathcal{N}_\lambda(Nt)} \left( 1 - \mathbf{1}_{\{R_0^N(t_n^{N,\lambda^-})=0\}} \prod_{p=1}^d \mathbb{I}_{(1-Z_1^N(t_n^{N,\lambda^-}), 1]}^{(V_n^p)} \right) \\ &\quad - \frac{1}{N} \sum_{n=1}^{\mathcal{N}_\lambda(Nt)} \sum_{p=1}^d \mathbb{I}_{(S_{0,1}^N(t_n^{N,\lambda^-}) - X_{0,1}^N(t_n^{N,1^-}), S_{0,1}^N(t_n^{N,\lambda^-}))}^{(V_n^p)} \\ &\quad - \frac{1}{N} \sum_{n=1}^{\mathcal{N}_\lambda(Nt)} \mathbf{1}_{\{R_0^N(t_n^{N,\lambda^-})=0\}} (F_{0,1,n}^N + F_{1,1,n}^N) \prod_{p=1}^d \mathbb{I}_{(1-Z_1^N(t_n^{N,\lambda^-}), 1]}^{(V_n^p)}. \end{aligned}$$

Since  $\bar{x}_{0,0}(t) + \bar{x}_{\cdot,1}(t) = 0$  and  $t$  is a point of differentiability, necessarily  $\dot{\bar{x}}_{\cdot,1}(t) = 0$  and thus

$$0 = \dot{\bar{x}}_{\cdot,1}(t) \tag{76a}$$

$$= \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \lim_{k \rightarrow \infty} X_{\cdot,1}^{N_k}(t + \epsilon) - X_{\cdot,1}^{N_k}(t) \tag{76b}$$

$$= \lambda d \bar{x}_{1,\cdot}(t) + \lambda - R_0(t) \tag{76c}$$

$$- \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \lim_{k \rightarrow \infty} \frac{1}{N_k} \sum_{n=\mathcal{N}_\lambda(N_k t)+1}^{\mathcal{N}_\lambda(N_k(t+\epsilon))} \mathbf{1}_{\{R_0^{N_k}(t_n^{N_k,\lambda^-})=0\}} (F_{0,1,n}^{N_k} + F_{1,1,n}^{N_k}) \prod_{p=1}^d \mathbb{I}_{(1-Z_1^{N_k}(t_n^{N_k,\lambda^-}), 1]}^{(V_n^p)}. \tag{76d}$$

The terms in (76c) are a direct application of Lemmas 1 and 4. Equation (76) also states that the limit in (76d) exists. Now, using (73)-(74) and (25), for any  $\epsilon > 0$

$$1 - \mathbf{1}_{\{R_j^{N_k}(t_n^{N_k,\lambda^-})=0\}} \prod_{p=1}^d \mathbb{I}_{(0, 1-Z_j^{N_k}(t_n^{N_k,\lambda^-})) \cup (1-Z_{j+1}^{N_k}(t_n^{N_k,\lambda^-}), 1]}^{(V_n^p)} \tag{77a}$$

$$\geq \sum_{i=0}^j F_{i,j,n}^{N_k} \geq \tag{77b}$$

$$1 - \mathbf{1}_{\{R_j^{N_k}(t_n^{N_k,\lambda^-})=0\}} \prod_{p=1}^d \mathbb{I}_{(0, 1-Z_j^{N_k}(t_n^{N_k,\lambda^-})) \cup (1-Z_{j+1}^{N_k}(t_n^{N_k,\lambda^-}), 1]}^{(V_n^p)} - \sum_{p=1}^d \sum_{\substack{I \subseteq \{1, \dots, d\}: \\ \|I\|=p}} \prod_{i=0}^j \mathbb{I}_{(0 \vee s_{i,j-1} - C\epsilon, s_{i,j} + C\epsilon]}^{(V_n^p)} \tag{77c}$$

for all  $k$  sufficiently large. Using both inequalities and that

$$\mathbb{I}_{(0, 1-Z_j^{N_k}(t_n^{N_k,\lambda^-})) \cup (1-Z_{j+1}^{N_k}(t_n^{N_k,\lambda^-}), 1]}^{(V_n^p)} \times \mathbb{I}_{(0, 1-Z_j^{N_k}(t_n^{N_k,\lambda^-}))}^{(V_n^p)} = \mathbb{I}_{(1-Z_{j+1}^{N_k}(t_n^{N_k,\lambda^-}), 1]}^{(V_n^p)}$$

for the term (76d) we obtain

$$\begin{aligned} R_0(t) - R_1(t) - \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \times O(\epsilon)\epsilon \\ \leq \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \lim_{k \rightarrow \infty} \frac{1}{N_k} \sum_{n=\mathcal{N}_\lambda(N_k t)+1}^{\mathcal{N}_\lambda(N_k(t+\epsilon))} \mathbf{1}_{\{R_0^{N_k}(t_n^{N_k,\lambda^-})=0\}} (F_{0,1,n}^{N_k} + F_{1,1,n}^{N_k}) \prod_{p=1}^d \mathbb{I}_{(1-Z_1^{N_k}(t_n^{N_k,\lambda^-}), 1]}^{(V_n^p)} \\ \leq R_0(t) - R_1(t). \end{aligned}$$

where the  $O(\epsilon)\epsilon$  term is obtained by applying Lemma 1 to the terms in the double sum of (77c), which gives

$$\lim_{k \rightarrow \infty} \frac{1}{N_k} \sum_{n=\mathcal{N}_\lambda(N_k t)+1}^{\mathcal{N}_\lambda(N_k(t+\epsilon))} \sum_{p=1}^d \sum_{I \subseteq \{1, \dots, d\}: \|I\|=p} \prod_{q \in I} \mathbb{I}_{\bigcup_{i=0}^q (0 \vee s_{i,j-1} - \epsilon C, s_{i,j} + \epsilon C]}^{(V_n^p)} = \sum_{p=1}^d \sum_{I \subseteq \{1, \dots, d\}: \|I\|=p} (2\epsilon C)^p \times \lambda \epsilon = O(\epsilon)\epsilon.$$

Thus, when  $\bar{x}_{0,0}(t) + \bar{x}_{\cdot,1}(t) = 0$  and  $t$  is a point of differentiability, we obtain

$$\lim_{\epsilon \downarrow 0} \frac{1}{\epsilon} \lim_{k \rightarrow \infty} \frac{1}{N_k} \sum_{n=\mathcal{N}_\lambda(N_k t)+1}^{\mathcal{N}_\lambda(N_k(t+\epsilon))} \mathbf{1}_{\{R_0^{N_k}(t_n^{N_k, \lambda^-})=0\}} \left( F_{0,1,n}^{N_k} + F_{1,1,n}^{N_k} \right) \prod_{p=1}^d \mathbb{I}_{(1-Z_1^{N_k}(t_n^{N_k, \lambda^-}), 1]}^{(V_n^p)} = R_0(t) - R_1(t)$$

and substituting this in (76) we obtain

$$0 = \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \lim_{k \rightarrow \infty} X_{\cdot,1}^{N_k}(t+\epsilon) - X_{\cdot,1}^{N_k}(t) = \lambda d\bar{x}_{1,\cdot}(t) + \lambda - 2R_0(t) + R_1(t),$$

which implies that  $R_1(t)$  exists and furthermore that (since  $R_1(t)$  is necessarily non-negative by definition) is given by

$$\begin{aligned} R_1(t) &= 0 \vee (2R_0(t) - \lambda d\bar{x}_{1,\cdot}(t) - \lambda) \mathbf{1}_{\{\bar{x}_{0,0}(t) + \bar{x}_{\cdot,1}(t) = 0\}} \\ &= 0 \vee \lambda (1 - 2d\bar{x}_{0,\cdot}(t)) - d\bar{x}_{1,\cdot}(t) \mathbf{1}_{\{\bar{x}_{0,0}(t) + \bar{x}_{\cdot,1}(t) = 0\}} \end{aligned}$$

where the last equation follows by substituting the expression of  $R_0(t)$  given in Lemma 8. This proves (34) when  $j = 1$ .

If  $j > 1$ , the same argument applies again. First, taking summations over the sample path expressions (19) and (20), we obtain

$$\begin{aligned} X_{\cdot,j}^N(t) &= X_{\cdot,j}^N(0) + \frac{1}{N} \sum_{n=1}^{\mathcal{N}_\lambda(Nt)} \sum_{p=1}^d \mathbb{I}_{(S_{i,j}^N(t_n^{N, \lambda^-}), S_{i,I}^N(t_n^{N, \lambda^-}))}^{(V_n^p)} \\ &\quad - \sum_{i=0}^{j-1} \frac{1}{N} \sum_{n=1}^{\mathcal{N}_\lambda(Nt)} \sum_{p=1}^d \mathbb{I}_{(S_{i,j}^N(t_n^{N, \lambda^-}) - X_{i,j}^N(t_n^{N, 1^-}), S_{i,j}^N(t_n^{N, \lambda^-}))}^{(V_n^p)} \\ &\quad + \frac{1}{N} \sum_{n=1}^{\mathcal{N}_\lambda(Nt)} \mathbf{1}_{\{R_{j-2}^N(t_n^{N, \lambda^-})=0\}} \left( \sum_{i=0}^{j-1} F_{i,j-1,n}^N \right) \prod_{p=1}^d \mathbb{I}_{(1-Z_{j-1}^N(t_n^{N, \lambda^-}), 1]}^{(V_n^p)} \\ &\quad - \frac{1}{N} \sum_{n=1}^{\mathcal{N}_\lambda(Nt)} \mathbf{1}_{\{R_{j-1}^N(t_n^{N, \lambda^-})=0\}} \left( \sum_{i=0}^j F_{i,j,n}^N \right) \prod_{p=1}^d \mathbb{I}_{(1-Z_j^N(t_n^{N, \lambda^-}), 1]}^{(V_n^p)}. \end{aligned}$$

Then, using (73) and noting that the term in (73e) can be bounded as in (77) we obtain

$$\begin{aligned} 0 &= \dot{\bar{x}}_{\cdot,j}(t) = \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \lim_{k \rightarrow \infty} X_{\cdot,j}^{N_k}(t+\epsilon) - X_{\cdot,j}^{N_k}(t) \\ &= \lambda d\bar{x}_{j,\cdot}(t) + (R_{j-2}(t) - R_{j-1}(t)) - (R_{j-1}(t) - R_j(t)), \end{aligned}$$

provided that  $\sum_{j'=0}^j \bar{x}_{\cdot,j'}(t) = 0$  and  $t$  is a point of differentiability. This condition inductively proves the existence of  $R_j(t)$  and implies that

$$R_j(t) = 0 \vee (2R_{j-1}(t) - R_{j-2}(t) - \lambda d\bar{x}_{j,\cdot}(t)) \mathbf{1}_{\{\sum_{j'=0}^j \bar{x}_{\cdot,j'}(t) = 0\}}.$$

Since we have already obtained an expression for  $R_0$  and  $R_1$ , we can derive an expression for  $R_2$  and so forth iteratively for all  $j$ . By iteratively substituting in previous equation the expressions of  $R_{j-1}$  and  $R_{j-2}$ , we obtain (34).

## C Proof of Lemma 10

As shown at the beginning of the proof of Lemma 8, if  $\sum_{i=0}^j \bar{x}_{\cdot,i}(t) > 0$ , then  $R_j^{N_k}(t_n^{N_k, \lambda^-}) > 0$  for all  $k$  sufficiently large, and therefore  $\Gamma_{i,j} = 0$ . Thus, let us assume in the following that  $\sum_{i=0}^j \bar{x}_{\cdot,i}(t) = 0$ . We give a proof when  $i \leq j$ ; when  $i = j + 1$ , the proof uses the same argument and is omitted.

So far we have assumed that  $\omega \in \mathcal{C}$  was fixed but for now let us explicit the dependence on  $\omega$  and treat quantities  $\bar{x}(t)$  and  $X^N(t)$  as random variables.

Let

$$\Gamma_{i,j}^{\epsilon,N}(t) \stackrel{\text{def}}{=} \frac{1}{\epsilon N} \sum_{n=\mathcal{N}_\lambda(Nt)+1}^{\mathcal{N}_\lambda(N(t+\epsilon))} \underbrace{\mathbf{1}_{\{R_j^{N_k}(t_n^{N_k, \lambda^-})=0\}} F_{i,j+1,n}^N \prod_{p=1}^d \mathbb{I}_{(1-Z_{j+1}^{N_k}(t_n^{N_k, \lambda^-}), 1]^{(V_n^p)}}}_{\stackrel{\text{def}}{=} Y_{n,N}} \quad (79)$$

For all  $n$ , the random variable  $Y_{n,N}$  is  $\mathcal{F}_n$ -measurable where  $\mathcal{F}_n \stackrel{\text{def}}{=} \{X^N(t_n^{N, \lambda^-}), V_n^1, \dots, V_n^d, U_n, W_n\}$ , and

$$\mathbb{E}[Y_{n,N} | \mathcal{F}_n \setminus W_n] = \mathbf{1}_{\{R_j^{N_k}(t_n^{N_k, \lambda^-})=0\}} \frac{X_{i,j+1}^N(t_n^{N_k, \lambda^-}) - M_{i,j+1,n}^N}{X_{\cdot,j+1}^N(t_n^{N_k, \lambda^-}) + \overline{M}_{j+1,n}^N - \underline{M}_{j+1,n}^N} \prod_{p=1}^d \mathbb{I}_{(1-Z_{j+1}^{N_k}(t_n^{N_k, \lambda^-}), 1]^{(V_n^p)}} \quad (80)$$

because we recall that the 0-1 random variable  $F_{i,j,n}^N$  is one if and only if  $(X_{\cdot,j}^N(t_n^{N_k, \lambda^-}) + \overline{M}_{j,n}^N - \underline{M}_{j,n}^N) W_n \in (\sum_{k=0}^{i-1} X_{k,j}^N(t_n^{N_k, \lambda^-}) - M_{k,j,n}^N, \sum_{k=0}^i X_{k,j}^N(t_n^{N_k, \lambda^-}) - M_{k,j,n}^N]$ , by definition (17), with  $W_n$  uniform over  $[0,1]$ ; the set  $\mathcal{F}_n \setminus W_n$  denotes the set  $\mathcal{F}_n$  with  $W_n$  removed.

Let  $Z_{n,N} \stackrel{\text{def}}{=} Y_{n,N} - \mathbb{E}[Y_{n,N} | \mathcal{F}_n \setminus W_n]$ . Then,  $\mathbb{E}[Z_{n,N} | \mathcal{F}_n \setminus W_n] = 0$  and  $|Z_{n,N}| \leq 2$ , and applying the Azuma–Hoeffding inequality, we get

$$\mathbb{P} \left( \frac{1}{N} \left| \sum_{n=1}^N Z_{n,N} \right| > \delta \right) \leq 2 \exp \left( -\frac{(N\delta)^2}{8N} \right) \quad (81)$$

for any  $\delta > 0$ . Since  $\sum_N \exp(-N\delta^2/8) < \infty$ , an application of the Borel–Cantelli lemma shows that  $\frac{1}{N} \sum_{n=1}^N Z_{n,N} \rightarrow 0$  almost surely. In particular,

$$\lim_{N \rightarrow \infty} \Gamma_{i,j}^{\epsilon,N}(t) - \frac{1}{N\epsilon} \sum_{n=\mathcal{N}_\lambda(Nt)+1}^{\mathcal{N}_\lambda(N(t+\epsilon))} \mathbb{E}[Y_{n,N} | \mathcal{F}_n \setminus W_n] = 0 \quad (82)$$

almost surely.

Now, we fix  $\omega \in \mathcal{C}$  and use (80) and Lemma 3 to obtain that for any  $\epsilon > 0$

$$\begin{aligned} & \lim_{k \rightarrow \infty} \frac{1}{\epsilon N_k} \sum_{n=\mathcal{N}_\lambda(N_k t)+1}^{\mathcal{N}_\lambda(N_k(t+\epsilon))} \mathbb{E}[Y_{n,N_k} | \mathcal{F}_n \setminus W_n] \\ & \leq \lim_{N_k \rightarrow \infty} \frac{1}{\epsilon N_k} \sum_{n=\mathcal{N}_\lambda(N_k t)+1}^{\mathcal{N}_\lambda(N_k(t+\epsilon))} \mathbf{1}_{\{R_j^{N_k}(t_n^{N_k, \lambda^-})=0\}} \frac{\bar{x}_{i,j+1}(t) + C\epsilon}{\bar{x}_{\cdot,j+1}(t) - C\epsilon} \prod_{p=1}^d \mathbb{I}_{(1-Z_{j+1}^{N_k}(t_n^{N_k, \lambda^-}), 1]^{(V_n^p)}} \end{aligned}$$

Replacing  $\epsilon$  by  $-\epsilon$  in the last fraction term, the previous inequality can be reversed and letting  $\epsilon \downarrow 0$ , we obtain

$$\lim_{\epsilon \downarrow 0} \lim_{k \rightarrow \infty} \frac{1}{\epsilon N_k} \sum_{n=\mathcal{N}_\lambda(N_k t)+1}^{\mathcal{N}_\lambda(N_k(t+\epsilon))} \mathbb{E}[Y_{n,N_k} | \mathcal{F}_n \setminus W_n] = R_j(t) \frac{\bar{x}_{i,j+1}(t)}{\bar{x}_{\cdot,j+1}(t)}. \quad (83)$$

Finally, (82) and (83) give (36).

## Acknowledgments

The authors would like to thank Sem Borst, Bruno Gaujal and the referees for their valuable comments and remarks.

## References

- [1] Azar Y, Broder AZ, Karlin AR, Upfal E (1999) Balanced allocations. *SIAM J. Comput.* 29(1):180–200, ISSN 0097-5397, URL <http://dx.doi.org/10.1137/S0097539795288490>.
- [2] Bramson M (1998) State space collapse with application to heavy traffic limits for multiclass queueing networks. *Queueing Syst. Theory Appl.* 30(1/2):89–148, ISSN 0257-0130, URL <http://dx.doi.org/10.1023/A:1019160803783>.
- [3] Bramson M, Lu Y, Prabhakar B (2010) Randomized load balancing with general service time distributions. *Proceedings of the ACM SIGMETRICS International Conference on Measurement and Modeling of Computer Systems*, 275–286, SIGMETRICS '10 (New York, NY, USA: ACM), ISBN 978-1-4503-0038-4, URL <http://dx.doi.org/10.1145/1811039.1811071>.
- [4] Bramson M, Lu Y, Prabhakar B (2012) Asymptotic independence of queues under randomized load balancing. *Queueing Syst. Theory Appl.* 71(3):247–292, ISSN 0257-0130, URL <http://dx.doi.org/10.1007/s11134-012-9311-0>.
- [5] Chen H, Ye HQ (2012) Asymptotic optimality of balanced routing. *Oper. Res.* 60(1):163–179, ISSN 0030-364X, URL <http://dx.doi.org/10.1287/opre.1110.0998>.
- [6] Dieker AB, Suk T (2015) Randomized longest-queue-first scheduling for large-scale buffered systems. *Advances in Applied Probability* 47(4):1015–1038, URL <http://dx.doi.org/10.1017/S0001867800048990>.
- [7] Gamarnik D, Tsitsiklis JN, Zubeldia M (2016) Delay, memory, and messaging tradeoffs in distributed service systems. *Proceedings of the 2016 ACM SIGMETRICS International Conference on Measurement and Modeling of Computer Science*, 1–12, SIGMETRICS '16 (New York, NY, USA: ACM), ISBN 978-1-4503-4266-7, URL <http://dx.doi.org/10.1145/2896377.2901478>.
- [8] Gamarnik D, Tsitsiklis JN, Zubeldia M (2018) Delay, memory, and messaging tradeoffs in distributed service systems. *Stochastic Systems* 8:1–54, URL <http://dx.doi.org/10.1214/17-SSY234>.
- [9] Gardner K, Harchol-Balter M, Scheller-Wolf A, Velednitsky M, Zbarsky S (2017) Redundancy-d: The power of d choices for redundancy. *Operations Research* 65(4):1078–1094, URL <http://dx.doi.org/10.1287/opre.2016.1582>.
- [10] Gast N, Van Houdt B (2017) A refined mean field approximation. *Proc. ACM Meas. Anal. Comput. Syst.* 1(2):33:1–33:28, ISSN 2476-1249, URL <http://dx.doi.org/10.1145/3154491>.
- [11] Gupta V, Walton N (2017) Load Balancing in the Non-Degenerate Slowdown Regime. *ArXiv e-prints* .
- [12] Lu Y, Xie Q, Kliot G, Geller A, Larus JR, Greenberg A (2011) Join-idle-queue: A novel load balancing algorithm for dynamically scalable web services. *Perform. Eval.* 68(11):1056–1071, ISSN 0166-5316, URL <http://dx.doi.org/10.1016/j.peva.2011.07.015>.
- [13] Luczak MJ, Norris JR (2013) Averaging over fast variables in the fluid limit for markov chains: Application to the supermarket model with memory. *Ann. Appl. Probab.* 23(3):957–986, URL <http://dx.doi.org/10.1214/12-AAP861>.
- [14] Maguluri ST, Srikant R, Ying L (2012) Heavy traffic optimal resource allocation algorithms for cloud computing clusters. *Proceedings of the 24th International Teletraffic Congress*, 25:1–25:8, ITC '12 (International Teletraffic Congress), ISBN 978-1-4503-1896-9, URL <http://dl.acm.org/citation.cfm?id=2414276.2414307>.

- [15] Mitzenmacher M (2001) The power of two choices in randomized load balancing. *IEEE Trans. Parallel Distrib. Syst.* 12(10):1094–1104, ISSN 1045-9219, URL <http://dx.doi.org/10.1109/71.963420>.
- [16] Mitzenmacher M, Prabhakar B, Shah D (2002) Load balancing with memory. *The 43rd Annual IEEE Symposium on Foundations of Computer Science, 2002. Proceedings.*, 799–808, ISSN 0272-5428, URL <http://dx.doi.org/10.1109/SFCS.2002.1182005>.
- [17] Mukherjee D, Borst SC, van Leeuwen JSH, Whiting PA (2016) Asymptotic Optimality of Power-of- $d$  Load Balancing in Large-Scale Systems. *ArXiv e-prints* .
- [18] Stolyar AL (2015) Pull-based load distribution in large-scale heterogeneous service systems. *Queueing Syst. Theory Appl.* 80(4):341–361, ISSN 0257-0130, URL <http://dx.doi.org/10.1007/s11134-015-9448-8>.
- [19] Tsitsiklis JN, Xu K (2012) On the power of (even a little) resource pooling. *Stoch. Syst.* 2(1):1–66, URL <http://dx.doi.org/10.1214/11-SSY033>.
- [20] van der Boor M, Borst SC, van Leeuwen JSH, Mukherjee D (2018) Scalable load balancing in networked systems: A survey of recent advances. *arXiv e-prints* arXiv:1806.05444.
- [21] Vvedenskaya N, Dobrushin R, Karpelevich R (1996) Queueing system with selection of the shortest of two queues: an asymptotic approach. *Problems of information transmission* 32:15–27.
- [22] Weber RR (1978) On the optimal assignment of customers to parallel servers. *Journal of Applied Probability* 15(2):406–413, URL <http://dx.doi.org/10.1017/S0021900200045678>.
- [23] Winston W (1977) Optimality of the shortest line discipline. *Journal of Applied Probability* 14(1):181–189, URL <http://dx.doi.org/10.1017/S0021900200104772>.
- [24] Xie Q, Dong X, Lu Y, Srikant R (2015) Power of  $d$  choices for large-scale bin packing: A loss model. *Proceedings of the 2015 ACM SIGMETRICS International Conference on Measurement and Modeling of Computer Systems*, 321–334, SIGMETRICS '15 (New York, NY, USA: ACM), ISBN 978-1-4503-3486-0, URL <http://dx.doi.org/10.1145/2745844.2745849>.
- [25] Ying L, Srikant R, Kang X (2017) The power of slightly more than one sample in randomized load balancing. *Mathematics of Operations Research* 42(3):692–722, URL <http://dx.doi.org/10.1287/moor.2016.0823>.